



مرکز تحقیقات اسلامی

اصفهان

گامی



عمران
علیه السلام

www. **Ghaemiyeh** .com
www. **Ghaemiyeh** .org
www. **Ghaemiyeh** .net
www. **Ghaemiyeh** .ir



مدیریت منابع اطلاعاتی وب

۱-۲



به کوشش
دکتر فاطمه منتظر
فرزانه شادان پور

بِسْمِ اللَّهِ الرَّحْمَنِ الرَّحِيمِ

مدیریت منابع اطلاعاتی وب

نویسنده:

غلامعلی منتظر

ناشر چاپی:

سازمان اسناد و کتابخانه ملی جمهوری اسلامی ایران

ناشر دیجیتال:

مرکز تحقیقات رایانه‌ای قائمیه اصفهان

فهرست

۵	فهرست
۱۴	مدیریت منابع اطلاعاتی وب
۱۴	مشخصات کتاب
۱۵	جلد ۱
۱۵	اشاره
۲۰	فهرست مطالب
۲۲	سخن نخست
۲۴	به جای مقدمه
۲۶	فصل اول: مبانی مدیریت و آرشيو وب
۲۶	اشاره
۲۷	چکیده
۲۸	بایگانی شبکه وب: مباحث و روش ها
۲۸	اشاره
۲۸	۱-مقدمه
۲۹	۲-حفاظت از میراث
۳۱	۱-۲-به اندازه کافی خوب نیست؟
۳۷	۱-۲-۲-خود بقای رسانه؟
۴۰	۱-۲-۳-یک وظیفه غیر ممکن
۴۶	۳-ویژگی های وب برای حفاظت
۴۷	۱-۳-کارديناليته وب
۵۰	۲-۳-وب به عنوان سیستم نشر فعال
۵۳	۳-۳-وب به عنوان یک محصول فرهنگی
۵۶	۴-روش های جدید برای رسانه جدید
۵۷	۱-۴-حفاظت از وب و زیر ساخت های اطلاعات
۵۹	۲-۴-فرآهم آوری
۶۱	۱-۲-۴-بایگانی جانبی سرویس گیرنده
۶۷	۲-۲-۴-بایگانی تراکشی
۷۰	۲-۳-۴-بایگانی سرور- جانبی
۷۳	۳-۴-سازماندهی و ذخیره سازی
۷۳	۱-۳-۴-نظام فایل های محلی به خدمت گرفته شده بایگانی ها
۷۳	توصیف
۷۵	توضیح
۷۸	استفاده ترجیحی
۷۸	ابزارها
۷۸	۲-۳-۴-بایگانی های مبتنی بر خدمت وب
۷۸	اشاره
۸۰	توصیف
۸۰	توضیح
۸۳	استفاده ترجیحی
۸۵	ابزارها

۸۵	۳-۴-بایگانی گیروپ
۸۵	توصیف
۸۷	توضیح
۸۷	مرجع
۸۷	اشاره
۸۷	۳-۴-خلاصه
۸۸	۴-۴-کیفیت و تمامیت (کامل بودن)
۹۴	۵-بررسی عمومی مراحل اولیه جاری
۹۵	۱-۵-بازبینی بایگانی
۹۷	۲-۵-حوزه (دامنه)
۹۷	۵-۲-۱-بایگانی مرکزی سایت
۹۷	۵-۲-۲-بایگانی مرکزی عنوان
۹۸	۵-۲-۳-بایگانی مرکزی حوزه ای
۱۰۰	۳-۵-روش های استفاده شده
۱۰۲	۶-نتیجه گیری
۱۰۳	منابع
۱۲۲	چکیده
۱۲۳	از آرشیو اینترنت تا آرشیو در اینترنت
۱۲۳	اشاره
۱۲۳	مقدمه
۱۲۴	پیشینه: جاب اینترنتی اولیه
۱۲۴	آغاز به کار آرشیو اینترنت
۱۲۸	ساختار پیوندی و روبات های نواری
۱۳۰	۱۹۹۸-حضور داده های آرشیوی بر روی (تقریباً) هر دسکتاپ
۱۳۰	۱۹۹۹: از نوار تا دیسک، یک خزشگر جدید و تصاویر متحرک
۱۳۲	۲۰۰۰: ایجاد مجموعه های موضوعی وب
۱۳۵	۲۰۰۱: دسترسی از طریق ماشین Wayback. آرشیو یازده سپتامبر
۱۳۹	۲۰۰۲: کتابخانه اسکندریه، کتابخانه سیار، و حق مؤلف
۱۴۲	۲۰۰۳: گسترش دستیابی ما به کتابخانه های ملی و مؤسسات آموزشی
۱۴۴	۲۰۰۴: آرشیو اروپا و پتاپاکس
۱۴۴	اشاره
۱۴۴	آینده
۱۴۴	منابع
۱۴۶	چکیده
۱۴۷	کاربرد وب و مطالعات مربوط به آن
۱۴۷	اشاره
۱۴۷	خلاصه
۱۴۷	اشاره
۱۴۸	۱-تحلیل محتوا
۱۵۰	۲-بررسی ها
۱۵۱	۳-تحلیل بلاغی

۱۵۳	تحلیل گفتنمان
۱۵۵	تحلیل دیداری
۱۵۷	قوم نگاری
۱۵۸	تحلیل شبکه
۱۵۹	ملاحظات اخلاقی
۱۶۰	نتیجه گیری
۱۶۰	منابع
۱۶۶	چکیده
۱۶۷	خصوصیات وب ایران: مریم پیروزمند
۱۶۷	اشاره
۱۶۷	۱. مقدمه
۱۷۰	۲. کارهای مرتبط در داخل و خارج
۱۷۲	۳. سامانه خودکار ارزیابی وب ایران
۱۷۸	۴. نتایج بدست آمده
۱۷۹	۴.۱. آمار وب گاه ها
۱۸۴	۴.۲. آمار صفحات
۱۸۸	۴.۳. آمار نوع صفحات
۱۹۱	۵. نتیجه گیری و کارهای آینده
۱۹۳	منابع
۱۹۳	اشاره
۱۹۶	چکیده
۱۹۷	آینده آرشیو وب
۱۹۷	اشاره
۱۹۷	خلاصه اجرایی
۱۹۷	بخش اول. نمای کلی از چهار سناریوی احتمالی برای آینده:
۱۹۷	اشاره
۱۹۸	مقدمه
۲۰۱	ساختن آینده
۲۰۲	سناریوها (طرح ها)
۲۰۲	اشاره
۲۰۲	سناریوی نیروانا
۲۰۴	سناریوی آپوکالیپس (آخر الزمان)
۲۰۶	سناریوی انفرادی
۲۰۶	سناریوی غبار آلود
۲۱۰	مرور آینده
۲۱۰	یادگیری از وب پویا
۲۱۲	مصورسازی
۲۱۲	چالش ها
۲۱۴	برنامه های کاربردی جست و جو همانند شکارچی
۲۱۵	تحلیل های شبکه اجتماعی
۲۱۷	مثال ها:

۲۲۲	سنجش های دگرساز
۲۲۵	وب نوشت (حاشیه نگاری) اجتماعی
۲۲۷	معماران جدید
۲۳۰	ماشین های اجتماعی
۲۳۲	شیکه های نقشه برداری
۲۳۳	علم وب
۲۳۴	اشاره
۲۳۵	درک تجربه به جای محتوا
۲۳۶	تحلیل وب معنایی و مجموعه داده های پیوند شده
۲۳۹	چالش های فعلی و آینده
۲۳۹	اشاره
۲۴۰	وب مجتمع: زندگی آرشیو وب
۲۴۰	چالش بلند مدت: دو چالش در این سؤال نهفته است، که هر دوی آن ها مستلزم مشارکت کنندگان و فعالان زیادی خواهد بود. نخست: ما باید دوباره در مورد اینکه چگونه اینترنت را ببینیم و مهندسی کنیم فکر کنیم، در حال حرکت از موجودیتی تک لایه یا پیوندهای ج
۲۴۲	وب در حال تغییر
۲۴۴	کاربرد های آرشیو ها و وب گاه ها
۲۴۷	متخصص وب
۲۴۸	وب دیداری
۲۵۲	وب همان گونه که بود
۲۵۴	ساختار وب
۲۵۴	ایده ها چگونه تکثیر می شوند
۲۵۵	وب غیر قانونی
۲۵۶	رد پای رقومی
۲۵۷	وب داده
۲۶۰	نتایج: مسیر پیش رو
۲۶۱	منابع
۲۶۷	فصل دوم: تجارب جهانی و مسائل بومی در آرشیو سازی وب
۲۶۷	اشاره
۲۶۸	چکیده
۲۶۹	آرشیو وب در دنیای وب ۲/۰ شعبه آرشیو وب و حفاظت رقومی کتابخانه ملی استرالیا
۲۶۹	اشاره
۲۶۹	مقدمه
۲۷۰	سه روش شناسی آرشیو
۲۷۶	جمع آوری فایل ها
۲۷۸	دستورالعمل های جمع آوری
۲۷۹	دستور عمل های آینده
۲۷۹	شرح حال مختصری از پدید آورنده
۲۸۰	چکیده
۲۸۱	آسیب شناسی زبان و خط فارسی در بازیابی اطلاعات: نگاهی به موتور های کاوش و پایگاه های برخط
۲۸۱	اشاره
۲۸۱	درآمدی بر مشخصه های زبان فارسی
۲۸۳	پیشینه پژوهش

۲۸۴	رسم الخط فارسی و بازیابی اطلاعات
۲۸۶	مسائل صرفی و بازیابی اطلاعات
۲۸۷	مسائل معنایی و بازیابی اطلاعات
۲۹۶	سخن پایانی: پیشنهادهایی در جهت بهبود وضعیت
۲۹۸	منابع
۳۰۲	چکیده
۳۰۳	ارزیابی کاربرد پذیری وبگاه نهاد کتابخانه های عمومی کشور
۳۰۳	اشاره
۳۰۳	مقدمه
۳۰۴	۲-روش پژوهش و توجیه رویی آن
۳۰۴	۳-شیوه گردآوری اطلاعات و تجزیه و تحلیل آن ها
۳۰۴	۴-تجزیه و تحلیل داده ها و ارائه یافته ها
۳۱۲	۵-نتیجه گیری
۳۱۲	۶-پیشنهادهای پژوهش
۳۱۲	منابع
۳۱۸	چکیده
۳۱۹	امکان سنجی برداش وبگاه ها در سازمان اسناد و کتابخانه ملی ایران
۳۱۹	اشاره
۳۱۹	مقدمه
۳۲۰	وبگاه ها؛ شناخته شده ترین منابع وب
۳۲۰	ارزیابی وبگاه ها
۳۲۰	استانداردها و طرحهای ایر داده ای سازماندهی منابع وب
۳۲۲	بردازش منابع الکترونیکی در سازمان اسناد و کتابخانه ملی ایران
۳۲۲	بیان مساله
۳۲۴	پیشینه پژوهش
۳۲۴	پیشینه پژوهش در ایران:
۳۲۴	پیشینه پژوهش در خارج از ایران:
۳۲۵	اهمیت پژوهش
۳۲۵	اهداف پژوهش
۳۲۷	برسش های اساسی
۳۲۷	تعاریف عملیاتی
۳۲۷	جامعه آماری
۳۲۷	روش پژوهش و ابزار گردآوری داده ها
۳۲۷	یافته های پژوهش و پاسخ به برسش های اساسی
۳۲۸	۱.استفاده از فهرست وب گاه ها
۳۲۸	۲.ورود اطلاعات وبگاه ها در وبگاه سازمان اسناد و کتابخانه ملی ایران
۳۳۵	بحث و نتیجه گیری
۳۳۷	پیشنهاد های برخاسته از پژوهش
۳۳۷	منابع
۳۴۱	چکیده
۳۴۲	مقدمه

۳۴۳	ایجاد آرشیو گزینشی منابع تحت وب
۳۴۳	اشاره
۳۴۴	تاریخچه
۳۴۶	آرشیو پاندورا در حال حاضر
۳۴۶	محدوده وظایف
۳۴۸	نیروی انسانی و افراد شاغل در پاندورا
۳۴۹	هزینه دستیابی به وب گاه ها و نشریات برخط
۳۴۹	محدوده برآورد هزینه
۳۵۱	روش شناسی انجام کار
۳۵۲	نحوه محاسبه هزینه ها
۳۵۲	هزینه های فراهم آوری منابع آرشیوی
۳۵۵	مقایسه با نوع چاپی
۳۵۶	هزینه فعالیت های خاص
۳۵۶	امکان کاهش هزینه ها
۳۵۷	نتیجه گیری
۳۵۸	قدردانی
۳۵۸	منابع
۳۶۰	چکیده
۳۶۱	بایگانی وب علمی در مقیاس کوچک
۳۶۱	اشاره
۳۶۱	۱-جرایب بایگانی علمی در مقیاس کوچک
۳۶۲	۲-بایگانی دیجیتال برای مطالعات زبان چینی
۳۶۲	اشاره
۳۶۴	۲-۱-گام های اولیه
۳۶۵	۲-۲-توسعه پایدار سازمانی
۳۶۶	۲-۳-سخت افزار
۳۶۷	۲-۴-نرم افزار
۳۶۷	۲-۵-فرآیندها
۳۶۹	۲-۶-سیاست گذاری و خط مشی مجموعه
۳۷۱	۲-۷-مشارکت
۳۷۲	۳-درس های آموخته شده: جمع بندی
۳۷۲	۴-منابع مفید
۳۷۶	چکیده
۳۷۷	بررسی و مقایسه قابلیت های قالبهای بونی مارک و مارک ۲۱ برای سازماندهی منابع اطلاعاتی وب
۳۷۷	اشاره
۳۷۷	۱-مقدمه و بیان مسئله
۳۷۸	۲-بررسی های پژوهش
۳۷۹	۳-پیشینه پژوهش
۳۸۰	۴-روش شناسی پژوهش
۳۸۰	۵-تجزیه و تحلیل یافته ها

۳۸۹	منابع
۳۹۲	چکیده
۳۹۳	سنجش رابط کاربر پایگاه های اطلاعاتی پیوسته، مجلات تمام متن فارسی
۳۹۳	اشاره
۳۹۳	مقدمه
۳۹۴	بیان مسأله و ضرورت پژوهش
۳۹۵	هدف پژوهش
۳۹۶	پیشینه ی پژوهش
۳۹۸	جمع بندی پیشینه
۳۹۹	روش شناسی پژوهش
۳۹۹	یافته های پژوهش
۴۰۶	بحث و نتیجه گیری
۴۰۷	پیشنهاد هایی برای پژوهش های آینده:
۴۰۷	منابع
۴۰۹	چکیده
۴۱۰	قانون واسپاری وب فرانسه: راهبردهایی برای گرده آوری دامنه ملی
۴۱۰	اشاره
۴۱۰	۱. قلمرو فرانسه
۴۱۰	۱.۱. تعریف حوزه قانون و اسپاری
۴۱۵	۲.۱. در حال حاضر کجا هستیم
۴۱۷	۲. طراحی خزش
۴۱۷	۱.۲. هدف چیست؟
۴۱۹	۲.۲. فهرست هسته
۴۲۴	۳.۲. تنظیمات خزش گر
۴۲۵	۱.۳.۲. حوزه
۴۲۷	۲.۳.۲. اولویت های خزش
۴۲۹	۳.۳.۲. سایر تنظیمات
۴۳۱	۴.۲. پروتکل حذف روایات ها
۴۳۳	۵.۲. برنامه ریزی با همکاری آرشیو اینترنت (IA)
۴۳۴	۳. خزش گر در حال کار
۴۳۴	۱.۳. آزمایش خزش ها
۴۳۴	۲.۳. ارتباط کتابخانه ملی فرانسه و آرشیو اینترنت در طول خزش
۴۳۵	۳.۳. «خزش تکه ای»
۴۳۵	۴. پیامدهای خزش
۴۳۵	۱.۴. اشکال اصلی
۴۳۸	۲.۴. توزیع هر سرآیند
۴۴۲	۳.۴. فایل های ویدئویی
۴۴۳	۴.۴. توزیع هر TLD
۴۴۶	۵.۴. Robots.txt
۴۴۷	۶. عمق خزش
۴۵۰	۷.۴. وب گاه های بزرگ

۴۵۰ ۱.۷.۴ دامنه ها
۴۵۲ ۲.۷.۴ دامنه های سطح دوم
۴۵۳ دنتیجه گیری
۴۵۵ تشکر و قدردانی
۴۵۶ منابع
۴۶۱ چکیده
۴۶۲ معرفی آرشوهای وب به عنوان یک خدمت جدید کتابخانه: تجربه کتابخانه ملی فرانسه
۴۶۲ اشاره
۴۶۳ چارچوب حقوقی
۴۶۵ دامنه
۴۶۷ ابزارها و روشهای جمع آوری
۴۶۸ کشف منبع: خدمات و ابزارهای دسترسی برای کاربران نهایی
۴۶۸ تعیین جای خدمت و محدودیت های دسترسی
۴۶۹ ابزارهای جستجو و دیدن
۴۷۴ مجموعه های ویژه
۴۷۸ کاربرد اطلاعات و ارقام
۴۷۸ تحلیل کمی
۴۸۰ تحلیل کیفی
۴۸۴ نظر سنجی ها
۴۸۷ راهبردهای به حساب آوردن آرشوهای وب به عنوان بخشی از کار روزانه کتابخانه
۴۸۸ راهبردهایی برای رسیدن به کاربران نهایی
۴۹۰ منابع
۴۹۱ چکیده
۴۹۲ یک سال آرشو وب گزینشی
۴۹۲ اشاره
۴۹۲ مقدمه
۴۹۳ ۱.۲ آرشو وب گزینشی در کتابخانه ملی زلاندنو
۴۹۳ ۲-۱ انگیزه
۴۹۳ ۲-۲ تاریخچه دروا هاروستینگ
۴۹۵ ۲-۳ نرم افزار گردآوری وب
۴۹۶ شکل ۱. فهرست WCT
۴۹۸ ۴.۲ اعضا و منابع
۵۰۰ ۲-۵ سطح گردآوری / درو
۵۰۳ گردآوری با WCT ویرایش ۱.۱
۵۰۳ ۱.۳ تجربه اولیه
۵۰۵ ۳.۲ مشکلات گردآوری
۵۰۵ کمیوهای تورق نرم افزار
۵۰۷ نسخه های میانجی کاربر
۵۰۷ اشتباه «در مرحله ایست درگیر شدن»
۵۰۷ وب گاه های بزرگ
۵۰۷ ۳.۳ خلاصه تجربه با ویرایش ۱.۱

۵۰۷	گردآوری با WCT یا ویرایش ۱.۲
۵۰۷	۴.۱. نرم افزار گردآوری تاریخ
۵۰۸	۲.۴. نرم افزار تورق (مروزرگر)
۵۰۸	۳.۴. نرم افزار هرس درخت تصمیم
۵۰۸	۴.۴. وب گاه های بزرگتر جمع آوری می شوند
۵۰۹	۵.۴. گسترش ذخیره دیجیتال با ارزش
۵۰۹	۶.۴. ارتباط
۵۰۹	۷.۴. فهرست گردش کار و دسترسی
۵۱۰	۵. گردآوری رویداد انتخابات هیئت محلی
۵۱۲	۶. نتیجه گیری
۵۱۲	منابع
۵۱۴	جلد ۲
۵۱۴	مشخصات کتاب
۵۱۵	اشاره
۵۱۸	فهرست مطالب
۵۲۲	سخن نخست
۵۲۶	فصل اول : مبانی مدیریت و آرشو وب
۵۲۶	اشاره
۵۲۸	آرشو اشیای داده‌ای با استفاده از فیدهای وب
۵۷۲	آرشو صفحات وب بر مبنای تحلیل دیداری و DIFF
۵۹۱	آرشو منابع ویدئویی وب
۶۱۷	استفاده از عملیات هوشمند نرم افزاری جهت ایجاد قابلیت تعامل پذیری در خدمات محتوایی و اطلاعاتی سازمانها
۶۴۴	تحلیل انسجام و مصورسازی در آرشو وب
۶۷۸	بایگانی وب پنهان
۷۰۵	بررسی تاثیر بستر نحوی بر میانگین پذیری استانداردهای فراداده ای گلمی در راستای بکارچه سازی نظامهای اطلاعاتی
۷۳۸	بهینه سازی کیفیت آرشوهای وب
۷۶۱	خرش هوشمند در برنامه های کاربردی وب
۷۸۳	دسته بندی مفهومی صفحات وب با استفاده از اتوماتای یادگیر توزیع شده
۷۹۸	DigiBoard: ابزار افزایش کارایی فعالیتهای پیچیده آرشو وب در کتابخانه کنگره
۸۲۰	رونوشت برداری از وبگاه ها

سرشناسه: منتظر، غلامعلی، 1348 -، گردآورنده

عنوان و نام پدیدآور: مدیریت منابع اطلاعاتی وب [کتاب] / به کوشش غلامعلی منتظر و فرزانه شادان پور.

مشخصات نشر: تهران: سازمان اسناد و کتابخانه ملی ایران، 1391.

مشخصات ظاهری: ج.2.

شابک: دوره 0-344-446-964-978 ؛ 150000 ریال: ج.1 1-978-343-446-964-978 ؛ ج.2 2-978-345-446-964-978 ؛

200000 ریال (ج.2، چاپ اول)

وضعیت فهرست نویسی: فایا

یادداشت: ج.2 (چاپ اول: 1391).

مندرجات: ج.1. مبانی و تجربه های جهانی .-ج.2. دیدگاه های فناورانه، اخلاقی و مدیریتی.

موضوع: وب--سایت ها--مدیریت

موضوع: منابع اطلاعاتی --مدیریت

موضوع: آرشپوسازی وب

شناسه افزوده: شادان پور، فرزانه، 1344-، گردآورنده

شناسه افزوده: سازمان اسناد و کتابخانه ملی جمهوری اسلامی ایران

رده بندی کنگره: 1391 4م8/TK5105/888

رده بندی دیویی: 005/72

شماره کتابشناسی ملی: 3077380

دسترسی و محل الکترونیکی: <http://dl.nlai.ir/UI/2fb77759-f3eb-4f7f-ad9d-2cc0b917ed1d/Catalogue.aspx>

خیراندیش دیجیتال: انجمن مددکاری امام زمان (عج) اصفهان

ویراستار کتاب: خانم شهناز محققیان

ص: 1

جلد 1

اشاره

سازمان اسناد و کتابخانه ملی

جمهوری اسلامی ایران

مدیریت منابع اطلاعاتی وب

جلد اول

مبانی و تجربه های جهانی

به کوشش:

دکتر غلامعلی منتظر

و

فرزانه شادان پور

زمستان 1391

ص: 2

فهرست نویسی پیش از انتشار کتابخانه ملی جمهوری اسلامی ایران

سرشناسه: منتظر، غلامعلی 1348 - ، گردآورنده

عنوان و نام پدیدآور: مدیریت منابع اطلاعاتی وب / به کوشش غلامعلی منتظر و فرزانه شادان پور.

مشخصات نشر: تهران: سازمان اسناد و کتابخانه ملی جمهوری اسلامی ایران 1391

مشخصات ظاهری: 2 ج.

شابک: دوره: 0 - 344 - 446 - 964 - 978؛ ج. 1: 3 - 343 - 446 - 964 - 978؛

مندرجات: ج. 1. مبانی و تجربه های جهانی - ج. 2. دیدگاه های فناورانه، اخلاقی و مدیریتی

موضوع: وب--سایت ها -- مدیریت

موضوع: منابع اطلاعاتی-- مدیریت

موضوع: وب-- آرشیو سازی

موضوع: شادان پور، فرزانه، 1344- گردآورنده

شناسه افزوده: شادان، پور فرزانه 1344 - ، گردآورنده

شناسه افزوده: سازمان اسناد و کتابخانه ملی جمهوری اسلامی ایران

رده بندی کنگره: 00 5 /72

رده بندی دیویی: TK5105/88884 1391

شماره کتابشناسی ملی: 3077380

سازمان اسناد و کتابخانه ملی جمهوری اسلامی ایران

عنوان: مدیریت منابع اطلاعاتی وب، جلد اول: مبانی و تجربه های جهانی

به کوشش: دکتر غلامعلی منتظر (دانشیار دانشگاه تربیت مدرس) و فرزانه شادان پور (مربی، سازمان اسناد و کتابخانه ملی جمهوری

اسلامی ایران)

ویراستار ادبی: آرزو تجلی (کارشناس ارشد جامعه شناسی، سازمان اسناد کتابخانه ملی جمهوری اسلامی ایران)

تنظیم و تصحیح: مهشید برجیان (کارشناس ارشد کتابداری و اطلاع رسانی، سازمان اسناد و کتابخانه ملی جمهوری اسلامی ایران)

ویراستار استنادی (مقالات تألیفی): فروزان رضایی نیا کارشناس کتابداری و اطلاع رسانی، سازمان اسناد و کتابخانه ملی جمهوری اسلامی ایران)

نمونه خوانی و اصلاحات: مهشید برجیان، فاطمه رمضانپور، آمنه هزارخانی، زهرا زاهدی، محمد رضا میقانی، ملیحه حاجی زاده مقدم

طراحی جلد و صفحه آرایی: شهره خوری

ناظر فنی چاپ: نصرت الله امیرآبادی

ناشر: سازمان اسناد و کتابخانه ملی جمهوری اسلامی ایران

شمارگان: 500 نسخه

بها: 15000 تومان

نشانی: تهران بزرگراه شهید حقانی (غرب به شرق)،

بعد از ایستگاه مترو بلوار کتابخانه ملی

تلفن فروشگاه 81623318-81623315-88941946 دورنگار: 88947496

وب سایت: www.nlai.ir

پست الکترونیک انتشارات: Publication@nlai.ir

ص: 3

سازمان استاد و کتابخانه ملی جمهوری اسلامی ایران

مدیریت منابع اطلاعاتی وب

جلد اول

مبانی و تجربه های جهانی

ص: 4

سخن نخست ... نه

به جای مقدمه ... 1

فصل اول: مبانی مدیریت و آرشیو وب ... 3

بایگانی شبکه وب : مباحث و روش ها: نوشته ژولین ماسانه / ترجمه فهیمه باب الحوائجی ... 4

از آرشیو اینترنت تا آرشیو در اینترنت: نوشته میشل کیمتون / ترجمه مرضیه هدایت ... 62

کاربرد وب و مطالعات مربوط به آن: نوشته استیو جونز گمیل جانسون / ترجمه سید مهدی طاهری، سید محمد موسوی ... 74

خصوصیات وب ایران: نوشته مریم پیروزمند ... 90

آینده آرشیو وب: نوشته اریک تی. مه یر، آرتور توماس، رالف شرودر، مؤسسه اینترنت آکسفورد / ترجمه رضا خانیپور، محبوبه قربانی ...

106

فصل دوم: تجارب جهانی و مسائل بومی در آرشیو سازی وب ... 149

آرشیو وب در دنیای وب 2/0 شعبه آرشیو وب و حفاظت رقومی کتابخانه ملی استرالیا: نوشته ادگار کروک / ترجمه مرجان هادیزاده ...

150

ص: 5

آسیب شناسی زبان و خط فارسی در بازیابی اطلاعات: نگاهی به موتورهای کاوش و پایگاه های برخط: نوشته شعله، ارسطو پور فاطمه
احمدی نسب ... 158

ارزیابی کاربردپذیری وبگاه نهاد کتابخانه های عمومی کشور نوشته صدیقه محمد اسماعیل، ماهرخ ناصحی اسکویی ... 176

امکان سنجی پردازش وبگاه ها در سازمان اسناد و کتابخانه ملی ایران: نوشته رضا خانی پور، محبوبه قربانی، سهیلا فعال ... 190

ایجاد آرشیو گزینشی منابع تحت وب، بررسی هزینه های مربوط به فراهم آوری منابع تحت وب در کتابخانه ملی استرالیا: نوشته مارگارت
فیلیس / ترجمه صدیقه محمد اسماعیل ... 204

بایگانی وب علمی در مقیاس کوچک DACHS نوشته هانو لشر / ترجمه حمزه علی نورمحمدی ... 218

بررسی و مقایسه قابلیت های یونی مارک و مارک 12 برای سازماندهی منابع اطلاعاتی وب نوشته رقیه حجازی مهرداد کوبی ... 230

سنجش رابط کاربر پایگاه های اطلاعاتی پیوسته مجلات تمام متن فارسی نوشته صدیقه جعفرزاده، معصومه پیروزفر، عبدالحسین فرج
پهلوی ... 242

قانون و اسپاری وب فرانسه: راهبردهایی برای گردآوری دامنه ملی نوشته فرانس لاس، فارگوس کلمنت کیوری برت وندلاند / ترجمه سودابه
نوذری ... 254

معرفی آرشیوهای وب به عنوان یک خدمت جدید کتابخانه: تجربه کتابخانه ملی فرانسه نوشته سارا اویری / ترجمه زهرا تهوری ... 286

یک سال آرشیو وب گزینشی با WCT در کتابخانه ملی زلاندنو نوشته گوردون پنیتر، سوزانا جو، وائیتا لا لا، گیلیان لی / ترجمه احترام
السادات کیانمهر ... 304

از ویژگی های قرون گذشته بی خبری بود و تمایز جدی عصر جدید نسبت به گذشته دسترسی آسان به اطلاعات است. بشر با از سر گذراندن سه موج و پارادایم، کشاورزی صنعت و اطلاعات امروز در قرن بیست و یکم پا در عصر انفجار اطلاعات نهاده است این امر فی نفسه نه مطلوب است نه مذموم، بلکه به نحوه مدیریت ما نسبت به اطلاعات باز می گردد.

بشر امروزی به دلیل رشد روزافزون علم و فناوری در شرایط هشدارآمیز عدم قطعیت بسر می برد و همین مدیریت و تصمیم گیری را با چالش جدی روبرو ساخته است. اگر اطلاعات درست مدیریت شود و در تصمیم گیری ها به موقع به کار آید، و از دو ویژگی صحت و سرعت برخوردار باشد، می تواند منشأ تصمیم های تحول آفرین شود. ویژگی دیگر این عصر ظهور و حضور همه جانبه اطلاعات دیجیتال است. دورانی فرا رسیده است که در آن بناست دانش مدون و تفکر مضبوط بشر علاوه بر کاغذ، و حتی بیش از آن، بر محمول «بیت» ها مسیر تولید، نشر و اشاعه، و مصرف را پیماید. هم اطلاعات تولید شده تحت وب و هم میزان استفاده از این اطلاعات با سرعت فزاینده ای رو به رشد است. کشور ما بنابر اطلاعات وثیق از حیث تعداد کاربران و میزان حضور و فعالیت آن ها در وب جایگاه نخست را در منطقه خاور میانه داراست. این روند رو به رشد، با نصب العین قرار دادن آرمان های بلند انقلاب اسلامی در ترویج تفکر رهایی بخش اسلام ولایت مدار، وظیفه خطیری بر دوش نهادها و دستگاه های مسئول تولید، سیاستگذاری و نشر محتوا در محیط وب قرار می دهد و آن انجام بررسی های علمی و مستند به منظور ابتنای سیاستگذاری ها و عملکردها بر مبانی صحیح و کارآمد و متناسب با نیازهای گوناگون کاربران در این محیط است. اما وجه دیگر، صیانت از این محتوا و انتقال آن به نسل های آینده است که با توجه به ناپایداری محتوای قرار گرفته بر اینترنت و فناوری پیشرفته ای که برای چنین امر خطیری لازم است از اهمیت مضاعفی برخوردار می شود.

سازمان اسناد و کتابخانه ملی جمهوری اسلامی ایران بنابر مأموریت خویش دایر بر صیانت از میراث فکری کشور و اشاعه آن، عزم راسخ داشته است که برای مدیریت منابع اطلاعاتی مهم و رو به رشد وبی نیز چاره اندیشی نماید؛ بنابر این در سال 1389 نخستین بار در کشور به تهیه ساز و کار لازم برای ایجاد آرشیو ملی وب همت گماشته است.

از دیگر سو سازمان با علم به این که مدیریت در این حوزه مشارکت همه صاحبان اندیشه در حوزه تولید، سازماندهی و اشاعه اطلاعات تحت وب را می طلبد، مصمم شد «نخستین کنفرانس ملی مدیریت منابع اطلاعاتی وب» را برگزار نماید تا اهل علم و فناوری در این مجمع با هم اندیشی و تضارب آراء همچون

گذشته این سازمان را یار و یاور باشند.

این اثر مجموعه ای است فراهم آمده از تلاش پژوهشگرانی که با وجود نبودن مباحث مطرح شده در محورهای موضوعی کنفرانس، به ارائه ثمره پژوهش های خود همت نمودند؛ که با برگزیده ای از مقالات ترجمه ای در این عرصه پژوهشی ادغام و به طبع رسیده است رجاء واثق دارم که با الطاف الهی از این پس مدیریت منابع اطلاعاتی وب، و آرشیو وب به طور خاص، موضوع پژوهش و ابتکار عمل اهل دانش و فناوری در کشورمان قرار خواهد گرفت و در این عرصه نیز فرزندان این مرز و بوم تجسم گفتار نغز رسول اعظم صلی الله علیه و آله خواهند بود که «علم اگر تا ثریا، برود مردانی از فارس بدان دست خواهند یافت».

اسحق صلاحی

رئیس کنفرانس و رئیس سازمان اسناد و کتابخانه ملی جمهوری اسلامی ایران

ص: 8

گسترش روزافزون اطلاعات در شبکه اینترنت و سادگی بارگذاری انواع داده‌ها بر وب جهان را با شکل جدیدی از تولید، انتشار و مصرف اطلاعات مواجه کرده است. تغییر جایگاه شهروندان جامعه از مصرف‌کننده صرف اطلاعات به مولد و ناشر اطلاعات و فارغ از ساز و کارهای موسوم، سبب ساز روابطی جدید در عرصه ارتباطات اجتماعی و فرهنگ شده است. از سویی حجم رو به تزاید داده‌ها و چرخه عمر کوتاه اطلاعات موجود در وب، موجب شده که «گردآوری»، «پالایش»، «سازماندهی»، «ذخیره سازی» و «اشاعه» آن‌ها در زمره مسائل پژوهشی در نهادهای علمی و نیز بخش‌های پژوهش و نوآوری شرکت‌ها قرار گیرد؛ ضمن اینکه حفظ و دسترسی پایدار به اطلاعات موجود در وب، که خود جزئی از میراث فکری ملت‌ها محسوب می‌شود، به دغدغه‌ای جدی برای سازمان‌ها متولی حفظ و اشاعه میراث فکری به ویژه کتابخانه‌های ملی، بدل شده است.

این حوزه در جهان موضوعی نسبتاً جدید است و پیشینه آن به کمتر از پانزده سال می‌رسد، لیکن با سرعتی شتابان در حال رشد است و محققان مختلفی را از زوایای مختلف فنی، حقوقی، اقتصادی و حتی اخلاقی به سوی خود جذب کرده که گواه آن نیز طیف وسیعی از مقاله‌ها کتاب‌ها و گزارش‌های سازمانی است که در طی چند سال اخیر در سطح جهانی منتشر شده است. به رغم این نکات، در ایران همچنان این زمینه، حوزه‌ای بکر و کمتر مورد توجه محسوب می‌شود و در طی سالهای اخیر کمتر تحقیق بدان پرداخته، لیکن تزاید اطلاعات فارسی بر روی وب و برنامه‌های ملی کشور مبنی بر توسعه کاربری‌های مختلف بر شبکه‌های اطلاعاتی (از جمله توسعه دولت الکترونیکی، یادگیری الکترونیکی و کتابخانه‌های دیجیتالی) لزوم توجه به این موضوع را بیش از پیش نمایان می‌سازد. به همین دلیل سازمان اسناد و کتابخانه ملی جمهوری اسلامی همزمان با برگزاری «نخستین کنفرانس ملی مدیریت منابع اطلاعاتی وب» در صدد برآمد تا این حوزه را هرچه بیشتر به متخصصان و پژوهشگران بازناساند. کتاب پیش رو حاصل همین نیت متولیان این موضوع مهم است.

این کتاب مجموعه‌ای قریب به 30 مقاله برگزیده از مهم‌ترین منابع علمی منتشر شده در جهان و نیز قریب به 15 مقاله برگزیده از صاحب نظران ایران است که در قالب دو جلد تقدیم

حضور خوانندگان ارجمند می شود. این مقالات در چهار موضوع اصلی به شرح زیر تقسیم شده اند:

• مبانی مدیریت و آرشیو وب

• تجارب جهانی و مسائل بومی در مدیریت و آرشیو وب

• مسائل فناوریانه

• مسائل اخلاقی و مدیریتی

بی گمان این مجموعه می توانست به افزودنی های دیگر (هم از منابع خارجی و هم از دیدگاه سایر متخصصان ایرانی) به اثری پربارتر بدل گردد لیک نخستین گامی است که در این حوزه برداشته شده و مطمئناً در مراحل بعدی با همت سایر اندیشمندان، ویراست هایی غنی تر از آن حاصل خواهد آمد. نگارنده امیدوار است این مجموعه به مثابه بذری باشد که در کشتزار ذهن پژوهشگران کاشته شده و ان شاء الله در آینده ای نه چندان دور به نهالی پر طراوت در عرصه علم و عمل در جامعه اسلامی مان مبدل گردد.

در پیدایی این اثر کسان بسیاری همراهی و همکاری داشته اند که مقدم بر همه اندیشمندانی است که متن هر مقاله به خامه دانش افزای آنان امکان وجود یافته است. از این رو نگارنده سپاس فروتنانه خود را نثار نگارندگان و مترجمان ارجمند این اثر می نماید. گردآوری، تنظیم و آماده سازی مطالب کتاب به همت خانم ها فرزانه شادان پور و مهشید برجیان بوده و ویراستاری آن را خانم آرزو تجلی بر عهده داشته اند. ویراستار استنادی مقالات تألیفی را سرکار خانم فروزان رضایی نیا به انجام رسانده اند و سرکار خانم دکتر میترا صمیمی زحمت چکیده نویسی شماری از مقالات را که فاقد چکیده بودند متقبل شدند. نمونه خوانی و اصلاحات اثر حاصل تلاش خانم ها مهشید برجیان فاطمه رمضانپور آهنگری، آمنه هزار خوانی، زهرا زاهدی، ملیحه حاجی زاده مقدم و آقای محمد رضا میقانی بوده است. ضمن اینکه زیبایی متن و صفحه آرایی آن مدیون حسن سلیقه سرکار خانم شهره خوری است. زحمات لیتوگرافی، چاپ و صحافی کتاب نیز بر عهده جناب آقای امیر آبادی بوده که بر خود فرض می داند از همه این بزرگواران صمیمانه تشکر کند.

بی گمان پدید آمدن این اثر به همت مسؤولان گرانمایه سازمان اسناد و کتابخانه ملی جمهوری اسلامی بوده است و نگارنده امیدوار است خداوند آنان را در مسیر خدمت به فرهنگ و دانش ایران اسلامی مورد تأیید قرار دهد.

اللهم وفقنا لما تحب وترضی

غلامعلی منتظر

تهران- بهمن ماه یک هزار و سیصد و نود و یک خورشیدی

ص: 2

فصل اول: مبانی مدیریت و آرشيو وب

اشاره

ص: 3

بسیاری از جنبه های اجتماعی، اتفاقی هستند یا به طور کلی درباره اینترنت و به ویژه درباره وب بازتاب یافته اند. محافظت از وب، به این دلیل ضرورتی فرهنگی و تاریخی است. اما وب، از نظام های انتشاراتی قبل نیز برای ضرورت یک بازبینی ریشه ای از عملکردهای حفاظتی مرسوم، متفاوت می باشد. مفهوم میراث جمعی مشترک، شامل هر مصنوع محصول انسانی می شود از بناهای تاریخی معماری گرفته تا کتاب های جدید قرن بیستم ولو اینکه با فعالیت هایی حفاظتی مرتبط باشند (مثل آن ها که به طور اصولی و اختیاری سازماندهی شده اند) و قبلاً نمایان شده اند. در عصر حاضر علت بایگانی، دلیلی است که بر جزئیات تأکید میکند، به گونه ای پایدار ما را به اجتناب از عمومیت گرایی هدایت میکند و به طور خاص و منحصر به فرد که حوادث و رویدادها را مورد بررسی قرار می دهد. در حقیقت، امکاناتی که وب برای انتشار ایجاد می کند، منبعی منحصر به فرد از محتوا را ارائه می دهد که استدلال بایگانی وب، معقولانه انجام و تصدیق شده است. با توسعه شبکه های وب، به طور چشمگیری، حجم آن چه که می تواند منتشر شود و همچنین تعداد «ناشران» بالقوه یا خالقان محتوا با تقلیل هزینه های انتشار به تقریباً هیچ افزایش یافته است. دلگرم کننده است که ببینیم که بسیاری از مؤسسه های (حفظ) میراث، در بایگانی وب در حال به کارگیری هستند. بررسی اخیر توسط گروه پژوهش کتابخانه (RLG 2006) نشان داد که 60 درصد اعضای مورد بررسی اشان، بایگانی وب را قسمتی از مأموریت خود پنداشته اند.

* بایگانی شبکه وب: مباحث و روش ها (1)

ژولین ماسانه (2) ترجمه: فهیمه باب الحوائجی (3)

1- مقدمه

محصولات فرهنگی گذشته همیشه نقش مهمی در اطلاع رسانی و خود ادراکی جامعه و ساختن آینده آن ایفا می نمایند. شبکه جهان گستر وب (به طور خلاصه وب) رسانه ای جامع و گذراست در جایی که با درکی عمیق، فرهنگ مدرن در جهت یافتن شکل طبیعی عبارات و اصطلاحات است. انتشار، مباحثه، ایجاد، کار و تبادل اجتماعی در یک درک عمیق: بسیاری از جنبه های اجتماعی، اتفاقی هستند یا به طور کلی درباره اینترنت و به ویژه درباره وب بازتاب یافته اند. محافظت از وب، به این دلیل ضرورتی فرهنگی و تاریخی است. اما وب، از نظام های انتشاراتی قبل نیز برای ضرورت یک بازبینی ریشه ای از عملکردهای حفاظتی مرسوم، متفاوت می باشد.

این فصل، بررسی موضوع های برخواسته از حفاظت از وب را ارائه می دهد و روش هایی که تا این تاریخ برای غلبه یافتن بر آن ها توسعه یافته است. ابتدا استدلال هایی را در برابر ضرورت و احتمالات بایگانی وب مطرح می کنیم. سپس، سعی می کنیم تفاوت های بسیار برجسته ای که وب را از دیگر

ص: 5

Web Archiving: Issues and Methods: in Masanes, Julien (ed.), Web Archiving. Berlin Heidelberg New - 1
York: Springer.pp.1-46

Julien Masané s -2

3- دانشیار کتابداری و اطلاع رسانی و دانشیار گروه کتابداری و اطلاع رسانی دانشگاه آزاد اسلامی واحد علوم و تحقیقات تهران

محصولات فرهنگی متمایز می کند، ارائه دهیم و دلالت های شان را برای محافظت ترسیم می نماییم.

این اساس وب را در بر می گیرد و آن را به عنوان نوعی سیستم نشر فعال و یک ابررسانه ویرایش اجتماعی و با یک محصول فرهنگی سراسری مورد بررسی قرار می گیرد. احتمالات و محدودیت های حفاظت هر یک از این جنبه های وب را مطرح می کنیم. سپس، رویکردهای روش های بررسی مهم را برای اکتساب سازماندهی و ذخیره سازی محتویات وب ارائه می کنیم.

فصل 2 و 4 و 5، جزئیات بیشتری درباره روش های بررسی و ابزارهایی برای اکتساب مضامین تهیه کند و فصل های 6-8 بر روی دسترسی استخراج اطلاعات و حفاظت از محتوای وب تمرکز کرده است. دو فصل آخر این کتاب بررسی هایی موردی را ارائه می کند. بایگانی اینترنتی که بزرگ ترین بایگانی وب در جهان است (فصل 9) و DACHS یک پژوهش که در جهت گزینش بایگانی وب (فصل 10) می باشد. این فصل مقدمه ای کلی برای این کتاب مورد بررسی قرار گیرد:

در نهایت، موارد اولیه را در این حوزه فراهم می کند و طبقه بندی بایگانی های وب را برای ترسیم وضعیت جاری حفاظت وب را پیشنهاد می نماید.

- میراث جامعه و وب

2-محافظت از میراث

مفهوم میراث جمعی مشترک، شامل هر مصنوع محصول انسانی می شود، از بناهای تاریخی معماری گرفته تا کتاب های جدید قرن بیستم، ولو اینکه با فعالیت هایی حفاظتی مرتبط باشند (مثل آن ها که به طور اصولی و اختیاری سازماندهی شده اند) و قبلاً نمایان شده اند. شکل، اهداف، و کارایی حفاظت از میراث، به طور قابل توجهی با زمان و رسانه های بررسی شده متفاوت اند و تلاش این مقاله برای خلاصه کردن این تکامل کافی نیست. اجازه دهید فقط به یاد آوریم که از آماده سازی روشنفکری مذهبی (طبق کتابخانه کاسیودوروس ویواریوم، ریچه 1998) (1) برای ایجاد مجموعه ای به عنوان نشانه های قدرت (ابتکارات موزه مدرن توسط مدیسیس (2) در فلورانس در اواخر قرن پانزدهم را ببینید) برای کنترل وضعیت نظام مند حفاظت از فرهنگ ملی (ابتکار سپرده قانونی فرانسیلر (3) را ببینید) استفاده شده است و از انگیزه های گوناگونی در جهت جمع آوری نظام مند و حفاظت از محصولات فرهنگی در تاریخ، ناشی شده است.

در عصر حاضر، بایگانی ها به طور کلی، تمایل زیادی به فراگیر شدن دارند (آسبورن 1999) (4). همان طور که مایک فیدراستون (5) اظهار می کند:

علت بایگانی، دلیلی است که بر جزئیات تأکید می کند، به گونه ای پایدار ما را به اجتناب از عمومیت گرایی هدایت می کند و به طور خاص و منحصر به فرد که حوادث و رویدادها را مورد بررسی قرار می دهد. تمرکز بر بایگانی مؤثر کالاها به طور فزاینده ای تغییر یافته است به سوی جزئیاتی با واقعیتهای

ص: 6

Medicis -2

Franceisler -3

osborn -4

Mike Featherstone -5

عظیم از زندگی روزمره دنیوی تمرکز یافته است (فیدرستون 2000) (1).

در حقیقت، امکاناتی که وب برای انتشار ایجاد می کند، منبعی منحصر به فرد از محتوا را ارائه می دهد که استدلال بایگانی تمایل به ستایش آن دارد. از این رو، می توانیم فرض کنیم که قانونی بودن بایگانی وب، معقولانه انجام و تصدیق شده است. با وجود این، محافظت از وب، مورد سؤال قرار گرفته و تاکنون مورد قبول همه واقع نشده است. استدلال ها در برابر بایگانی وب می تواند در سه مقوله دسته بندی شود. آن ها که مبتنی بر محتوای یافت شده در وب هستند، آن ها که تصور می کنند که وب خود - محافظ است، و آن ها که فرض می کنند بایگانی وب امکان پذیر نمی باشد.

2-1- به اندازه کافی خوب نیست؟

نخستین مقوله، بحث هایی را درباره کیفیت محتوای وب در بر دارد که گفته می شود فرض بر عدم مطابقت معیارهای مورد نیاز برای محافظت می باشد. این موقعیت، مدت زیادی توسط برخی متخصصان نشر جهانی (ناشران و کتابداران) حفظ شده است و در جهت تهدید بزرگی که توسط این رسانه جدید برای بقا اعمال می شود همراه می شود. معمولاً با تأکیدی درباره میزان گسترده اطلاعات وب و فقدان دانش درباره روش های بایگانی وب و هزینه های آن همراه شده است.

منافع این موقعیت، از انتقال سیستم نشر پیوسته و آن هایی که بر حفاظت و بازدهی صنعت نشر پیوسته ادامه می دهند آگاهی دارد. اما آن ها از مرزهای آن چه که محافظت شده است را اجتناب می ورزند همان اندازه که شبکه وب، محدودیت آن چه را که انتشار یافته است را گسترش می دهد. معادلات اقتصادی تولیدات فیزیکی حامل دانش (ادواری ها، کتاب، و مانند آن) میراث انقلاب گوتنبرگ هستند که باید طبق این دیدگاه بر استقرار محدودیت ها برای کاری که باید حفاظت شود ادامه دهند. حتی زمانی که این معادلات عمیقاً تعدیل شده باشند. از نظر تاریخی، این حقیقت که چه چیزی می تواند منتشر شود، به واسطه هزینه های فیزیکی محدود می شود (که شامل تولید و نقل و انتقال، ذخیره سازی، و هزینه های اداره است) و فیلترسازی را ایجاد می کند، برای آن که نظام های نشر بیشتر از پنج قرن است که این کار را انجام می دهند. اما این مورد مناسبی نیست و نسبتاً میراث تعادل ثابتی از قرن پانزدهم است که نقض شده است. توسعه شبکه های وب، به طور چشمگیری، حجم آن چه که می تواند منتشر شود و همچنین تعداد «ناشران» بالقوه یا خالقان محتوا با تقلیل هزینه های انتشار به تقریباً هیچ افزایش یافته است. مباحث درباره ارزیابی کیفیت، به ناچار ذهنی، در واقع با پنهان نمودن مذاکرات واقعی درباره گسترش جو انتشار می یابد.

اگر چه رشد ادواری ها در پایان قرن نوزدهم، قابل مقایسه با وسعت آن در انقلاب جاری نمی باشد، برخی خصوصیات (مثل در همکرد نوع نشر با وضعیت جسمانی و فکری) را به اشتراک گذاشته و عکس العمل های مشابهی را به وجود آورده است. گاهی اوقات برای جامعه کتابخانه برای مثال برای پذیرش این نوع انتشار در قفسه هایشان و همچنین در قلب شان، در بر گرفته شده است. همان طور که

ص: 7

فایت - شایب (2000) (1) برای موردی در فرانسه نشان داده است، رفتار توصیفی خاص که در سطح عنوان به آن نیاز داشت، توسط این جامعه مورد اغماض قرار گرفت و بخش کاملاً جدید مدیریت اطلاعات در کنار کتابخانه هارا، به همین دلیل، به وجود آورد (مستندسازی، نمایه سازی، نوشته های علمی). مباحث بر روی بایگانی وب، برخی شباهت هایی را بر حسب این رویدادهای فرعی، به اشتراک می گذارد که اگر در همان حالت باقی بماند، دیده خواهد شد.

فیلترسازی، اگر چه به مدت طولانی نیازی به تخصیص تولیدات فیزیکی منابع کار آمد دانش را ندارد، به طور کامل از بین نرفته است. بلکه از نقشی مرکزی به یک نقشی پیرامونی تغییر می یابد و باز هم در برخی فضاها مورد نیاز است (برای مثال، اعتبار سنجی علمی) و به شکل های جدید تجربه می شود (مثل ویکی پدیا، اسلش دات، بوگس فر (2))

چنان چه اکسل برونز (3) توضیح می دهد:

عکس العمل ناگهانی این فعل و انفعالات و وسایل ارتباط جمعی مشارکتی بسیار محسوس است. اگر کسی بتواند یا حداقل توانایی بالقوه آن را، [مثل یک] یک ناشر دارد، چه اثراتی بر روی مؤسسه های انتشاراتی موجود خواهد داشت؟ اگر اطلاعات موجود در وب بتواند با اطلاعات وسیع متنوع به آسانی ارتباط (4) داشته باشند، چه اثراتی بر روی چارچوب های (5) انتشاراتی سنتی می تواند داشته باشند؟ اگر توانایی بالقوه برای مخاطبان وب برای مشارکت در تولید و ارزیابی تعاملی محتوا وجود داشته باشد، برای ایجاد نقش های تولید کننده و مصرف کننده در وسایل ارتباط جمعی چه اتفاقی می افتد؟ (برانز، 2005).

در زمینه حفاظت این مسئله به طور جدی مورد توجه قرار گرفته است. یک موضوع قطعی است و آن مدینه فاضله ای است که امیدوار باشیم که تعداد کمی از کتابداران فیلتر سازی ناشران را در مقیاس وب جهانی جایگزین نمایند حتی اگر آن ها سنت مدیدی در انتخاب محتوا داشته باشند، این کار را در محیطی ساختار یافته تر انجام می دهند که چندین مرتبه در ابعاد کوچک تر داشته باشد. اگر چه این کار هنوز احتمالی است و برای جامعه ای که خوب تعریف شده و اهداف کوچک مفید است (فصل درباره انتخاب روش شناختی ها و فصل 10 را درباره داجز (6) پژوهشی برگرفته درباره بایگانی وب، پژوهش انجام شده و هم چنین بروگر (7) 2005) را ببینید) به کارگیری این موارد به عنوان یک ساز و کار جهانی برای بایگانی، وب واقع گرایانه نیست اما این حقیقت که گزینش دستی محتوا برای وسعت وب مقیاس گذاری نشده است دلیلی برای نپذیرفتن بایگانی وب نیست این فقط دلیل خوبی برای بررسی مجدد موضوع گزینش و کیفیت در این محیط است.

آیا می تواند بر اساس یک ارزیابی کیفی توزیعی در سطح بالا و جامع باشد؟ این ارزیابی، به طور

ص: 8

Fayet-schibe -1

impact bogosphere, Slashdo -2

Axel Bruns -3

link -4

format -5

DACHS -6

ضمنی در دو سطح ساخته می شود:

کاربران وب به وسیله دسترسی بر محتوا، خالقان (محتوا) با پیوند دادن شکل محتوای صفحه های شان (ما در اینجا به قضاوتی که توسط خود خالقان قبل از اینکه محتوای شان را به صورت برخط در وب بگذارند در نظر نمی گیریم، که اگر به عنوان یک معیار انتخاب استفاده شود، به معنای بایگانی هر چیز است) همچنین، می تواند به طور ضمنی توسط افزایش انتخاب کنندگان فعال ایجاد گردد.

اجازه دهید ابتدا دسترسی کاربران را بررسی کنیم. گسترش جو انتشار پیوسته تحت چیزی که ظرفیت اقتصادی برای چاپ فیزیکی مجاز دانسته است، نتایج دیگری را در بر می گیرد: افت مکانیکی در میانگین تعداد خوانندگان هر واحد از محتوای منتشر شده برخی صفحه ها حتی نه تنها توسط هیچ انسانی خوانده نمی شوند بلکه توسط هیچ رباتی هم نمایه نمی شود. بوفخاد و ویونات (2003) (1)، استفاده از سیاهه های مربوط و پرونده های سرور یک وبگاه بزرگ دانشگاهی را نشان داده اند که 5 درصد صفحه ها فقط توسط ربات ها قابل دسترسی بودند و به 25 درصد آن ها هرگز دستیابی نداشتند، بدان معنا که توزیع دسترسی به محتوای پیوسته بسیار طولانی مدت ارائه خواهد شد.

اما این تکامل در نشر مدرن، کاملاً جدید نیست. رشد و درجه بالایی از تخصصی شدن انتشارات ادواری تقریباً الگوی مشابه دسترسی را نشان می دهد. آیا این استدلالی برای عدم حفظ ادواری هاست؟ در بیش تر کشور ها، نظام های سپرده قانونی، به طور مستقل، از انتشارات آن چه را که مورد استفاده قرار می گیرند، را نگهداری می کنند. این بی تکلیفی علائق خوانندگان آینده پیش بینی می کند.

مطمئناً برای حفاظت کنندگان در ارزیابی مفید بودن محتوای پیوسته برای نمایش و تلاش برای پیش بینی برای آینده، تا زمانی که برای جوامع کاربران از پیش تعریف شده باشد، امکان پذیر است.

الگوهای دسترسی همچنین می توانند برای اداره نظام های بایگانی جهانی استفاده شوند: در مورد بایگانی وب اصلی، تاکنون، مجموعه بایگانی اینترنت توسط الکسا (2) اهدا شده است که از الگوی دستیابی برای تعیین عمق خزش برای هر سایت استفاده کرده است (فصل 9 کیمپتون) (3) و دیگران (2006) را ببینید). همچنین، می تواند به وسیله پرس و جو های فرستاده شده برای موتور جست و جو اجرا گردد (پانندی و اولتسون، 2005) (4). اما سؤال کلیدی برای بایگانی های وب این است: چگونه به این اطلاعات دست یابیم و کدام مرز مورد استفاده قرار می گیرد؟

ترافیک اطلاعات، معمولاً وجود ندارد و موتورهای جست و جو از ابتکار الکسا پیروی می کنند و آن را از میلیون ها نوار ابزار نصب شده در مرورگرهایی که از اطلاعات شناوری (5) کاربران به آن ها می گذرند به دست می آورد.

مؤسسه های بایگانی از کجا می توانند به آن [اطلاعات] دسترسی پیدا کنند چنان چه آن ها خودشان جست و جوی ویژه ای را ارائه ندهند؟ مرزهای آن چه باید باشند؟ آیا باید در صفحه یا در سطح سایت

ص: 9

Kimpton et al -3

Pandey and Olston -4

navigation -5

به کار برده شود؟ (الکسا در سطح سایت از آن استفاده می کند) آیا عمق خزش را فقط در سطح نخست هر سایت محدود می کند (که این بدان معناست که حداقل در سطح اول هر سایت در تمام موارد اشغال خواهد بود)؟

حتی اگر این معیار، مباحث اجرایی عملی زیادی را نمودار سازد، مزیت بردن به عنوان حمل کننده برای تمرکز بر بایگانی، ستانده های میلیون ها کاربر - نه اجتماعی کوچک - را دارد که به خوبی با مدل انتشار جامعه وب تطبیق یافته است.

معیار دیگر، سطح اهمیتی است که توسط درجه پیوند درونی یک صفحه (یا یک سایت) اندازه گیری می شود. این مسئله استدلال شده است (میسائز، 2002) (1) که این تعادل مناسب در محیطی فرامتن، با درجه ای از عمومیت است که ویژگی های انتشارات سنتی را مشخص می کند و عملاً مزیتش در قابلیت استفاده برای استخراج ماتریس های پیوندی وب است (پیچ) (2) و همکاران،

1998؛ ایتبول (3) و همکاران، 2003، 2002؛ پاستور - ساتوراس و و اسپگنانی (4) (2004). روش دیگری در انبوه کردن ارزیابی کیفی ایجاد شده است، البته نه توسط کاربران، بلکه به وسیله ایجاد کنندگان صفحه (و پیوندها) این مدل ارزیابی کیفی توزیع شده به خوبی با طبیعت انتشار بر روی اینترنت و به طور عملی با امکان اجرا مطابقت دارد.

سرانجام، ممکن است توسط مشارکت های بیشتر، وظیفه گزینش مطالب برای بایگانی کردن مقیاس را بالا برد. این کار می تواند با در بر گیری مؤسسه های بیشتری در انجام و تسهیل توسط ایجاد خدمات بایگانی انجام شود که با بخش فنی مسئله سروکار دارد. این کار به وسیله بایگانی خدمات در بایگانی اینترنتی پیشنهاد شده است که در سال 2006 شروع شده و عبارت است از توانمند سازی تنظیم و مدیریت آسان برای کتابخانه و آرشیو که نمی تواند در زیر ساختهای عملکردی مورد نیاز برای بایگانی وب سرمایه گذاری کند.

تحول ممکن دیگر با عمومیت بخشیدن به این توانایی برای هر کاربر وب در مشارکت فعال می باشد؛ البته چنانچه بخواهند در بایگانی وب شرکت کنند. انگیزه اصلی کاربران در این مورد، سازماندهی حافظه وب شخصی برای امکان بازگشت به مرجع بعدی برای محتوای با ثبات و، کاوش و سازماندهی آن به عنوان راهی برای مبارزه با علائم «گم شدن در فضای مجازی» است چندین بررسی در مورد کاربران نشان داده است که حفظ نشانه های محتوای مشاهده شده برای بسیاری از کاربران ضروری است (تیوان 2004) (5) البته آن ها از روش های مناسبی نیز استفاده می کنند (6) (جونز 2003، 2001). در بایگانی وب شخصی، دنبال کردن پیشینه کاربر بر روی وب می تواند برای یک سازمان دهی شخصی و محوریت زمانی

ص: 10

Masanes -1

Page -2

Abiteboul -3

Pastor-Satorras and Vespignani -4

Teevan -5

Jones -6

در حافظه وب میسر شود (رکیموتو 1999) (1) دامیس و همکاران (2003) (2)، رینگر و همکاران، (2003) (3). خدمات پیوسته متعددی (Furl, My Yahoo) قبلاً در بایگانی وب شخصی در سطح صفحه ارائه شده که با قابلیت علامت زدن ترکیب شده است. خدمات بایگانی های هانزو (4)، حوزه گسترش یافته (مضمون، تمام سایت) و همچنین در هم و بر همی قابلیت های بایگانی، با ابزارها و خدمات دیگر (مثل بلاگ ها، مرورگرها، و مانند آن) از طریق باز کردن ای پی آی (5) را اجازه می دهد که آن با یک بایگانی سرویس گیرنده با قابلیت های P2P توسعه بیشتری خواهد یافت؛ و به طور چشمگیری امکانات برای کاربران در ثبت تجربیات شان از وب را به عنوان بخشی از زندگی رقومی توسعه خواهد داد (فریمن و گلنتر 1996) (6)، گمل و همکاران، (2002). در مورد استفاده بالقوه از مخازن کاربران در یک بایگانی وب نظیر به نظیر نیز می توانید مانتراتزیس آرگون (2004) (7) را ببینید.

دیده شده است اگر این گسترش و دموکراسی شدن نقش بایگانی بتواند مانند تفسیر و سازماندهی اطلاعات توسعه یابد موجب پیشرفت علامت گذاری (8) گولدر و هامبرمن 2005 (9) و نظام های بلاگ کردن می شود (هالاویس 2004) (10)، ؛ برونز 2005 (11). که در این صورت کمک ارزشمندی خواهد بود و برای محافظت سازمان هایی که می توانند نظارت طولانی مدت بر این محتواها داشته باشند، درونداد [خوبی] خواهد بود.

همان طور که دیدیم استدلال ها در برابر بایگانی وب بر اساس کیفیتی است که درباره فرضیاتی بر پا شده است مانند (1) کیفیت محتوا تحت فضای سنتی محتوای ویرایش شده به طور مرسوم، کافی و مناسب نیست؛ و (2) فقط گزینش دستی و تک به تک که توسط حفاظت کنندگان ایجاد شده، می تواند جایگزین فقدان فیلترسازی ناشران شود (روشی که نمی تواند درست به اندازه مقیاس وب باشد، در حالی که همه با (فیلپز 2005) (12)، موافق اند. این دو استدلال، فقدان درک اساس توزیعی وب را نشان می دهند و اینکه چگونه می تواند وسیله نفوذ سازماندهی حافظه اش در مقیاس بزرگ باشد.

2-2-1- خود بقای رسانه؟

دومین مقوله استدلال ها اظهار می دارد که وب رسانه ای خود بقاست. در این دیدگاه، منابعی که برای حفاظت شدن مناسب هستند، بر روی سرورها نگهداری می شوند. بقیه به اراده به وجود آورنده اصلی، ناپدید خواهد شد. از آن جا که نوع اول استدلال درباره کیفیت تقریباً در مجموعه برنامه های جهانی

ص: 11

Rekimoto -1

Dumais et al -2

Ringel et al -3

Hanzo -4

API -5

Freeman and Gelernter -6

Orgun, Mantratzis -7

Tagging -8

Golder and Huberman -9

Halavais -10

Bruns -11

Phillips -12

یافت شده است، این‌ها بیشترین طرفداران را در علوم کامپیوتر جهانی به دست آورده‌اند. اگر چه به شدت در روزهای اول پشتیبانی شده است، باید بگوییم که همانطور که زمان می‌گذرد و محتوا از وب محو می‌شود این مسئله کمتر مورد [چالش] است. اسناد مطالعاتی زیادی بر طبیعت بی‌دوام منابع وب، این ادعا را که که وب رسانه‌ای خود بقاست دچار شکست می‌کند. برای مثال، برای مروری بر ادبیات موضوع (کهلر 2004) (1) و (اسپاینلیس 2003) (2) را ببینید. این مطالعات بر دسترس پذیری منابع با همان URL تأکید می‌کنند و نه تغییرات بالقوه‌ای که می‌تواند متحمل شود. مطالعات نشان می‌دهند که میانگین نیم عمر هر صفحه وب (مدت زمانی که نیمی از صفحه‌ها ناپدید خواهد شد)، فقط دو سال است. این بررسی‌ها بر روی دسترس پذیری منابع در URL مشابه متمرکز هستند و تغییرات بالقوه‌ای وجود ندارد که آن‌ها متحملش شوند. همچنین برخی، مضمون‌ها را مورد تحقیق و بررسی قرار داده و میزان تغییر را ارزیابی می‌کنند. چو و گارسیا - مولینا (2000) (3)، نیم عمره 5 روزه‌ای را برای میانگین صفحه‌های وب کشف کردند. فترلی و همکارانش (2003) (4) نشان دادند چگونه این میزان تغییر با اندازه و موقعیت محتوا مرتبط می‌باشد.

دلایل بسیار زیادی برای تمایل منابع به ناپدید شدن از وب وجود دارد. نخست، محدودیت زمانی مجاز شدن دامنه نام (معمولاً 1-3 سال) است که به وسیله طراحی هر فضای وب در یک انتقال و شرایط غیر ایمن واقع می‌شود.

دوم، توان الکتریکی پایدار، پهنای باند و سرورهایی که به پشتیبانی انتشارات نیاز دارند - همان طور که در مقابل ماهیت خارجی انتشار صورت می‌گیرد. اما حتی وقتی که نامیدن فضا و منابع نشر، مصون هستند، سازماندهی و طراحی اطلاعات می‌توانند نقش مهمی را در حالت ارتجاعی منابع روی سرورها ایفا کند (برنرز لی 1998) (5). همان طور که برنرز مخترع وب ادعا کرده است:

«اصلاً استدلالی در نظریه برای افراد وجود ندارد تا URL‌ها را تغییر دهند (یا اسناد نگهداری شده را متوقف کنند) اما میلیون‌ها استدلال در عمل وجود دارد (برنرز، 1998).

تغییر افراد، سازماندهی داخلی، طرح‌ها، فناوری‌های سرور وب، عملیات نام دادن، و مانند آن می‌تواند ناشی از بازسازی و گاهی فقدان اطلاعات باشد.

از این دیدگاه، سبک رشد نظام مدیریت محتوا (6) در انتشار، برداشت‌های گمراه‌کننده در برقراری نظم هنگام بحران - چنانچه نظام مدیریت محتوا آورد چون معمولاً یک سبک ساختاری از اطلاعات یکپارچه و اغلب قابلیت‌های بایگانی کردن را دارا می‌باشد. مسئله این است که آن‌ها به لایه‌های دیگر وابستگی به نرم افزار اضافه می‌شوند (نرم افزار نظام مدیریت محتوا) چون استاندارد سازی در این حوزه وجود ندارد. معماری‌های اطلاعات بر اساس نظام مدیریت محتوا ثابت شده است که خنک است تا زمانی که نظام

ص: 12

Koehler -1

Spinellis -2

Cho and Garcia-Molina -3

Fetterly etal -4

Berners-Lee -5

CMS -6

مدیریت محتوا تغییر نکند، یعنی که خیلی طولانی نباشد.

اما آیا طراحی اطلاعات دستی است یا توسط سیستم انجام می شود. وب رسانه ای خود بقا نیست و نخواهد شد. مهم ترین دلیل آن مغایرت فعالیت های انتشار و حفاظت می باشد. انتشار، به معنای ایجاد تازگی است حتی زمانی که در مطالب کهنه هزینه شده باشد و (برای مثال در یک فضای نام گذاری مشابه یا کتاب های جدید و قدیمی باید در مخزن ناشران مشابه، با هم قرار گیرند).

تجربه ثابت می کند که انگیزه حفاظت، در میان تولید کنندگان محتوا کافی نیست و آن ها را برای حفاظت وابسته می سازد. در واقع، مرحله نخست حفاظت، اجبار در انجام آن توسط انواع مختلف سازماندهی، اجرای توسط اهداف متفاوت، انگیزه ها و حتی روش متفاوت است. وب، به عنوان زیر ساخت اطلاعاتی نمی تواند به طور اساسی مشکلات سازمانی را حل کند. از این رو، بایگانی کردن وب به عنوان فعالیتی مستقل از انتشار مورد نیاز است.

2-1-3-یک وظیفه غیر ممکن

سومین مقوله استدلال ها در مقابل بایگانی وب، از سوی افرادی مطرح می شود که نیاز به بایگانی وب را تصدیق می کنند، اما درباره امکان انجام آن شبهه دارند.

تردیدها، یا در مورد اندازه وب است یا در موارد دیگر (تأکید بر خصوصی سازی، خاصیت، روشن فکر گرایی و موانع حق مؤلف) که بایگانی وب را به چالش می کشاند.

نخستین جنبه، منتسب به بیکرانی وب است که باید در رابطه با هزینه های ذخیره و ظرفیت ابزارهای خودکار برای گردآوری حجم زیاد اطلاعات مورد بررسی قرار گیرد. خطوط DSL فعلی و ظرفیت پردازش کامپیوترهای شخصی، خزش روزانه میلیون ها صفحه را امکان پذیر می کند. مقیاس میانگین بایگانی وب، در تناسب با مقیاس خود وب می باشد. حتی اگر تخمین دقیق آن مشکل باشد (داهن 1)، 2000؛ آگهه 2000 (2)؛ دوپرا و فاینبرگ (3)، 2004) از منابع مختلف در می یابیم (4) که اندازه وب سطحی، به طور متداول در دامنه ده ها بیلیون صفحه است و این اطلاعات به شکل سایر نظام های اطلاعاتی پیچیده وب که - نمی تواند خزش کند (وب پنهان) - قابل دستیابی می باشد، البته به اندازه یا دو مرتبه بزرگ تر است.

بایگانی وب سطحی ثابت شده است طی یک دهه کامل توسط بایگانی اینترنتی، سازمانی کوچک با سرمایه گذاری خصوصی کوچک شدنی است (کاهل 5)، 2002، 1997).

دلیل این امر این است که این میزان مشابهی از محتوا، به وجود آورندگان، ارزش قابل توجهی را برای ایجاد، حفظ و نگهداری و دسترسی بالا پرداخت می کنند. ذخیره سازی فقط قسمت کمی از

ص: 13

4- منابع به اندازه نمایه موتورهای کاوش مستند شده اند (یاهو ادعا می کند که 20 بیلیون صفحه را نمایه می کند، گوگل می گوید که بیشتر نمایه می کند (بتل، 2005) در یک نگاه کلی اندازه بایگانی اینترنت 10 بیلیون صفحه است)، مطالعات اخیر بر اساس روش شناختی های نمونه گیری است (گلی و سیگنورینی، 2005).

هزینه های انتشار وب را امروز در بر می گیرد. بر عکس، بایگانی اینترنتی فقط برای ذخیره سازی، با استفاده از فشرده سازی برای مثال خزش گر توسط الکسا اهدا شده است) و دسترسی ها هزینه را پرداخت می کند و مورد دوم پرداخت برای هر واحد محتواست، که بسیار کوچک تر از چیزی است که سرور اصلی می پردازد. این نتایج در میزبان کردن یک کپی برداری کاملاً گسترده از وب در مؤسسه ای واحد (کوچک) به طور محسوس ممکن است.

جنبه دوم، نگرانی های خصوصی سازی، مالکیت معنوی و موانع حق مؤلف است فقط توجه داشته باشید که وب یک برنامه کاربردی انتشار غیر تجاری در اینترنت نیست. ارتباطات پنهانی برای رخ دادن در وب تصور نشده اند، اما درباره برنامه های کاربردی ارتباطات (مانند ایمیل و انتقال پیام) زمانی که این کار انجام می شود (لیوگ و فیشر (1)، 2003) همیشه احتمالات برای حفاظت از آن ها (که به طور وسیع می شود) به وسیله ورود به سیستم و اسم رمز وجود دارد. از این رو، فضاهای حفاظت شده به عنوان بخشی از وب عمومی مورد بررسی قرار نگرفته اند و بنابراین نباید در بایگانی های عمومی حفاظت شوند. این طرح طبیعی از جو خصوصی / عمومی در اینترنت به وسیله روشی که خزش گرها اجرا می کنند، تقویت می شود (به وسیله دنبال کردن پیوندها) به این معناست که صفحه ها و سایت ها به داشتن درجه معینی از پیوند درونی برای کشف شدن و تصرف شدن نیاز دارند. بقیه، اجزای غیر متصل وب هستند (برادر و همکارانش (2)، 2000) ، که به طور طبیعی از خزش گرها حذف می شوند. سایتی می تواند از این استفاده کرده و مرزهای بیشتری را برای شامل شدن در مجموعه (بیش از یک لینک درونی) برای محدود کردن تصرف بخش های مرئی تر تنظیم نماید.

در رابطه با وضعیت قانونی بایگانی وب، به وضوح، موقعیت های گوناگونی در هر کشور وجود دارد و این یک فضای در حال نمو می باشد. کشف این جنبه ها، فراتر از دامنه این کتاب است که در کتاب چارلزورث (2003) (3) به آن ها اشاره شده است. توجه داشته باشید که محتوای منتشر شده در وب غیر تجاری است چه توسط تبلیغات بر روی سایت ها یا بوسیله اشتراک پرداخت شود.

برای تمام موارد، بایگانی های وب، حتی با دسترسی پیوسته باید شرایط غیر رقابتی را با وبگاه های اصلی پیدا کنند و این کار می تواند در خصوص محدودیت های دسترسی به محتوا انجام شود (برای مثال همان طور که توسط تولیدکننده در متن فایل روبات ها گفته شده است). داشتن یک دوره ممنوعیت، قابلیت های کمتری را نشان می دهد (جست و جوی سایت و تعاملات پیچیده) و همچنین عملکردهای سطح پایین (سرعت دسترسی به محتوا). بنابراین، استفاده از بایگانی وب برای دسترسی به محتوا، زمانی انجام می شود که دسترسی اصلی امکان پذیر نباشد و درآمد بازدهی داشته باشد. در این صورت، برای ناشر اصلی استفاده از بایگانی وب تهدیدی محسوب نمی شود (این موضوع را در لایمن 2002 (4) ببینید). در مقابل، بایگانی وب می تواند به طور قابل توجهی برای به وجود آوردن سایت، حفظ بار محتوای منسوخ

ص: 14

Lueg and Fisher -1

.Broder et al -2

Charlesworth -3

Lyman -4

(قدیمی) را کم کند امکان تمرکز بر روی محتوای فعلی را می دهد. حتی در این موقعیت، نویسندگان و ناشران ممکن است درخواست کنند که مطالب شان از بایگانی های قابل دسترس عموم برداشته شود. درخواست همچنین می تواند توسط شخص سوم به دلایل مختلف انجام شود. چگونه بایگانی های عمومی وب به این درخواست ها پاسخ خواهند داد؟

توصیه هایی در این زمینه در ایالات متحده پیشنهاد شده است. جدول 1-1 را ببینید، (آبویس 2002) (1).

عکس

بایگانی شبکه وب: مباحث و روش ها ۱۵

(قدیمی) را کم کند امکان تمرکز بر روی محتوای فعلی را می دهد. حتی در این موقعیت، نویسندگان و ناشران ممکن است درخواست کنند که مطالبشان از بایگانی های قابل دسترس عموم برداشته شود. درخواست، همچنین می تواند توسط شخص سوم به دلایل مختلف انجام شود. چگونه بایگانی های عمومی وب به این درخواست ها پاسخ خواهند داد؟
توصیه هایی در این زمینه در ایالات متحده پیشنهاد شده است. جدول 1-1 را ببینید، (آبویس 2002).

جدول ۱.۱

درخواست	توصیه
درخواست مدیر وب توسط یک وبگاه خصوصی (غیردولتی) به دلایل خصوصی سازی، افترا یا خجالت	<p>۱- آرشیویست ها باید یک سایت به روش سلف سرویس ایجاد کنند و مالکان سایت می توانند مطالبشان را با استفاده از آن حذف کنند که براساس استفاده از معیار پروتکل رویوت متن می باشد (با این قابل متنی می توان، میزان دسترسی موتور جست و جوگر به محتوای یک سایت را کنترل کرد).</p> <p>۲- درخواست کنندگان ممکن است بخواهند با دلیل مدرک مالکیت شان را با تغییر با افزودن یک پروتکل رویوت متن، در سایتشان اثبات کنند.</p> <p>۳- این کار به بایگان ها اجازه می دهد تا مطمئن شوند که مطالب گذشته بیش از این جمع آوری نخواهند شد یا قابل دسترسی نیستند.</p> <p>۴- این درخواست، عمومیت نخواهد داشت از این رو آرشیویست ها باید کپی های تمام درخواست های حذف شده را نگه دارند.</p>
درخواست های پاک شده شخص سوم براساس بیانیه حق مؤلف هزاره رقمی سال ۱۹۹۸ می باشد.	<p>۱- آرشیویست ها باید در تلاش برای بررسی اعتبار شکایت ها با کنترل آنها باشند که آیا صفحه ها اصلی ثبت شده اند و اگر مناسب است بر روی درخواست نسبت به سایت اصلی از نظر قانونی اعمال حکم می شود.</p> <p>۲- اگر شکایت معتبر باشد، آرشیویست ها باید موافقت نمایند.</p> <p>۳- آرشیویست ها برای ایجاد درخواست های عمومی از بیانیه حق مؤلف هزاره رقمی، از طریق اثرات ناامیدکننده^۱ و اخطار به جست و جوکنندگان زمانی که صفحه ها درخواست شده حذف شده اند، تلاش خواهند کرد.</p> <p>۴- بایگان ها به مدیران سایت های مذکور از طریق ایمیل، اخطار خواهند داد.</p>
درخواست های پاک شده شخص سوم براساس شکایت های غیر از بیانیه حق مؤلف هزاره رقمی شخصی و (شامل علائم تجاری و رازهای تجاری در تولید).	<p>۱- آرشیویست ها برای بررسی اعتبار شکایت با کنترل این که آیا صفحه ها اصلی ثبت شده اند، تلاش خواهند کرد، تا اگر مناسب است بر روی درخواست در رابطه با سایت اصلی از نظر قانونی اعمال حکم کنند.</p> <p>۲- اگر صفحه ها اصلی پاک شده اند و آرشیویست ها تعیین کند که پاک کردن آنها از سرورهای عمومی، مناسب است، آرشیویست ها صفحه ها را از سرورهای عمومی شان بر می دارند.</p> <p>۳- آرشیویست ها برای ایجاد این درخواست های عمومی به وسیله اثرات ناامیدکننده و اخطار به جست و جوکنندگان وقتی صفحه های درخواست شده برداشته شده باشد، تلاش می کنند.</p> <p>۴- آرشیویست ها به مدیران سایت های مذکور از طریق ایمیل، اخطار خواهند داد.</p>

جدول 1.1

ص: 15

Ubois -1

درخواست‌های پاک شده شخص سوم براساس اعتراض به محتوای بحث انگیز باشد (مثل مباحث سیاسی، مذهبی و عقاید دیگر)	همان‌طور که در مجموعه قوانین کتابخانه بیل آ طبقت ذکر شده کتابخانه‌ها باید مطالب و اطلاعات را تهیه کنند و تمام نقطه نظرات درباره مباحث جاری و تاریخی ارائه دهند. مطالب نباید ممنوعیت (انتشار) داشته باشد یا به علت عدم تصویب یا عدم رضایت طرفداران پاک شود. از این رو، آرشیویست‌ها نباید به‌طور کلی به این درخواست‌ها عکس عمل نشان دهند.
درخواست‌های پاک شده شخص سوم براساس اعتراض به افشاء داده‌های شخصی که به‌طور محرمانه تهیه شده است.	گاهی اوقات، داده‌های محرمانه افشاء شده به‌وسیله یک طرف به طرف دیگر ممکن است در نهایت توسط یک شخص سوم عمومیت پیدا کند. برای مثال، اطلاعات پزشکی که به‌طور محرمانه تهیه شده، گاهی اوقات عمومیت پیدا می‌کند. وقتی عملکردهای شرکت‌های بیمه یا عملکردهای پزشکی متوقف می‌شود، این درخواست‌ها به‌طور کلی به معنای درخواست‌هایی تلقی می‌شوند که توسط نویسندگان یا ناشران داده‌های اصلی ایجاد می‌شود.
درخواست توسط دولت	آرشیویست‌ها بهترین تاثیر مرتبط با قبول قابلیت کاربرد قرار صادره از دادگاه، به‌کار می‌برند. بعلاوه، همان‌طور که طبق قوانین کتابخانه بیل در ذکر شد کتابخانه‌ها باید سانسور عقاید در اجرای تعهداتشان برای تهیه اطلاعات و آگاهی حقیقی را به چالش در آورند.
درخواست‌های دیگر و شکایت‌ها، مباحث حقوقی ایجاد مجدد و نگاه‌ها براساس تغییر مالکیت می‌باشد.	درخواست‌های دیگر و شکایت‌ها، مباحث حقوقی اصولی را در بر می‌گیرد. اینها براساس مورد به مورد توسط بایگانی و مشاوران آنها انجام می‌شوند. کنترل و ایجاد مجدد و نگاه‌ها براساس تغییر مالکیت می‌باشد. این توصیه‌ها می‌توانند با محیط‌های قانونی دیگر وفق داده شوند تا جایی که استفاده مجدد از ساز و کارهای عملکردی مهم ایجاد گردند (ارتباطات از مالک سایت از طریق استفاده گسترده از استاندارد متن روبوتس و همترازی روی آن چه که بر سایت اصلی در ادعای سوم انجام شده است).

نیاز برای درک بهتر همزیستی بین به‌وجود آورندگان سایت و بایگانی‌های وب وجود دارد تا آنجا که بتوانند با عنایت به حقوق ایجاد کننده اطمینان پیدا کنند که از حافظه می‌تواند محافظت شود. اما این نیز بخشی از فرآیند تکامل رسانه وب است.

در مجموع، استدلال‌ها در برابر ضرورت و همچنین امکانات بایگانی وب است. جای تعجب نیست که، از نظر ما، در مغایرت با نقش مرکزی وب در خلق فرهنگ و انتشار آن و همچنین براساس طبیعت مطلقش، ایجاد شده است. فصل ۲ بینش بیشتری درباره چگونگی اهمیت بایگانی‌های وب برای پژوهش در بسیاری از حوزه‌ها را فراهم می‌کند.

در اینجا، سعی می‌کنیم نشان دهیم که در صورتی که چالش‌های جدی و مهمی را برای عملکردهای سنتی مطرح کردیم، بایگانی وب امکان‌پذیر می‌شود و یکی از موارد اصلی در برنامه حفاظت از میراث فرهنگی امروز است.

۳- ویژگی‌های وب برای حفاظت

وب خصوصیات مهمی دارد که هر تلاش حفاظتی باید درباره آن انجام شود. ما آنها را در این بخش در

نیاز برای درک بهتر همزیستی بین به‌وجود آورندگان سایت و بایگانی‌های وب وجود دارد تا آنجا که بتوانند با عنایت به حقوق ایجاد کننده اطمینان پیدا کنند که از حافظه می‌تواند محافظت شود. اما این نیز بخشی از فرآیند تکامل رسانه وب است.

در مجموع، استدلال‌ها در برابر ضرورت و همچنین امکانات بایگانی وب است. جای تعجب نیست که از نظر ما، در مغایرت با نقش مرکزی وب در خلق فرهنگ و انتشار آن و همچنین بر اساس طبیعت مطلقش، ایجاد شده است. فصل ۲ بینش بیشتری درباره چگونگی

اهمیت بایگانی های وب برای پژوهش در بسیاری از حوزه ها را فراهم می کند.

در اینجا، سعی می کنیم نشان دهیم که در صورتی که چالش های جدی و مهمی را برای عملکردهای سنتی مطرح کردیم، بایگانی وب امکان پذیر می شود و یکی از موارد اصلی در برنامه حفاظت از میراث فرهنگی امروز است.

3-ویژگی های وب برای حفاظت

وب خصوصیات مهمی دارد که هر تلاش حفاظتی باید درباره آن انجام شود. ما آن ها را در این بخش در

ص: 16

زوایای مختلف مورد بررسی قرار می دهیم. ابتدا، کاردینالیتی وب است یعنی اینکه چه تعداد نمونه از هر قسمت محتوا موجود است. و دوم اینکه وب به عنوان یک نظام انتشاراتی فعال و آخرین مورد وب به عنوان یک محصول فرهنگی جهانی، با طبیعت مافوق رسانه ای و طبیعت نشر آزاد بودن آن مورد بررسی قرار گرفته است.

3-1-3- کاردینالیتی وب

نخستین سؤال که برای حفاظت از محصولات فرهنگی عنوان شده است، کاردینالیتی بودن آن است. تعداد مواردی که هر اثر در حال توزیع شدن است. بایگانی ها و موزه ها معمولاً، با محصولات منحصر به فردی سر و کار دارند و حتی اگر در برخی موارد چندین قالب وجود دارد که کپی یا نشانه هایی از یک پیکر تراشی واحد، نقاشی یا اثر عکاسی است.

بر عکس، کتابخانه ها تقریباً موارد غیر منحصر به فرد را در مجموعه چاپی نگهداری می کنند (حفاظت از نسخه خطی، از این دیدگاه نزدیک به عملکرد بایگانی می باشد) منحصر به فرد بودن دارای اهمیت اجتماعی و نمادین عمیق است (بنجامین 1963) (1). همچنین، اثر بسیار زیاد آشکاری بر عملکردهای حفاظتی دارد. کتابخانه ها همیشه یک فرصت ثانویه برای یافتن کتاب های چاپ شده بعد از انتشارشان دارند.

چنین تخمین زده می شود بیشتر از 20 میلیون کتاب برای 30/000 ویرایش، بین 1455 و 1501 به چاپ رسیده است (فبوره و مارتین، 1976) که به این معنا که به طور میانگین نخستین دوره کاردینالیتی، متجاوز از 650 بوده است. کاردینالیتی مستلزم این است که حفاظت با تأخیر معینی بعد از انتشار رخ دهد، همان طور که کپی های متعدد برای یک دوره زمانی حتی در غیاب حفاظت فعال باقی می مانند. همچنین، یک سطح طبیعی از افزونگی یک ویژگی در یک نظام را موجب می شود که کتابخانه متفقاً انجام می دهند. با استفاده از داده های یکی از بزرگ ترین پایگاه داده های کتابشناختی (worldcat) لایوه و شانفلد (2005) (2) سه ردیف توزیع کاردینالیتی اثر منتشر شده در کتابخانه ها را کشف نمودند که از Worldcat استفاده کرده است (تقریباً 20/000 در آمریکای شمالی): 37 درصد فقط یکبار نگهداشته شده اند، 30 درصد، 2-0 بار و 33 درصد بیشتر از 5 بار نگاه داشته شده اند.

زمان و افزونگی (تکرار اطلاعات میان فایل های گوناگون) دو مزیت قابل توجه از یک چشم انداز حفاظتی هستند که یکدیگر را تقویت می کنند. آن ها همیشه وجود ندارند. تولید مجدد نسخه های خطی در رابطه نقایص آن برای قرن ها قبل از اختراع چاپ اقدام شده بود، از این رو، اکنون حتی زمانی که چندین نسخه (واقعاً تعدادی) گوناگون وجود دارد. کتابداران بزرگ ترین کتابخانه قدیمی اسکندریه (3) را تشکیل می دهند که از کپی برداری نسخه های خطی استفاده می کردند که به دورن شهرها انتقال یافته بودند اما آن ها

ص: 17

Benjamin-1

Lavoie and Schonfeld-2

Alexandria-3

نسخه اصلی را نگهداری می کردند (کن فوراً، 1989).

ترجمه و جمع آوری، تفسیر و توضیحات و حاشیه نویسی، غالب اوقات بنیاد و پایه اصلی برای تولید مجدد متن به جای حفاظت قابل اعتماد بوده اند که برای زیان و ضررهای اجتناب ناپذیر اضافه شده بود که مستلزم کپی برداری دستی بودند. بیشتر کپی برداری های اصولی از متون، اغلب به دلایل خارجی ایجاد می شدند مانند وقتی متون یونانی، اساساً در موقعیت اختراع یک نوشته جدید، (حروف کوچک) در دوران امپراطوری رم در قرن نهم کپی برداری شده اند، تثبیت و انتقال آن ها به شکلی که ما امروزه می شناسیم خواهد بود.

کپی برداری آینده به طور قابل توجهی، شرایط را در این رابطه تغییر می داد که آن محتوا را به حالت تثبیت کرد زمانی که توزیع گسترده تر آن را مجاز کرد (آیزنشتاین 1979 (1)؛ فبوره و مارتین (1976) (2). همچنین، با افزایش قابل توجهی از کاردینالیته آثار، راندمان حفاظت را بدون سابقه کرد. در جایی بر آورد شده است که یکی از 40 اثر شناخته شده از دوران قدیم، حفظ شده است (و کمتر اگر آثار ناشناخته را در نظر بگیریم). راندمان حفاظت به بیشتر از یکی از دو تا در قرن هفدهم در فرانسه و نزدیک 80 درصد یک قرن بعد از آن (استیوالز 1965) (3) بالا رفته است و برای یک مؤسسه واحد، کتابخانه سلطنتی نیز بعد از تقویت سپرده گذاری قانونی توسط فرانسیس ل (4) در سال 1537 (استیوالز، 1961؛ بالیه، 1988) (5).

امروزه، حفاظت از کارهای چاپ شده، در بیشتر کشورها به کارآیی و رشد مؤثری دست یافته است؛ از نقطه نظر عملی و سازمانی که با ثبات مطالب چاپ شده و همچنین کار دینالیته، مجاز شده اند.

هر آن چه که بود، کاردینالیته در محصولات فرهنگی حداقل از ایجاد تا دسترسی یکپارچه می باشد. این تنها مورد در وب نیست. کاردینالیته محتوایی وب ساده نیست، بلکه هر یک (چند جزئی) است. همان طور که منبع محتوا معمولاً یک سرور منحصر به فرد است شخص می تواند به طور محسوس نماید که کاردینالیته اش مانند آثار هنری و نسخه های خطی، یکی می باشد. در واقع، همان آسیب پذیری را نشان می دهد، حتی توسط این حقیقت که محتوا به تولید کننده خود وابسته است، افزایش می یابد. از طرفی دیگر، دسترسی و همچنین کپی های محتوای وب می تواند از نظر مجازی، نامحدود باشد. این اختلاف میان دو کاردینالیته های وب، ما را به سمت این مفهوم مهم از منبع وب هدایت می کند. هر منبع، یک منبع منحصر به فرد (سرور وب) و یک شناسه منحصر به فرد دارد، اما می تواند از نظر مجازی به طور نامحدود و با درجه های گوناگون برای برنامه ریزی های گوناگون تولید شود.

از دیدگاه حفاظتی، هر منبع دو خصوصیت مهم دارد:

نخست اینکه به طور دائمی به منشأ انحصاری اش برای موجودیت وابسته است. این کار یک تفاوت قابل توجهی با انتشار ایجاد می کند، جایی که مدیران انتشار، فقط یک بار به آن نیاز پیدا می کنند و بعد

ص: 18

Eisenstein -1

Febvre and Martin -2

Estivals -3

François 1er -4

از آن، کتاب‌ها به وجود می‌آیند. دوم اینکه سرورهای وب می‌توانند محتوا را برای هر نوع منبع، مناسب سازند و آن را در هر زمانی برای یو.آر.ال مشابه متفاوت می‌سازد. وب از این دیدگاه، یک ظرف محتوی فایل‌های ثابت نیست، اما یک جعبه سیاه با منابعی است که کاربران فقط نمونه‌هایی را به دست می‌آورند.

همان‌طور که کریشنامرتی (1) و رکسفورد (2) درباره پروتکل وب توضیح داده‌اند:

یک روش برای ادراک پروتکل، این تصور است که منبع سرور حاوی جعبه‌های سیاه با نمایش منابعی باشد که توسط یو.آر.ال‌ها معنا شده‌اند. منبع سرور اصلی، شیوه درخواست برای منبع مشخص شده را به وسیله یو.آر.ال درخواست می‌کند. دریافت مشترک از خواندن یک منبع از یک فایل و نوشتن پاسخ برگشت به سرویس‌گیرنده، دور از دید جعبه سیاه، مجزا و مختصر شده است. این نگرش، مفهوم یک منبع را عمومیت می‌بخشد و آن را از پاسخ ارسال شده به سرویس‌گیرنده تفکیک می‌نماید. درخواست‌های مختلف برای یو.آر.ال‌های مشابه، می‌تواند ناشی از پاسخ‌های متفاوت باشد و به عوامل مختلفی بستگی دارد. فیلدهای بالایی درخواست، زمان درخواست با تغییرات برای منابعی که ممکن است رخ دهند (کریشنامرتی و رکسفورد، 2001) حفاظت وب، منابع را طبق کاردینالیتی دوگانه (به ظاهر مهم و در واقع درست) مورد بررسی قرار می‌دهد و این کار مستلزم چندین استنباط است. نخست اینکه چون از نظر مجازی تعداد نامحدود کپی برداری می‌تواند به آسانی ایجاد شود، شخص می‌تواند ادارک گمراه‌کننده‌ای داشته باشد که بایگانی فعال وب برای حفاظت مورد نیاز نیست. از این رو، تعدد نمونه‌ها، به طور گسترده مخفی‌اند و به یک منبع تکی بستگی دارند. هر زمانی که لازم باشد. (سرور) کد می‌تواند برچیده و روزآمد شود از این رو برای یک بایگانی فعال مورد نیاز می‌باشد.

استنباط دوم این است که بایگانی‌های وب می‌تواند فقط برخی موارد در منابع را به طور بالقوه درجات گوناگونی از میان آن‌ها را به تصرف در آورد (3). این مورد زمانی رخ می‌دهد که محتوا برای یک مرورگر خاص یک زمان معین یا یک موقعیت جغرافیایی معین یا زمانی که محتوا با هر کار بر وفق داده شده است، مناسب گردد. همچنین، در بخش بعد خواهیم دید که وب در حقیقت یک سیستم نشر فعال است و از این رو، تفاوت پاسخ‌ها در واقع جنبه‌ای مهم برای بررسی است، به خصوص وقتی که بایگانی انجام می‌شود.

3-2-وب به عنوان سیستم نشر فعال

وب، نوعی برنامه کاربردی نشر اصلی در اینترنت است همچنین، به طور اساسی شامل ترکیب سه

ص: 19

Krishnamurthy -1

Rexford -2

3- توسعه پویای صفحه‌ها برای ایجاد یگانگی در طراحی و معماری نیز در کل سایت استفاده شده است (دستگاه‌های شناوری و مانند آن). استفاده از تمپلیت‌ها به طور یکسان نگاه کردن به صفحه‌ها را آسان کرده و تغییر طراحی توسط تمپلیت‌ها را آسان‌تر از صفحه‌های انفرادی می‌کنند برآورد شده است که تمپلیت‌های مبتنی بر صفحه‌ها 40 درصد تا 50 درصد از صفحه‌ها را نمایش می‌دهند (جیسون و دیگران 2005).

استاندارد می باشد: 1) URL (لی - برنرز، 1994) که فضای نام گذاری شده را برای یک شیء تعریف می کند (1)؛ 2) HTTP (فیلدینگ و همکارانش (2)، 1994) پروتکل تعامل سرویس دهنده - سرویس گیرنده را با استفاده از فرآیندها در هسته اش تعریف می کند؛ و 3) HTML (برنرز-لی و کونولی (3)، 1995)، نوعی (4) SGML DTD (تعریف نوع داده) که ارائه صفحه آرایی در مرورگرها را معین می نماید. اجرای این سه استاندارد، هر کامپیوتر متصل به اینترنت را برای وارد شدن به سیستم نشر، قادر می سازد. شبکه سرویس دهندگان وب، نوعی سیستم اطلاعاتی منحصر به فرد را تشکیل می دهد که می تواند در هر حالتی برای تولید، روزآمد شدن، و نشر محتوا در حالتی که کامپیوترهای جدید اجازه می دهند، به کار رود.

در مقایسه با وسایل انتشاراتی دیگر، انقلاب در نشر، گسترش امکانات در تمام جهات ممکن برای تولید، سازماندهی، دستیابی، و ارائه محتوا را نشان می دهد. برای مثال، پیوندها را مورد بررسی قرار دهید:

شخص می تواند استدلال کند این فقط شکل جدیدی از یک ارجاع است که پیش از زمانی که برای نخستین بار نوشته شده به وجود آمده است (5). اما حقیقت این است که روشی که به وسیله قطعه قطعه کردن محتوا به تکه های نشانی پذیر کوچک تر و توجه کلی به برنامه ریزی خاصیت انتقال از طریق شناوری دسترسی به محتوایی که تغییرات عمیقی را در نوشتن و همچنین در خواندن پیدا کرده است، به دلیل تغییرات وب در روش ارجاع قابل تعقیب قانونی است (آرسث (6)، 1997؛ لندو (7)، 1997؛ بولتر (8)، 2001).

این حقیقت، که محتوا تنها بر روی سیستم و با دقت بیشتری بر روی سرورهای ناشران موجود است، به انتشار دائم ایجاد کننده بستگی دارد. یک کتاب می تواند بعد از ترک چاپخانه به طور مستقل از ناشر باقی بماند، اما محتوای وب هیچ موجودیتی فراتر از سرور اصلی اش نخواهد داشت (به استثنای ساز و کارهای حافظه با سرعت بالای ناپایدار (9) هافمن و بیومونت، 2005). انتشار دائم، کنترل واضح را گسترش می دهد که ایجادکنندگان بر روی محتوا دارند. آن ها می توانند، با وب، در هر زمانی تغییر کنند، به روز شوند، و در زمان واقعی مواردی را از انتشار پاک کنند. علاوه بر این، تولیدکنندگان وب، از نظام اطلاع رسانی وب (10) استفاده می کنند که بتواند اطلاعات را از هر نوع نظام اطلاعاتی موجود (پایگاه داده ها، مخزن اسناد، برنامه های کاربردی، و مانند آن) ترکیب، مجتمع و سازماندهی مجدد کند. از این رو، وب یک فضای اطلاعاتی ثابت نیست، بلکه یک فضای نشر فعال است که ناشی از اثرات یک مجموعه آمیخته شده از نظام های اطلاعاتی فعال می باشد.

ص: 20

1- این استاندارد مهم ترین در میان سه مبتکر وب است (لی - برنرز و فیشتی، 2000؛ گیلز کیلیو، 2000؛ چنان چه وب را در موقعیت دسترسی جهانی قرار داده است که کلاً منبع سندی قابل دسترس در اینترنت است.

Fielding -2

Connolly -3

SGML DTD -4

5- برای مقایسه استنادهای علمی سنتی و این که چگونه می تواند برای ارزیابی علمی استفاده شود اینگورسن (1998) را ببینید بچورن بورن و اینگورسن (2001) تحلیل انتقادی آن در توال (2001)، توال و هریس (2004) و توال (2006).

Aarseth -6

Landow -7

Bolter -8

Hofmann and Beaumont -9

(Web information systems (WIS -10

از این رو، بایگانی وب نخست به جدا کردن محتوا از نشر ثابت ایجادکنندگان اصلی اش نیاز دارد، و دوم اینکه باید مطمئن شود که محتوا می تواند از عدم پذیرش و تکامل جاری وب، عدول کند.

قبلاً به کپی برداری و بایگانی محتوا در یک زیر ساخت مجزا نیاز بود (مطالب زیر، (1) فصل 3 و روشه، 2006 را ببینید). مورد آخر مستلزم حفاظت فعال از محتوای وب (فصل 8، دی (2) 2006 را ببینید) برای رفع وابستگی از اجزای سیستم های گوناگون (پروتکل ها، به فرمت های دیجیتال، برنامه های کاربردی، و نظیر آن) و اجتناب از منسوخ بودن اصول فنی آن هاست. حفاظت از وب، این نیازها را در کل برای حفاظت از اصول فنی فعال با اشیای رقومی به اشتراک گذاشته است، اما جداسازی از ایجاد کننده نشر دائم، برای حفاظت از وب مشخص است.

اما رفع هر گونه وابستگی از سرور اصلی، مستلزم این است که از قابلیت های گوناگون و شیوه تعاملی وب، بایگانی وب بتواند فقط تعداد کمی را حفظ کند. هزینه هایی برای جداسازی از شبکه اصلی نظام های اطلاعاتی وب وجود دارد.

قابلیت هایی که بر بخش سرویس گیرنده ها اجرا می شوند، آن هایی هستند که شخص می تواند به طور معقولانه به حفظ آن امیدوار باشد. دامنه قابلیت هایی در کد صفحه و کد فایل مربوط جاسازی شده اند که به میل سرویس گیرنده اجرا می شوند و بیشتر اوقات بر روی نسخه های بایگانی قابل اجرا هستند؛ اما این قابلیت ها که توسط کد و یا با اطلاعات سرور تهیه شده اند جاسازی نمی شوند. جنبه های سندی مطالب اصلی هستند که گم شده اند (مانند انواع تعاملات که شخص می تواند بر روی یک ویدئو ثبت نماید)، اما این کار فقط می تواند برای تعداد محدودی از صفحه ها و نقطه نظرات معین و شرایط خاص انجام شود (کریستنسن - دالسگاد (3) ، 2001 ؛ بروگر، 2005) (4).

3-3-وب به عنوان یک محصول فرهنگی

علاوه بر یک نظام نشر فعال، وب یک فضای اطلاعاتی با مشخصات خاص است. واژه وب در این مضمون، یک محصول فرهنگی رقومی گسترده را برگزیده است (لایمن و کاهله، 1998) که می تواند با حقایق زیر مشخص شود:

ص: 21

Roche -1

Day -2

Christensen-Dalsgaard -3

4- نقطه نظر طراح وبگاه در این مورد نیز جالب است. در دابری و دیگران (2002). جلیس هودج پیشنهاد می کند، سایت هایی که او در حال طراحی است بایگانی شوند: - درخواست برای پروپوزال؛ - بیان هدف و دلیل استفاده؛ - توصیف استفاده از مضمون (مثال های مورد نیاز)؛ - توصیف استفاده کنندگان واقعی و مورد هدف؛ - جایگزین های ثابتی که به اندازه کافی دیدگاه و احساس را احاطه می کنند؛ - مثال هایی از چند راه مهم در سایت؛ - توصیف فناوری هایی که به کار برده شده و یا حمایت می کنند؛ و - هر ماژول مرتبط مثل پویا نمایی های فلش، فیلم ها، PDF ها، و مانند آن.

- از هر محلی که متصل به اینترنت است قابل انتشار و دسترسی (معمولاً مجانی) می باشد؛

- به عنوان یک فوق رسانه با استفاده از پیوندهای مستقیم و قابل تعقیب قانونی بین قطعات محتوا ساخت

یافته است (1)؛

- نه تنها شامل متن است، بلکه شامل ترکیبی از تصاویر، صداها، و محتوای متنی نیز می باشد؛ و

- ناشی از یک تألیف و تصنیف باز و توزیعی می باشد (2).

اگر چه وب، این کارها را به طور بسیار گسترده انجام می دهد، از اشکال قبلی انتشار هم می نماید (3) (کراستون و ویلیامز 1997) (4)؛ اریکسون و آیهلستروم، 2000؛ شفرد و پولانی، 2000). همچنین، مواردی جدید را ابداع می کند. برای مثال، بلاگ ها با تلفیق ساده گسترده ای منتشر می شوند (حتی مهارت های فنی که برای سایت های عادی لازم است و بیش از اینها مورد نیاز نیستند). و یک مدیریت مرجع قدرتمند (شامل مراجع معکوس یا آگاه سازی از استنادها با استفاده از برگشت پینگ) و سهولت در به روزرسانی، افزودن توضیحات و حذف محتوا، همه اینها ناشی از یک نشر آزاد و توضیحات شخصی به وسیله ده ها میلیون نفر بوده است (5).

این خصوصیات وب، به عنوان یک فوق رسانه توزیعی، به طور آشکار و ثابت در یک مقیاس سراسری، تألیف و فراهم شده است که بایگانی وب بتواند فقط، به حفاظت از جنبه های محدود شده محصولات فرهنگی موجود و بزرگ تر دست یابد.

اتصالات درونی محتوا، یک کیفیت مهم در وب است که زمانی که بایگانی می گردد، مباحثی از آن ها ایجاد می شود، اما به عنوان یک نتیجه عمومی، نشان داده است که بایگانی همیشه، بهترین نوع گزینش (به گزینی) را در برنمی گیرد، حتی اگر در موارد انتخاب دستی و سایت به سایت باشد. از این مباحث برای بایگانی های بزرگ و وسیع اجتناب می شود، تا جایی که ممکن است، قطع تداوم اطلاعاتی که وب ارائه می دهد (لایمن و دیگران 1998). یا برای تعریف یک هدف تحلیلی خاص برای زمینه یابی تصمیم های گزینش انجام می شود (بروگر، 2005).

اما در عمل، اجرای خزش، در ارتباط با موارد اولویت دار و سیاست گزینش دستی، اجزای بایگانی وب را در بر می گیرد که همیشه فقط یک برش در فضا و زمان وب اصلی خواهند بود.

چگونه این نمونه های معنی دار و جایگزین را در وب بزرگ تر ایجاد می کنید؟ چه استنباطی هایی

ص: 22

1- آبرون و مک کرلی (2003) نشان می دهند که یک سوم پیوندهای مستخرج از بیلیون ها صفحه از نقطه جهت مشابه می خزند، یک سوم از عرض، بالا یا پایین در سلسله مراتب جهت ها از همان سایت و یک سوم از سایت خارجی.

2- نویسندگی محدود به تعدادی از افراد نیست بلکه از طریق ده ها و صدها فرد توزیع شده است. مثلاً تخمین زده شده است که در مورد فرانسه، گره های نشر، شخص یا ساختاری که منتشر می کند (ویراستاران نه نویسندگان) با انتشار در وب سه برابر اندازه توسعه یافته است: در حدود 5000 ناشر یا ساختار دهنده نشر به پنج میلیون سایت و سایت شخصی (منبع: انجمن فرانسوی دسترسی به اینترنت). این موارد

شامل وب نوشت ها نیست.

3- این در مورد چاپ نیز بود که برای مدت طولانی از دست نوشته ها و سازماندهی صفحه قبل از اختراع چاپ تبعیت می کرد (فبوره و مارتین، 1976).

Crowston and Williams -4

5- در مورد حفاظت از بلاگ ها انتلیج (2004) را ببینید.

در درک از آینده کاری که وب انجام خواهد داد دارید؟ تمام این سؤال‌ها باید زمانی که بایگانی و را به کار گرفته اید مورد بررسی قرار بگیرند. حتی تعریف کاری که «وب اصلی» انجام می‌دهد مجموعه تجربیات استفاده‌کنندگان از وب یا مجموع برنامه‌ریزی‌های محتواس‌ت که باید از قبل در نظر بگیریم که بر مجموعه‌ای از محتواهای ثابت بیشتر در نظر گرفته می‌شوند.

خصوصیات دیگری که برای اندیشه و سازماندهی مجدد عملکردهای حفاظت سنتی باید انجام طبیعت نویسنده‌گی باز وب است. در حقیقت، این موضوع فیلتر کردن و حفاظت ساختار را بسیار مشکل می‌سازد که بر اساس ناشران و نویسندگان است. آن‌ها به اندازه بسیار قابل توجهی بر روی وب می‌باشند و برای شناسایی و ثبت مشکل‌اند. گاهی اوقات، اطلاعات تألیف و تصنیف، بر روی سایت موجود است، گاهی اوقات هم نیست، و گاهی اوقات در روش قابل اعتمادی نیز قابل دسترسی نیست. تنها اطلاعات ثبت شده (در یک روش کاملاً کنترل نشده و آزاد)، اطلاعات درباره کسانی که نام دامنه را برای مدیرتی DNS مجاز می‌کنند. اگر چه این اطلاعات مقادیر بزرگی برای کامل کردن مطالب بایگانی شده در وب هستند، و به طور قطع، برای تفسیر و استفاده مستقیم آسان نیست.

به عنوان یک محصول فرهنگی، وب، سبک متفاوتی از سازماندهی اطلاعات و الگوهای ساختاری متفاوت را برای استفاده در حفاظت ارائه می‌دهد. نشانه‌های پیوند محتواها و شناوری کاربران از ساختارهای طبیعی که بایگانی‌ها باید بیشتر وقت‌ها برای سازماندهی موارد نمونه‌های جمع‌آوری شده استفاده کنند. از این رو، خصوصیات وب به تغییر شکل و روش‌های حفاظتی عمیق نیاز دارد. رویکرد کل‌نگری برای بایگانی وب، آمادگی بیشتر برای انطباق با خصوصیات وب است، اما هر نوع بایگانی وب باید آن‌ها را در هسته روش‌های مشارکت دهد.

4- روش‌های جدید برای رسانه جدید

کتابخانه‌ها آرشیوها (بایگانی)، و موزه‌ها، روش‌های بسیار کارآمدی را با موضوعات مورد علاقه‌شان تطبیق داده‌اند که نقش مهمی در ساخت حافظه اجتماعی ایفا می‌کند. اگر چه باید بیشتر فراگرفته شوند و بتواند برای حفاظت از وب مجدداً استفاده‌گردند. ماهیت وب و کیفیت‌های مورد نیاز، همان‌طور که دیده شد برای بررسی مجدد و تطبیق عملکردهای حفاظتی به ارث برده شده از این سنت طولانی در حفاظت از محصولات فرهنگی، فیزیکی می‌باشد. این بخش، یک بررسی عمومی از روش‌های جدید و رویکردهایی را نشان خواهد داد که باید برای حفاظت از وب مورد استفاده قرار گیرند (فصل 3 تا 8 مباحث را با جزئیات تهیه می‌کنند).

قبل از شروع مبحث روش شناختی، باید درباره استقرار بایگانی وب در زیرساخت های اطلاعات (1) به طور کلی، و به ویژه در اینترنت پرسش هایی مطرح شود.

بورگمن (2000)، تعریف کتابخانه رقومی جهانی را مطرح کرد (pp47sq) و تفاوت میان نگرش تکاملی و انقلابی در فناوری اطلاعات را توضیح داد:

«نگرش انقلابی، کتابخانه هایی است که رقومی هستند و با پایگاه های داده ها به وسیله شبکه های کامپیوتری پیوند یافته اند و در کل، می توانند آرایه ای از خدمات را فراهم نمایند که کتابخانه ها را از ریشه برخواهند کند. نگرش تکاملی این است که کتابخانه های رقومی سازمان هایی هستند که به تهیه محتوای و خدمات در شکل های گوناگون و فقط به عنوان پیش نیازهای مؤسسه ها ادامه خواهند یافت و سرانجام کتابخانه های تکاملی خواهند بود که امروزه وجود دارند» (همان، ص 48).

او، تعریف میان گذر از تکامل - انقلابی را پیشنهاد کرده است: که «کتابخانه های رقومی یک گسترش، افزایش و یکپارچه سازی در نظام های بازیابی اطلاعات و مؤسسه های با اطلاعات چندگانه و کتابخانه ها که فقط یکی باشد را بیان می کند. محدوده قابلیت های کتابخانه های رقومی نه تنها شامل بازیابی اطلاعات است، بلکه ایجاد و استفاده از این اطلاعات می باشد».

موقعیت برای بایگانی های وب در جهتی که آن ها برای موجود شدن مورد بررسی قرار می گیرند و قبلاً فضای اطلاعات را ساخته اند، متفاوت است. همچنین، به طور آشکار قابل دستیابی اند. از روی سنجش و اندیشه، در این فضا، فقط به دروازه بان ها نیاز نداریم. همان طور که محدودیت های دسترسی فیزیکی وجود ندارند. در این زمینه، نقش بایگانی وب، در سازماندهی اطلاعات زیاد است.

کتابخانه های فیزیکی باید در هر دو سازماندهی فیزیکی و فکری اشیا ایجاد شوند و این طیف وسیعی از امکانات و انتخاب ها را میسر می سازد. همچنین، در حالی که دسترسی فیزیکی به محتوا را مدیریت می کنند، نقش میانجی گری اجتناب ناپذیری دارند. کتابخانه های رقومی، این نقش میانجی گری را به وسیله ایجاد محیط های همکاری و دانش وابسته به متن تحت تابع بنیانی جست و جو و دستیابی را گسترش می دهد (لاگوز و دیگران (2)، 2005).

بایگانی های وب در قسمت های مربوط به خود محتوایی بارگذاری شده با روابط تعبیه شده و قابل

ص: 24

1- مفهوم زیر ساخت به طور کل در استار و روهتلر (1994) با ابعاد مختلف تعریف شده است: - جاسازی (در هم جا دادن)، - شفافیت؛ - به دست آوردن یا دامنه (ساختار به یک حادثه عواید یا عمل یک جانبه رسیده است؛ - به عنوان بخشی از عضویت فراگرفته شده است (اعضای جدید برای عضو شدن نیاز به یک آشنایی دارند)؛ - پیوندهایی با کنوانسیون ها؛ - بر پایه پیاده سازی ساخته شوند؛ و - در حالت تفکیک قابل رؤیت باشند. این مفهوم در مضمون ساختارهای اطلاعاتی در بورگمن، (2000، 2003) بحث شده است

Lagoze et al -2

تعقیب قانونی و ساختارهای اطلاعاتی غنی ایجاد شده که توسط میلیون ها نفر در سراسر جهان ویرایش می شود وقتی بایگانی های سنتی و کتابخانه ها نگرش سازمانی، شخصی و ابزارها را در این محتوا ایجاد می کنند (مدخل های شخصی و وب سازان و مانند آن)، فقط در این ویرایش جهانی وب شرکت می کنند. این کیفیت را به عنوان متخصصان حوزه تقلیل نخواهد داد اما آن ها را در یک تلاش سازمانی بزرگ تر قرار خواهد داد.

به عنوان بایگانی، وب آن ها مسئولیت های بیشتری دارند چون محتوا و مفهوم را تصرف کرده و تحت سیطره خواهند داشت و می توانند آزمایش هایی را در بهبود نقش منحصر به فرد قدیمی خود در سازمان دهندگان اطلاعات داشته باشند همچنین می توانند فقط برای تثبیت و محافظت از نمونه های شخصی از یک محصول فرهنگی جاودانی بزرگ تر به دست آیند.

این کار می تواند قانونی باشد زمانی که طبق خط مشی گزینش مناسب با نیاز جامعه کاربران برپا یا توسط اهداف پژوهشی معین شده اداره شود لچر (1) ، 2006؛ میزانس، 2006 ب). اما هزینه ها و محدودیت ها و همچنین امکانات فنی برای بایگانی هر دو در یک مقیاس بزرگ تر و در یک مسیر بی طرف نیز نیازمند بررسی روش های جایگزین در بایگانی وب می باشند این جایگزین در نقش خود معتدل تر ولی در دامنه جاه طلب تر است. نقش سازمان دهندگان اطلاعات برای به تصرف در آوردن محدود است و برای ساختار اصلی ایجاد شده توسط ویرایش میلیون ها نفر در سراسر جهان درست است.

در مشکل رسیدن به جامعیت، همچنین در بخش قبلی دیدیم که شخص می تواند حداقل تلاش در بی طرفی را برای تصرف محتوا با پیروی از ماهیت جمع آوری و توزیعی وب برای راهنمایی در تصرف و گسترش آن، تا جایی که ممکن است داشته باشد از این رو تلاش بر روی کمیت و موضوع مقیاس گذاری بر روی منابع فنی می باشد. این رویکرد، چندین کتابخانه ملی برای حوزه ملی و بایگانی اینترنتی در یک مقیاس جهانی داشته است.

هیچ یک از این ابتکارهای عملیاتی نمی توانند از نظر عمق و کیفیت محتوای بایگانی شده به تنهایی توسعه یابد. تلاش های گوناگونی به عنوان بخشی از یک بایگانی جهانی مورد بررسی قرار خواهد گرفت البته وقتی که اتصالات درونی میان بایگانی وب به عنوان اتصالات درونی بین سرورهای نشر از طریق وب سازماندهی شوند.

تنها با این کار، کاربران قادر به نفوذ به تمام این تلاش ها خواهند بود و به بهترین حافظه ممکن وب منتج خواهد شد. در این جهت، هر چه شراکت بیشتر مؤسسه ها و افراد مختلف وجود داشته باشد، بهتر می توانند مکمل یکدیگر باشند و زوایا، عمق و کیفیت های مختلف بایگانی های متفاوتی را ارائه دهند. اما این کار نیازمند این است که آن ها در برخی نقاط از یک شبکه بایگانی وب بزرگ تر، همکاری کنند. چنین شبکه ای باید بایگانی وب را پیوند دهد، به طوری که با یکدیگر نوعی فضای شناوری جهانی مانند خود وب را شکل دهند این امر فقط در صورتی ممکن است که آن ها در یک مسیر نزدیک به وب اصلی ساخته

ص: 25

شوند و آزادانه قابل دسترسی باشند کنسرسیوم (شرکت) حفاظت از اینترنت بین المللی (1) بر روی ایجاد و تنظیم زمینه هایی برای ایجاد کننده به وسیله توسعه استانداردها و ابزارهایی که ساخت این نوع بایگانی را تسهیل می کند عملکردهایی داشته است (برخی از آن ها در پایان این بخش توصیف خواهند شد).

دسترسی باز در خصوص مقررات و خط مشی است و در این برهه از زمان به عنوان یک بحث آزاد باقی می ماند.

بایگانی وب، به طور اختصاصی یا به عنوان کل می تواند در ایجاد زیر ساخت های اینترنتی مناسب باشد. آن ها از پروتکل ها و استانداردهای مشابه برای سازماندهی اطلاعات و ایجاد دسترسی به آن استفاده کنند. وب می تواند به طور طبیعی آن ها را در برگرد چون آن ها کاملاً با آن سازگار هستند (2). از نقطه نظر زیربنایی بایگانی وب می تواند به آسانی موقعیتی را به عنوان مکمل ایجاد زیر ساخت های اینترنتی پیدا نماید. آن ها در حال فراهم کردن حافظه وب هستند که خود بخشی از وب است و اثر شدید منفی ماهیت ضروری ناپایدار انتشار وب را محدود می نماید.

شخص می تواند با چشم پوشی از این نقش ناراضی باشد شرایط این کار در نظر نگرفتن ارزش طبیعت توزیع گراو جامع بودن این رسانه است که آن را توجیه می کند.

4-2-فراهم آوری

اصطلاح فراهم آوری برای معانی فنی گوناگونی به کار می رود، مانند رسیدن محتوا به درون بایگانی. این اصطلاح شامل تصرف پیوسته و غیر پیوسته تحویل محتوا می شود که فرآیند انتخاب را پوشش می دهد و نه فرآیند درج و توسعه فردا را.

از دیدگاه فنی این مرحله از تعامل با تولید کنندگان و در صورت سنتی برای مؤسسه میراث حافظه، هر چیزی می تواند با شرایط در بایگانی وب جزئی می باشد. زیرا هیچ روش واحدی برای فنون انتشار گسترده وب کافی نیست گسترده سازی دامنه سازندگان و افزایش اندازه محتوا به درجه ای معین با خودکار سازی که در محیط وب ممکن می شود متعادل می شود. از این رو، موانع اصلی اکتساب ابزارهایی است که باید بر آن غلبه یابند و عدم توانایی پروتکل HTTP در ایجاد کپی دسته ای از محتوای سرور می باشد. سرورهای HTTP فقط می توانند تا زمانی که URI درخواست نماید، فایل به فایل تحویل دهند. این کار موجب کشف گذرگاه فردی برای هر فایل را در یکی از مباحث کلیدی در بایگانی وب

می شود.

در این بخش سه نوع روش فراهم آوری را بررسی می کنیم. چرا سه روش؟ چون فرآیند جمع آوری می تواند یا به عنوان یک سرویس گیرنده دور افتاده، در نزدیکی به خروجی سرور انجام شود یا به وسیله دسترسی مستقیم به فایل های سرور صورت گیرد (تصویر 1). گزینه نخست با خزنده بایگانی یا ماشین

ص: 26

2- فرد می تواند بحث کند که یک بازگشت قهقراپی بالقوه در اینجا وجود دارد در مورد آن چه که با آن مخالفت می شود که آرشیوهای وب باید از بایگانی کردن سایر آرشیوهای وب اجتناب ورزند و خود را با وب زنده (موجود) محدود کنند.

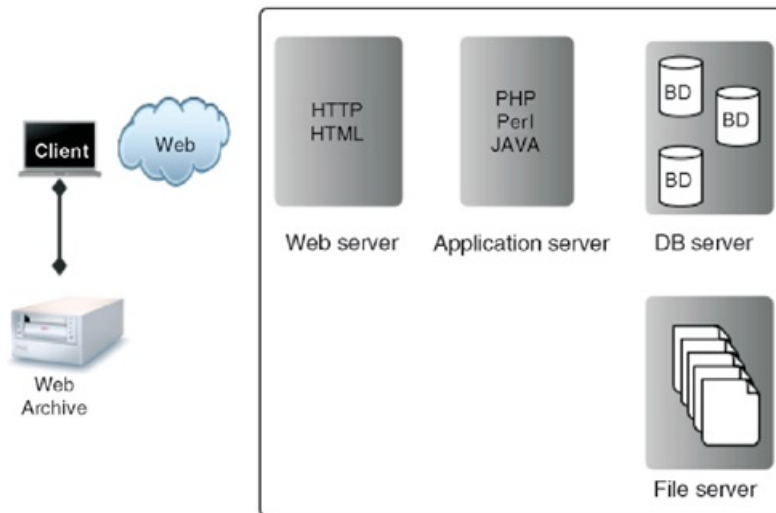
کپی کننده وبگاه انجام می شود، مشتق شده و در فناوری موتور جست و جوی سازگار و استنتاج می شود یک ابزار قدرتمند برای تصرف موقعیت سرویس گیرنده را فراهم می کند. روشه (2006)، توصیفی با جزئیات این ابزارها و کاربردها برای بایگانی وب ارائه می دهد. در این مقاله فقط یک بررسی عمومی از این فناوری را ارائه خواهیم داد که برای ارزیابی در هر یک از موارد می تواند به کار برده شود. چون خزش گر، برای سرور وب یک سرویس گیرنده مانند دیگران است از اصطلاح «بایگانی جانبی سرویس گیرنده» برای این روش فراهم آوری استفاده می کنیم. بسته به ساختار پسین وب و سطح تعامل با سرویس گیرنده خزش گر ها می توانند یا یک وبگاه کامل و یا بخشی از اجزای آن را به تنهایی تصرف کنند. جزء باقی مانده برای خزش گرها غیر قابل دسترسی است در اصطلاح شناسی موتور جست و جو «وب عمیق» یا وب پنهان نامیده شده است.

این اصطلاح شناسی را تصدیق خواهیم کرد تا زمانی که مشخص شود که تعیین حدود وب مخفی، به طور محض فنی بوده و به طور دائم به عنوان خزنده ها توانایی شان را در یافتن راهی برای اسناد، بهبود بخشید. دو روش متناوب برای گردآوری محتوا باقی می ماند، حتی اگر بسیار کم به کار برده شوند و حتی هم چنان تحت تحقیق و بررسی مانده باشند هر دو نیاز به فعالیت از جانب سرور دارند که نه تنها به یک مجوز نیاز دارد بلکه همچنین به یک مشارکت فعال انتشار دهنده سایت برای استفاده شدن نیز نیاز دارد. اولی، بر اساس کاربران سایت است که مسیر هدایت را توسعه می دهد و شناسایی محتوای سایت را برای بایگانی آن انجام می دهد و چون بر اساس ثبت تراکنش های ایجاد شده بین کاربران سایت و سرور می باشد آن را بایگانی تراکنشی می نامیم.

دومی، شامل بایگانی به طور مستقیم از ناشر اجزای گوناگون سیستم اطلاعات وب و انتقال آن ها به یک شکل بایگانی شده می باشد بنابراین بایگانی جانبی سرور نامیده می شود. این روش های متناوب، سخت تر از بایگانی جانبی سرویس گیرنده می باشند چون همانطور که در بالا ذکر شد، نه تنها به یک مشارکت فعال از جانب سازندگان نیاز دارند بلکه باید بر مبنای مورد به مورد اجرا گردند اما حتی اگر افزایش مقیاس نداشته باشند می توانند در مواردی به کار برده شوند؛ برای مثال در جایی که خزنده برای تصرف دقیق موفق نمی شود و زمانی که محتوا در آن کمتر به کار رفته باشند.

4-2-1-بایگانی جانبی سرویس گیرنده

این روش مهمی در فراهم آوری است در هر دو، به علت سادگی، قابلیت مقیاس پذیری، و تطبیق با یک محیط سرور - سرویس گیرنده می باشد (شکل 1-1 را ببینید). خزنده ها با آن چه که روش معمول در دستیابی به وب است سازگار شده اند و این کار بایگانی هر سایت را مجاز می سازد که آزادانه به وب باز شبکه داخلی یا خارجی دسترسی داشته باشند تا زمانی که خزنده به یک اجازه مناسب دست یابد.



تصویر ۱- آرشیو سمت کاربر (Client-side): آرشیو وب در تقابل با آرشیو کاربر است تا محتوا را از سرور وب جمع‌آوری کند. سرور وب می‌تواند محتوا را از سرورهای متعدد و متنوعی فراهم کند (درخواست‌ها، پایگاه داده‌ها، فایل‌های سرور).

این روش، نه تنها موقعیت مشابهی را که کاربران عادی وب نیاز دارند، اتخاذ می‌کند؛ بلکه شکل تعامل‌هایش با سرورها را تقلید می‌کند. خزش‌گرها از صفحه‌ها هسته شروع می‌شوند، آنها را تجزیه می‌کنند، پیوندها را برداشت می‌کنند، و سند پیوندی را واکنشی می‌کنند؛ سپس آنها این پردازش با سند واکنشی شده را تکرار می‌نمایند و تا زمانی که پیوندهایی را کشف کنند^۱ و سند را درون حوزه تعیین شده پیدا کنند. این پردازش مورد نیاز است، چون HTTP فرماتی را که باید فهرست کامل سند قابل دسترس بر روی سرور را باید بازگرداند که برای مثال بر خلاف FTP است. از این رو، هر صفحه باید به وسیله پیوندی از صفحه‌ها دیگر «کشف» شود.

فناوری خزش، برای اهداف نمایه‌سازی توسعه یافته است^۲. به کار بردن آن برای بایگانی وب، برخلاف اینکه جنبه‌های بیشتری از این فناوری را دوباره استفاده می‌کند، تغییراتی را برای آن ایجاد می‌کند.

نخست اینکه، خزش‌گرهای بایگانی باید برای واکنشی تمام فایل‌ها تلاش کنند هرچه فرمتشان برای بایگانی یک مدل کامل از سایت‌ها باشد برخلاف خزش‌گرهای موتور کاوش که معمولاً فقط فایل‌هایی را واکنشی می‌کند که آنها بتوانند فهرست نمایند. خزش‌گرهای موتور کاوش، برای مثال اغلب از انتقال از

۱. برای نظرات اخیر درباره فناوری خزشگر پلنت و دیگران (۲۰۰۴) و چاکرلبارتی (۲۰۰۲) را ببینید.
۲. برای نظر اخیر درباره توسعه موتور کاوش تجاری سونتریج (۱۹۹۷) را ببینید.

تصویر ۱- آرشیو سمت کاربر (client-side): آرشیو وب در تقابل با آرشیو کاربر است تا محتوا را از سرور وب جمع‌آوری کند. سرور وب می‌تواند محتوا را از سرورهای متعدد و متنوعی فراهم کند (درخواست‌ها، پایگاه داده‌ها، فایل‌های سرور).

این روش، نه تنها موقعیت مشابهی را که کاربران عادی وب نیاز دارند، اتخاذ می‌کند؛ بلکه شکل تعامل‌هایش با سرورها را تقلید می‌کند. خزش‌گرها از صفحه‌ها هسته شروع می‌شوند، آنها را تجزیه می‌کنند، پیوندها را برداشت می‌کنند، و سند پیوندی را واکنشی می‌کنند؛

سپس آن‌ها این پردازش با سند واکنشی شده را تکرار می‌نمایند و تا زمانی که پیوندهایی را کشف کنند (1) و سند را درون حوزه تعیین شده پیدا کنند این پردازش مورد نیاز است چون HTTP فرمانی را که باید فهرست کامل سند قابل دسترس بر روی سرور را باید بازگرداند که برای مثال بر خلاف FTP است. از این رو هر صفحه باید به وسیله پیوندی از صفحه‌ها دیگر «کشف» شود.

فناوری خزش، برای اهداف نمایه‌سازی توسعه یافته است (2). به کار بردن آن برای بایگانی وب، بر خلاف اینکه جنبه‌های بیشتری از این فناوری را دوباره استفاده می‌کند تغییراتی را برای آن ایجاد می‌کند.

نخست این که خزش گره‌های بایگانی باید برای واکنشی تمام فایل‌ها تلاش کنند هر چه فرمت‌شان برای بایگانی یک مدل کامل از سایت‌ها باشد بر خلاف خزش گره‌های موتور کاوش که معمولاً فقط فایل‌هایی را واکنشی می‌کند که آن‌ها بتوانند فهرست نمایند. خزش گره‌های موتور کاوش برای مثال اغلب از انتقال از

ص: 28

1- برای نظرات اخیر درباره فناوری خزشگر پلنت و دیگران (2004) و چاکرابارتی (2002) را ببینید.

2- برای نظر اخیر درباره توسعه موتور کاوش تجاری سونریج (1997) را ببینید.

فایل های کاربردی و ویدئویی بزرگ چشم پوشی می کند. بارگذاری این نوع فایل ها می تواند تفاوت قابل توجهی را در رابطه با زمان و پهنای باند مورد نیاز برای خزش در کل سایت ایجاد نماید.

تفاوت دوم، با مدیریت موقتی خزش ها مرتبط می باشد. برای اجتناب از اضافه بار سرورهای وب، خزش گر ها از قوانین مطلوب و معتبری استفاده می کنند (روشه، 2006). این مستلزم این است که تصرف یک وب می تواند در طی چندین دقیقه منتهای مراتب چندین ساعت و گاهی اوقات چندین روز طول بکشد. یک محاسبه ساده نشان می دهد زمانی را که در خصوص یک تأخیر 3 ثانیه ای بین دو درخواست، بیشتر از 3 روز برای بایگانی یک سایت با 100/000 صفحه طول خواهد کشید. این تأخیر، مبحث ثبات موقت تصرف را که دستخوش تغییراتی در طول زمان می شود، افزایش می دهد. اگر صفحه نمایه برای مثال در طی تصرف تغییر یابد روش بایگانی با آخرین بایگانی که با صفحه ها بایگانی پیوند داده شده است سازگار نمی باشد.

این یک مورد برای خزش گر های بایگانی است چون خزشگر برای تهیه محتوا فرض شده است البته نه فقط در جهت هدایت محتواها خزش گر های موتور کاوش عادت دارند که به صفحه ها زنده اشاره کنند. به این معنی که مضمون ابرمتن برای آن ها، یکی است که توسط سرور اصلی تهیه شده است (که البته در سراسر صفحه ها دارای ثبات است و به هنگام سازی شده است). بر عکس، خزش گر های بایگانی باید محتوا را به طور کلی تصرف کنند که انجام خواهند داد با وابستگی یا بدون وابستگی داخلی، به عنوان فقط مضمون هایی برای شناوری و تفسیر.

این امر به مراتب نتایج قابل دسترسی نسبت به خط مشی خزش گرها دارد. چون مطلوبیت برای سرورها حکم می کند که همیشه عملیات خیلی محدودی برای خزش داشته است، خزش گر های موتور های کاوش از اولویت خزش در سطح پهنای نخست استفاده کرده اند با برخی تغییرات اساسی که با خزش در بهترین صفحه های اول، هدفمند شده است (چو و همکاران، 1998؛ ناجون و هیدان 2001) (1)؛ ناجورک و واینر 2001 (2)؛ کاستیلو و همکاران، 2004؛ بازا - یس و کاستیلو (2005) (3).

اتخاذ این خط مشی، همچنین، روش برای به حداقل رسانیدن اثرات شدید تله های روباتیک (4) بر روی کل خزش با قرار گیری بیرون از خزش بر روی تعداد زیادی از سایت های مختلف می باشد.

اما این راهبرد زمان بندی خزش، در دسر افزایش تفاوت موقتی خزش ها در سطح سایت را دارد. بنابراین برای اتخاذ خزش های بایگانی در اولویت نخست یک سایت پیشنهاد شده است. (5) اما برای خزش های مقیاس بزرگ، هنوز برای بهینه سازی بازده خزنده با اطمینان از منابعی که در حداکثر ظرفیت خود استفاده شده اند، ضروری است تأخیر آشکار بین درخواست ها و منابع قابل دسترس خزش

ص: 29

Najork and Heydon -1

Wiener -2

Baeza-Yates and Castillo -3

4- فاصله ای عمده در یک سیستم پردازش اطلاعات که به منظور جمع آوری تغییر یا خراب کردن آتی اطلاعات به وجود آمده است
5- به طور موردی، برای ملزومات درون بحث شده است (میسانس، 2004)؛ در مورد سیاست های الگوریتم زمان بندی خزش که در سایت به عنوان بعدی عمودی آمیخته می شود منابع زیر را ببینید (Castillo et al. 2004, Baeza-Yates and Castillo (2005)).

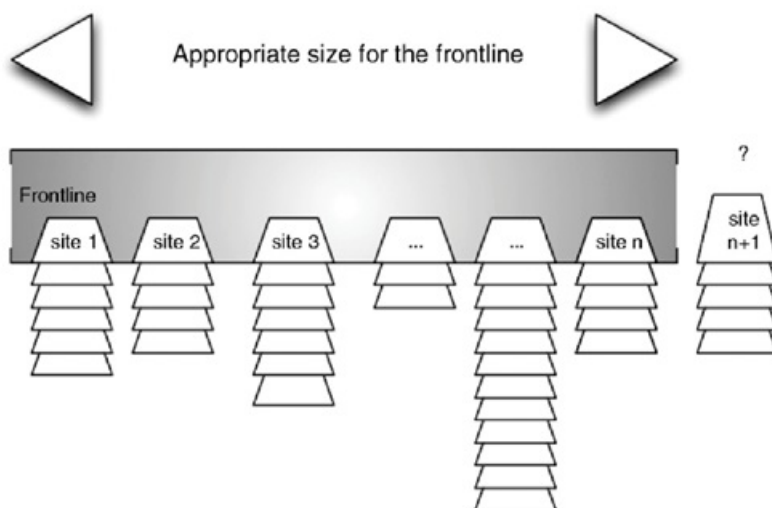
موجود است، شخص باید تعداد مطلوبی از سایت‌ها را برای شروع در زمان مشابه با ایجاد فرکانس‌های مطمئن جست و جو کند که به وسیله قوانین مطلوب بدون هیچ تأخیر غیر ضروری بین درخواست‌ها واقع شده‌اند. شکل 2-1 خط مقدم یک خزش‌گر را نشان می‌دهد که اندازه‌ای مناسب با

تخصیص بهینه منابع خزش دارد.

عکس

۳۰ مدیریت منابع اطلاعاتی وب

موجود است، شخص باید تعداد مطلوبی از سایت‌ها را برای شروع در زمان مشابه با ایجاد فرکانس‌های درخواست‌های مطمئن جست‌وجو کند که به وسیله قوانین مطلوب، بدون هیچ تأخیر غیر ضروری بین درخواست‌ها واقع شده‌اند. شکل 2-1 خط مقدم یک خزش‌گر را نشان می‌دهد که اندازه‌ای مناسب با تخصیص بهینه منابع خزش دارد.



تصویر 1-2. Frontline شامل سایت‌هایی است که باید با یک خزشگر یکسان و به‌طور همزمان خزش شوند. اندازه آن (n) در حد اَپتیمم است و اگر بین درخواست‌ها تأخیری رخ دهد، فقط به‌وسیله قوانین ساده‌ای محدود می‌شود و منابع خزش همچنان درگیر و مشغول خواهند بود. اگر n+1 سایت خزش شوند، محدودیت منابع خزش نوع تأخیر اضافی و عدم ربط موقت را به کار خواهند گرفت. اگر n-1 سایت خزش شوند، منابع بدون استفاده خواهند ماند.

محدودیت‌هایی برای هر آنچه که بتوان با استفاده از این روش بایگانی کرد وجود دارد. بیش از همه طی برداشت پیوند و برخی در طی بازیابی از طریق واسط HTTP رخ می‌دهد. دلیل مورد قبلی این حقیقت است که URL استخراج شده، به‌طور نامساعد شکل گرفته است یا از پارامترهای پیچیده استفاده کرده‌اند، یا به‌سختی برای تجزیه کردن URL از فایل آغازگر یا فایل اجرا یا حتی کد HTML استفاده کرده است. دومی می‌تواند به‌علت تجدید مسیرها، مذاکره محتوا، اجازه، پاسخ‌های تدریجی (کند)، اندازه نهایی، اتصالات TCP استثناهایی، پاسخ‌های سرور نامعتبر، و مانند آن باشد. با استفاده از این نوع ابزارها، فراهم‌آوری مقیاس بزرگی از متوا در یک مسیر کل نگ، که البته این طور نیست، مجاز می‌شود.

تصویر 1-2 Frontline شامل سایت هایی است که باید با یک خزشگر یکسان و به طور همزمان خزش شوند. اندازه آن (n) در حد اپتیمم است و اگر بین درخواست ها تأخیری رخ دهد فقط به وسیله قوانین ساده ای محدود می شود و منابع خزش همچنان درگیر و مشغول خواهند بود. اگر n+1 سایت خزش شوند، محدودیت منابع خزش نوع تأخیر اضافی و عدم ربط موقت را به کار خواهند گرفت. اگر n-1 سایت خزش شوند، منابع بدون استفاده خواهند ماند.

محدودیت هایی برای هر آن چه که بتوان با استفاده از این روش بایگانی کرد وجود دارد. بیش از همه طی برداشت پیوند و برخی در طی بازیابی از طریق واسط HTTP رخ می دهد. دلیل مورد قبلی این حقیقت است که URL استخراج شده به طور نامساعد شکل گرفته است یا از پارامترهای پیچیده استفاده کرده اند یا به سختی برای تجزیه کردن URL از فایل آغازگر یا فایل اجرا یا حتی کد HTML استفاده کرده است. دومی می تواند به علت تجدید مسیرهها، مذاکره محتوا، اجازه، پاسخ های تدریجی (کند)، اندازه نهایی، اتصالات TCP استثنایایی، پاسخ های سرور نامعتبر، و مانند آن باشد.

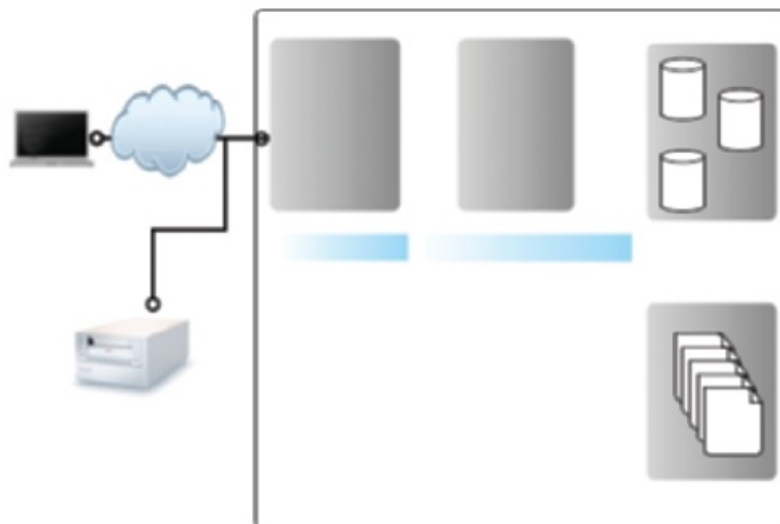
با استفاده از این نوع ابزارها، فراهمآوری مقیاس بزرگی از متوا در یک مسیر کل نگ، که البته این طور نیست. مجاز می شود.

بایگانی تراکنشی (شکل 3-1 را ببینید) به وسیله فیچ (2003) (1) پیشنهاد شده است که شامل تصرف و بایگانی «تمام پاسخ های متمایز اساسی که توسط یک وبگاه تولید شده است و در رابطه با محتوایشان و چگونگی تولید می باشد». این کار در سیستم پرش صفحه (2) با استفاده از یک فیلتر درون ورودی سرور وب (درخواست) جریان و جریان خروجی (پاسخ) اجرا شده است. این توابع عملیاتی اکنون بر روی برخی از نظام های مدیریت محتوای وب مانند Vignette TM قابل دسترس هستند.

عکس

۴-۲-۲- بایگانی تراکنشی

بایگانی تراکنشی (شکل ۳-۱ را ببینید) به وسیله فیچ (۲۰۰۳) پیشنهاد شده است که شامل تصرف و بایگانی «تمام پاسخ‌های متمایز اساسی که توسط یک وبگاه تولید شده است و در رابطه با محتوایشان و چگونگی تولید می‌باشد». این کار در سیستم پرش صفحه^۱ با استفاده از یک فیلتر درون ورودی سرور وب (درخواست) جریان و جریان خروجی (پاسخ) اجرا شده است. این توابع عملیاتی اکنون بر روی برخی از نظام‌های مدیریت محتوای وب مانند Vignette TM قابل دسترس هستند.



شکل ۳-۱ بایگانی تراکنشی

جفت‌های درخواست/ پاسخ منحصر به فرد ذخیره و بایگانی شده‌اند از این رو، ایجاد یک بایگانی کامل از تمام محتوا برای یک سایت مشخص پیش‌بینی شده است. درخواست‌ها فقط با اندکی تفاوت (غیر مادی) به صورت منحصر به فرد مورد بررسی قرار می‌گیرند، به استثنای محاسبه مجموع مقابله‌ای قسمتی از کد، که آنها را به صورت رمز در آورده است، کیفیت دقیق اینها چگونه می‌تواند با تعداد زیادی از روش‌های شخصی‌سازی محتوا، که واضح نمی‌باشد، منطبق گردد. این نوع بایگانی وب می‌تواند فواید پیگردی و ثبت هر برنامه‌ریزی ممکن در محتوا را ثابت کند.

1. Fitch
2. <http://www.projectComputing.com/products/pageVault>

شکل 3-1 بایگانی تراکنشی

جفت‌های درخواست/ پاسخ منحصر به فرد ذخیره و بایگانی شده‌اند از این رو، ایجاد یک بایگانی کامل از تمام محتوا برای یک سایت مشخص پیش‌بینی شده است. درخواست‌ها فقط با اندکی تفاوت («غیر مادی») به صورت منحصر به فرد مورد بررسی قرار می‌گیرند، به استثنای محاسبه مجموع مقابله‌ای قسمتی از کد، که آن‌ها را به صورت رمز در آورده است،

کیفیت دقیق اینها چگونه می‌تواند با تعداد زیادی از روش‌های شخصی‌سازی محتوا، که واضح نمی‌باشد، منطبق گردد.

این نوع بایگانی وب می تواند فواید پیگردی و ثبت هر برنامه ریزی ممکن در محتوا را ثابت کند.

ص: 31

Fitch -1

<http://www.projectComputing.com/products/page Vault> -2

محتوایی که هرگز دیده نشده، بایگانی نخواهد شد (همان طور که ذکر شد، بوفخواد و وینناد (2003) (1)، برآورد کردند که 25 درصد صفحه ها از یک وبگاه بزرگ علمی، هرگز قابل دستیابی نیستند). اما محتوای وب پنهان تا زمانی که به دستیابی برسد ثبت خواهد شد و این یک مزیت مهم است.

محدودیت مهم این روش این است که باید با موافقت و همکاری مالک، سرور، اجرا شود. از این رو، برای بایگانی داخلی وب نشان داده شده است. و این مسئله مزیت توانا بودن ثبت دقیق از چیزی را که و زمانی را که دیده است دارا می باشد. برای بایگانی شرکت یا سازمانی، اغلب به وسیله جواب گویی قانونی، برانگیخته می شود. این کار، می تواند یک مزیت باشد. حتی برای ترکیب با اطلاعات از سرور ثبت وقایع درباره کسی که محتوا را دیده است، امکان پذیر باشد. به طور واضح، هر چه را که می تواند به عنوان یک مزیت برای بایگانی وب داخلی دیده شود ممکن است برای یک بایگانی عمومی مشکل باشد، چون می تواند تأکید بر محرمانگی جدی را بالا ببرد. اما به هر حال محتوای قابل استفاده نمی باشد.

3-2-4-بایگانی سرور - جانبی

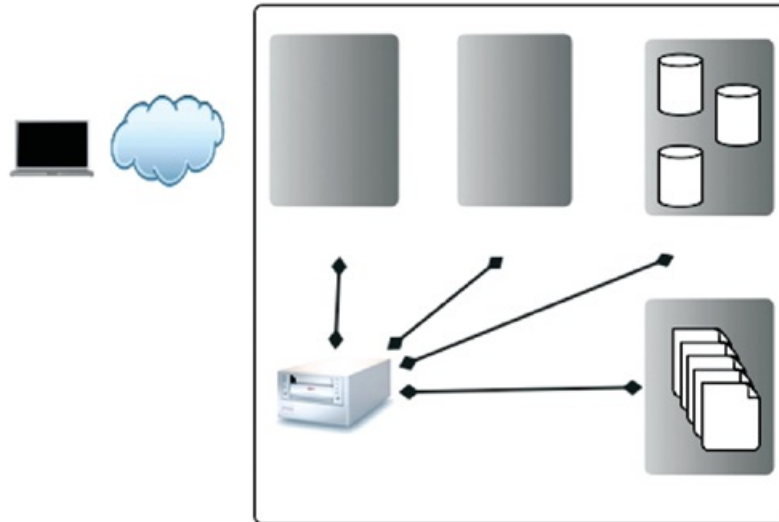
نوع آخر از روش فراهم آوری بایگانی وب، کپی مستقیم فایل ها از سرور، بدون استفاده از واسطه های HTTP است. این روش مانند روشی قبلی می تواند فقط با همکاری مالکان سایت استفاده گردد (شکل 4-1).

عکس

محتوایی که هرگز دیده نشده، بایگانی نخواهد شد (همان‌طور که ذکر شد، بوفخواد و وینناد^۱، ۲۰۰۳ برآورد کردند که ۲۵ درصد صفحه‌ها از یک وبگاه بزرگ علمی، هرگز قابل دستیابی نیستند). اما محتوای وب پنهان تا زمانی که به دستیابی برسد، ثبت خواهد شد و این یک مزیت مهم است. محدودیت مهم این روش، این است که باید با موافقت و همکاری مالک سرور، اجرا شود. از این رو، برای بایگانی داخلی وب نشان داده شده است. و این مسئله مزیت توانا بودن ثبت دقیق از چیزی را که و زمانی را که دیده است دارا می‌باشد. برای بایگانی شرکت یا سازمانی، اغلب به وسیله جوابگویی قانونی، برانگیخته می‌شود. این کار، می‌تواند یک مزیت باشد. حتی برای ترکیب با اطلاعات از سرور ثبت وقایع درباره کسی که محتوا را دیده است، امکان‌پذیر باشد. به‌طور واضح، هر چه را که می‌تواند به‌عنوان یک مزیت برای بایگانی وب داخلی دیده شود، ممکن است برای یک بایگانی عمومی مشکل باشد، چون می‌تواند تأکید بر محرمانگی جدی را بالا ببرد. اما به هر حال محتوا قابل استفاده نمی‌باشد.

۴-۲-۳-بایگانی سرور- جانبی

نوع آخر از روش فراهم آوری بایگانی وب، کپی مستقیم فایل‌ها از سرور، بدون استفاده از واسطه‌های HTTP است. این روش مانند روشی قبلی می‌تواند فقط با همکاری مالکان سایت استفاده گردد (شکل ۴-۱).



شکل ۴-۱. آرشیو Server - side: قطعه‌های مختلف اطلاعات به‌طور مستقیم از سرورها آرشیو می‌شوند. تولید نمونه کار محتوای آرشیو شده و یک نمونه پشتیبان (backup) از فایل‌ها چالش اصلی این روش است.

1. Boufkhad and Viennot

اگر چه بسیار ساده به نظر می رسد، در واقع مشکل زیادی را برای ایجاد محتوای کپی شده قابل استفاده افزایش می دهد و حتی در مورد فایل های ایستای HTML، شخص می تواند به زحمت در محتوا از طریق پیوندهای مطلق به عنوان نام دامنه هدایت شود که در بایگانی متفاوت خواهند بود. اما بیشترین مشکل ناشی از محتوای توسعه یافته پویاست که تکه های به هم پیوسته محتوا از منابع گوناگون (قالب ها (1) و پایگاه های داده) است که توسط درخواست های کاربرد در حالت پرواز در ارتفاع کم توسعه یافته است. کپی برداری فایل های پایگاه داده ها، قالب ها، و فایل های آغاز گر به این معنی نیست که آن برای تولید مجدد محتوا از بایگانی آسان خواهد بود. بر عکس یک وظیفه چالش برانگیز خواهد بود چون نیازمند توانا شدن برای اجرای در محیط مشابه، با پارامترهای مشابه در بایگانی است. در واقع، وقتی امکان پذیر شد، محتوای توسعه یافته پویا در شکل نهایی اش بهتر حفظ می شود، معمولاً در فایل های HTML مسطح بهتر حفظ می شود (برای مثال این موردی برای بیشتر CMSها، بلاگ ها و ویکی ها (2) صورت می گیرد).

اما گاهی اوقات این کار مشکل است، حتی برای خزش گرهای یافتن مسیر برای برخی اسناد یک وبگاه و فایل هایی که می تواند از طریق یک تعامل پیچیده به دستیابی برسند، غیر ممکن است (مانند ارسال یک پرس و جو و یک فرم یا پنجره) که به سختی توسط خزش گرها تصرف خواهد شد این بخش وب، وب پنهان یا عمیق نامیده شده است (برگمان، 2001؛ چانگ و همکاران 2004) که بزرگ تر از وب سطحی است و (همچنین به طور عمومی وب قابل نامیده شده است). (3)

در این مورد، بایگانی جانبی سرور می تواند یک راه حل باشد، همان طور که در بالا ذکر شد، به مشارکت فعال مجری سایت نیاز دارد. بیشتر از یک پشتیبانی ساده است که دسترسی به محتوا را در نمونه های اصلی اش تضمین نمی کند، آن بر توانا بودن برای «نمایش» مجدد سایت در محیط بایگانی را دلالت می کند. این کاهش وابستگی به پایگاه داده ها و اجرای فایل های آغاز گر جانبی سرور به قدری که امکان پذیر است را نشان می دهد. این کار می تواند به وسیله استخراج اطلاعات ساخت یافته محتوی در پایگاه داده ها و انتقال آن درون XML انجام شود. نوعی معماری اطلاعات دروازه اسناد نامیده شده است که شامل اسناد غیر وب با آن هایی که به وسیله کاتالوگی می تواند مانند اینها بایگانی شود، قابل دستیابی می باشد. این کار، برای چندین سایت انجام شده است که به مقوله سایت های پنهان به وسیله کتابشناسی ملی فرانسه (4) وابسته می باشد (فصل 5 را ببینید).

این کار فقط در چارچوب واسپاری قانونی امکان پذیر می باشد که در فرانسه مانند بسیاری از کشورهای دیگر به کار گرفته می شود. حقیقت این است که وب پنهان همچنین غالب اوقات دارای محتوای بسیار غنی با این نوع معماری اطلاعات است که انبوه زیادی از محتوای از قبل موجود بر روی وب منتشر کرده است. عمومیت این نوع معماری اطلاعات، بایگانی سایت سرور را می سازد، روشی که به جایی که می تواند به کار رود توجه دارد.

ص: 33

templatestemplates -1

wikis -2

3- این اصطلاح برای قسمتی از وب مشخص شده است که می تواند توسط خزش گرهای نمایه شود (لورنس و گیلز، 1998، 1999)

Bibliothèque nationale de France -4

همان طور که قبلاً دیده شد کپی برداری از یک وبگاه، یک وظیفه غیر پیش پا افتاده است. آن در واقع بر ایجاد مجدد یک سیستم اطلاعاتی اشاره می کند که برای کاربران قابل دسترسی خواهد بود. همان طور که آنتونیول و همکارانش (1) آن را در وبگاه قرار دادند «وبگاه ممکن است به سادگی یک فایل واحد یا یکی از پیچیده ترین مجموعه های محصولات مصنوعی نرم افزاری مشترک باشد که تا به حال درک شده است». به طور مطلوب، بایگانی می توانست در اصل متناظر (هم ریخت) (ساختار سلسله مراتبی مشابه، نام گذاری فایل ها، سازوکار پیونددهی، چارچوب) باشد، اما به دلایل عملی اینگونه نیست. همان طور که در بخش قبلی دیده شده اکتساب سایت ها در حقیقت موارد یک تغییر شکل فایل را به طور مؤثر کاهش می دهد.

چالش بیشتر، در ایجاد مجدد سیستم های اطلاعاتی مشابه می باشد. نظام اطلاعاتی وب معماری اطلاعاتی پیچیده را نشان می دهد که به سیستم های عملیاتی خاص پیکربندهای سرور و محیط کاربردی وابسته هستند که در بیشتر موارد حتی برای ایجاد مجدد استفاده از (فایل) چرکنویس برای طراحان و مدیران مشکل است به همین دلیل است آرشیویست ها مجبور به اتخاذ راهبرد تبدیل یا تغییر شکل می شوند. این تبدیل ها می توانند اثر شدیدی بر ساز و کارهای آدرس دهی و پیوندی، فرمت ها، و همچنین تغییر خود شیء تأثیر داشته باشند.

تاکنون سه راهبرد برای ساخت بایگانی وب اتخاذ شده است. راهبرد نخست برای ایجاد کپی محلی از فایل های سایت و هدایت از طریق این کپی در یک مشابه برای مثال بر روی وب می باشد.

راهبرد دوم، برای اجرای یک سرور وب با به کار بردن مضمون در یک محیط برای مرورگرهای کاربران است راهبرد دوم راه اندازی سرور وب و به کارگیری محتوا در این محیط برای مرورگرهای کاربران است. سومین راهبرد، سازماندهی مجدد اسناد منطبق بر نام گذاری متفاوت (غیر وب)، آدرس دهی منتقل کردن است. بخش های زیر موافقان و مخالفان این رهنمودهای متفاوت و همچنین موارد استفاده را نشان می دهد.

1-3-4- نظام فایل های محلی به خدمت گرفته شده بایگانی ها

توصیف

این نوع بایگانی (شکل 5-1) بر اساس احتمالاتی است که مشخصه های URL از پیشوند سیستم فایل محلی، «فایل» (2) استفاده کند در یک شمای URL برای کپی و دسترسی به فایل های محلی از وبگاه اصلی استفاده کند مانند:

`Http://www.example.org/example.HTML`

`File://Users/archire2005/eample.org/example.HTML`

این، استفاده از سیستم فایل محلی را برای شناسایی از طریق مواد بایگانی شده در وب، قادر می سازد.

Antoniol et al -1

file -2

همچنین به استفاده از شکل جزئی (وابسته) از URL که نه تنها از پیشوند دوری می کند بلکه همچنین از نام سرور و مسیر هدف نیز اجتناب می کند.

مرورگرهای استاندارد می توانند به طور مستقیم باز شوند (به طور مثال بدون سرور وب) و چنین فایل های ذخیره شده را محلی نمایند و تا زمانی که پیوندها اسناد وابسته هستند، شناسایی در بایگانی مشابه مانند همانی است که در سایت اصلی خواهد بود، به جز فقط در قسمت آدرس دهی (1) مرورگر وقتی در جست و جوی پیشوند URL هستند. (در اینجا «فایل» به جای HTTP می باشد).

توضیح

فایده اصلی این راهبرد، تسهیل دستیابی به بایگانی با مسیره‌ی ساختار وبگاه اصلی به سوی نظام بایگانی فایل است. با استفاده از مرورگر استاندارد و سیستم فایل مجاز به اجتناب از هزینه های بالا که همراه است با اجرای دسترسی به بایگانی مبتنی بر سرور. بنابراین، حتی گروهی با مهارت های فنی فناورانه بسیار اساسی می توانند این نوع بایگانی را بر پا و اجرا کنند. ولی محدودیت هایی در این رویکرد وجود دارد. از یک نگرش محافظه کارانه نقص اصلی این است که تبدیل های گوناگون فایل های اصلی، مورد نیاز می باشند. از این رو، صداقت محض برای فایل های اصلی، نمی تواند مورد ملاحظه قرار گیرد مگر به وسیله سندسازی با تغییرات دقیق که در فایل اصلی به کار می رود یا به وسیله نگهداری یک کپی از اصل. انتقال محتوا در دو سطح در رویکرد بایگانی «FS» محلی مورد نیاز می باشد.

نخست، به علت تفاوت در نام قرار داد بین URL و نظام فایل (تصویب و ذخیره کاراکترها، رهایی از قوانین، حساسیت موردی)، نام اشیاء ممکن است به تغییراتی نیاز داشته باشد (فصل 1 - B را برای نشان دادن جزئیات بیشتر این تغییرات ببینید). در مورد جایی که صفحه طبق پارامترها پرس و جو شده و به طور چشمگیری تولید شده است، یک نام حتی برای ایجاد شدن در صفحه نتیجه، پارامترهایی را برای اطمینان از منحصر به فرد بودن صفحه بایگانی شده تأمین می کند.

دوم اینکه، پیوندهایی مطلق باید در پیوندهای مرتبط در خودکد صفحه منتقل شوند تا شناسایی مبتنی بر نظام فایل را امکان پذیر کنند. حتی اگر این کار بتواند به سادگی توسط تغییر شکل URL اصلی درون یک توضیح در کد، سندیت داده شود این دستکاری منبع را نشان می دهد.

از نظر عملی، اشکالات اصلی ناشی از خود نظام فایل می باشد که یک معماری اطلاعات متفاوت قابل توجه نسبت به وب دارد. نخست اینکه، سازماندهی بایگانی باید در سازماندهی سلسله مراتبی نظام های فایل مناسب باشد، با این حال، یک بایگانی نه تنها از ترکیب سایت ها، بلکه همچنین گروه هایی در سایت ها (مجموعه ها) و نسخه هایی از سایت ها تشکیل شده است. مسیره‌ی این سازماندهی برای یک ساختار سلسله مراتبی، بدون تغییر و انتخاب روی نمی دهد. چگونه سایت ها در یک حالتی که زمان پایدار، یک مبحث مهم برای بررسی می باشد، با یکدیگر در یک گروه قرار می گیرند. نام های مجموعه

ص: 35

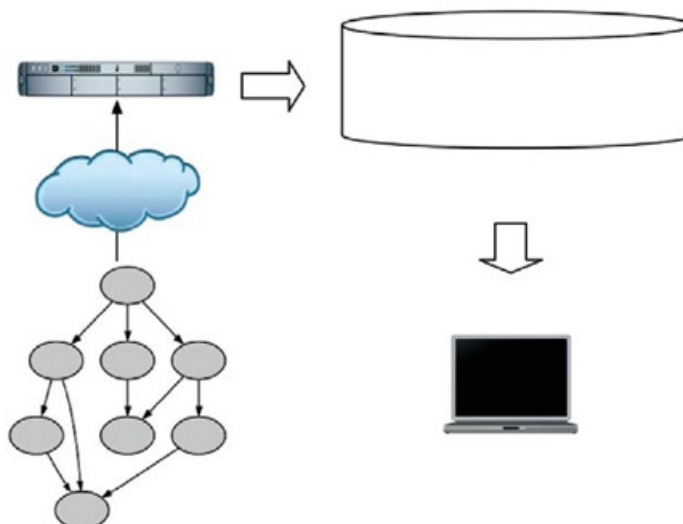
باید دائمی باشند و زمان گروه بندی نیز باید برای ضبط فرکانس ها، تطبیق یافته باشد. در مسیر تمام این مباحث، تصمیم گیری کامل باید پیشاپیش انجام شده باشد. آن ها بر روی چگونگی انتخاب ساختار که در برابر توسعه مجموعه مقاومت می کنند، تحت فشار قرار می گیرند. سازماندهی زمان در اجرای هر جمله از یک برنامه برای اهداف اشکال زدایی (از یک مدل به مدل دیگر) سایت نیز یک مبحث مهم برای هر یک از لایه های نرم افزار است که باید به سیستم فایل با استاندارد بالا افزوده شود. این لایه باید قادر باشد حداقل مدل های متفاوت سایت ها را که به تاریخ شان بستگی دارد با یکدیگر متصل نماید (نسخه پردازی) و این کار را با یک واسط کاربر مناسب برای شناسایی از طریق زمان به سادگی ارائه می دهد. این کار اغلب با استفاده از مدیریت خارجی یک پایگاه داده در سایت ها و ضبط اطلاعات و ابزارهایی برای ایجاد صفحه های نمایش میانجی با فهرستی از تاریخ اجرا می شود که در آن سند بایگانی شده است.

محدودیت دیگر این رویکرد، ناشی تعداد عظیم فایل های بایگانی وب است که باید جا به جا شوند. این مسئله عادی است که بایگانی هایی با بیلیون ها صفحه را ببینیم. این شکل به محدودیت های ظرفیت نظام های فایل جاری می رسد و دسترسی می یابد حتی وقتی یک نظام فایل بتواند این مقدار فایل را جا به جا کند، عملکرد می تواند تحت تأثیر قرار گیرد. برای کم کردن باری که بر روی نظام فایل قرار گرفته، بایگانی با مقیاس بزرگ برای فایل های مخزن استفاده شده است. البته این کار، ارتباط مستقیم در نام گذاری و پیوند دادن را که رویکرد بایگانی نظام فایل محلی برای اتخاذ رویکرد دوم ارائه می دهد را می شکند و بایگانی مبتنی بر سرور به کار رفته در وب برای تحویل محتوا از این فایل مخزن است.

شکل 1-5. بایگانی سیستم فایل محلی. فایل اصلی مورد خزش قرار می گیرد و صفحه ها و سایر فایل ها به صورت انفرادی روی سیستم فایل بایگانی ذخیره می شوند. دست یابی بوسیله شناسایی مستقیم در نظام فایل انجام می شود.

عکس

باید دائمی باشند و زمان گروه‌بندی نیز باید برای ضبط فرکانس‌ها، تطبیق یافته باشد. در مسیر تمام این مباحث، تصمیم‌گیری کامل باید پیشاپیش انجام شده باشد. آنها بر روی چگونگی انتخاب ساختار که در برابر توسعه مجموعه مقاومت می‌کنند، تحت فشار قرار می‌گیرند. سازماندهی زمان در اجرای هر جمله از یک برنامه برای اهداف اشکال‌زدایی (از یک مدل به مدل دیگر سایت) نیز یک مبحث مهم برای هر یک از لایه‌های نرم‌افزار است که باید به سیستم فایل با استاندارد بالا افزوده شود. این لایه باید قادر باشد حداقل مدل‌های متفاوت سایت‌ها را که به تاریخ‌شان بستگی دارد با یکدیگر متصل نماید (نسخه‌برداری) و این کار را با یک واسط کاربر مناسب برای شناسایی از طریق زمان به سادگی ارائه می‌دهد. این کار اغلب با استفاده از مدیریت خارجی یک پایگاه داده در سایت‌ها و ضبط اطلاعات و ابزارهایی برای ایجاد صفحه‌های نمایش میانجی با فهرستی از تاریخ اجرا می‌شود که در آن سند بایگانی شده است. محدودیت دیگر این رویکرد، ناشی تعداد عظیم فایل‌های بایگانی وب است که باید جابه‌جا شوند. این مسئله عادی است که بایگانی‌هایی با بیلیون‌ها صفحه را ببینیم. این شکل به محدودیت‌های ظرفیت نظام‌های فایل جاری می‌رسد و دسترسی می‌باید حتی وقتی یک نظام فایل بتواند این مقدار فایل را جابه‌جا کند، عملکرد می‌تواند تحت تأثیر قرار گیرد. برای کم کردن باری که بر روی نظام فایل قرار گرفته، بایگانی با مقیاس بزرگ برای فایل‌های مخزن استفاده شده است. البته این کار، ارتباط مستقیم در نام‌گذاری و پیوند دادن را که رویکرد بایگانی نظام فایل محلی برای اتخاذ رویکرد دوم ارائه می‌دهد را می‌شکند و بایگانی مبتنی بر سرور به‌کار رفته در وب برای تحویل محتوا از این فایل مخزن است.



شکل ۱-۵. بایگانی سیستم فایل محلی. فایل اصلی مورد خزش قرار می‌گیرد و صفحه‌ها و سایر فایل‌ها به صورت انفرادی روی سیستم فایل بایگانی ذخیره می‌شوند. دستیابی بوسیله شناسایی مستقیم در نظام فایل انجام می‌شود.

استفاده ترجیحی

این روش برای بایگانی سایت شرکت ها یا سازمانی و بایگانی که در مقیاس کوچکی توسعه یافته، توصیه شده است. بسته به استفاده از این بایگانی، صحت مبحث باید به دقت مورد بررسی قرار گیرد مخصوصاً برای بایگانی سازمانی برای بایگانی با مقیاس کوچک، تعادل بین سختی برای سازماندهی مداوم مجموع فایل ها و سادگی دسترسی ایجاد شده توسط این رویکرد باید بر اساس «مورد به مورد» ارزیابی گردد. برای بایگان های وب با مقیاس بزرگ و متوسط از این روش باید اجتناب شود.

ابزارها

این راهبرد برای بایگانی وب در مقیاس کوچک و متوسط ساده ترین رویکرد است که با بسیاری از ابزار های موجود مانند HTTrack قابل دسترس است.

2-3-4-بایگانی های مبتنی بر خدمت وب

اشاره

*بایگانی های مبتنی بر خدمت وب (1)

با وجود تقاضاهای بیشتر، این گزینه تطبیق بهتری را برای نام گذاری منبع و ساختار اسناد ایجاد می کند (شکل 6-1). همچنین، برای اجتناب از محدودیت های اندازه نظام فایل که برای بایگانی وب در مقیاس بزرگ بحرانی است، مجاز می کند.

عکس

استفاده ترجیحی

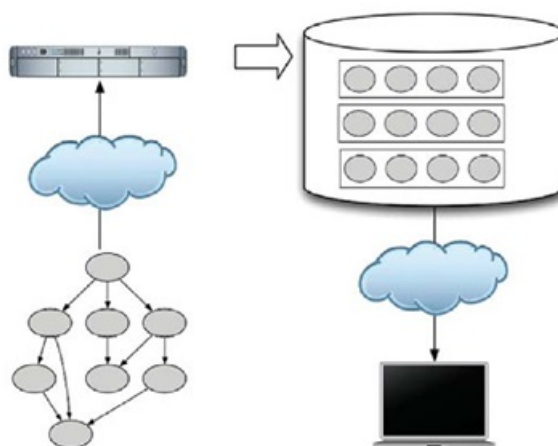
این روش برای بایگانی سایت شرکت ها یا سازمانی و بایگانی که در مقیاس کوچکی توسعه یافته، توصیه شده است. بسته به استفاده از این بایگانی، صحت مبحث باید به دقت مورد بررسی قرار گیرد مخصوصاً برای بایگانی سازمانی. برای بایگانی با مقیاس کوچک، تعادل بین سختی برای سازماندهی مداوم مجموع فایل ها و سادگی دسترسی ایجاد شده توسط این رویکرد باید براساس «مورد به مورد» ارزیابی گردد. برای بایگان‌های وب با مقیاس بزرگ و متوسط از این روش باید اجتناب شود.

ابزارها

این راهبرد برای بایگانی وب در مقیاس کوچک و متوسط ساده ترین رویکرد است که با بسیاری از ابزار های موجود مانند HTTrack قابل دسترس است.

۴-۳-۲-بایگانی های مبتنی بر خدمت وب^۱

با وجود تقاضاهای بیشتر، این گزینه، تطبیق بهتری را برای نام گذاری منبع و ساختار اسناد ایجاد می کند (شکل ۶-۱). همچنین، برای اجتناب از محدودیت های اندازه نظام فایل که برای بایگانی وب در مقیاس بزرگ بحرانی است، مجاز می کند.



شکل ۶-۱. مدل مبتنی بر خدمت وب. سایت اصلی مورد خزش قرار می گیرد و پاسخ هادر مخزن بدون تغییر ذخیره می شوند (فایل های) که اجتناب از پیمان نامه مسیریابی نام گذاری فایل های نظام و تغییر ساختار پیوند را اجازه می دهند. دستیابی نیازمند یک سرور وب است که محتوا را در مخازن واگشی کند و آن را به عنوان یک پاسخ به کاربر بفرستد.

1. Web-served archive

شکل ۶-۱. مدل مبتنی بر خدمت وب. سایت اصلی مورد خزش قرار می گیرد و پاسخ ها در مخزن بدون تغییر ذخیره می شوند (فایل های) که اجتناب از پیمان نامه مسیریابی نام گذاری فایل های نظام و تغییر ساختار پیوند را اجازه می دهند. دست یابی نیازمند یک سرور وب است که محتوا را در مخازن واگشی کند و آن را به عنوان یک پاسخ به کاربر بفرستد.

این روش، بر اساس بایگانی پاسخ است (در مقایسه با اولی که بر اساس بایگانی فایل است). پاسخ ها از سرور اصلی (منبع) بدون تغییر در فایل های مخزن (1) WARC ذخیره می شوند که اجازه می دهد تا پشت سر آخرین تا کاربر بایگانی با یک سرور HTTP خدمت رسانی را انجام دهد.

پیشینه های یک فایل WARC (فایل آرشیوی مبتنی بر خدمت) رشته ای از فایل های وب را ثبت می کند هر صفحه به وسیله یک (برچسب) سرآمد که به طور مختصر محتوای حاصل شده و طولش توصیف کند، پیشی می گیرد. به علاوه، محتوای اولیه ثبت شده و WARC (فایل آرشیوی مبتنی بر خدمت) محتوای ثانویه را نیز در بر می گیرد مانند ابر داده و نقل و انتقالات فایل اصلی (منبع). اندازه یک فایل WARC (فایل آرشیوی مبتنی بر خدمت) می تواند تا حدود صدها مگابایت فرق داشته باشد. هر پیشینه یک نقطه شروع دارد که دسترسی مستقیم به پیشینه های انفرادی را بدون بارگیری و تجزیه تمام فایل های WARC (فایل آرشیوی مبتنی بر خدمت) انجام می دهد. نقاط شروع پیشینه های انفرادی در یک نمایه مرتب شده توسط URL ذخیره می شوند. از این رو، به سرعت برای پیشینه های انفرادی که بر اساس URL شان خارج از یک مجموع فایل WARC (فایل آرشیوی مبتنی بر خدمت) هستند، استخراج می شوند؛ که برای دسترسی در حالت شناوری تطبیق یافته اند. سپس ثبت هایی برای سرور وب تصویب می شوند که آن ها را سرویس گیرنده تهیه کرده است.

نگهداری و حفاظت از طرح نام گذاری شمای اصلی منبع (شامل پارامترهایی در صفحه های پویا)، شناوری در سایت را مجاز می کند همان طور که خز شده است کاربر بایگانی می تواند از تمام سیرهای پیروی شده توسط خزشگر، بار دیگر پیمایش نماید (حرکت کند).

توضیح

مزیت عمده استفاده از مخزن های WARC امکانات غلبه یافتن بر محدودیت سیستم فایل ذخیره سازی در رابطه با اندازه (تعداد کمی از فایل های انفرادی در نهایت در نظام فایل بایگانی ذخیره می شوند) و فضای نام (نام گذاری انفرادی فایل های وب می تواند حفظ و نگهداری شود) می باشد. دستیابی به بایگانی اینترنت از طریق Wayback Machin (که دسترسی به 500tb از مجموعه وب را می دهد) نشان می دهد که این رویکرد به نسبت ثابت افزایش داشته است نه مانند دیگران. اشکال این رویکرد این است که دسترسی مستقیم به فایل های ذخیره شده غیر ممکن است. دو لایه اضافی استفاده شده برای دسترسی به محتوا ضروری هستند:

نظام نمایه فایل WARC و سرور وب. این دو لایه، پیچیدگی برجسته ای ندارند، اما برای دستیابی به محیط به اجرا نیاز دارد که می تواند برای راه اندازی و نگهداری در سازمان های کوچک سخت باشد. این میانجی گر می تواند همچنین مشکلاتی را برای انتقال محتوا افزایش دهد چون نیاز دارد که ساز و کارهای

پیوند به طور مناسب از محیط وب زنده به محیط بایگانی، مسیره‌ی نمایندگی (فرض می‌کنیم که پیوندهای اصلی بدون تغییر در بایگانی نگهداری می‌شوند که این مزیت مهم این روش است). این کار می‌تواند در سطح نمایش صفحه و در سطح نماینده انجام شود.

گزینه نخست در افزودن به صفحه فرستاده شده به مرورگر کاربر بایگانی را در بر می‌گیرد که متن سند این کار را خواهد کرد، در خزش تفسیر مجدد پیوندها در صفحه ای به نقطه ای در بایگانی انجام خواهد شد (یا آن‌ها را در یک شکل مرتبط تفسیر می‌کند) بایگانی اینترنت برای مثال این کار را با پیروی از کد جاوا-اسکریپت (1) هر صفحه فرستاده شده برای کاربران پیوست شده انجام می‌دهد.

```
<>>SCRIPT language="Javascript<
```

```
<!--
```

```
FILE ARCHIVED ON 20050308085053 AND RETRIEVED FROM THE//
```

```
.INTERNET ARCHIVE ON 20060514055212//
```

```
JAVASCRIPT APPENDED BY WAYBACK MACHINE, COPYRIGHT//
```

```
.INTERNET ARCHIVE
```

```
ALL OTHER CONTENT MAY ALSO BE PROTECTED BY COPYRIGHT//
```

```
.U.S.C(17
```

```
)).SECTION 108(a)(3//
```

```
!";var sWayBackCGI ="http://web.archive.org/web/20050308085053
```

```
)function xLateUrl(aCollection, sProp
```

```
var i = 0( )
```

```
++)for(i = 0; i < aCollection.length; i
```

```
if (aCollection[i][sProp].indexOf("mailto:") == -1
```

```
)aCollection[i][sProp].indexOf("javascript:") == -1
```

```
];aCollection[i][sProp] = sWayBackCGI + aCollection[i][sProp
```

```
var i = 0 }
```

```
++)for(i = 0; i < aCollection.length; i  
  
if (aCollection[i][sProp].indexOf('mailto:') == -1  
  
)aCollection[i][sProp].indexOf('javascript:') == -1  
  
];aCollection[i][sProp] = sWayBackCGI + aCollection[i][sProp  
  
var i = 0 }
```

```
++)for(i = 0; i < aCollection.length; i
```

```
if (aCollection[i][sProp].indexOf('mailto:') == -1
```

ص: 39

```
)aCollection[i][sProp].indexOf('javascript:') == -1
```

```
];aCollection[i][sProp] = sWayBackCGI + aCollection[i][sProp
```

```
}
```

```
");if (document.links) xLateUrl(document.links, "href
```

Web Archiving: Issues and Methods 35 1

```
");if (document.images) xLateUr (document.images, "src
```

```
");if (document.embeds) xLateUrl(document.embeds, "src
```

```
+if (document.body document.body.background( document.body.background = sWayBackCGI
```

سند. بدنه. زمینه

```
//->
```

```
>SCRIPT</
```

```
>HTML</
```

مشکل این روش این است که برخی پیوندها (جاسازی شده در متن سند) [\(1\)](#) تفسیر نخواهد شد و از این رو، در نقطه ای در وبگاه منبع می مانند. در برخی موارد، تفسیر کد صفحه، برخی رفتارها را فعال می کند، مانند تغییر مسیر حتی قبل از اینکه کد پیوست شده به عنوان مرورگر جدید تفسیر گردد. منتظر نباشید تا سند کامل برای تفسیر به دست آید و آن را نمایش دهد.

با استفاده از نماینده ای که تمام درخواست ها از مرورگر کاربر به بایگانی تغییر مسیر می دهد، بسیار مؤثرتر و کارآمد می باشد همانطور که مسیر دهی بعد از تفسیر پیوند رخ می دهد توسط تعامل کاربر (با کلیک کردن) و مرورگر انجام می شود و کد را برای تولید درخواست مناسب تفسیر می کند (HTML سند متن سرویس گیرنده جانبی فرمت های دیگر) این کارآمدترین به عنوان ظرفیت مرورگر اصلی برای تفسیر مجموعه کدهای استاندارد برای کاری که معمولاً بر روی وب انجام می شود. این رویکرد، نیاز به ایجاد یک نماینده دارد که به یگانی تغییر مسیر دهد و به یک پارامتر که یک مرورگر از آن استفاده کند تا بتواند تقاضاهای زیادی برای یک محیط بایگانی پیوسته داشته باشد. استفاده از اتصال مرورگر برای مدیریت انتقال شکل باز به محیط نماینده می تواند این کار را برای کاربران نهایی آسان سازد.

استفاده ترجیحی

این روش برای بایگانی مقیاس بزرگ و متوسط و همچنین بایگانی های کوچک تر مناسب است که در زمینه حفاظت از صحت محتوا

هستند. چون این روش ها ، پاسخ ها را از سرور منبع همانطور که از سرویس گیرنده دریافت شده است ذخیره می کنند بدون تغییر و تبدیل. در واقع، این روش پایداری

ص: 40

1- ساختار الگو گونه برای نمایش ترتیب حوادث: Script

بیشتری را از روش های دیگر فراهم می کند. چون بر سازماندهی فایل محلی بستگی دارد، همچنین مناسب برای انتقال و همچنین تحویل محتواست.

ابزارها

این روش، به یک زیر ساخت دسترسی و به یک خزنده بایگانی (مانند HERITRIX) و یک نظام نمایه سازی برای فایل های WARC نیاز دارد. IIPC زنجیره کاملی از ابزارهایی دارد که برای ایجاد این توابع عملیاتی توسعه یافته است.

3-3-4-بایگانی غیروب

توصیف

در این رویکرد که در شکل 1-7 ترسیم شده است، اسنادی که بر روی وب بوده و از متن ابر متن استخراج شده و در یک سبک متفاوت برای دسترسی منطقی و/یا فرمت سازمان دهی مجدد می شوند.

این می تواند موردی باشد که وقتی یک مجموعه اسناد از وب برداشته می شوند، از دسترسی منطقی مبتنی بر پیوند به یک پیوند مبتنی بر فهرست دوباره سازماندهی می شوند.

عکس

بیشتری را از روش های دیگر فراهم می کند. چون بر سازماندهی فایل محلی بستگی دارد، همچنین مناسب برای انتقال و همچنین تحویل محتواست.

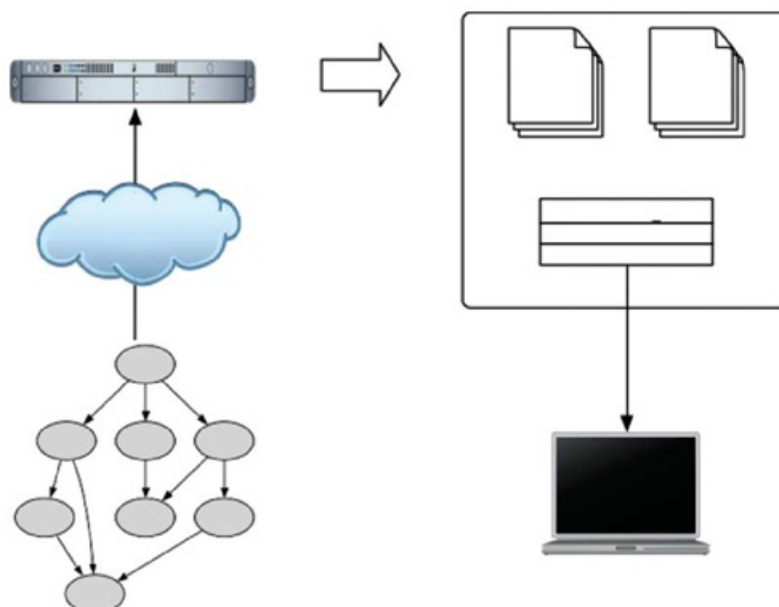
ابزارها

این روش، به یک زیر ساخت دسترسی و به یک خزنده بایگانی (مانند HERITRIX) و یک نظام نمایه سازی برای فایل های WARC نیاز دارد. IIPC زنجیره کاملی از ابزارهایی دارد که برای ایجاد این توابع عملیاتی توسعه یافته است.

۴-۳-۳- بایگانی غیروب

توصیف

در این رویکرد که در شکل ۷-۱ ترسیم شده است، اسنادی که بر روی وب بوده و از متن ابر متن استخراج شده و در یک سبک متفاوت برای دسترسی منطقی و/یا فرمت سازماندهی مجدد می شوند. این می تواند موردی باشد که وقتی یک مجموعه اسناد از وب برداشته می شوند، از دسترسی منطقی مبتنی بر پیوند به یک پیوند مبتنی بر فهرست دوباره سازماندهی می شوند.



شکل ۷-۱. اسناد از سایت اصلی در بایگانی سازماندهی مجدد می شوند، ساختار غیر وبی را پیروی می کنند، به طور موردی از فهرستی که دسترسی به اسناد انفرادی را فراهم می کند، استفاده می کند.

شکل ۱-۷. اسناد از سایت اصلی در بایگانه سازماندهی مجدد می شوند، ساختار غیر وبی را پیروی می کنند، به طور موردی از فهرستی که دسترسی به اسناد انفرادی را فراهم می کند، استفاده می کند.

همچنین، این موردی است که وقتی یک صفحه یا حتی یک وبگاه کامل به فرمت PDF تغییر شکل یافته است. (1 Adobe Acrobat) این توابع عملیاتی را دارد و می تواند تمام یک وبگاه را در یک سند واحد PDF تغییر شکل دهد. در این مورد سند به طور مجازی چاپ شده است که یک انتقال بی حرکت و یک صفحه یاد داشت مانند سازماندهی را شامل می شود، حتی اگر پیوند دهی بتواند هنوز با استفاده از شمای داخلی نام گذاری مناسب کار کند.

توضیح

این رویکرد اساساً برای دریافت موضوع هایی که در اصل ایجاد شده اند، به طور مستقل از وب سازماندهی شده اند، ایجاد شده است. این مورد برای مجموعه های بزرگی از کتاب های رقومی به عنوان نمونه، مقاله ها، موزیک، و ویدئوهای ساخته شده که بر روی وب موجود هستند، ولی سازماندهی اصلی شان ابرمتی نمی باشد چون مبتنی بر کاتالوگ هستند.

این رویکرد می تواند در این مورد برای چسبیدن به معماری اطلاعات اصلی و بایگانی این مجموعه ها با یکدیگر با کاتالوگ های ادغام شده در کاتالوگ بایگانی، مرجع باشد. فرض شده است که محتوای ابر متن بدون ربط پنداشته شده است و می تواند منتشر شود.

برای مثال این مورد در طرح انبارهای الکترونیکی kb در نیوزیلند جایی که انتشارات علمی الزویر (2) اجرا شده است که در یک نظام مبتنی بر فهرست بایگانی شده است. حقیقت این است که الزویر به این مطالب دسترسی داشته است و به عنوان یک مجموعه کمکی برای محتوای خودش مورد بررسی قرار گرفته و به عنوان انتشار علمی سنتی ساخت یافته است.

مرجع

اشاره

این روش برای مجموعه های محتوا و نه ساخت یافته در حالت وب دلالت می کند.

4-3-4 خلاصه

جدول 2-1- انواع گوناگونی از بایگانی های وب، موارد مرجع های مورد استفاده، ابزارهای و مزیت های و عدم مزیت ها را به طور خلاصه شرح داده است.

ص: 42

1- نرم افزاری از محصول شرکت Adobe برای ساخت و تهیه فایل های (PDF)

Elsevier -2

بایگانی شبکه وب: مباحث و روش ها ۴۳

جدول ۱-۲. خلاصه ای از انواع بایگانی وب

نوع بایگانی	نظام فایل محلی	مبتنی بر خدمت	غیروب
توصیف	تمام پیوند به پیوندهای مرتبط (ربطی) تبدیل شده‌اند. و شناوری ابر متن، به‌طور قطع بر روی نظام‌فایل محلی انجام شده است.	سرور یک وب، برای دسترسی از طریق هر سندی که به‌کار رفته است نصب شده است و شناوری ابر متن به بایگانی اصلی بسته شده است.	اسناد از محتوای ابر متن اصلی استخراج شده و مجدداً برای یک منطق مختلف سازماندهی شده‌اند.
استفاده مرجع	بایگانی سایت منفرد و بایگانی در مقیاس متوسط و کوچک	بایگانی در مقیاس کوچک و متوسط	بایگانی مجموعه‌های خاص (غیروب)
ابزارها	کپی کننده وبگاه (مانند HTTrack)	خزنده بایگانی (مانند Heritrix) و نظام نمایه‌سازی برای فایل‌های WARC	بستگی به ساختار محتوای نهایی دارد.
مزیت‌ها	برای اجرا ساده است	صحت، قابلیت مقیاس پذیری	قادر به ایجاد یکپارچگی در فهرست‌های سنتی یا سایر معماری‌های اطلاعاتی محلی
معایب	افزایش مقیاس ندارد. نیاز به نامگذاری مجدد و محدود شدن سازماندهی مجدد محتوا برای شناوری ابر متن‌ها است. نیاز به مدیریت در سطح نظام فایل در مجموعه بایگانی و نگارش موارد دارد.	اجرا در غیاب نرم افزار یکپارچه‌سازی سخت می‌باشد (این کار ممکن است در آینده تغییر کند).	فقدان ساختار ابر متن. فقط می‌تواند برای اسناد مجزا و غیروب به‌کار برده شود.

۴-۴- کیفیت و تمامیت (کامل بودن)

به‌طور کلی، کیفیت می‌تواند در یک حالت وظیفه‌ای (متناسب برای استفاده خاص) یا در یک حالت هدفمند (منطبق با ویژگی‌های اندازه‌گیری) تعریف شود. اصطلاح کیفیت، برای مجموعه فرهنگی در مضمون‌های مختلف به‌کار برده شده است. شخص می‌تواند از آن برای کنترل وضعیت حفاظت، کامل نمودن موارد یا مجموعه‌ها، سطح محتواهای هوشمند و علمی و غیره استفاده شود. در هر مورد، آن با مقیاس مطلوب کامل در یک ناحیه خاص و (حفاظت فیزیکی، پوشش یک قلمرو، صحت انتخاب) مرتبط می‌باشد.

برای بایگانی‌های وب همان‌طور که دیده شده، بیشترین نقایص ناشی از سختی برای جمع‌آوری محتوا از طریق واسط HTTP (بخش قبلی در مورد بایگانی سرویس گیرنده جانبی را ببینید) و سختی تحویل در یک حالت منسجم از نتایج محتوا (بخش «سازماندهی و ذخیره‌سازی» را ببینید) می‌باشد. از این رو، کیفیت بایگانی وب در این فصل به‌عنوان ۱) تکامل مطالب (فایل‌های پیوندی) بایگانی شده درون فضای پیرامون هدف، و ۲) توانایی برای تحویل شکل اصلی سایت، به‌ویژه در رابطه با هدایت و تعامل با کاربران (میزانس، ۲۰۰۵) مورد بررسی قرار گرفته است.

جدول ۱-۲. خلاصه ای از انواع بایگانی وب

۴-۴- کیفیت و تمامیت (کامل بودن)

به طور کلی، کیفیت می‌تواند در یک حالت وظیفه‌ای (متناسب برای استفاده خاص) یا در یک حالت هدفمند (منطبق با ویژگی‌های

اندازه‌گیری) تعریف شود. اصطلاح کیفیت، برای مجموعه فرهنگی در مضمون‌های مختلف به کار برده شده است. شخص می‌تواند از آن برای کنترل وضعیت حفاظت، کامل نمودن موارد یا مجموعه‌ها، سطح محتواهای هوشمند و علمی و غیره استفاده شود. در هر مورد، آن با مقیاس مطلوب کامل در یک ناحیه خاص و (حفاظت فیزیکی پوشش یک قلمرو، صحت انتخاب) مرتبط می‌باشد.

برای بایگانی‌های وب همانطور که دیده شده، بیشترین نقایص ناشی از سختی برای جمع‌آوری محتوا از طریق واسط HTTP (بخش قبلی در مورد بایگانی سرویس گیرنده جانبی را ببینید) و سختی تحویل در یک حالت منسجم از نتایج محتوا (بخش «سازماندهی و ذخیره سازی» را ببینید) می‌باشد.

از این رو، کیفیت بایگانی وب در این فصل به عنوان 1) تکامل مطالب (فایل‌های پیوندی) بایگانی شده درون فضای پیرامون هدف و 2) توانایی برای تحویل شکل اصلی سایت به ویژه در رابطه با هدایت و تعامل با کاربران (میزانس، 2005) مورد بررسی قرار گرفته است.

تکامل می تواند به طور افقی با تعدادی از نقاط ورودی مربوط که در فضای پیرامون طراحی شده، با نمایش هندسی ارزیابی شود و از نظر برنامه کاربردی عمودی با تعدادی از گره های پیوندی می توان که از این نقطه ورودی دریافت شده اند اندازه گیری می شود. معمولا نقاط ورودی صفحه های اصلی سایت هستند و پیوندها می توانند کاربران را به هر یک از نقاط ورودی هدایت کنند (سایت دیگر) یا به عناصر همان سایت بفرستند. این مورد برای بایگانی مبتنی بر سایت است.

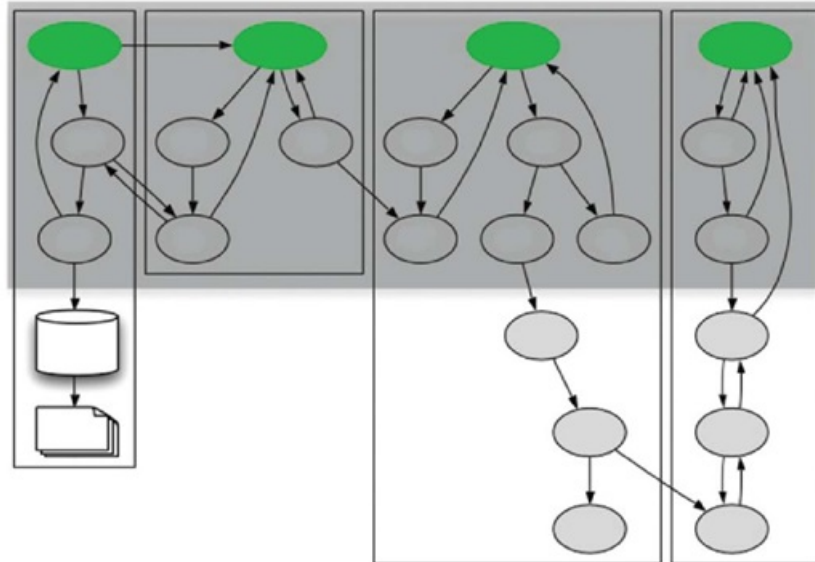
به هر حال، در برخی موارد، گنجایش برنامه کاربردی عمودی برای جا سازی عناصر (برای مثال تصاویر همراه با یک صفحه) محدود شده است و مجموعه فقط به طور افقی با نادیده گرفتن سطح سایت، سازماندهی شده است. برای مثال این مورد برای خزش موضوع های خالص در جایی که صفحه ها را در بر نمی گیرد بر اساس تعلق شان به سایت می باشد، اما فقط بر روی ارتباط شان با موضوع است.

به طور مطلوب، بایگانی وب باید به طور عمودی همچنین افقی کامل شود. اما در عمل، برای دسترسی سخت است و اولویت هایی باید قرار داده شوند. بایگانی زمانی که کامل شدن افقی به کامل شدن عمودی ترجیح داده شود، «بسیط» نامیده می شود.

عکس

تکامل می‌تواند به‌طور افقی با تعدادی از نقاط ورودی مربوط که در فضای پیرامون طراحی شده، با نمایش هندسی، ارزیابی شود و از نظر برنامه کاربردی عمودی، با تعدادی از گره‌های پیوندی می‌توان که از این نقطه ورودی دریافت شده‌اند، اندازه‌گیری می‌شود. معمولاً، نقاط ورودی، صفحه‌های اصلی سایت هستند و پیوندها می‌توانند کاربران را به هریک از نقاط ورودی هدایت کنند (سایت دیگر) یا به عناصر همان سایت بفرستند. این مورد برای بایگانی مبتنی بر سایت است.

بمهر حال، در برخی موارد، گنجایش برنامه کاربردی عمودی برای جا سازی عناصر (برای مثال تصاویر همراه با یک صفحه) محدود شده است و مجموعه فقط به‌طور افقی با نادیده گرفتن سطح سایت، سازماندهی شده است. برای مثال، این مورد برای خزش موضوع‌های خالص در جایی که صفحه‌ها را در بر نمی‌گیرد بر اساس تعلقشان به سایت می‌باشد، اما فقط بر روی ارتباطشان با موضوع است. به‌طور مطلوب، بایگانی وب باید به‌طور عمودی همچنین افقی کامل شود. اما در عمل، برای دسترسی سخت است و اولویت‌هایی باید قرار داده شوند. بایگانی زمانی که کامل شدن افقی به کامل شدن عمودی ترجیح داده شود، «بسیط» نامیده می‌شود.



شکل 1-8 مجموعه‌های بسیط شامل بیشتر سایت‌هاست ولی فقط در سطحی بایگانی می‌شود. فقط محتوا در ناحیه سایه‌دار بایگانی می‌شود. صفحه‌ها عمیق در سلسله مراتب (سی 6، سی 7، سی 8، دی 3، دی 3، دی 5) و نیز محتوایی که در پایگاه داده پنهان شده (وب پنهان) تصرف خواهند شد.

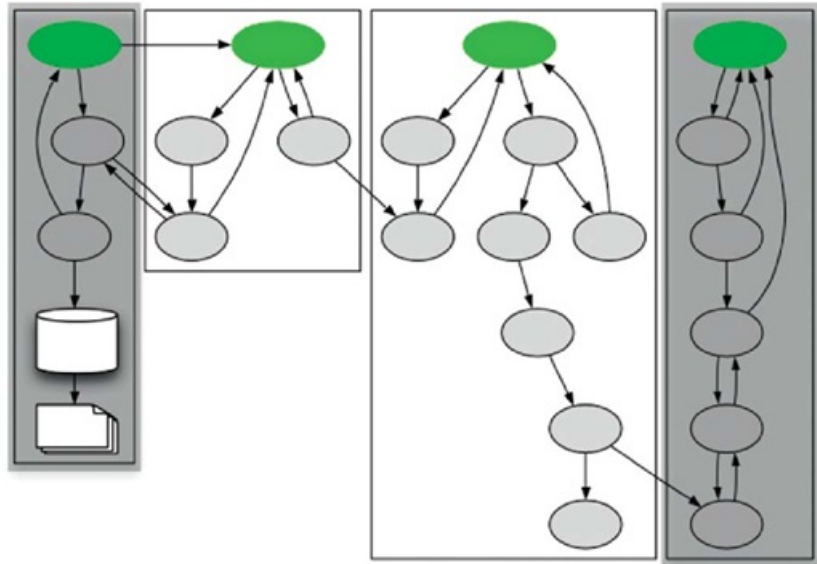
شکل 1-8. مجموعه‌های بسیط شامل بیشتر سایت‌هاست ولی فقط در سطحی بایگانی می‌شود. فقط محتوا در ناحیه سایه دار بایگانی می‌شود. صفحه‌ها عمیق در سلسله مراتب (سی 6، سی 7، سی 8، دی 3، دی 3، دی 5) و نیز محتوایی که پایگاه داده پنهان شده (وب پنهان) تصرف خواهند شد.

برای مثال، این مورد برای مجموعه بایگانی اینترنت است که به وسیله الکسا (همچنین برنر، 1997؛ کیمپتون و همکاران، 2006) ارائه شده است، خزشگر الکسا یک رویکرد خزش سطح اول را استفاده می کند و عمق خزش را برای یک سایت بر طبق ترافیک اندازه گیری شده برای این سایت وفق می دهد.

بر عکس، بایگانی «متمرکز» نامیده می شود زمانی است که تکامل عمودی به کامل شدن افقی ترجیح داده می شود (شکل 9-1 را ببینید).

عکس

برای مثال، این مورد برای مجموعه بایگانی اینترنت است که به وسیله الکسا (همچنین برنر، ۱۹۹۷؛ کیمپتون و همکاران، ۲۰۰۶) ارائه شده است، خزشگر الکسا یک رویکرد خزش سطح اول را استفاده می‌کند و عمق خزش را برای یک سایت بر طبق ترافیک اندازه‌گیری شده برای این سایت وفق می‌دهد. برعکس، بایگانی «متمرکز» نامیده می‌شود زمانی است که تکامل عمودی به کامل شدن افقی ترجیح داده می‌شود (شکل ۹-۱ را ببینید).



شکل ۹-۱. بایگانی بسیط سایت‌های کمتری مورد خزش قرار می‌گیرند ولی خزش در عمق انجام می‌شود. ولی سایت الف و دی مورد بایگانی قرار می‌گیرند ولی در حالت کلی شامل بخش وب پنهان سایت الف.

برای مثال، در این مورد، وقتی اولویت نخست - سایت برای خزش گرها استفاده شده یا وقتی که یک بازبینی دستی، در جایی که مورد نیاز است، بر روی آن انجام می‌شود، بایگانی تکمیلی انجام می‌شود. بایگانی متمرکز حتی متقاضی بیشتری برای سایت‌های وب پنهان دارد (همچنین، سایت‌های وب عمیق نامیده شده است) جایی که دسترسی به تمام محتوا با خزش گرها ممکن نیست.

۵ - بررسی عمومی مراحل اولیه جاری

بایگانی وب می‌تواند در چندین روش طبقه‌بندی شود. در این بخش، روش‌های مهم را بررسی خواهیم کرد و این فرصت را برای ارائه برخی مراحل اولیه بایگانی وب و مقایسه با رویکردهای گوناگون نشان می‌دهیم.

شکل 1-9. بایگانی بسیط سایت‌های کمتری مورد خزش قرار می‌گیرند ولی خزش در عمق انجام می‌شود. ولی سایت الف و دی مورد بایگانی قرار می‌گیرند ولی در حالت کلی شامل بخش وب پنهان سایت الف.

برای مثال، در این مورد، وقتی اولویت نخست - سایت برای خزش گرها استفاده شده یا وقتی که یک بازبینی دستی، در جایی که مورد نیاز است، بر روی آن انجام می‌شود، بایگانی تکمیلی انجام می‌شود. بایگانی متمرکز حتی متقاضی بیشتری برای سایت‌های وب پنهان دارد (همچنین، سایت‌های وب عمیق نامیده شده است) جایی که دسترسی به تمام محتوا با خزش گرها ممکن نیست.

5- بررسی عمومی مراحل اولیه جاری

بایگانی وب می تواند در چندین روش طبقه بندی شود. در این بخش روش های مهم را بررسی خواهیم کرد و این فرصت را برای ارائه برخی مراحل اولیه بایگانی وب و مقایسه با رویکردهای گوناگون نشان می دهیم.

ص: 45

نوع ایجاد سازماندهی و میزبانی از بایگانی، نخستین معیار برای طبقه بندی بایگانی های وب می باشد. برخی دسترسی عمومی به مجموعه های شان را فراهم می کنند (WA عمومی) برخی این کار را نمی کنند (WA خصوصی (یا مخفی)).

در میان بایگانی های عمومی وب، برخی دسترسی پیوسته را فراهم می کنند برخی دسترسی به سایت را در اتاق های مطالعه فراهم می کنند (بایگانی وب عمومی پیوسته و بایگانی وب عمومی غیر پیوسته). همچنین برخی و در بیشتر موارد در ابتدا مجموعه غیر رقومی را مدیریت می کنند (بایگانی وب هیبرید) و بالاخره برخی در حالت سرمایه گذاری یا بدون منفعت (بایگانی وب غیر تجاری) با در نظر گرفتن اینکه برخی شرکت های تجاری هستند (بایگانی وب تجاری).

مؤسسه های میراث سنتی (کتابخانه ها، بایگانی ها، موزه ها) که مجموعه هایشان را برای وب گسترش داده اند، با یکدیگر بیشترین قسمت مقوله بایگانی وب هیبرید عمومی را تشکیل داده اند. کتابخانه های ملی چندین کشور متعلق به این مقوله هستند (سوئد و استرالیا که نخستین بار به 1996 بر می گردد و اکنون کشورهای زیادی هستند) (1).

بایگانی های ملی، منطقه ای و شهری همچنین شروع به بایگانی وبگاه های مجاز دولتی و محلی نموده اند (2). سازمانی بر روی شکل های جدید صنعتی مانند v2 مستقر در روتردام هلند کار می کند که در حال یکپارچه سازی شبکه در یک انعکاس و عملکردی عمومی برای حفاظت از رسانه ناپایدار می باشد (فوکونیر و فرومه 2004) (3).

تمام این بایگانی ها می توانند به عنوان آرشیو وب عمومی هیبرید رده بندی شوند. همان طور که محتوای وب را در یک مضمون بزرگ تر از مجموعه ها مجتمع می کنند. بیشتر آن ها فقط دسترسی غیر پیوسته برای اسناد را در حال حاضر فراهم می کنند.

از میان آن ها، کتابخانه اسکندریه در مصر یکی از معدود دسترسی های پیوسته برای مجموعه بایگانی وب است (معکوس نمودن بایگانی اینترنتی) و یک مثال از بایگانی وب عمومی هیبرید غیر تجاری پیوسته می باشد.

نفوذپذیری اینترنت همچنین ضرورت برخی از انواع جدید سازمان های بایگانی را مجاز می داند که فقط مجموعه رقومی و تهیه دسترسی پیوسته را نگه می دارد که به عنوان بایگانی وب پیوسته غیر تجاری عمومی، رده بندی خواهند شد.

بایگانی اینترنت، مثال مهمی در این مقوله می باشد (فصل 9 را ببینید) (کیمپتون و همکارانش، 2006).

ص: 46

1- با پیروی از آن ها، چندین کتابخانه ملی بایگانی وب را شروع کرده و برنامه هایی را اجرا می کنند (این سیاهه جامع نیست): در اروپا: فنلاند، دانمارک، نروژ، ایسلند فرانسه، جمهوری چک، اسلونی، ایتالیا و یونان؛ در آسیا ژاپن چین و سنگاپور و کتابخانه کنگره در آمریکا با پیروی از آن ها، چندین کتابخانه ملی بایگانی وب را شروع کرده و برنامه هایی را اجرا می کنند (این سیاهه جامع نیست): در اروپا: فنلاند، دانمارک، نروژ، ایسلند فرانسه، جمهوری چک، اسلونی، ایتالیا و یونان؛ در آسیا ژاپن چین و سنگاپور و کتابخانه کنگره در آمریکا.

2- آرشیوهای ملی استرالیا (آرشیوهای ملی استرالیا، 2001) بریتانیا (براون، 2006)؛ کانادا ایالات متحده آمریکا (کارلین، 2006) بایگانی وب نظام مند را شروع کرده اند. هم چنین طرح شهر Antwerp DAVID را ببینید (بودرس و آیتده، 2002).

برخی شرکت های تجاری مجموعه های بزرگی از محتوای وب عمومی را بایگانی می کنند مانند گوگل با نهانگاهش (1) و بایگانی هانزو (2). مثال هایی از بایگانی وب تجاری عمومی پیوسته هستند.

سرانجام، بسیاری از سازمان ها، بایگانی وب داخلی را برای اهداف شخصی توسعه می دهند که به عنوان بایگانی وب خصوصی (مخفی) طبقه بندی خواهند شد. منحصر به فرد بودن نوع دسترسی (پیوسته یا غیر پیوسته) و همچنین وضعیت تجاری در اینجا ارتباط کمتری دارد چون این بایگانی ها فقط برای استفاده خصوصی هستند.

2-5-2- حوزه (دامنه)

روش مفید دیگر برای دسته بندی بایگانی های وب، بررسی حوزه ای است که آن ها اتخاذ می کنند. بایگانی های وب می توانند هر یک، سایت، سر فصل ها یا متمرکز بر حوزه باشند.

5-2-1- بایگانی مرکزی سایت

این نوع بایگانی، بر روی یک سایت مشخص متمرکز شده است که تقریباً به وسیله و برای ایجاد کننده سایت اجرا شده است. این حوزه بندی، از این رو، تقریباً برای بایگانی وب خصوصی استفاده شده است. مثلاً بسیاری از شرکت ها، مسئول تمام محتواهایی هستند که منتشر می کنند و باید مطمئن باشند که می توانند به روش های قدیمی تر سایت و بنوشت ها مراجعه کنند. این نوع بایگانی ترجیحاً از کپی کنندگان سایت و برخی تهیه کنندگان خدمات اینترنتی استفاده می کند که برای این نوع بایگانی داخلی مناسب، به وجود آمده اند (3).

5-2-2- بایگانی مرکزی عنوان

بایگانی های وب، عمومی و عمومی تر شده اند، اغلب به وسیله نیازهای پژوهشی مستقیم اجرا می شوند تا موقعی که کار بر روی یک فیلد مشخص و انعکاس بر روی وب انجام می شود، دانش پژوهان زیادی با طبیعت بی دوام انتشار وب روبه رو می شوند، جایی که طول عمر وبگاه برای بازبینی علمی (تکذیب نیازمند دستیابی به همان داده است) و همچنین برای ارجاع با دوام، نامناسب می باشد.

به این دلیل است که طرح های مختلف اغلب در کتابخانه های دانشگاه میزبانی می شوند و تحت حفاظت مطالب اولیه برای پژوهش قرار می گیرند، مانند بایگانی رومی برای مطالعات چینی در دانشگاه هایدلبرگ (4) در آلمان یا آرچیپل (5) برای تحلیل سایت های سیاسی هلندی در دانشگاه گرونینگن (6) در هلند

ص: 47

Cashe -1

Hanzo -2

3- برای نمونه hanzoarchives.com را ببینید.

4- Digital Archive for Chinese Studies (DACHS) at Heidelberg University

5- Archipol

6- Groningen

(ورمن 1) و همکاران، 2002) این پروژه ها نه تنها جهت دهی یک عنوان را به اشتراک می گذارند، بلکه از

یک شبکه آگاهی دهنده نیز استفاده می کنند.

از طرفی دیگر، پژوهشگران که تغذیه های صحیح و به روز

طرح های متمرکز دیگری در کتابخانه ها از طریق جست و جوی فعال و بایگانی وبگاه گزینشی (2) انجام شده است و مانند طرح مینروا (3) از کتابخانه کنگره (اشنایدر و همکاران 2003) یا بایگانی وب انتخابات فرانسه که توسط کتابشناسی ملی فرانسه ایجاد شده است (میزانس ، 2005). در مقایسه با رویکرد قبلی مبتنی بر عنوان متمرکز، کشف سایت ها به طور طبیعی به عنوان محصول فرعی از فعالیت پژوهشی ایجاد نشده است. و به عنوان یک فعالیت خاص نیاز به بررسی دارد.

بالاخره، برخی طرح ها که با این مقوله مرتبط هستند. از خزش عنوان برای کشف و ضبط محتوا مرتبط با مطلب مشابه استفاده. می کنند (چاکرابارتی (4) و همکاران، 1999؛ برگمارک 2002 (5)؛؛ برگمارک و همکاران، 2002؛ کیواین (6) و همکاران، 2004). کشف و فیلتر کردن خودکار با استفاده از فن سنتی خزش در ترکیب ارزیابی سطح صفحه در محتوای متنی انجام شده است که گاهی اوقات با کاویدن ساختارهای پیوندی آمیخته می شود. مجاورت با عنوان، می تواند از مجموعه ای از نوشتجات یا از باز خورد کاربر، فرا گرفته شود. اگر چه محتمل است، این ناحیه باز هم نیاز به پژوهش برای به کار برده شدن در بایگانی دارد.

3-2-5-بایگانی مرکزی حوزه ای

ساخت بایگانی می تواند بر اساس محل محتوا نیز انجام شود و به این ترتیب، نوع سوم در بایگانی وب حوزه ای را مشخص می کند. در اینجا کلمه «حوزه» (7) در جهت واژه شبکه یا به وسیله پسوند در جهت واژه ملی استفاده شده است که یک معیار ترکیبی برای سایت های هدف از یک کشور خاص می باشد (8).

نظام ملی حوزه ای یک گزینش ساده و قابل تعقیب قانونی در محتوا را بر اساس نام های دامنه اجازه می دهد. این حقیقتی است که نام های دامنه حتی برای سطوح بالایی حوزه ای توسط نمایندگی رسمی آی.سی.آن (9) مدیریت شده است، که واقعاً از قوانین در رابطه با نام دادن مشخصه عملیاتی و سازمان ها پیروی نمی کنند، ولی بیشتر سنت ها را پیروی می کنند (لیو و آلبرت 1999) (10). همچنین می توانید در مورد تکامل بر روی نام های اینترنتی به (کوهلر، 1999) نگاه کنید. از این رو، شخص می تواند انواع پسوند های

ص: 48

Voerman -1

Electoral Web sites -2

Minerva -3

Chakrabarti -4

Bergmark -5

Qin -6

Domain -7

8- برای بحث در مورد روش ممکن مرزبندی کردن فضای اینترنت ملی آرویدسون و دیگران (2000)؛ ابایت بول و دیگران (2002)؛ لامپوس و دیگران (2004) را ببینید، برای مطالعاتی در مورد ویژگی های فضای اینترنت ملی بیزایتس و دیگران (2005 الف، 2005 ب) و گومس و سیلوا (2003) را ببینید.

ICANN -9

Liu and Albitz -10

عملیاتی یا عمومی مانند Com و edu و انواع سیستم اطلاعاتی جغرافیایی مانند (1) (ch.jp) و انواعی را در دامنه سطح نخست تشخیص دهد و که (اغلب دامنه سطح بالا نامیده می شود). دامنه های سطح بالای جغرافیایی اغلب بخش های فرعی عملیاتی دارند و (مانند asso.fr gob.mex). که به این معناست که دامنه سطح دوم (2) همچنین، در همان روش مدیریت خواهد شد. استثناها برای عملکردهایی مانند TLDها وجود دارد که بخش های فرعی جغرافیایی دیگر را دارا می باشند. (به وسیله ایالت ها). به هر حال، توجه داشته باشید که تمام این قسمت ها فضای حوزه اینترنتی است که توسط نمایندگی مدیریت شده اند (3). هر هویت تحت فرمان آن ها می تواند یک سیاست خاص را در رابطه با تخصیص و کنترل فضاهای شان به کار ببرد بنابراین رسیدن به بهره برداری از SLD, TLD برای گزینش بایگانی وب، در هر مورد طبق ارزیابی این سیاست بستگی دارد (مثلاً org.com.gTLD) به وسیله تمام انواع سازمان ها استفاده شده اند نه تنها توسط سازمان های تجاری برای com. و سازمان های بدون فایده برای org. چون محدودیتی برای ثبت نام وجود ندارد) به علاوه، برخی هویت ها تحت فرمان مدیریت TLD، سیاستشان را طبق زمان تغییر می دهند.

org.net با داشتن محدودیت ها قبل از سال 1996 استفاده شده اند و frTLD به طور قابل توجهی محدودیت ها را در سال 2005 کاهش داده است).

مزیت بزرگ دیگری که در اینجا باید ذکر کنیم، آوردن معیارهایی است که بتواند به طور خودکار توسط خزشگرها آشکار شوند مانند نام های دامنه پروژه های دیگری در واقع، رویکرد مرکزی دامنه را اجرا می نمایند. برخی، بر روی یک دامنه عمومی مانند (کروس و و همکاران، 2003؛ کارلین، 2004 و یا edu (لایل، 2004) متمرکز شده اند.

برخی از دامنه های ملی مانند ایالت کالتیورارو (4) در سال 1997 که توسط کتابخانه سلطنتی سوئد ایجاد شد، استفاده می کنند (آرویدسون 2000) که SeTLD و همچنین صفحه ها سوئدی پیوند یافته از آن و در مانند com. واقع شده است را پوشش می دهد.

3-5- روش های استفاده شده

طرح ها می توانند همچنین به طور قابل توجهی نسبت به رویکرد روش های بررسی که برای کشف اکتساب و توصیف محتوا در بر می گیرند و متفاوت باشند. یکی از تفاوت های مهم این است که محدوده ها در سراسر این مراحل، از دست به جای پردازش خودکار استفاده می کنند. اگر چه سادگی ظاهری این تضاد باید به عنوان پردازش خودکار متعادل می شود و در چندین سطح رخ می دهد (ضبط، استفاده از موتور جست و جو برای کشف دستی و غیره (میزانس، 2006 ب)، این مسئله، همچنین می تواند بایگانی وب را بر طبق این تضاد رده بندی نماید که به طور مستقیم اثر شدیدی بر روی قابلیت مقیاس پذیری و

ص: 49

1- ایزو 3166 دو حرف اول نام کشورها را اجازه می دهد به جز uk که باید gb باشد، و همچنین به جز این که اخیراً به سه حرف برای هر منطقه گسترش یافته است مثل حوزه cat در کاتولینا در اسپانیا.

SLD-2

3- در مورد حوزه دولتی و مفاهیم آن میولر (2002) را ببینید.

Kulturarw-4

همان طور که توانستیم پیش بینی کنیم، خود کار سازی این وظایف، در پایین آوردن فاحش هزینه های هر دستیابی به سایت توانا هستند (1). به طور مطلوب یک اپراتور تنها با اجرای خزش می تواند میلیون ها صفحه را کشف و بارگذاری کند. نظر به این که نمایه سازی تمام متن کمکی برای یافتن قدرتمند قابل مقایسه ای را فراهم میکند که اگر در برخی موارد برای فهرست کردن عالی نباشد، بنابراین می بینیم که بار دیگر اینجا چقدر خود کار سازی به طور چشمگیری هزینه ها را کاهش می دهد، چون می تواند بر روی یک مقیاس بزرگ برده شود (استک، 2005) (هال گریسون، 2006).

متأسفانه، خود کار سازی، محدودیت هایی دارد و بررسی دستی باید در برخی موارد انجام شود. برای مثال کشف می تواند به طور دستی با اتوماتیک انجام شود. وقتی کار به طور دستی انجام شود می تواند یک فعالیت خاص یا یک محصول فرعی از فعالیت های دیگر باشد مانند DACHS (لچر، 2006) و نمایش بایگانی های وب (وثرمان (2) و دیگران، 2002) Achipol این نوع رویکرد معمولاً برای بایگانی مرکزی عنوان در انجام می شود. اگر چه خزش عنوان به طور مؤثر برای کشف مطلب سایت های یا صفحه های مربوط آزمایش می گردد. ابزارهای خود کار می توانند به طور مطمئن (نه در این زمان) در مقایسه با یک شبکه متخصصان، مرجع هایی را برای بهترین مطالب که از آن مطلع هستند، فراهم نمایند.

به هر حال، فقدان حوزه تخصصی و عدم درک تنها معایب خزش گرها نیستند. این مسئله مورد توجه است که تأخیر برای یافتن سایت های جدید مورد نیاز است. زمان بسیار زیادی برای یافتن سایت های کلی گرفته می شود. وقتی وارد سایت های زودگذر (بی دوام) می شود، برای مثال مرتبط با یک رویداد، تاخیری می تواند بسیار طولانی باشد و آن ها را بایگانی کند این تفاوت به وسیله (میزانس، 2005) با یک مقایسه سایت های کشف شده توسط خزشگر الکسا و قابلیت دسترسی و روزانه به بایگانی اینترنت و سایت های مرتبط با انتخابات فرانسه در سال 2002 که توسط گروهی از کتابداران مرجع انجام شد و به وسیله کتابخانه ملی فرانسه بایگانی گردید، مورد بررسی قرار گرفته است.

این بررسی، مزیت های آشکاری را برای گزینش فعالیت دستی در مجموعه های مبتنی بر ربط، برای کشف به موقع و در تمرکز عمیق نشان می دهد.

رده بندی بایگانی وب بر طبق روش بررسی اش همچنین می تواند در یک شاخه مالی انجام شود. بر تراز پردازش اتوماتیک و دستی دو حالت، شخص می تواند برای مثال نوع منبع استفاده شده برای کشف، تناوب جست و جو، و ضبط، سطح کیفیت بازبینی ایجاد شده و دانه دانه بودن آیتم های بایگانی شده (سایت ها، صفحه ها)، و مانند آن را مورد بررسی قرار می دهد.

ص: 50

1- فیلیپز (2005) تخمین های بسیار مفید و دارای جزئیات در خصوص زمان و هزینه های فرآیند دستی خود کار سازی سایت ها برای یکی از قدیمی ترین آرشیوهای وب موجود را فراهم کرده است. تخمین های زمانی به صورت ذیل هستند: - تشخیص و انتخاب 30 دقیقه - گردآوری، اطمینان از کیفیت، و موارد آرشیوی: 210 دقیقه - فهرست نویسی 81 دقیقه

به هر حال، این امر که بیشتر بایگانی های وب به دو مدل عمده تمایل دارند، تمایز اصلی است که آیا انتخاب به طور دستی یا غیر آن انجام شده است. یکی مدل های خزش گر های کلی است که معمولاً مبتنی بر مرکز (دامنه های ملی یا دامنه های عمومی) است یا به صورت آزاد (بایگانی اینترنت)، و دیگری مدل گزینش فردی تعداد محدودی از جست و جوها یا نقاط ورودی است که به طور دستی انجام می شود (معمولاً سایت ها). تمایز در بیشتر در رویکرد روش شناختی آن ها به ندرت دیده شده یا هیچ یک برای طبقه بندی آن ها مورد استفاده قرار نمی گیرد.

6- نتیجه گیری

وب فقط 15 سال است که وجود دارد و می توان گفت که نگهداری و حفاظت از حافظه اش، به طور نسبی در مقایسه با رسانه های دیگر از ابتدا آغاز شده بود. (1) اما فقط مراحل لازم اولیه برای حفاظت آن ایجاد شده است. وضعیت حفاظت جاری، به تعداد بسیار کمی از مؤسسه ها وابسته است و پوشش زیادی را حاصل نشده است.

نقش ها و مسئولیت ها برای بیشتر سهامداران بسیار واضح و آشکار نیست و توانایی در حمایت از بسیاری از مجموعه های شاخص به وجود نیامده است. و ما هنوز در دوره ای هستیم که هیچ گسیختگی فناورانه از زمان آغاز وب، رخ نداده است.

مرورگرهای جاری با یکدیگر با تعداد محدودی که برنامه های کامپیوتری متصل می توانند فرمت های زیادی را جابه جا کنند که می تواند بر روی وب یافت شود.

اما این موقعیت، تا ابد طول نخواهد کشید و حفاظت از وب با یک چالش جدی روبه رو خواهد شد وقتی تغییرات مهم فناوری در وب رخ می دهد (که ممکن نیست مانند اینها بعداً دیده شود).

بنابراین، دلگرم کننده است که بینیم که بسیاری از مؤسسه های (حفظ) میراث، در بایگانی وب در حال به کارگیری هستند. بررسی اخیر توسط گروه پژوهش کتابخانه (RLG2006) نشان داد که 60 درصد اعضای مورد بررسی شان، بایگانی وب را قسمتی از مأموریت خود پنداشته اند (RLG2006) که بسیار دلگرم کننده است.

امیدواریم که بازنمون های ایجاد شده در این فصل، مباحث مهم و روش ها، متفقاً با منطق و دلیل، به آن ها و دیگران برای مشارکت در این تلاش گروهی، کمک خواهد کرد.

ص: 51

- Aarseth, E. J. (1997). *Cybertext: perspectives on ergodic literature*. Baltimore, MD: Johns Hopkins University Press . 1
- Abiteboul, S., Cobena, G., Masanès, J., Sedrati, G. (2002). A first experience in archiving the French Web. Paper presented at the Proceedings of the 6th European Conference on Research and Advanced Technology for Digital Libraries . 2
- Abiteboul, S., Preda, M., Cobena, G. (2003). Adaptive on-line page importance computation. Paper presented at the Proceedings of the twelfth international conference on World Wide Web . 3
- Antoniol, G., Canfora, G., Cimitile, A., De Lucia, A. (1999). Websites: files, programs or database. Paper presented at the 1st International Workshop on Web Site Evolution, Atlanta, USA . 4
- Arvidson, A., Persson, K., Mannerheim, J. (2000). The Kulturarw3 project – The Royal Swedish Web Archive – An example of "complete" collection of web pages. Paper presented at the 66th IFLA – International Federation of Library Associations and Institutions, Jerusalem . 5
- Baeza-Yates, R., Castillo, C. (2005). Characteristics of the Web of Spain. *Cybermetrics*, 9 . 6
- Baeza-Yates, R., Castillo, C., Efthimiadis, E. (2005a). Characterization of national Web domains . 7
- Baeza-Yates, R. A., Castillo, C., Marin, M., Rodriguez, A. (2005b). Crawling a country: better strategies than breadth-first for Web page ordering. Paper presented at the WWW 05: Proceedings of the 14th international conference on World Wide Web, Chiba, Japan . 8
- Balayé, S. (1988). *La Bibliothèque nationale, des origines à 1800 (Histoire des idées et critique littéraire; vol. 262)*. Genève: Droz . 9
- Battelle, J. (2005). Google Announces New Index Size, Shifts Focus from Counting. <http://battellemedia.com/archives/001889.php> . 10
- Benjamin, W. (1963). *Das Kunstwerk im Zeitalter seiner technischen Reproduzierbarkeit; drei Studien zur Kunstsoziologie*. [Frankfurt am Main]: Suhrkamp . 11
- Bergman, M. I. K. (2001). The deep Web: Surfacing hidden value. *The Journal of Electronic* . 13

- Bergmark, D. (2002). Collection synthesis. Paper presented at the 2nd ACM/IEEE-CS joint conference .14
on Digital libraries, Portland, USA
- Bergmark, D., Lagoze, C., Sbityakov, A. (2002). Focused crawls, tunneling, and digital libraries. Paper .15
presented at the 6th European Conference on Research and Advanced Technology for Digital Libraries,
Roma, Italy
- Berners-Lee, T. Connolly, D. (1995). Hypertext Markup Language – 2.0. RFC,1866 .16
- Berners-Lee, T. (1994). Universal Resource Identifiers in WWW, A Unifying Syntax for the Expression .17
of Names and Addresses of Objects on the Network as used in the World- Wide Web. RFC 1630
- Berners-Lee, T. (1998). Cool URIs don't change. [http://www.w3.org/Provider/ Style/ URI.html](http://www.w3.org/Provider/Style/URI.html) .18
- and ultimate destiny of the World Wide Web by its inventor (1st pbk. ed.). New York: HarperCollins .19
- .Björneborn, L. Ingwersen, P. (2001). Perspective of webometrics .20
Scientometrics, 50(1), 65–82 .21
- Bolter, J. D. (2001). Writing space: Computers, hypertext, and the remediation of print (2nd ed.). .22
Mahwah, NJ: Lawrence Erlbaum Associates
- Borgman, C. L. (2000). From Gutenberg to the global information infrastructure: access to information .23
in the networked world (Digital libraries and electronic publishing). Cambridge, MA: MIT
- Borgman, C. L. (2003). The Invisible Library: Paradox of the Global Information Infrastructure. Library .24
Trends, 51(4), 652–674
- Boudrez, P. Eynde, V. D., Sofie. (2002). Archiving Websites .25
- Boufkhad, Y. Viennot, L. (2003). The Observable Web. RR .26
- Boyko, A. (2004). Test Bed Taxonomy. IIPC Reports, 16 .27
- Broder, A., Kumar, R., Maghoul, F., Raghavan, P., Rajagopalan, S., Stata, R., et al. (2000). Graph .28
structure in the web. Paper presented at the 9th International World Wide Web Conference (WWW9),
Amsterdam, Netherlands

Brown, A. (2006). Archiving the Web: A guide for information management professionals. Library . 29
.Assn Pub

Brügger, N. (2005). Archiving Websites, general considerations and strategies. A arhus, Denmark: . 30
Center for Internet Research

ص: 53

- Bruns, A. (2005). *Gatewatching: Collaborative online news production* (Digital formations, v. 26). .31
New York: P. Lang
- Burner, M. (1997). *Crawling towards Eternity Building An Archive of The World Wide Web*. New .32
Architect, 5
- Canfora, L. (1989). *The vanished library* (Hellenistic Culture and Society; 7). Berkeley: University of .33
California Press
- Canfora, L. (1996). *Les bibliothèques anciennes et l'histoire des textes*. In M. Baratin, C. Jacob (Eds.), .34
Le pouvoir des bibliothèques: la mémoire des livres en Occident. (pp. 338 p). Paris: A. Michel
- Carlin, J. W. (2004). *Harvest of agency public websites*. NARA Bulletin, 2005-02 .35
- Castells, M. (1996). *The rise of the network society*. Malden, MA: Blackwell .36
- Castillo, C., Marin, M., Rodriguez, A., Baeza-Yates, R. A. (2004). *Scheduling Algorithms for Web* .37
Crawling
- Chakrabarti, S. (2002). *Mining the Web: discovering knowledge from hypertext data*. San Francisco, .38
CA: Morgan Kaufmann
- Chakrabarti, S., Berg, M. V. D., Dom, B. (1999). *Focused crawling: A new approach to topic-specific* .39
Web resource discovery. *Computer Networks* (Amsterdam, Netherlands: 1999), 31, 1623-1640
- Chang, K. C.-C., He, B., Li, C., Patel, M., Zhang, Z. (2004). *Structured* .40
,databases on the web: observations and implications. *SIGMOD Record* .41
61-70 ,(3)33 .42
- Charlesworth, A. (2003). *Legal issues relating to the archiving of Internet resources in the UK, EU, USA* .43
and Australia
- Cho, J., Garcia-Molina, H. (2000). *The evolution of the web and implications for an Incremental* .44
Crawler. Paper presented at the Proceedings of the 26th International Conference on Very Large Data Bases
- Cho, J., Garcia-Molina, H., Page, L. (1998). *Efficient Crawling Through url ordering*. *Computer* .45
Networks and Isdn Systems, 30, 161-172

Christensen-Dalsgaard, B. (2001). Archive experience, not data. Paper presented at the Preserving the .46
Present for the Future – Strategies for the Internet, The Royal Library, Copenhagen, Denmark

Crowston, K., Williams, M. (1997). Reproduced and emergent genres of communication .47

ص: 54

- on the World-Wide Web. Paper presented at the 30th Annual Hawaii International Conference on System Sciences (HICSS-30), Wailea, USA
- Cruse, P., Eckman, C., Kunze, J. (2003). Web-based government information: Evaluating solutions for .48 capture, curation, and preservation. An Andrew W. Mellon funded initiative of the California Digital Library
- Dahn, M. (2000). Counting Angels on a Pinhead: Critically Interpreting Web Size Estimates. Online, .49
January/February, 35-40
- Day, M. (2006). The long-term preservation of Web content. In J. Masanè s (Ed.), Web archiving. .50
Berlin Heidelberg New York: Springer
- Dikaiakos, M. D. (2004). Intermediary infrastructures for the World Wide web. Computer Networks, .51
45(4), 421-47
- Dobra, A., Fienberg, S. E. (2004). How Large Is the WorldWide Web?. In M. Levene, A. Poulouvasilis .52
(Eds.), Web dynamics web dynamics – adapting to change in content, size, topology and use. (pp. 23-44).
Berlin Heidelberg New York: Springer
- Dubberly, H., Forlizzi, J., Hodge, C., Laurel, B., Lyman, P., Meggs, P. B., et al. (2002). Archiving .53
experience design, a virtual roundtable discussion. LOOP: AIGA Journal of Interaction Design Education,
Number 6
- Dumais, S. T., Cutrell, E., Cadiz, J. J., Jancke, G., Sarin, R., et al. (2003). Stuff I've seen: A system for .54
personal information retrieval and re-use. Toronto, Canada
- Egghe, L. (2000). New informetric aspects of the Internet: some reflections – many problems. Journal of .55
Information Science, 26(5), 329-335
- Eisenstein, E. L. (1979). The printing press as an agent of change: Communications and cultural .56
transformations in early modern Europe. Cambridge [Eng.]; New York: Cambridge University Press
- (Entlich, R. (2004). Blog Today, Gone Tomorrow? Preservation of Weblogs. RLG DigiNews, 8(4 .57
- Eriksen, L. B. Ihlström, C. (2000). Evolution of the web news genre – The slow move beyond the print .58
metaphor. Paper presented at the 33rd Hawaii International Conference on System Sciences (HICSS-33),
Hawaii, USA

Estivals, R. (1961). Le dé pôt lé gal sous l'Ancien Ré gime, de 1537 a 1791. Paris: M. Rivière .59

Estivals, R. (1965). La statistique bibliographique de la France sous la monarchie au .60

ص: 55

- Fauconnier, S. Frommé, R. (2004). Capturing unstable media, summary of research .61
- Fayet-Scribe, S. (2000). Histoire de la documentation en France: Culture, science, et technologie de l'information, 1895-1937 (CNRS histoire). Paris: CNRS .62
- Featherstone, M. (2000). Archiving cultures. *British Journal of Sociology*, 51(1) .63
- Febvre, L.P.V. Martin, H. J. (1976). The coming of the book: The impact of printing 1450-1800 ([New ed.] ed.). London: NLB .64
- Fielding, R. T., Gettys, J., Mogul, J., Nielsen, H. F., Masinter, L., J, P., et al. (1999). Hypertext Transfer Protocol - HTTP/1.1. RFC, 2616 .65
- Fitch, K. (2003). Web site archiving: An approach to recording every materially different response produced by a website. Paper presented at the AusWeb (2003): The Ninth Australina World Wide Web Conference, Sanctuary Cove, Australia .66
- Florescu, D., Levy, A., Mendelzon, A. (1998). Database techniques for the World- Wide Web: A survey. *SIGMOD Record* 27, 59-74 .67
- Freeman, E. Gelernter, D. (1996). Lifestreams: A storage model for personal data. *SIGMOD Record*, 25(1), 80-86 .68
- Gemmell, J., Bell, G., Lueder, R., Drucker, S., Wong, C. (2002). MyLifeBits: fulfilling the Memex vision. Juan-les-Pins, France .69
- Gibson, D., Punera, K., Tomkins, A. (2005). The volume and evolution of web page templates. Paper presented at the WWW'05 14th international conference on World Wide Web, Chiba, Japan .70
- Gillies, J. Cailliau, R. (2000). How the Web was born: The story of the World Wide Web. Oxford: Oxford University Press .71
- Golder, S. Huberman, B. A. (2005). The Structure of Collaborative Tagging Systems 73. Gomes, D. .72
- Silva, M. J. (2003). A Characterization of the Portuguese Web. Paper presented at the 3rd Workshop on Web Archives (IWA'03), Trondheim, Norway .73

Gulli, A. Signorini, A. (2005). The indexable web is more than 11.5 billion pages. Chiba, Japan .74

Halavais, A. (2004). Tracking Ideas in the Blogosphere .75

Hallgrímsson, T. (2006). Access and finding aids or web archives. In J. Masanè s (Ed.), Web archiving. .76
Berlin Heidelberg New York: Springer

Hine, C. (2000). Virtual ethnography. London; Thousand Oaks, CA: Sage .77

ص: 56

- Hofmann, M. Beaumont, L. R. (2005). Content networking: Architecture, protocols, and practice (The Morgan Kaufmann Series in Networking). Amsterdam; Boston: Morgan Kaufmann
- Ingwersen, P. (1998). The calculation of web impact factors. *Journal of Documentation*, 54(2) (Documentation, 54(2) .80
- Jones, S. Johnson, C. (2006). Web Use and Web Studies. In J. Masanè s (Ed.), *Web Archiving*. Berlin Heidelberg New York: Springer .81
- Jones, W., Bruce, H., Dumais, S. (2001). Keeping found things found on the web. Atlanta, GA, USA .82
- Jones, W., Bruce, H., Dumais, S. (2003). How do people get back to information on the Web? How can they do it better? Paper presented at the IFIP INTERACT'03 .83
- Kahle, B. (1997). Preserving the Internet. *Scientific American*, 397, 82–84 .85
- Kahle, B. (2002). The Internet Archive. *RLG DigiNews*, 6(3) .86
- Kimpton, M., Braggs, M., Ubois, J. (2006). Year by Year: From an Archive of the Internet to an Archive on the Internet. In J. Masanè s (Ed.), *Web Archiving*. Berlin Heidelberg New York: Springer .87
- Koehler, W. (1999). Unraveling the ISSUES, ACTORS, ALPHABET SOUP of the Great Domain) Name Debates. *Searcher*, 7(5) .88
- Koehler, W. (2004). A longitudinal study of Web pages continued: a consideration of document persistence. *Information Research*, 9(2) .89
- Krishnamurthy, B. Rexford, J. (2001). Web protocols and practice: HTTP/1.1, networking protocols, caching, and traffic measurement. Boston, MA: Addison–Wesley 91. Lagoze, C., Dean B. K., Sandy, P., Jesurogaili, S. (2005). What Is a Digital Library Anymore, Anyway? Beyond Search and Access in the NSDL. *D-Lib Magazine*, 11-11 92. Lampos, C., Eirinaki, M., Jevtuchova, D., Vazirgiannis, M. (2004). Archiving the Greek Web. Paper presented at the 4th International Web Archiving Workshop (IWA'04), (Bath (UK
- Landow, G. P. (1997). *Hypertext 2.0* (Rev., amplified ed.). Baltimore: Johns Hopkins University Press .93
- Lavoie, B. F. Schonfeld, R. C. (2005). The systemwide print book collection. Paper presented at the CNI .94

- Lawrence, S. Giles, C. L. (1998). Searching the Web. *Science*, 281, 175 .95
- Lawrence, S. Giles, C. L. (1999). Accessibility of Information on the Web. *Nature*, 400, 107–109 .96
- Lecher, H. E. (2004). Informant networks, alarm systems, and research contributors. Selection and ingest process for the Digital Archive for Chinese Studies. Paper presented at the Archiving Web Resources Conference – Issues for Cultural Heritage Institutions, NLA, Canberra, Australia .97
- Lecher, H. E. (2006). Academic Web archiving: DACHS. In J. Masanè s (Ed.), *Web archiving*. Berlin Heidelberg New York: Springer .98
- Levy, P. (1997). *Collective intelligence: Mankind's emerging world in cyberspace*. Cambridge, MA: Perseus Books .99
- Liu, C. Albitz, P. (1999). *DNS BIND (3rd ed.)*. O'Reilly Associates .100
- Lueg, C. Fisher, D. (2003). *From Usenet to CoWebs: Interacting with social information spaces (Computer supported cooperative work)*. Berlin Heidelberg London New York: Springer .101
- Lyle, J. A. (2004). Sampling the Umich.edu Domain. Paper presented at the 4th International Web Archiving Workshop (IWA'04), Bath (UK) .102
- Lyman, P. (2002). Archiving the World Wide Web. In CLIR (Ed.), *Building a national strategy for preservation: issues in digital media archiving*. Council on Library and Information Resources and the Library of Congress .103
- Lyman, P. Kahle, B. (1998). Archiving digital cultural artifacts. *D-Lib Magazine* .105
- Mantratzis, C. Orgun, M. (2004). Towards a peer2peer world-wide-web for the broadband-enabled user community .106
- Masanè s, J. (2002). Towards continuous Web archiving: First results and an agenda for the future. *D-Lib Magazine*, 8(12) .107
- Masanè s, J. (2004). Site-first priority: Implementing the frontline .108
- Masanè s, J. (2005). Web archiving methods and approaches: A comparative study. *Library Trends* .109

Masanès, J. (2006a). Collecting the hidden web. In J. Masanès (Ed.), Web archiving. Berlin . 110
Heidelberg New York: Springer

Masanès, J. (2006b). Selection for Web Archives. In J. Masanès (Ed.), Web archiving. Berlin . 111
Heidelberg New York: Springer

- Mohr, G., Kimpton, M., Stack, M., Ranitovic, I. (2004). Introduction to Heritrix, an archival quality web crawler. Paper presented at the 4th International Web Archiving Workshop (IWAW'04), Bath (UK) .112
- Mueller, M. (2002). Ruling the root: Internet governance and the taming of cyberspace. Cambridge, MA: MIT .114
- Najork, M. Heydon, A. (2001). High-performance Web crawling. SRC Research Report .115
- Najork, M. Wiener, J. (2001). Breadth-first search crawling yields high-quality pages. Paper presented at the 10th World Wide Web Conference (WWW'10), Hong Kong .116
- National Archives of Australia. (2001). Archiving Web resources: A policy for keeping records of web-based activity in the Commonwealth Government .117
- Osborn, T. (1999). The ordinariness of the archive. *History of the human sciences*, 12(2) .118
- Page, L., Brin, S., Motwani, R. Winograd, T. (1998). The Pagerank citation ranking: Bringing order to the Web, 17 .119
- Pandey, S. Olston, C. (2005). User-centric Web crawling. Chiba, Japan .120
- Pant, G., Srinivasan, P. Menczer, F. (2004). Crawling the Web. In M. Levene, A. Poulouvasilis (Eds.), *Web Dynamics*. (pp. 153-178). Berlin Heidelberg New York: Springer .121
- Pastor-Satorras, R. Vespignani, A. (2004). Evolution and structure of the Internet: A statistical physics approach. Cambridge, UK; New York: Cambridge University Press .122
- Phillips, M. E. (2005). Selective archiving of Web Resources: A study of acquisition costs at the National Library of Australia. *RLG DigiNews*, 9(3) .123
- Qin, J., Zhou, Y. Chau, M. (2004). Building domain-specific web collections for scientific digital libraries: A meta-search enhanced focused crawling method. Tuscon, AZ, USA .124
- Rekimoto, J. (1999). Time-machine computing: A time-centric approach for the information environment. Paper presented at the 12th annual ACM symposium on User interface software and technology, Asheville, North Carolina, USA .126

Riché , P. (1996). La bibliothè que et la formation de la culture mé dié vale. In M. Baratin, C. Jacob) .127
. (Eds.), Le pouvoir des bibliothè ques: la mé moire des livres en Occident (p

ص: 59

- Ringel, M., Cutrell, E., Dumais, S., Horvitz, E. (2003). Milestones in Time: The Value of Landmarks . 128
in Retrieving Information from Personal Stores. Paper presented at the IFIP INTERACT '03
- ?RLG. (2006). Web Archiving Program. <http://www.rlg.org/en/page.php> .129
- Page-ID=399 .130
- Roche, X. (2006). Copying web sites. In J. Masanè s (Ed.), Web Archiving. Berlin Heidelberg New . 131
York: Springer
- Rosenfeld, L. Morville, P. (2002). Information architecture for the World Wide Web (2nd ed.). . 132
Cambridge, MA: O'Reilly
- Scharl, A. (2000). Evolutionary Web development (Applied computing). Berlin Heidelberg New . 133
York: Springer
- Shepherd, M. Polanyi, L. (2000). Genre in Digital Documents. Paper presented at the Proceedings of . 134
the 33rd Hawaii International Conference on System Sciences - vol. 3
- Sonnenreich, W. (1997). A History of Search Engines. <http://www.wiley.com/legacy/compbooks/sonnenreich/history.html> . 135
- Spinellis, D. (2003). The decay and failures of web references. Communications of ACM, 46(1), 71- . 136
77
- Stack, M. (2005). Full Text Search of Web Archive Collections. Paper presented at the IWWAW'05, . 137
Vienna, Austria
- Star, S. L. Ruhleder, K. (1994). Steps towards an ecology of infrastructure: Complex problems in . 138
design and access for large-scale collaborative systems. Chapel Hill, NC, United States
- Teevan, J. (2004). How people re-find Information when the Web changes. AIM- 2004-012 .139
- Thelwall, M. (2001). Extracting macroscopic information from Web links. Journal of the American . 140
Society for Information Science and Technology, 52(13), 1157-1168 141. Thelwall, M. (2006).
Interpreting social science link analysis research: A theoretical framework. Journal of American Society of

Thelwall, M. Harries, G. (2004). Do the websites of higher rated scholars have .142

ص: 60

significantly more online impact? *Journal of the American Society for Information Science and Technology*,
55(2), 149-59

Thelwall, M. Vaughan, L. (2004). A fair history of the Web? Examining country balance in the Internet .143
archive. *Library Information Science Research*, 26(2), 162- 176

Ubois, J. (2002). The Oakland archive policy. Recommendations for managing removal requests and .144
preserving archival integrity

Voerman, G., Keyzer, A., Hollander, F. D., Druiven, H. (2002). Archiving the Web: Political Party) .145
Web sites in the Netherlands. *European Political Science*, 2 (1

پیشرفت در رایانه و ارتباطات این امکان را فراهم ساخته که ذخیره کتاب، صفحه موسیقی، فیلم، بسته های نرم افزاری و تمام صفحات عمومی وب که تاکنون ساخته شده اند هزینه - سودمندی داشته باشند و دسترسی به این مجموعه ها از طریق اینترنت برای همه، از جوان تا پیر، در سرتاسر دنیا فراهم شود. برای سال های آینده رسالت آرشیو مشخص است: ایجاد کتابخانه ای جهانی که تمام دانش به سهولت در اختیار هر مرد و زن و کودک در سرتاسر جهان قرار گیرد. در مقاله حاضر بحث آرشیو اینترنت و آرشیو در محیط اینترنت، دسترس پذیری دراز مدت تمام دانش برای همه افراد جهان پیشرفت های فن آوران ایجاد آرشیو اینترنتی مورد بررسی قرار می گیرد. سپس برنامه ها و پروژه های نمونه آرشیو اینترنتی در دنیا آرشیو اروپا و پتاباکس معرفی می شوند.

اشاره

*از آرشیو اینترنت تا آرشیو در اینترنت (1)

میشل کیمتون (2) | جف یوبویس (3)

ترجمه : مرضیه هدایت (4)

مقدمه

«آرشیو اینترنت» (5)، از آغاز کار خود در 1995، تاکنون، هدف درازمدت فراهم کردن دسترسی جهانی به تمام دانش در طول عمر را دنبال کرده است.

در 10 سال گذشته، آرشیو، در برنامه های کشورهای مختلف، برای کاربران متعدد و گوناگون، شرکت داشته و برای حفظ صفحات وبی، کتاب، موسیقی، نرم افزار، و تصاویر پویا فعالیت کرده تا آن ها را از طریق اینترنت دسترس پذیر سازد.

در حال حاضر، ماشین Wayback آرشیو، روزانه به 70,000 بازدید کننده، و در هر ثانیه به 200 درخواست پاسخ می دهد. مجموعه 600 ترابایتی آن شامل 50 بیلیون صفحه وب، 30000 جلد کتاب 36000 قطعه موسیقی، و 15000 قطعه فیلم است.

حضور صدها مشترک این امر را امکان پذیر ساخته است. به واسطه پیشرفت های فنی در صنعت رایانه،

ص: 63

Year – by – Year: From an Archive of the Internet to an Archive on the Internet: in Masanes, Julien (ed.), – 1

.Web Archiving. Berlin Heidelberg New York: Springer.pp.201-212

Michele Kimpton –2

Jeff Ubois –3

4- عضو هیئت علمی سازمان اسناد و کتابخانه ملی ایران

Internet Archive –5

امکانات این سیستم هر 12 - 18 ماه دو برابر می شود. در جامعه آرشیوی 10 سال زمان زیادی نیست، اما زمان مناسب برای خلق یک عامل 100 برابر پیشرفت است: از سال 1997 تاکنون قیمت دیسک خام بیش از 99 درصد کاهش یافته است؛ یعنی از 180 دلار به 50 سنت برای هر گیگابایت رسیده و بیش از 25 میلیون ارتباط پهنای باند فقط در ایالات متحده اضافه شده است.

با توجه به اینکه میزان پیشرفت رایانه بر اساس ذخیره و بازیابی است، می توان آن چه را که راج ردی (1) از دانشگاه کارنگی ملون (2) با عبارت «دستیابی جهانی به تمام دانش» بیان کرده است، محتمل دانست.

پیشینه: چاپ اینترنتی اولیه

از اواسط دهه 1980 معلوم بود که تغییری در دنیای چاپ الکترونیکی در حال ظهور است. در نیمه اول دهه 1990، حرکت روزنامه ها از شکل پایگاه های بسته با مالکیت انحصاری به شکل اینترنتی شروع شد و نظریه اینترنت به عنوان کتابخانه شروع به شکل گیری کرد. نظام های چاپی اینترنتی از قبیل (3) WAIS و Gopher، به عنوان متمم های صفحات وبی دیده شدند؛ استعاره اینترنت به عنوان کتابی با صفحات وبی که فهرست مندرجات آن توسط سرورهای Gopher و نمایه آن توسط سرورهای WAIS تهیه می شد، به عنوان جایگزینی برای استعاره اینترنت به عنوان «فرشاهراه» محسوب شد.

شروع به کار خدمات آلتا ویستا (4) در دسامبر 1995 ثابت کرد که تمام صفحات موجود در وب را می توان مجموعه ای یگانه دانست که قابل نمایه شدن و جست و جو کردن بر روی شبکه برای تمام کاربران است؛ اما معلوم نبود که این صفحات را چگونه باید حفظ و نگهداری کرد.

آغاز به کار آرشیو اینترنت

آرشیو اینترنت، به طور رسمی با همکاری بروس گیلین (5) و بروستر کال (6)، در آوریل 1996 شروع شد. در آن زمان، پیوندهای شکسته (404 درگاه) یک مشکل روبه رشد بود و روشن بود که بیشتر صفحات وبی عمر کوتاهی دارند. برای این مشکل راه حلی مورد نیاز بود و نظامی برای آرشیو صفحات وبی قبل از اینکه پاک شوند یک راهکار ضروری به نظر می رسید.

این مسئله، به تصمیم گیری راجع به طرح اولیه در «آرشیو» درباره سیاستگذاری مجموعه ها منجر شد؛ مانند حریص بودن نسبت به جمع آوری مطالبی که در معرض خطر از بین رفتن بودند، و شکار فرصت ها برای جمع آوری و رقومی کردن مواردی از گذشته مثل پستینگ های یوزنت (7).

با این حال، هنوز در 1996، در جامعه اینترنتی اهمیت از دست رفتن صفحات وبی مشکل چندان حساسیت

ص: 64

Raj Reddy -1

Carnegie Mellon University -2

(Wide Area Information Servers (WAIS -3

Alta Vista -4

Bruce Gilliat -5

Brewster Kahle -6

Usenet -7

برانگیز نبود. چون وب سابقه تاریخی چندانی نداشت توضیح فواید صفحات از دست رفته مشکل بود.

برای نشان دادن ارزش بالقوه چنین صفحاتی، آرشیو، با مؤسسه اسمیتسونین (1) در واشنگتن دی.سی، برای جمع آوری نسخه فوری (2) وبگاه های تمام نامزدهای ریاست جمهوری 1996 (3) همکاری کرد. ابزار انجام این پروژه چندان مورد رضایت نبود. این ابزار به طور اساسی عبارت بودند از: تجهیزات رایانه های شخصی که بر پایه ضبط وب سایت ها از طریق دنبال کردن پیوندها از صفحه اصلی عمل می کردند

این داده ها، در نهایت، به آرشیو ریاست جمهوری اسمیتسونین منتقل شد که در حال حاضر شامل

صفحاتی از 5 حزب سیاسی و کاندیداهای پرشمار ریاست جمهوری از بیل کلینتون گرفته تا پت بوکانن (4) است. بسیاری از سایت ها در این آرشیو با حذف کاندیدها بسته شدند.

بر اساس این موفقیت، کتابخانه کنگره به آرشیو مأموریت داد تا مجموعه پیوسته متمرکزی از انتخابات سال 2000 ایجاد کند و این درخواست را نیز برای انتخابات 2002 تجدید کرد.

همچنین، در سال 1996، آرشیو، ارتباطش را با اینترنت الکسا شروع کرد. اینترنت الکسا، مؤسسه ای انتفاعی است که خزشگری و آرشیوسازی وب را در نوامبر شروع کرد تا داده های پیشنهادی نوار ابزار مرورگر درباره سایت های دیده شده را Plug-in (وصلینه) کند و بر اساس داده های جمع آوری شده از سایر کاربران و پیشنهاد درباره صفحات مرتبط - که ممکن است مورد توجه باشند - عمل می کرد.

دو پیشرفت دیگر از سال 1996 نیز قابل ذکرند. نخستین پیشرفت، فناوریانه است. در 1996، هنوز ارجحیت نوار نسبت به دیسک از نظر قیمت قابل توجه بود و آرشیو اولین نسل زیرساختار را با استفاده از رویات های ذخیره نواری ساخت که با ADIC 50 شروع شد. با وجود همکاری های سخاوتمندانه فروشندگان اصلی، در نهایت ثابت شد این کار قابل دفاع نیست. نیازهای دستیابی که توسط کاربران آرشیو مطرح می شدند بیش از حد شدید و زمان بازایی خیلی کند بود. همانگونه که بروس گیلیت به طنز می گوید: گرفتن یک صفحه می تواند در چند ثانیه باشد... یا چند روز بعد.

دومین پیشرفت به مسائل قانونی مربوط می شد. مسائل قانونی جمع آوری صفحات وبی با الگوریتم حریصانه و ارائه آن ها بر پایه «صرف نظر کردن» (5) آن گونه که آرشیو در آن سال شروع به انجامش کرد - روشن نبود. بهبود پروتکل ارتقای رویات های محدودیت متنی متعلق به آلتا ویستا گامی مهم بود چون «پیش فرض ها» را تغییر داد - صاحبان صفحات وبی که از گذاشتن آن ها در نمایه موتور جست و جو یا در آرشیو ابا داشتند، روشی ساده برای صرف نظر کردن یا انتقال صفحاتی که صاحب آن بودند (اثبات آن از طریق توان آن ها برای مناسب سازی راهنمای پایه ای وبگاه مورد نظر بود) داشتند. رویات های متنی برای انتقال موارد مورد نظر برای آن ها که صاحب صفحات در آرشیو بودند راه حلی ارائه داد، ولی نتوانست مشکل انتقال صفحات متعلق به دیگران را رفع کند. این موارد انتقالی، در کنفرانس کوچکی

ص: 65

Snapshots -2

<http://movie0.archive.org/96-Elections/index.htm> -3

Pat Buchanan -4

Opt-out -5

توضیح داده شد، که توسط آرشیو اینترنت در 2002، در یوسی برکلی (1) برگزار شد (یوبویس 2002).

ساختار پیوندی و روبات های نواری

جمع آوری صفحات وبی، داده های پیوندی و «رد پایهای استفاده ای (2)» به عنوان نمونه، فرصت هایی که میلیون ها کاربر وب با رفتن از صفحه ای به صفحه ای ایجاد کرده بودند؛ در اوایل 1997 با معرفی نوار ابزار الکسا، یک مرورگر Plug-in از طریق تهیه اطلاعات روی سایتی که در حال دیده شدن بود و نیز پیشنهادهایی برای سایت های مرتبط به کاربران کمک می کرد تا به ناوبری در وب پردازند.

داده های پیوندی و ردپاهای استفاده ای که توسط الکسا جمع آوری می شدند، به عنوان نظام فیلترینگ اشتراکی عمل می کرد و صفحاتی را برجسته می کرد که جامعه اینترنتی دارای بیشترین ارزش می دانست. پیوندها و کلیک ها بر اساس ارزش صفحه داده شده امتیاز می گرفتند.

توان مشخص کردن ارزش یک صفحه، به طور خودکار، تا حد زیادی با یک طرح اساسی دیگر مرتبط بود. در 1997، تعدادی از بزرگ ترین و موفق ترین کتابخانه های وبی، فهرست سایت هایی مانند یاهو بودند. اما در مورد نحوه سنجش راهکارهای فهرست نویسی دستنامه ای که ممکن بود زمانش بگذرد، اطمینانی نبود. آیا احتمال حذف فهرست نویسی دستنامه ای و پذیرش روند انتخاب های قطعی با دیدگاه «جمع آوری کامل» که با فراداده ایجاد شده توسط کاربر به جای فهرست ترکیب می شد، امکان پذیر بود؟

به نظر می رسید پاسخ مثبت باشد. به همین دلیل، الکسا بر اساس ثبت میزان استفاده، به خزش در صفحات پرداخت. سنجش آن از طریق داده های جمع آوری شده توسط نوار ابزار مرورگر انجام می شد. صفحاتی که بسیار دیده شده بودند، نخستین هایی بودند که برایشان نسخه پشتیبان تهیه می شد.

خزشگر الکسا، طوری تنظیم شده بود که هر 8 هفته یک نسخه فوری از وب تهیه کند، و هنوز این برنامه اجرا می شود، اگر چه اندازه هر خزش از یک ترابایت در سال 1997 به 100 ترابایت در 2004 رسیده است.

مورد دیگری که آرشیو در سال 1997 با آن روبه رو شد این بود که آیا به ذخیره روی نوار می توان اعتماد کرد یا ذخیره روی دیسک؟ از نظر قیمت هنوز نوار بر دیسک ارجح بود، ولی دستیابی کند بود.

همانگونه که در مقاله گری و شنوی (1999) (3) آمده «نسبت قیمت نوار، دیسک و RAM تقریباً 10:1 است، به این معنی که ذخیره روی دیسک 10 بار گران تر از ذخیره روی نوار و ذخیره روی RAM 100 بار گران تر از ذخیره روی دیسک است».

اما وقتی هزینه دستیابی را با یارد آهنی بسنجیم، دیسک واقعاً بسیار ارزان تر است؛ «هزینه آرشیو نواری در مقایسه با هر ترابایت ذخیره روی دیسک نصف است، ولی دستیابی آسان به داده ها از طریق نوار ممکن نیست. هزینه هر دستیابی از طریق نوار، به صورت تصادفی، حدود 100 هزار بار بیشتر است (دستیابی از طریق دیسک: 100 دستیابی در ثانیه با هزینه 1 دلار در مقابل 10/000 دستیابی در ثانیه از طریق نوار با هزینه 10000 دلار است)» (گری و شنوی 1999).

UC Berkeley -1

Usage trails -2

Gray and Shenoy -3

1998: حضور داده های آرشیوی بر روی (تقریباً) هر دسکتاپ

پس از آن که در 1996، نت اسکپ (1) در دسترس عموم قرار گرفت؛ در 1998، اینترنت و مشاغل وابسته به آن در یک تغییر جهانی در اولویت های سرمایه گذاری مورد توجه قرار گرفت. همانطور که بیلیون ها دلار در بازارهای عمومی، سرمایه گذاری های مخاطره آمیز و راه اندازی اینترنت هزینه می شد، شماره صفحات در آرشیو هر 3 - 6 ماه دو برابر و تعداد کاربران اینترنت نیز هر چند ماه دو برابر می شد. به تدریج دستیابی به یک دغدغه تبدیل می شد.

در تلاش برای ایجاد دسترسی به موجودی و تأسیس آرشیو و اینترنت الکسا، به عنوان بخشی از زیر ساختار اینترنت، الکسا قراردادهایی با مرورگر های مایکروسافت و نت اسکپ بست. این امر به معنای حضور الکسا در 90 درصد رایانه های دنیا بود و کاربران چه می دانستند چه نه، الکسا بود که دسترسی به داده های تهیه شده توسط آرشیو اینترنت را به آن ها می داد.

برای آرشیو، ارائه داده ها برای ده ها میلیون کاربر تأثیر شدیدی بر زیر ساختار نوار مدار آن داشت.

در پایان 1998 دو امر مشخص بود:

- به دلیل میزان زیاد تقاضا رفتن از نوار به دیسک نیاز مبرمی بود. همان طور که تقاضای دسترس افزایش می یافت توان رویات های نواری برای پاسخگویی مورد تردید قرار می گرفت.

- سیاستگذاری های مجموعه های دستی گران تر از دیسک بود، یعنی آرشیو سازی بر پایه فهرست نویسی دستی، به ویژه وبگاه ها، بسیار گران تر از آرشیو سازی تمام سایت های قابل دسترسی بر اساس داده های جمع آوری شده از کاربران نهایی توسط الکسا بود.

1999: از نوار تا دیسک، یک خزشگر جدید و تصاویر متحرک

موفقیت تجاری خدمات الکسا، که به واسطه حضورش در هر رایانه شخصی مرتبط با اینترنت، بخش مشخصی بود، باعث شد آمازون، در سال 1999، الکسا را خریداری کند. این امر، در نهایت، به تغییراتی در ساختار هر دو سازمان منجر شد.

یک پیشرفت فناورانه مهم در سال 1999، ابداع خزشگر جدید توسط اندی جوئل (2) بود. این خزشگر جدید، برای جمع آوری داده های وبی توانایی بیشتری داشت و از طریق ماشین های چندگانه قابل مدیریت بود. این خزشگر، به افزایش توانایی الکسا در فیلتر کردن 16 بیلیون URL و 4 بیلیون خزشگری و گسترش وسعت و عمق خزش خود پردازد.

همچنین، تصمیماتی درباره فرمت فایل ARC اخذ شد، که برای ذخیره صفحات وبی استفاده می شد. ویژه سازی فرمت فایل ARC - که برای نیازهای متعددی طراحی شده بود - از آغاز، در سال 1996 توسط مایک برنر (3) و بروستر کال توسعه یافت (برنر و کال 1996). نیازهای مذکور عبارت اند از:

• فایل باید خود شمول باشد: فایل باید اجازه بدهد اشیای متراکم، بدون استفاده از فایل نمایه همراه،

Netscape -1

Andy Jewel -2

Mike Burner -3

● فرمت باید برای پذیرش فایل های بازیابی شده از طریق پروتکل های شبکه ای متنوع مانند FTP ، HTTP ، اخبار، گوفر، و پست الکترونیکی قابل گسترش باشد؛

● فایل باید قابلیت «جریان داشتن» (1) داشته باشد: فایل باید بتواند فایل های چندگانه آرشیو را در جریان

داده ها به هم پیوند دهد؛

● رکورد باید با یکبار نوشتن ماندگار: باشد یکپارچگی فایل نباید به ایجاد بعدی نمایه درون فایلی محتوا وابسته باشد.

در عین حال که خاص سازی به نمایه بیرونی محتوا و شیء - جبرانی نیاز ندارد، چنین نمایه ای به میزان زیادی قابلیت بازیابی اشیای ذخیره شده در این فرمت را افزایش می دهد. در حال حاضر، آرشیو چنین نمایه هایی را فراهم کرده و به دنبال استاندارد سازی فرمت های آن ها از طریق کنسرسیوم حفاظت بین المللی اینترنتی (2) است.

در 999، حرکت از صفحات وبی، به سایر انواع داده ها نیز شروع شد. در این سال، هزینه ذخیره آن قدر کاهش یافت که آرشیو توانست به گردآوری تصاویر متحرک بپردازد. از طریق شراکت با ریک پری لینگر (3)، از آرشیو Prelinger، پروژه دیجیتالی کردن 1000 فیلم (با هزینه حداکثر 160000 دلار) و نیز آرشیو کردن اخبار تلویزیون در پایان سال شروع شد.

2000: ایجاد مجموعه های موضوعی وب

در سال 2000، آرشیو به سطحی از ثبات فناوری رسید. پذیرش داده های خزشگری روش معمول شد و مهاجرت از نوار به دیسک سپری مجموعه گشت.

عکس

شناسایی و باز شوند؛

- فرمت باید برای پذیرش فایل‌های بازبایی شده از طریق پروتکل‌های شبکه‌ای متنوع مانند HTTP، FTP، اخبار، گوفر، و پست الکترونیکی قابل گسترش باشد؛
- فایل باید قابلیت «جریان داشتن»^۱ داشته باشد: فایل باید بتواند فایل‌های چندگانه آرشیو را در جریان داده‌ها به هم پیوند دهد؛
- رکورد باید با یکبار نوشتن ماندگار باشد: یکپارچگی فایل نباید به ایجاد بعدی نمایه درون فایلی محتوا وابسته باشد.

در عین حال که خاص‌سازی به نمایه بیرونی محتوا و شیء - جبرانی نیاز ندارد، چنین نمایه‌ای به میزان زیادی قابلیت بازبایی اشیای ذخیره شده در این فرمت را افزایش می‌دهد. در حال حاضر، آرشیو، چنین نمایه‌هایی را فراهم کرده و به دنبال استانداردسازی فرمت‌های آنها از طریق کنسرسیوم حفاظت بین‌المللی اینترنتی^۲ است.

در ۹۹۹، حرکت از صفحات وبی، به سایر انواع داده‌ها نیز شروع شد. در این سال، هزینه ذخیره آنقدر کاهش یافت که آرشیو توانست به گردآوری تصاویر متحرک بپردازد. از طریق شراکت با ریک پری لینگر^۳، از آرشیو Prelinger، پروژه دیجیتالی کردن ۱۰۰۰ فیلم (با هزینه حداکثر ۱۶۰۰۰۰ دلار) و نیز آرشیو کردن اخبار تلویزیون در پایان سال شروع شد.

۲۰۰۰: ایجاد مجموعه‌های موضوعی وب

در سال ۲۰۰۰، آرشیو به سطحی از ثبات فناوری رسید. پذیرش داده‌های خزشگری روش معمول شد و مهاجرت از نوار به دیسک سپری گشت.

جدول ۱. مجموعه‌های آرشیو اینترنت در مارس ۲۰۰۰

اندازه	واحد	مجموعه
۱۳/۸ ترابایت	۱ بیلیون صفحه	وب (۱۹۹۶ تا ماه ۳ سال ۲۰۰۰)
۰/۰۵ ترابایت	۵۰/۰۰۰ سایت	FTP (۱۹۹۶)
۰/۵۹۲ ترابایت	۱۶ میلیون پستینگ	Usenet (۱۹۹۶ - ۱۹۹۸)

در سال ۲۰۰۰، انتخابات دیگری در ایالات متحده برگزار شد و این بار تمام انتخاب‌کنندگان دسترسی اینترنتی داشتند. از نظر سیاسی، روشن بود که حضور در اینترنت برای برنده شدن حیاتی است و این امر تمرکز بر برخط بودن سیاسی را افزایش داد. آرشیو اینترنت، با کتابخانه کنگره برای گردآوری سایت‌های سیاسی شریک شد.

این اقدام، اولین پروژه آرشیو با کتابخانه کنگره بود و برای بسیاری از کارکنان آن به معنی انتقال از

1. Streamable

2. International Internet Preservation Consortium (IIPC) <http://netpreserve.org>

3. Rick Prelinger

جدول 1. مجموعه‌های آرشیوی اینترنت در مارس 2000

در سال 2000، انتخابات دیگری در ایالات متحده برگزار شد و این بار تمام انتخاب‌کنندگان دسترسی اینترنتی داشتند. از نظر سیاسی، روشن بود که حضور در اینترنت برای برنده شدن حیاتی است و این امر تمرکز بر برخط بودن سیاسی را افزایش داد. آرشیو اینترنت با کتابخانه کنگره برای گردآوری سایت‌های سیاسی شریک شد.

این اقدام، اولین پروژه آرشیو با کتابخانه کنگره بود و برای بسیاری از کارکنان آن به معنی انتقال از

Streamable -1

International Internet Preservation Consortium (IIPC) <http://netpreserve.org> -2

Rick Prelinger -3

یک پروژه تجربی به یک مؤسسه ثابت بود.

فکر ایجاد دسترسی به آثار ناپایدار نگهداری شده باعث حرکت و جنبش آئی و آرشیو تصاویر متحرک (1) در سال 2000 تأسیس شد. در حال حاضر، این آرشیو، با پروانه Creative Commons، شامل فیلم‌هایی از آرشیو Prelinger مجموعه ای بیش از 1900 فیلم ناپایدار (2) (پیام‌های بازرگانی، آموزشی، صنعتی و آماتور) است. این مجموعه، در حال حاضر دارای بیش از 10 درصد تمام فیلم‌های ناپایدار تولید شده در سالهای 1927 و 1987 در آمریکاست و یکی از کامل‌ترین و متنوع‌ترین مجموعه فیلم‌های ژانرهای است که تعداد زیادی از آن‌ها نگهداری نشده‌اند.

2001: دسترسی از طریق ماشین Wayback: آرشیو یازده سپتامبر

طی یک بازه زمانی یکساله از مارس 2000 تا مارس 2001، آرشیو، اندازه موجودی خود را 3 برابر کرد و به بیش از 40 ترابایت رساند. در این دوره، آرشیو هر ماه تقریباً 10 ترابایت رشد می‌کرد.

عکس

یک پروژه تجربی به یک مؤسسه ثابت بود.

فکر ایجاد دسترسی به آثار ناپایدار نگهداری شده باعث حرکت و جنبش آئی و آرشیو تصاویر متحرک^۱ در سال ۲۰۰۰ تأسیس شد. در حال حاضر، این آرشیو، با پروانه Creative Commons، شامل فیلم‌هایی از آرشیو Prelinger، مجموعه‌ای بیش از ۱۹۰۰ فیلم ناپایدار^۲ (پیام‌های بازرگانی، آموزشی، صنعتی، و آماتور) است. این مجموعه، در حال حاضر، دارای بیش از ۱۰ درصد تمام فیلم‌های ناپایدار تولید شده در سال‌های ۱۹۲۷ و ۱۹۸۷ در آمریکاست و یکی از کامل‌ترین و متنوع‌ترین مجموعه فیلم‌های ژانرهای است که تعداد زیادی از آنها نگهداری نشده‌اند.

۲۰۰۱: دسترسی از طریق ماشین Wayback: آرشیو یازده سپتامبر

طی یک بازه زمانی یکساله از مارس ۲۰۰۰ تا مارس ۲۰۰۱، آرشیو، اندازه موجودی خود را ۳ برابر کرد و به بیش از ۴۰ ترابایت رساند. در این دوره، آرشیو، هر ماه تقریباً ۱۰ ترابایت رشد می‌کرد.

جدول ۲. مجموعه‌های آرشیو در مارس ۲۰۰۱

اندازه	واحد	مجموعه
۴۰ ترابایت	۴ بیلیون صفحه	۱۹۹۶- ماه ۳ سال ۲۰۰۱
۲ ترابایت	۲۰۰ میلیون صفحه	آرشیو انتخابات ۲۰۰۰
۰/۵ ترابایت	۱۶ میلیون پستینگ	۱۹۹۶-۱۹۹۸، ۲۰۰۰- ماه ۳ سال ۲۰۰۱
۰/۵ ترابایت	۳۶۰ فیلم	فیلم‌های آرشیو: حدود ۱۹۰۳- حدود ۱۹۷۳
<۰/۱ ترابایت	۵۰۰۰ صفحه	Arpanet: اسناد تاریخی

اما سال ۲۰۰۱ برای بسیاری از سازمان‌های دارای فناوری بالا در محدوده سانفرانسیسکو، سال سختی بود. سقوط بازار سرمایه، واگذاری صدها شرکت محلی، و حمله به مرکز تجارت جهانی در نیویورک بر کارهای آرشیو اثر گذاشت. به‌ویژه، از دست رفتن مشاغل «های‌تک»^۳ در واقعه ۱۱ سپتامبر باعث تمرکز بر آن شد. همانگونه که برای توان آرشیو نیز در تهیه تصاویر متحرک و پاسخ به این وقایع حکم نوعی آزمون را داشت. در اوایل ۲۰۰۱، شاید مهم‌ترین پرسش پیش روی آرشیو این بود که چگونه به بهترین وجه، دسترسی به مجموعه را فراهم کند. داده‌های بسیاری، به‌طور مستقیم، از طریق خدمات الکسا در اختیار عموم قرار می‌گرفت، ولی دسترسی مستقیم به مجموعه‌ها هنوز نیازمند مهارت‌های برنامه‌ریزی یونیکس^۴ بود.

برنامه‌ریزان الکسا، در قراردادی با آرشیو، برنامه ماشین Wayback را ساختند که خدمات دسترسی

1. Moving Images Archive
2. ephemeral
3. high tech
4. Unix

جدول ۲. مجموعه‌های آرشیو در مارس ۲۰۰۱

اما سال ۲۰۰۱ برای بسیاری از سازمان‌های دارای فناوری بالا در محدوده سانفرانسیسکو، سال سختی بود سقوط بازار سرمایه، واگذاری صدها شرکت محلی، و حمله به مرکز تجارت جهانی در نیویورک بر کارهای آرشیو اثر گذاشت. به‌ویژه، از دست رفتن مشاغل «های‌تک»⁽³⁾ در واقعه ۱۱ سپتامبر باعث تمرکز بر آن شد. همانگونه که برای توان آرشیو نیز در تهیه تصاویر متحرک و پاسخ به این وقایع حکم نوعی آزمون را داشت. در اوایل ۲۰۰۱، شاید مهم‌ترین پرسش پیش روی آرشیو این بود که چگونه به بهترین وجه، دسترسی به مجموعه را فراهم کند. داده‌های بسیاری، به‌طور مستقیم، از طریق خدمات الکسا در اختیار عموم قرار می‌گرفت، ولی دسترسی مستقیم به مجموعه

ها هنوز نیازمند مهارت های برنامه ریزی یونیکس (4) بود.

برنامه ریزان الکسا، در قراردادی با آرشیو ، برنامه ماشین Wayback را ساختند که خدمات دسترسی

ص: 69

Moving Images Archive -1

ephemeral -2

high tech -3

Unix -4

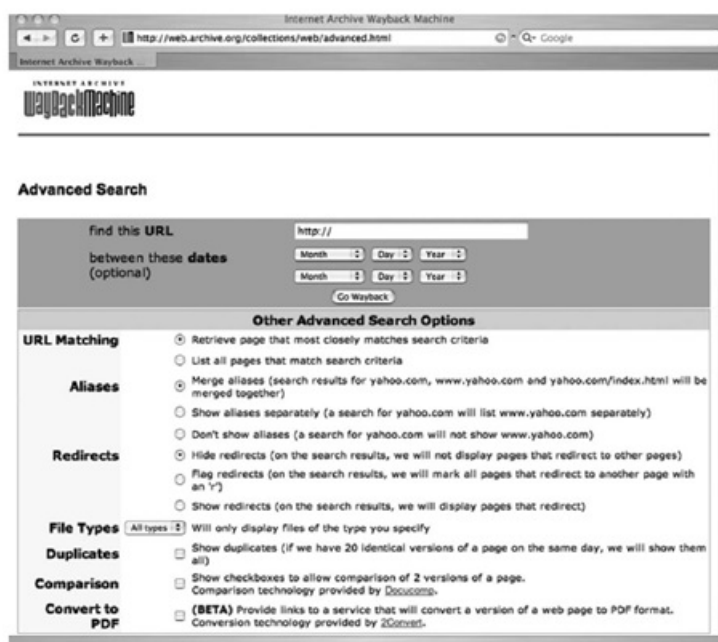
به محتوای آرشیو را بر پایه URL ها می داد 24 اکتبر 2001، ماشین Wayback به کار افتاد و دسترسی به بیش از 10 بیلیون صفحه وبی آرشیو شده و 100 ترابایت داده ممکن شد.

در آن زمان، داده ها روی Hewlett Packar ذخیره می شد و سرورهای uslab.com از سیستم های عامل Linux و FreeBSD استفاده می کرد. هر رایانه، حدود 512 مگابایت حافظه داشت و به طور کلی بیش از 300 گیگابایت روی دیسک های IDE بود.

عکس

۷۰ مدیریت منابع اطلاعاتی وب

به محتوای آرشیو را بر پایه URL ها می داد. ۲۴ اکتبر ۲۰۰۱، ماشین Wayback به کار افتاد و دسترسی به بیش از ۱۰ بیلیون صفحه وبی آرشیو شده و ۱۰۰ ترابایت داده ممکن شد.
در آن زمان، داده ها، روی Hewlett Packar ذخیره می شد و سرورهای uslab.com از سیستم های عامل Linux و FreeBSD استفاده می کرد. هر رایانه، حدود ۵۱۲ مگابایت حافظه داشت و به طور کلی بیش از ۳۰۰ گیگابایت روی دیسک های IDE بود.



تصویر ۹.۱. گزینه های فعلی جستجو برای Wayback machine

پروژه مهم دیگر سال ۲۰۰۱، آرشیو ۱۱ سپتامبر^۱ بود. با همکاری کتابخانه کنگره، آرشیو، تصاویری از بیش از ۳۰۰,۰۰۰ وبگاه منتخب را از ۱۱ سپتامبر ۲۰۰۱ تا ۱ دسامبر ۲۰۰۱ و صدها ساعت پخش اخبار را گردآوری کرد.

۲۰۰۲: کتابخانه اسکندریه^۲، کتابخانه سیار، و حق مؤلف

در سال ۲۰۰۲، آرشیو، ۵ برنامه مهم دیگر یعنی افزایش مجموعه ها، گزینه های دستیابی به این مجموعه ها، همکاری با سایر سازمان ها و تعیین راهکارها را بر عهده داشت.

1. September 11 Archive
2. The library of Alexandria

تصویر 1. 9. گزینه های فعلی جستجو برای Wayback machine

پروژه مهم دیگر سال 2001، آرشیو 11 سپتامبر (1) بود. با همکاری کتابخانه کنگره، آرشیو، تصاویری از بیش از 30,000 وبگاه منتخب را از 11 سپتامبر 2001 تا 1 دسامبر 2001 و صدها ساعت پخش اخبار را گردآوری کرد.

2002: کتابخانه اسکندریه، کتابخانه سیار، و حق مؤلف

*کتابخانه اسکندریه، کتابخانه سیار، و حق مؤلف (2)

در سال 2002، آرشیو، 5 برنامه مهم دیگر یعنی افزایش مجموعه ها، گزینه های دستیابی به این مجموعه ها، همکاری با سایر سازمان ها و تعیین راهکارها را بر عهده داشت.

ص: 70

September 11 Archive -1

The library of Alexandria -2

نخستین و بزرگ ترین پروژه ایجاد سایت قرینه در کتابخانه اسکندریه در مصر بود. سرورها و بیش از 100 ترابایت داده، که بیش از 5 میلیون دلار ارزش داشتند به مصر برده شد و برای افتتاح کتابخانه در ماه آوریل نصب شد.

دومین پروژه مهم، ایجاد کتابخانه سیار اینترنتی (1) بود که برای نمایش چگونگی ترکیب اسکن های الکترونیکی کتاب ها، چاپ بر اساس فناوری مورد تقاضا، و اینکه ارتباط شبکه ای ماهواره ای چگونه می تواند مناسب کتابخانه ای هزار جلدی در پشت یک ون باشد، طراحی شد پروژه یک میلیون کتاب (2) Million Books Project، در تابستان 2002، با شراکت کارنگی ملون شروع شد. هدف این بود که حداقل یک میلیون کتاب دیجیتالی شود و به صورت رایگان روی اینترنت قرار گیرد. با تشویق سوزان مبارک (3)، آرشیو، ساخت یک کتابخانه سیار را در ایالات متحده شروع کرد و با همکاری دیگران، در هند و کنیا، نیز نمونه هایی از آن ایجاد شد.

سومین پروژه مهم مربوط به سیاست گذاری بود. 30 سپتامبر 2002، در تلاشی برای افزایش آگاهی عمومی درباره اهمیت موارد راهکاری حق مؤلف و کتابخانه سیار اینترنتی، آرشیو به سفری در سر تاسر کشور و چاپ و توزیع کتاب های رایگان مبادرت کرد کتابخانه سیار را در حیاط ساختمان دادگاه عالی ایالات متحده پارک کرد و به چاپ کتاب پرداخت؛ جایی که قضات، در 9 اکتبر مباحثه میان الدرد (4) و اشکرافت (5) را شنیدند این حادثه مهمی بود که بر اساس آن تصمیم گرفته شد چه تعداد کتاب باید بخشی از کتابخانه دیجیتالی کتابخانه سیار و سایر کتابخانه های دیجیتالی در ایالات متحده باشد. متأسفانه، الدرد شکست خورد و تصمیم حق مؤلف مؤثر واقع افتاد؛ ولی پروژه کتابخانه سیار شکوفا شد و در نهایت، به صورت کتابخانه ای غیر انتفاعی درآمد. اخیراً، دولت هند، ساخت 25 کتابخانه سیار برای استفاده در سر تاسر کشور هند را شروع کرده است.

چهارمین حوزه فعالیت شامل ایجاد اولین آرشیو مجموعه کتاب و موسیقی است در ژوئن اولین مجموعه های کتاب، به صورت پیوسته در اختیار گذاشته شد در آگوست آرشیو موسیقی زنده، شامل مجموعه ای از اجراهای کنسرت - که به طور قانونی قابل بارگیری بودند - به صورت پیوسته درآمد.

پنجمین پروژه مهم، تأسیس کتابخانه دیجیتالی بین المللی کودکان (6)، با شرکت دانشگاه مری لند (7)، بود که از سوی کتابخانه کنگره، NSF، IMLS بنیاد کال / اوستین (8)، شرکت سیستم های ادوب (9)، بنیاد مرکل (10)، و اکتاوو (11) پشتیبانی شد. ICDL، بر ماهیت ذاتی اینترنت، مبتنی بر فراهم آوری دسترسی مستقیم و جهانی به محتوای کیفی برای کودکان متمرکز بوده و هست.

ص: 71

Internet Bookmobile -1

MBP)، -2

Suzanne Mubarak -3

Eldred -4

Ashcroft -5

International Children's Digital Library -6

University of Maryland -7

Kahle/Austin Foundation -8

.Adobe System Inc -9

Markle Foundation -10

Octavo -11

در پایان سال 2002، آرشیو، برای مراقبت از یکپارچگی آرشیوهای دیجیتالی از طریق استاندارد سازی معیارها، و جلوگیری از جابه جا یا غیر قابل دسترس شدن مواد تلاش کرد. در نشست در یوسی برکلی، نمایندگانی از آرشیو با سایر کتابداران دیجیتالی، به منظور تکمیل Okland Archive Policy ملاقات کردند که روند انتقال مواد را بر اساس قانون یا بر اساس خواست مالکان سایت و سایرین با جزئیات بیان می کرد.

2003: گسترش دستیابی ما به کتابخانه های ملی و مؤسسات آموزشی

در سال 2003، آرشیو، رسیدن به کتابخانه های ملی و مؤسسات آموزشی در سرتاسر دنیا را ادامه داد. آرشیو، با همکاری کنسرسیوم حفاظت بین المللی اینترنتی (IIPC) از نزدیک با سازمان های شریک روی استانداردهای جدید و یک خزشگر جدید منبع باز شروع به کار کردن کرد.

در جولای 2003، آرشیو، به کنسرسیوم حفاظت بین المللی اینترنتی برای آغاز به کار کمک کرد. گروهی متشکل از 12 کتابخانه ملی برای توسعه استانداردها، ابزار و راهکارهایی برای تهیه حفاظت و دسترس پذیری دانش و اطلاعات از اینترنت برای نسل های آینده در همه جا، ارتقای تبادل جهانی و ارتباطات بین المللی تلاش می کردند. برای انجام این رسالت، IIPC برای رسیدن به اهداف زیر کار می کند:

- رسیدن به مجموعه ای غنی از محتوای اینترنتی از سرتاسر دنیا که به گونه ای حفاظت شوند که بتوانند آرشیو و محافظت شده و در لحظه قابل دسترسی باشند؛

- توسعه و استفاده از ابزار عادی فناوری ها و استانداردهایی که ایجاد آرشیوهای اینترنتی را ممکن می سازند؛

- تشویق و حمایت از کتابخانه های ملی در همه جا برای آرشیو اینترنتی و حفاظت.

IPC در کتابخانه ملی فرانسه با 12 مؤسسه همکار مجاز شروع به کار کرد. اعضا موافقت کردند به اتفاق هزینه ها را پرداخت کرده و در برنامه ها و کارگروه ها برای رسیدن به اهداف مذکور شرکت کنند. موافقت اولیه برای 3 سال بود. طی این پروژه اعضا محدود به مؤسسات مجاز بودند.

در سال 2003، آرشیو بودجه قابل توجهی از خارج یعنی از سازمان های دیگر از جمله بنیادهای هیولیت واسلون (1) دریافت کرد و شروع به کار بر مجموعه های خاص کرد. کاهش های بعدی هزینه ذخیره روی دیسک و پهنای باند اینترنت آرشیو را به عرضه دائم «پهنای باند نامحدود، برای همیشه و رایگان» برای سازمان ها و افراد با مواد دیجیتالی، راهبر شد.

این امر به شراکت با Etree منجر شد. سازمانی که داوطلبانه در سال 1998 ایجاد شد تا تجارت آزاد و قانونی کنسرت های موسیقی زنده را امکان پذیر سازد. حاصل همکاری با Etree این شد که آرشیو e اکنون میزبان بیش از 15000 کنسرت موسیقی زنده است.

آرشیو، برای حمایت از نیازهای رو به رشد هم ذخیره و هم پهنای باند، یک مرکز جدید در سانفرانسیسکو باز کرد. این مرکز داده جدید، از طریق پیوند 1 Gbps به اینترنت متصل است و بیش از 1500 رایانه شخصی را که از لینوکس استفاده می کنند، میزبانی می کند.

اشاره

*آرشیو اروپا و پتاباکس (1) (2)

آرشیو، در سال 2004 شروع به انتقال داده ها به سومین نسل سخت افزاری خود، موسوم به پتاباکس، کرد. طرح پتاباکس، بر پایه سخت افزاری rack-mounted و سیستم عامل لینوکس، ذخیره RAID برای هر ترابایت به مبلغ تقریبی 2000 دلار یا 2 میلیون دلار برای هر پتابایت پیشنهاد کرد.

نخستین نصب این طرح جدید در آمستردام در آرشیو تازه تشکیل شده اروپا بود، مؤسسه ای که قرار بود پاسخگوی نیازهای جامعه اروپا باشد و آرشیو اینترنت با سایر شرکای اروپایی در اولین سال تأسیس، آن را حمایت می کرد. نصب آن در آمستردام برای ایجاد قرینه ای برای مجموعه های اسکندریه و سانفرانسیسکو است. ایجاد شبکه ای از مؤسسات مستقل در سرتاسر دنیا که هر یک قادر است به طور مستقل عمل کند به پیشگیری از نابودی های فاجعه آمیز اطلاعات کمک خواهد کرد.

همچنین، در سال 2004، کنسرسیوم حفاظت بین المللی اینترنتی (3 Heritrix)، یعنی خزشگر وی منبع باز، قابل توسعه با قابلیت تغییر اندازه وی، با کیفیت آرشیوی و بر پایه جاوا را شروع کرد.

در مسیر توسعه مجموعه در سال 2004، از طریق استخدام کارکنان بیشتر، تکمیل پروژه های اسکن کتاب و فیلم و دریافت داده از سایر مؤسسه ها، گام های مهمی رو به جلو برداشته شد.

آینده

پیشرفت در رایانه و ارتباطات این امکان را فراهم ساخته که ذخیره کتاب، صفحه موسیقی، فیلم، بسته های نرم افزاری و تمام صفحات عمومی وب که تاکنون ساخته شده اند، هزینه - سودمندی داشته باشند و دسترسی به این مجموعه ها از طریق اینترنت برای همه از جوان تا پیر در سرتاسر دنیا فراهم شود.

همانگونه که در اعلامیه حقوق بشر، ماده 19 آمده است: «هر انسانی حق آزادی بیان و عقیده را دارد. این حق شامل آزادی داشتن باور و عقیده بدون نگرانی از مداخله و مزاحمت و حق جست و جو، دریافت و انتشار اطلاعات و نظرها از طریق هر رسانه ای بدون ملاحظات مرزی است».

برای سال های آینده، رسالت آرشیو مشخص است: ایجاد کتابخانه ای جهانی که تمام دانش را به سهولت در اختیار هر مرد و زن و کودک در سرتاسر جهان قرار می دهد.

Burner, M. Kahle, B. (1996). The ARC File Format

Gray, J. Shenoy, P. (1999). Rules of Thumb in Data Engineering. Microsoft

Technical Report, MS-TR-99-100

Masanès, J. (2006). Web archiving: issues and methods. In J. Masanès (Ed.), Web Archiving. Springer,
Berlin Heidelberg New York

Ubois, J. (2002). The Oakland Archive Policy. Recommendations for Managing Removal Requests and
Preserving Archival Integrity

ص: 73

European Archive -1

Petabox -2

Heretrix -3

پیشرفت در رایانه و ارتباطات این امکان را فراهم ساخته که ذخیره کتاب، صفحه موسیقی، فیلم، بسته های نرم افزاری و تمام صفحه های عمومی وب که تاکنون ساخته شده اند، هزینه - سودمندی داشته باشند و دسترسی به این مجموعه ها از طریق اینترنت برای همه، از جوان تا پیر، در سر تاسر دنیا فراهم شود. برای سال های آینده، رسالت آرشیو مشخص است: ایجاد کتابخانه ای جهانی که تمام دانش را به سهولت در اختیار هر مرد و زن و کودک در سر تاسر جهان قرار دهد. در مقاله حاضر بحث آرشیو اینترنت و آرشیو در محیط اینترنت، دسترس پذیری دراز مدت تمام دانش برای همه افراد جهان پیشرفت های فناورانه ایجاد آرشیو اینترنتی مورد بررسی قرار می گیرد. سپس برنامه ها و پروژه های نمونه آرشیو اینترنتی در دنیا آرشیو اروپا و پتاباکس معرفی می شوند.

اشاره

*کاربرد وب و مطالعات مربوط به آن (1)

استیو جونز (2) | کمیل جانسون (3) (دانشگاه ایلینویز شیکاگو)

ترجمه: دکتر سید مهدی طاهری (4) | سید محمد موسوی (5)

خلاصه

اشاره

در سال 2002، مرکز کتابخانه رایانه ای پیوسته (OCLC) برآورد کرد که بیش از سه میلیون وبگاه در شبکه وب جهان گستر در دسترس عموم قرار دارد (اونیل و همکاران 2003، پاراگراف 9). در دنیای مادی اطلاعات، این تعداد وبگاه نزدیک به 14 تا 28 میلیون خواهد بود که با حجم کتاب های موجود در بزرگ ترین کتابخانه های دنیا برابری یا حتی از تعداد آن تجاوز می کند. این مجموعه گسترده اطلاعات، مجموعه ای بی پایان و دسترس پذیر از داده های دیداری و شنیداری را در اختیار پژوهشگران علاقه مند به مطالعه در حوزه فعالیت های پیوسته قرار می دهد. حجم خالص منابع دسترس پذیر بر روی وب همواره چالش هایی را به وجود آورده است چه چیزی را انتخاب و چگونه مطالعه کنیم؛ با این حال، شبکه وب ثابت کرده است این حجم اطلاعات برای کاربران دلهره آورتر هم شده است. در این مقاله، برخی رویکردهای روش شناختی را که پژوهشگران برای مطالعه در حوزه وب به کار گرفته اند، مرور خواهیم کرد.

ص: 75

Web Use and Web Studies: in Masanes, Julien (ed.), Web Archiving. Berlin Heidelberg New York: – 1
Springer. pp.55-67

Steve Jones –2

Camille Johnson –3

4- عضو هیئت علمی پژوهشگاه علوم و فرهنگ اسلامی

5- کارشناس ارشد کتابداری و اطلاع رسانی پژوهشگاه علوم و فرهنگ اسلامی

هدف از این کار، طبقه‌بندی جامع روش‌شناسی‌ها نیست به جای آن امیدواریم با درک روش‌های به کار گرفته شده در مطالعه وب و مطالعه کاربرد آن، افرادی که درصدد آرشو کردن و نگهداری وب هستند بتوانند نیازهای جامعه علمی را هرچه بهتر درک کنند.

شبکه وب عبارت است از طیف گسترده‌ای از انواع مواد که تنوع آن‌ها را با دو بُعد می‌توان بهتر درک کرد. نخست اینکه خود وب یک رسانه است نه محتوا برای روشن تر شدن مطلب باید گفت که وب هم رسانه‌ای است که محتوا را از طریق پروتکل‌های متعدد (نظیر HTTP) منتقل می‌کند و هم «ظرفی» برای محتواست که محتوا را شکل داده و به کاربران خود ارائه می‌دهد. هر چند ارائه محتوا به تشخیص کاربر وابسته است و گذشته از آن توسط ابزارهای دیداری مورد استفاده (مرورگرها و برنامه‌های کاربردی دیگر) شکل می‌گیرد. به عبارت دیگر، اگر چه محتوا ممکن است صفحه وب یکسانی باشد، دو کاربری که از مرورگرهای متفاوت، یا از تنظیمات متفاوتی از یک مرورگر استفاده می‌کنند، ممکن است در نهایت صفحه‌های متفاوتی را ببینند. دیگر اینکه، برخلاف دنیای مواد آنالوگ، همان تعریف رسانه برای ذخیره‌سازی وب مورد تردید است. وبگاه‌ها را می‌توان به صورت محلی ذخیره یا پنهان کرد، اطلاعات آن‌ها را منعکس نمود، یا ممکن است مجازی یا موقتی باشند همان‌طور که این ویژگی‌ها برای وب کم‌ها نیز صادق است.

به طور کلی، پژوهش درباره اینترنت و وب، توسط پژوهشگرانی با تخصص‌های مختلف از جمله متخصصان رشته‌های زبان‌شناسی، روزنامه‌نگاری، علوم سیاسی، مدیریت بازرگانی، جغرافیا، تبلیغات، ارتباطات، هنر، و دیگر رشته‌ها انجام می‌شود. از آن‌جا که با طیف وسیعی از سنت‌های آموزشی سر و کار داریم، انواع موادی که پژوهشگران از آن‌ها بهره می‌گیرند بسیار متنوع است و به همان اندازه نیز در روش‌های پژوهش تنوع مشاهده می‌شود. آن‌چه به عنوان تحلیل «قوم‌نگاری» متون وب، برای پژوهشگر بازاریابی به حساب می‌آید، ممکن است از دیدگاه پژوهشگر رشته ارتباطات تحلیل محتوای کیفی تعبیر گردد. بنابراین هدف دیگر این مقاله نه تنها ارائه مقوله‌های ثابت محتوای وب یا طرح‌هایی برای رویکردهای روش‌شناختی است - آن‌گونه که برای مطالعات وب به کار می‌روند - بلکه ارائه طیفی از تعابیر و کاربردهای این روش‌ها نیز مورد نظر است، به نحوی که مفیدترین روش‌ها را برای پژوهشگران حوزه وب به اثبات رساند.

1- تحلیل محتوا

تحلیل محتوا یکی از رایج‌ترین روش‌های مطالعه بر روی حوزه وب است. پژوهشگری که از روش تحلیل محتوا استفاده می‌کند، محتوای وب را - هم محتوای متنی و هم تصاویر - بر اساس معیارهای خاص رمزگذاری کرده و آن‌ها را درون مقوله‌ها یا موضوع‌های مرتبط قرار می‌دهد؛ به عبارت دیگر، این نوع تحلیل بیشتر بررسی محتوای وب است تا کاربران وب. تحلیل محتوا در میان مطالعات وب، به عنوان ابزار مقایسه مورد استفاده قرار گرفته است به نحوی که این ابزار به پژوهشگر اجازه می‌دهد تا مقایسه‌های معنی‌داری از محتوا میان متون مشابه وب انجام دهد. در بررسی وبگاه‌های سازمان‌های ضد جهانی سازی

از روش تحلیل محتوا استفاده شد تا نشان داده شود که آیا میان این سایت ها انسجام پیام و مقصود وجود دارد (ون الست و، والگریو 2002)؟ محتوا، بر اساس کارکرد اصلی خود در چهار حوزه رمزگذاری و مرتب شده است ارائه اطلاعات درباره سازمان مسائل ضد جهانی سازی افزایش تعامل با مؤسسان گروه و اعضای دیگر و افزایش بسیج افراد در مورد مسائلی چون اهدای پول که در یک فراخوان اینترنتی اعلام شده است و الست و والگریو، از طریق تحلیل خود به این نتیجه رسیدند که در حقیقت سازمان های ضد جهانی سازی پیوسته، عموماً به وسیله وبگاه های خود و از روش های مشابهی برای آگاه نمودن و مشغول کردن اعضای خود استفاده می کردند.

مطالعه تطبیقی دیگری به تحلیل محتوای وبگاه های مربوط به ایستگاه های رادیویی پرداخت، تا چگونگی پاسخگویی صنعت تجاری رادیو را به اندازه دسترس پذیری وب جهان گستر برای ارتقای ایستگاه های خود مشخص کند. این ارزیابی با تجزیه و تحلیل انواع اطلاعات کاربر - محور، که در وبگاه های ایستگاه هایی رادیویی فراهم شده بود برای مثال نقشه های ترافیکی و لاگ های برنامه ها، استفاده از وبگاه ها به عنوان ابزارهایی برای بهبود ایستگاه خود (مانند اطلاعات رقابتی و بایوس های DG)، و ترکیب ویژگی های تعاملی آن ها (نظیر نشانی ایمیل ها برای اعضای ایستگاه ها و نظرسنجی های شنوندگان) به دست آمد (پیتز و هارمز 2003). این مطالعه در حوزه وب منبع اطلاعاتی ارزشمندی را برای ارزیابی کاربردهای کنونی و بالقوه وب در اختیار ایستگاه های رادیویی قرار داد.

همچنین، روش تحلیل محتوا، برای مطالعه اثرات تغییرات خط مشی های سازمانی بر وبگاه های مرتبط به کار گرفته شده است. برای مثال، خط مشی های آموزشی، اهمیت الحاق فناوری اطلاعات و ارتباطات (ICT) به مأموریت های سازمانی مدارس راهنمایی را افزایش دادند، مطالعه ای، پیشرفت های صورت گرفته برای نیل به این هدف را همخوان با طرح ریزی انجام شده توسط شبکه ملی یادگیری در بریتانیا از طریق تحلیل محتوا در وبگاه های 150 مدرسه راهنمایی مورد ارزیابی قرار داد (هسکت و سولیوان 1999). تصاویر و متون رمزگذاری شدند، و به شناسایی وبگاه ها در یکی از مقوله های روشی زیر پرداختند: فعال، دانش آموز-محور، و منفعل. پژوهشگران از طریق این مقوله ها و تحلیل پرونده های پژوهشی مدارس، رابطه ای را بین تعهد به یکپارچگی فناوری اطلاعات و ارتباطات (ICT) و سرمایه های سازمانی درک نمودند؛ هر قدر سرمایه اقتصادی و اجتماعی مدرسه بیشتر باشد احتمال اینکه محتوای وبگاه آن منعکس کننده گرایش های مثبت و فعال الحاق فناوری اطلاعات و ارتباطات با مأموریت های مدرسه باشد بیشتر خواهد بود.

کمیسیون بازرگانی فدرال امریکا، از روش تحلیل محتوا، برای ارزیابی پیاده سازی برنامه های حفظ امنیت اطلاعات و حریم خصوصی پیوسته بر روی وبگاه های بازرگانی استفاده کرده است (Milne and Culnan). مطالعه انجام شده از سوی این کمیسیون، دارای تحلیل دراز مدت چهار نظر سنجی وب با گستره زمانی 1998 تا 2001 بود از معیارهایی نظیر حفظ یا حذف اعلان ها، با توجه به افشای اطلاعات بازدیدکنندگان با اشخاص ثالث و استفاده از «کوکی ها»، مجموعه ای از اطلاعات کاربران به غیر از نشانی ایمیل ها، و به کارگیری امنیت اطلاعاتی در سایت ها مورد استفاده قرار گرفت تا مشخص شود که آیا

وبگاه ها ملزومات شیوه های استفاده بی طرفانه از اطلاعات را تأمین می نمایند؟ مقایسه نتایج به دست آمده از تحلیل انجام شده در سال 2001 با تحلیل های سال های پیشین نشان می دهد، در میان وبگاه های تجاری که به گردآوری اطلاعات شخصی کاربران خود دست می زنند، روند افزایشی در اجرای شیوه های استفاده بی طرفانه از اطلاعات وجود داشته است.

2- بررسی ها

بررسی های مربوط به وب را می توان به دوروش متمایز تعریف کرد: بررسی هایی که بر روی وب انجام می گیرد، و بررسی های منتشر شده در وب. بررسی های نوع نخست که از تداول بیشتری برخوردار است به وسیله پژوهشگرانی انجام شده است که مطالعه درباره وبگاه ها را به منظور دسترس پذیری برای جمعیت علاقه مند ارسال می کنند. در بیشتر موارد، هدف از این نوع بررسی ها، گردآوری اطلاعاتی درباره موارد استفاده کاربران از اینترنت (1) است (مانند بررسی کاربران وب که در مرکز گرافیک، دیداری سازی، و کاربرد پذیری (The [Graphic, Visualization, and Usability Center](#) (2) انجام گرفته است (1998). GVU، از سال 1994، مطالعات مداومی درباره کاربرد اینترنت انجام داده است. این مرکز، شرکت کنندگانی را از طریق ثبت گروه های خبری مرتبط با اینترنت، تابلوهای تبلیغاتی در رسانه های خبری، وبگاه های موتورهای کاوش و اطلاعاتی های موجود در رسانه های ناپیوسته نظیر مجلات و روزنامه ها به کار می گیرد. چندین بررسی متمرکز توسط GVU از سال 1994 تولید شد که از آن میان می توان به بررسی سرشماری عمومی، بررسی مرتبط با سرشماری فناوری- که اطلاعاتی را نظیر سرعت ارتباط کاربران و مرورگر مورد استفاده گردآوری کرد، بررسی مربوط به گرایش کاربران به سوی حریم خصوصی و امنیت اینترنتی، بررسی عمومی رایانه ها، کاربرد اینترنت و وب، و نیز بررسی جست و جوی محصولات و فعالیت های مربوط به خرید به صورت پیوسته اشاره کرد. جامعه آماری، معمولاً بالای 1000 نفر [بودند که بررسی مربوط به سرشماری عمومی مشتمل بر 5000 شرکت کننده بود.

سازمان ها نیز بررسی های مربوط به وب را ابزار مفیدی برای ارزیابی کیفیت تجربیات کاربرانی می دانستند که از وبگاه ها یا کارکرد های وب بهره می برند. چنین مطالعه ای، بررسی بازدیدکنندگان حدود 450 وبگاه موزه ها را هدایت می کند به نحوی که هم اطلاعات مربوط به سرشماری، و هم پاسخ های مربوط به کیفیت و قابلیت استفاده از این وبگاه ها را گردآوری می نماید (Sarraf 1999). بسیاری از موزه ها به تدریج وب گاه هایی را تهیه کردند تا ابزاری برای ارائه اطلاعات درباره مجموعه های موجود در موزه ها و نیز ایجاد ارتباطی تعاملی و قابل دسترس برای بازدیدکنندگان کنونی و پیشین باشد.

ص: 78

1- با دسترسی 59 درصد از افراد بالغ آمریکایی به اینترنت در سال 2002، پژوهشگران در آمریکا نیز بررسی های مبتنی بر وب را به جای انواع بررسی های سنتی از قبیل بررسی های تلفنی و کاغذی درباره موضوعاتی به جز کاربردهای اینترنت و نگرش های موجود مد نظر قرار دادند به نظر می رسد پژوهشگران علاقه مند به گردآوری هنگام واکنش ها نسبت به وقایع اخیر به ویژه به بررسی های اینترنتی هستند، به نحوی که این علاقه مندی در بررسی های انجام گرفته درباره اثرات حملات تروریستی نیویورک در 11 سپتامبر سال 2001 بر شرکت کنندگان مشهود است (برای مثال لی و همکاران 2003) درست 9 روز پس از این حملات، گروهی از پژوهشگران توانستند بررسی از طریق ارزیابی وب-پایه پاسخ های روان شناختی آمریکایی ها نسبت به این واقعه تلخ را توزیع کنند.

2-، (GVU)

به منظور ارائه بازخوردی جامع برای موزه های مورد بحث به پاسخ های کیفی و کمی نیاز بود. کارزول و ونکاتش (2002)، با استفاده از بررسی وب محور توانستند پاسخ های ارزیابی بیش از 500 دانشجوی مقطع کارشناسی را که در کلاس های پیوسته غیر همزمان شرکت داشتند درخواست کنند. آن ها حمایتی نسبی را در مورد فرضیه های خود یافتند مبنی بر اینکه پذیرش، و در آینده قصد استفاده از فناوری در دوره ای پیوسته و غیر همزمان می تواند به صورت مثبت متأثر از گرایش ها و درک دانشجویان از فناوری باشد. در تمامی موارد از جمله GVU، مزایای استفاده از بررسی مبتنی بر وب آشکار بود: جمعیت های علاقه مند در بررسی ها همان کاربران وب بودند، و داده های گردآوری شده با فعالیت های وب ارتباط مستقیم داشت.

3- تحلیل بلاغی

*تحلیل بلاغی (1)

زیر مجموعه دیگری از پژوهش در حوزه وب نگاهی جدی به متون وب دارد تا راهکارهای متقاعد کننده ای را از طریق تحلیل بلاغی مشخص کند. همان طور که وارنیک توصیف می کند:

روش انتقادی بلاغی این مسئله را در نظر دارد که چگونه متن به برخی عناصر موجودیت می بخشد- در حالی که با برخی دیگر مخالفت می کند، چگونه ساختارهای گزارشی تجربه مرورگر را به روش های مشخص پیکربندی می کنند، و اینکه چگونه گفتمان، تمایلات و عادات ذهنی مخاطب خود را نقش می دهد.

وی، در بررسی خود به وبگاه های سیاسی در جریان مبارزات انتخاباتی رئیس جمهوری آمریکا در سال 1996 به ویژه راهکارهای بلاغی وبگاه های مقلد پرداخت - وبگاه هایی که از طراحی و محتوای وبگاه های قانونی مربوط به مبارزات انتخاباتی تقلید می کردند تا در صورت امکان از اعتبار آن سایت ها بکاهند و تا حدی آن ها را مورد تمسخر قرار دهند. بر اساس نتیجه گیری، وارنیک چنین سایت هایی بر گزارش های نادرست در دولت، تحریک مشارکت سیاسی از طریق ویژگی های تعاملی - نظیر شکایاتی که هرگز به حزب مورد نظرشان تسلیم نشد و اشاراتی درباره سوابق سیاسی و جنایی - که نامزدها بر آن ها متکی هستند تا خوانندگان شان را به پذیرش بینش بلاغی خود مجاب کنند - تکیه می کنند. در نهایت، وارنیک، این راهکارها را ریاکارانه می داند؛ زیرا از رفتارهای غیر اخلاقی مشابه بسیاری که ظاهراً از سوی نامزدهای مورد نظر به کار گرفته می شود بهره می گیرند.

از تحلیل بلاغی وبگاه هایی که برای شناسایی اجتماعات بلاغی چنین استفاده می گردد که گروه هایی از وب برای شکل گیری جهان بینی خود استفاده می کنند که ممکن است همسو با جریانات اخیر باشد یا نباشد کروبر، به بحث و تبادل نظر در این مسئله می پردازد که گروهی از وب گاه های شکل گرفته توسط مادرانی با گرایش های فمینیستی، خط مشی [خود] را به مخاطبان القا می کند تا به طرز مؤثری در برابر تصورات منفی فمینیستی در رابطه با مادر بودن مقاومت ورزند. زنان، توانستند از طریق متون و تصاویر موجود در وب گاه های شان دوباره به بحث و تبادل نظر در خصوص فمینیسم پردازند و برای درک معنای مادری به عنوان بخش پر نفوذ و حتی لازم فمینیست بودن به بحث و جدل می پردازند.

بررسی دیگر، نگاهی به موارد استفاده از وبگاه ها از سوی گروه های نفرت دارد، مانند شوالیه های

کوکلاکس کلان (1) و اتحاد ملی سازمان نئونازیسم که مقاصد ترغیبی را دنبال می کردند (Duffy 2003). در این مورد، دافی، از روش انتقادی بلاغی به نام تحلیل موضوعی فانتزی استفاده کرد که بر «نظریه همگرایی نمادین» استوار است. همانطور که دافی مشخص کرد، نظریه همگرایی نمادین «نظریه کلی علم معانی است که در آن، گروه ها تصوراتی را درباره گروه و گروه های خارج شکل می دهند و به اشتراک می گذارند، و در نتیجه هویت مشترکی را به وجود می آورند» (صفحه 293) (2) نویسنده، به وسیله تحلیل متون برگزیده موجود در در وبگاه هر گروه توانست گزارش های بلاغی متعددی را شناسایی کند که «ادعای انصاف و عدالت»، «نظم طبیعی و رستاخیز انسان» و «ساکنان اولیه زمین نامیده شدند تا مفهوم جدیدی را به نژاد انتقال دهد» از آن جمله هستند (صفحه 295-305).

4- تحلیل گفتمان

*تحلیل گفتمان (3)

مک کوئیل (2000)، تحلیل گفتمان را استعمال آن «در تمامی اشکال کاربردی زبان و اشکال نوشتاری» می داند، با استناد بر این طرز تفکر که «ارتباط در قالب نوشتار و گفتار رخ می دهد که با موقعیت اجتماعی خاص، موضوع ها، و انواع شرکت کنندگان سازگار می شود» (صفحه 494). تحلیل گفتمانی متون وب، موقعیت اجتماعی - فرهنگی وبگاه ها را در نظر می گیرد؛ یعنی به ساختار معنی محلی آن ها از طریق عناصر گفتاری و دیداری می پردازد. دگربار، متون وب متعلق به گروه های نفرت ظاهر می شوند؛ چنان که بیلینگ (2001) قراردادهای زبانی طنز را همانطور که از سوی وبگاه های «طنز» متعلق به گروه کوکلوکس کلان دیده می شود کشف کرد. او پی برد که در این وبگاه ها از تکذیب کنندگان برای هشدار علیه تفسیر طنزهای نژادپرستانه استفاده می کردند این طنزهای نژاد پرستانه به اقدامات خشونت آمیز فرمان می دادند در حالی که محتوای آن ها برای افراد «شوخی طبع» تعبیری طنزگونه داشت (صفحه 274). هر چند نویسنده به این نتیجه رسید که بخش عمده محتوای اینگونه سایت ها طنزگونه نبوده بلکه از جنس حقیقت است و فراگفتمانی را که او مطرح می سازد «انکار می کند که طنز طنز است» (صفحه 278). متون و تصاویر نژادپرستانه که هم خوانندگان این «طنزها» با گرایش های نژادپرستانه و هم مقاصد اخلاقی اینگونه طنزها را مخاطب قرار می داد شواهدی از گفتمان اهریمنی تری را فراهم می نمود، گفتمانی که بررسی بیلینگ را به تعاملات اجتماعی اوباش لینچ (4) تشبیه می کند (صفحه 287).

گفتمان در سطح کلان نیز در وب مورد تجزیه و تحلیل قرار گرفته است، و تأکید آن بر وبگاه هایی است که صدای دولت و ملت ها هستند. در بررسی پورسل و کورداس از وبگاه دولت اسلوونی، آن ها با متون موجود در این وبگاه به عنوان واکنشی به خطاهای مشاهده شده در معرفی اسلوونی به عنوان دولت بالکان - منطقه ای که گرفتار جنگ داخلی شد و بنابراین از سوی دولت های غربی به عنوان دولتی

ص: 80

Ku Klux Klan -1

2- فاس (1996)، نظریه همگرایی نمادین را بر اساس دو ادعای اصلی توصیف می کند: «ارتباط واقعیت را می سازد» و «نه تنها واقعیت را برای افراد می سازد بلکه معانی افراد را از نشانه ها می تواند با هم یکی کند تا واقعیت مشترکی را برای شرکت کنندگان به وجود آورند» (صفحه 122).

Discourse analysis -3

نامطلوب به خاطر اقدامات تروریستی و سرمایه گذاری های تجاری دیگر محکوم شد - برخورد کردند. پورسل و کورداس اظهار می دارند که اسلوانی تلاش می کند تا هویت از دست رفته خود را از طریق متون و تصاویر در وبگاه دولتی احیا کند که نوعی راهکار بلاغی محسوب می شود. اما بر اساس نتیجه گیری آن ها این سایت، سرانجام بخشی از تلاشی مستمر برای مذاکره موقعیت [اسلوانی] در «گفتمان جهانی» است.

شندلر (1998)، بیان کرده است که علاوه بر وبگاه های سازمانی، صفحه های شخصی نیز ممکن است شکلی از گفتمان باشند. شندلر در تحلیل خود از متون وب تولید شده توسط نوجوانان و لوز از طریق مصاحبه با سازندگان سایت ها، روش های بسیاری را بررسی کرد که بیان گر تصور نوجوانان از مخاطبان شان بود و درباره مرز بین فضاهای عمومی و خصوصی زندگی شان صحبت می کردند. بسیاری از نویسندگان جوان وب بیان کردند که سعی دارند تا در سایت های خود با مخاطبانی که احتمالاً علاقه های خود را به اشتراک می گذارند ارتباط برقرار کنند در حالی که افراد دیگر وبگاه را چیزی توصیف می کنند که صرفاً برای خودشان ایجاد کرده اند برخی، دیگر انگیزه ای پیدا می کنند تا بخشی از این اینترنت «پهناور» باشند و عقایدشان را درباره زندگی با دیگران به اشتراک بگذارند.

5- تحلیل دیداری

*تحلیل دیداری (1)

به دلیل قابلیت های چندرسانه ای شبکه جهانی وب، برخی پژوهشگران در صدد برآمده اند تا تحلیل های سنتی از متون وب گفتاری را با تحلیل متون دیداری تکمیل کنند یا بر آن ها غلبه کنند. ایجادکنندگان وبگاه ها نیز از گرافیک ها و تصاویر خلاقانه فراوان با چاشنی سلیقه در وبگاه ها بهره می گیرند، مانند آن چه در JenniCam.org به کار رفته است. از ویژگی های این وبگاه شخصی دوربین 24 ساعته ای با چشم انداز اتاق خواب جنی رینگلی نویسنده وب است. جیمروگلو (1999)، از سایت جنی کم (2) به عنوان بررسی موردی برای کشف مفهوم علم فرمانشی ترکیبی از ماشین و ارگانسیم هاروی (3) استفاده می کند (صفحه 149). همانند سایبورگ جیمروگلو مشاهده کرد که مرزهای بین بدنه جنی و فناوری از طریق دوربین پیوسته او محو می شود، مثلاً تصویر جنی که در مقابل صفحه نمایش رایانه اش نشسته است. سایت وب کم او تصویر جنی در پشت رایانه اش به نمادی برای آن تبدیل شده است، به طوری که جسمش ذوب و تبدیل به صفحه کلید می شود» (صفحه 441). برای بررسی معنای جنی کم، به عنوان موضوعی جنسیتی و دیداری از نظریه فمینیستی فیلم استفاده شد. تعهد رینگلی در ارائه نگاهی اجمالی به زندگی واقعی او در برگزیده لحظاتی بود که از جلوی وب کم به دلیل برهنگی یا داشتن فعالیت جنسی کنار می رفت. جیمروگلو، مشکل منحصر به فرد این تصاویر را تشخیص داده است: جنی موفق شد تا مرزهای سنتی اختصاص یافته به اندام زن را بین فضاهای عمومی و خصوصی بشکند، حال آن که با انتقاد برخی فمینیست ها روبه رو شد. آن ها بر این باور هستند که جنی در دام عینیت بخشی اندام زن، که در وب شایع

ص: 81

Visual analysis -1

JenniCam -2

Haraway -3

است، گرفتار شده است. در هر دو مورد «اندام جنی به عنوان کانون معنا، محل کمال، و ریشه معنی واحدی از جنی کم عمل می کند» (صفحه 449).

پژوهشگران دیگر، از تحلیل دیداری برای ارزیابی کاربرد اینترنت از سوی رسانه های خبری برای ارائه و توزیع تصویر استفاده کرده اند. فر پارکز (2001)، تحلیل انتقادی را از تصاویر مربوط به بحران نسل کشی رواندا در سال 1994 ارائه می دهند که رسانه های خبری آمریکا از طریق تلویزیون و اینترنت در دسترس عموم قرار دادند. پوشش زمینی این بحران، که در آن دوربین ها تصاویری را از پناهندگان و حشت زده گرفته بودند و امدادگران غربی سفید پوست به آن ها یاری می رساندند، و نیز تصاویر ماهواره ای از مردمی که دسته جمعی از رواندا می گریختند، با مخالفت روبه رو شدند، زیرا از قربانیان این بحران اختیار را سلب کرده بود. [آن ها] موقعیت سیاسی این واقعه را بیش از اندازه ساده انگاشته بودند و بیننده را از مردم درگیر و مسائل موجود دور نگه می داشتند. فر و پارکز، می افزایند که استفاده از وبگاه ها برای ارائه و پخش تصاویر رواندا بخشی از روند گسترده تر سلطه و مالکیت فناوری های دیداری از سوی رسانه های خبری در غرب است. در واقع، محتوای تولید شده به مصرف مخاطبان غربی یعنی کاربران اصلی اینترنت می رسد. فر و پارکز، به این نتیجه رسیدند که انتخاب تصاویر از سوی رسانه های خبری در پوشش خبری خود از بحران رواندا به همراه پخش گسترده تصاویر از طریق فناوری های رسانه ای غرب دست به دادند تا «آشفتگی» موجود آمریکا را از طریق فرهنگ و سیاست آفریقا تقویت کنند (صفحه 42).

دو مطالعه یادشده، توجه خاصی به تصاویر مورد استفاده در وبگاه ها، به منظور انتقال تفاسیر شان مبذول داشته اند. اما تحلیل های متون دیداری و گفتاری، جدا از منحصر به فرد بودن دو جانبه شان، معمولاً توسط پژوهشگران تلفیقی می گردند. همان گونه که در مطالعه هسکچ و سالیوان، درباره هویت بخشی مدارس از طریق تصاویر و متن ها در وبگاه ها، یا مطالعه و ارنیک درباره وبگاه های مقلد مربوط به مبارزات انتخاباتی و استفاده آن ها از متون گفتاری و غیر گفتاری برای پیشبرد دستور جلسات سیاسی ویژه مشاهده گردید. طراحی دیداری وبگاه ها، از قبیل استفاده از تصاویر و متن ها، و گذاشتن آن ها در صفحه وبگاه، نمونه دیگری از تلفیق این روش هاست (2000 Rivett). بررسی موردی وبگاه «نیوبیتل» ولکس واگن (1)، در مورد صفحه آرایی برگزیده، را که در صفحه اصلی این سایت قرار داشت مورد بحث قرار داد، و نویسنده سایت آن ها را با قراردادهای طراحی مجلات ناپیوسته مقایسه نمود. به علاوه، این سایت، با استفاده از یک پیش زمینه سفید، «هویت دیداری» پیوسته شرکت را منتقل می کند (صفحه 50). تحلیل دیداری وبگاه نیوبیتل مکمل تحلیل متنی است. در صورتی که داستان «حمله بیگانگان» در سرتاسر صفحه هایش به چشم می خورد. سپس، سایت ولکس واگن با طرح دیداری سایت دانشگاهی دانیل شندلر مقایسه گردید. سایت دانشگاهی دانیل شندلر، که مملو از متن بود بر هدف خود به عنوان منبع اطلاعات تأکید می کرد. در این سایت از تصویر کلاسوری به عنوان تصویر زمینه این سایت استفاده شد که بر ماهیت آموزشی آن تأکید می کند.

ص: 82

*قوم نگاری (1)

قوم نگاری، اغلب به عنوان «توصیف و تفسیر گروه یا نظام فرهنگی یا اجتماعی» درک می شود (Creswell 1998P58). تعریف موضوع در پژوهش های قوم نگاری به فرآیند پیچیده ای تبدیل شده است. زیرا در خلال مطالعه فرهنگ های اینترنت و وب انجام می گیرد. جنبه های عملی قوم نگاری، پیش از اینترنت، مطالعه فرهنگ های تک محله ای بوده است که از لحاظ جغرافیایی متمرکز شده اند. با وجود این انسجام پیوسته و پراکندگی فیزیکی ناپیوسته شرکت کنندگان در فرهنگ های اینترنتی مستلزم روشی چند مکانه برای قوم نگاری است. مفهوم جامعه چند مکانه نیز باید در موقعیت های شرکت کنندگان جامعه پیوسته به کار برده شود و ارائه توصیفی تیره از جامعه پیوسته مستلزم روشی متنوع تر از مشاهده ساده گفتمان گروه درون اتاق گفت و گو یا تحلیل کردن محتوای وبگاه های مربوط به آن جامعه است. در عوض، برخی پژوهشگران قوم نگار وب استدلال کرده اند که نیازمند روشی کلی نگر تر هستند، روشی که پژوهشگر تمامی سایت های مربوط به مشارکت و تجربه اعضای جامعه پژوهش را هم به صورت پیوسته و هم ناپیوسته بررسی کند (Howard 2002, Miller 2000, Slater).

میلر و اسلتر (2000) نمونه سودمندی را برای اجرای الگوی چند مکانه از قوم نگاری اینترنتی در مطالعه خود درباره اهالی ترینیداد (2) و اینترنت ارائه کردند. پژوهشگران تالارهای گفت و گو، وب گاه ها، و نیز خانه ها و کافی نت های مردم را بازدید کردند تا از روش های تلاقی اینترنت با زندگی های سیاسی، اقتصادی، و مذهبی «ترینی ها» آگاهی یابند. هر سایت، ارائه دهنده چشم انداز منحصر به فردی از آن فرهنگ است. اتاق های گفت و گو مکان هایی هستند که بینشی را از «ترینی» بودن به صورت پیوسته به دست می دهند. وب گاه ها، به عنوان نمایندگان ترینیداد به صورت پیوسته بودند (صفحه 13). میلر و اسلتر، تحلیل محتوایی کیفی را به کار گرفتند تا درباره وب گاه ها به بحث و گفت و گو پردازند. در حالی که، تلفیق نمادهای ملی نظیر پرچم و تصاویر نقشه ها را که با هویت های شخصی نویسندگان وب سراسر سایت های تحلیل شده اند - متذکر می شوند. افزون بر آن، آن ها به بحث درباره بازتاب دستور جلسات سیاسی در وبگاه ها، با سود آور ترین جاذبه توریستی ترینیداد، یعنی جشنواره کارنیوال، پرداختند که به محور اصلی در این سایت ها تبدیل شده بودند. پژوهشگران علاقه مند به جوامع طرفدار اینترنت، قوم نگاری را روشی سودمند برای مطالعات شان می دانند. بلوستین (2002)، تحلیل های انجام گرفته از وب گاه های طرفدار و انجمن های مربوط را در ترسیم خود از تعصب، و نمایش تلویزیونی «بافی»، قاتل خون آشام» گنجانند. او با استفاده از تحلیل های خود از متون وب و اطلاعات گردآوری شده از مصاحبه های شخصی و بازدید از اتاق های طرفداران توانست به نتیجه گیری هایی برسد که چگونه طرفداران بافی (که اغلب نوجوان بودند) بین جلوه فانتزی و جلوه های سحرآمیز نمایش تمایز قائل می شوند تا تخیل و «تفریح» جوانی را حفظ کنند، مادامی که به فضای جدی و ترسناک ابهام اخلاقی موجود در بزرگسالی وارد می شوند (صفحه 440).

ص: 83

Ethnography - 1

2- ترینیداد بزرگ ترین جزیره کشور ترینیداد و توباگو است. این جزیره جنوبی ترین جزیره دریای کارائیب است که تنها 11 کیلومتر با کرانه های ونزوئلا فاصله دارد. ترینیداد 4,768 کیلومتر مربع مساحت دارد.

*تحلیل شبکه (1)

به عقیده گارتون و همکارانش (1997)، «شبکه اجتماعی، به مجموعه افراد (یا سازمان ها یا نهادهای اجتماعی دیگر) اطلاق می شود که از طریق مجموعه ای از روابط اجتماعی نظیر دوستی، همکاری یا تبادل اطلاعات به هم مرتبط می گردند» (پاراگراف 2). هنگامی که این روابط اجتماعی با استفاده یا ایجاد وبگاه ها بنا برقرار یا ایجاد می گردند، فرایندها وسیله ای را برای مطالعه الگوهای ارتباطی میان اعضای شبکه ارائه می دهند. پارک (2003)، این نوع تحلیل در پژوهش وب را تحلیل فرایندها می داند که وبگاه ها را نمایانگر افراد، گروه ها، سازمان ها، و دولت ها و ملت ها، و فرایندهای بین سایت ها را «پیوند[های] رابطه ای» تصور می کنند (صص 50-51). برای مثال، هالویس (2000)، به منظور درک بهتر این موضوع که آیا وب به درستی شبکه «جهان گستر» سایت ها را رواج می دهد یا اینکه مرزهای ملی در حال پیوستن به مرزهای اینترنت هستند؟، به مطالعه فرایندهای موجود در وب گاه ها پرداخت. او با مطالعه 4000 وبگاه به بررسی پیوندهای آن وبگاه ها با صفحه های خارج از آن وبگاه ها پرداخت که از هویت میزبان ملی وبگاه هایی که به آن ها پیوند می یابند، پیروی می نمایند. هالویس، با استفاده از این فرآیند اینگونه نتیجه گرفت که در حقیقت بیشتر وبگاه های موجود مورد بررسی اصولاً به سایت های درون فرهنگ های ملی مرتبط می شوند. با وجود این، در مقایسه با انواع دیگر فناوری های اطلاعاتی مانند پست زمینی و تلویزیون او بر این باور بود که اینترنت بین المللی ترین فناوری اطلاعاتی است، و نسبت به دیگر رسانه ها در داخل و خارج از آمریکا در سراسر مرزهای بین المللی از تعداد مراجعات بیشتری برخوردار است (صفحه 23).

همان طور که پیشتر ذکر شد، تحلیل شبکه در سطح سازمانی، همزمان با تحلیل محتوایی در مطالعه وبگاه های ضدجهانی سازی، برای درک بهتر این مسئله مورد استفاده قرار گرفت که این سایت ها تا چه حد با هم تلفیق شده اند (Van Alst and Walgrave 2002). پیوندهای موجود بین 17 سایت، با استفاده از نرم افزار دیداری سازی شبکه (<http://vlado.fmf.uni-pajek>) مورد تجزیه و تحلیل قرار گرفت. نرم افزار یاد شده، نقشه گرافیکی پیوندهای میان وبگاه ها را ایجاد می کرد. سایت های بسیار مرتبط و دارای ارجاعات بسیار با گردآوری پیوندها در موقعیت خودشان، و خروج از آن ها به آسانی شناسایی شدند. حال آن که، سایت های پراکنده تنها به وسیله یک یا دو پیوند با دیگر سایت های شبکه پیوند می یابند، دور افتاده به نظر می رسند. ون الست و والگریو دریافتند که سایت های ضد جهانی سازی تا اندازه ای یکپارچه شده اند، اما یادآور می شوند که ارزیابی ماهیت روابط بین سازمان هایی که صرفاً مبتنی بر موجودیت پیوندهای بین سایت ها می باشند، دشوار است.

در موارد دیگر، فرایندها به منظور کشف شبکه های اجتماعی و اطلاعاتی افراد پیوسته مورد تجزیه و تحلیل قرار گرفته اند بیشتر مرورگرهای وب، مانند اینترنت اکسپلورر (3) و Netscape Navigator، مجهز

به کار کرد «تاریخچه (1)» هستند که نشانی های اینترنتی را که شخص در هنگام بازدید از سایت های شبکه به آن ها دسترسی داشته است، ذخیره می کنند. تاچر و گرینبرگ (1997)، از این کار کرد برای تحلیل رفتارهای مرورگرانه 23 نفر طی یک دوره 6 هفته ای استفاده کردند. آن ها به دنبال الگوهایی در بازدید صفحه های وب بودند. این دو پژوهشگر، متوجه شدند که تقریباً دو سوم بازدیدهای صفحه های وب از سوی آزمودنی ها، صفحه هایی هستند که قبلاً بازدید شده بودند (صفحه 112). همچنین، تعداد صفحه هایی که اغلب توسط مخاطبان دوباره بازدید می شوند، نسبتاً کم بودند. تاچر و گرینبرگ، فرض می کنند که کارکرد «بازگشت (2)» در مرورگرهای وب، که 30 درصد از فعالیت های مرورگریانه آزمودنی ها را تشکیل می دهد، ممکن است عامل مؤثری برای علاقه کاربران اینترنت برای بازدید دوباره صفحه ای واحد در وب باشد (صفحه 131).

علاوه بر شبکه های اجتماعی مطالعه شبکه های مدارک، به دلیل بزرگی شبکه وب رو به فزونی رفته است، و این بدان معناست که اشکال گوناگون تحلیل علمی از شبکه اجتماعی امکان پذیر است. تحلیل شبکه های مدارک می تواند فرصت هایی را برای فراتحلیل محتوا به وجود آورد. هنزینگر و لورنس (2004)، درباره روش های نمونه برداری از صفحه های وب به بحث و گفت و گو می پردازند تا «به طور خودکار نمونه وسیعی از علایق و فعالیت های انجام گرفته در این دنیا را ... به وسیله تجزیه و تحلیل ساختار پیوند وب و چگونگی روی هم انباشته شدن پیوندها در طول زمان تجزیه و تحلیل کنند» (صفحه 5186). آیزن و همکارانش (2004)، بر این ادعا هستند که استفاده از داده ها در وبگاهی با ترافیک بالا می تواند اطلاعاتی را درباره وقایع خارجی و احساسات ناگهانی در دسترس عموم قرار دهد که ممکن نیست به تنهایی از تحلیل محتوا و ساختار پیوند دسترس پذیر شود (صفحه 5254).

8-ملاحظات اخلاقی

اهمیت نیاز به پژوهش با رعایت اصول اخلاقی برای پژوهشگران اینترنت پوشیده نیست. اما خط مشی ها و ابزارهای پژوهش اخلاقی به روشنی مشخص نیستند. از دیدگاه پژوهشگران وب، و به ویژه آن هایی که مواد آرشیوی وب را استخراج می کنند، ملاحظات اخلاقی به صورت ضرورت در جریان پژوهش نمود می یابند. برخی اطلاعات موجود در وب و اطلاعات در دست آرشیو، که محرمانه به شمار می آیند، ممکن است سهواً دسترس پذیر شده، و سپس توسط موتورهای کاوشی نظیر گوگل نمایه و ذخیره شوند. دیدگاه های مربوط به امکان استفاده از اطلاعات خصوصی، که به صورت عمومی آرشیو شده اند، متعدد و خارج از حوصله این مقاله می باشند. با وجود این، باید خاطرنشان کرد که تمامیت موتورهای کاوش در جست و جوی وب می تواند منجر به ایجاد آرشیوهای اطلاعاتی شود که کاربران عموماً نه قصد فاش کردن آن ها را دارند و نه آن ها را آرشیو کرده اند. در حالی خوش بینانه، برآوردی سرعتی که موتورهای کاوش وب را نمایه سازی می نمایند، نشان داد که سرعت نمایه سازی و آرشیوسازی گوگل از ایجاد صفحه های جدید

ص: 85

وب فراتر رفته است (Whelan 2004). احتمالاً مسائلی که پژوهشگران حوزه مطالعات اینترنت مدت ها به بحث و تبادل نظر گذاشته اند، و نیز مطالعات انجام شده بر روی ارتباطات با واسطه رایانه، با توجه به تقابل داده های عمومی با داده های خصوصی، در آینده بسیار نزدیک در میان پژوهشگران حوزه مطالعات وب و حرفه مندان حوزه آرشیو وب و دیگر مواد الکترونیکی از اهمیت بالایی برخوردار خواهند شد.

9- نتیجه گیری

پژوهش های حوزه وب به وضوح سودمندی خود را نشان داده اند. این پژوهش ها بدون آن که در سایه تحلیل های متنی جوامع مبتنی بر متن قرار داشته باشند، ثابت کرده اند که مفاهیمی نظیر جامعه، فرهنگ، رفتار و ساختارهای معنایی می توانند در اینترنت به طور مؤثری بررسی شوند. در حالی که غنا و گوناگونی پژوهش های وب نوید بخش است، هنوز مسائل زیادی وجود دارد که باید کشف شوند.

میراث تحلیل متنی به جامانده از ارتباطات با واسطه رایانه در سوگیری شدید نسبت به تحلیل های زبانی در مطالعات حوزه وب مشهود است. «مطالعات مفقود» به مطالعاتی اطلاق می شود که وب را مانند چشم اندازی چندرسانه ای بررسی می کند، و توصیف می کند که چگونه تصاویر و صدا برای از میان برداشتن ارتباطات گفتاری مورد استفاده قرار می گیرند تا بر موانع زبانی در رسانه جهانی فائق آیند. به استثنای اشاره به کاربرد صدا در وبگاه های ترنیدادی، در مطالعه قوم نگاری میلر و اسلتر (2000)، وب گاه های شنیداری کاملاً از بخش کنونی پژوهش وب غایب بودند. احتمال دارد که پروژه های آینده به کاربردهای بلاغی صدا در وب از طریق تحلیل های وب گاه های تجاری یا سیاسی پردازند. دیگر مطالعاتی که به حوزه چند رسانه ای ها گرایش دارند توانستند روش هایی را که از صدا و تصاویر برای ایجاد هویت افراد موجود در اینترنت از افراد حقیقی و سازمان ها گرفته تا دولت های ملی استفاده می شدند، بررسی کنند.

به منظور افزایش تعداد پژوهشگران این حوزه دسترسی به همه انواع محتوای وب مربوط به تمامی دوره های زمانی، به ویژه بازآفرینی پیوندهای درون اشیا محتوایی و میان وبگاه ها الزامی است. به علاوه، باید مرورگرها و ابزارهای دیگری در دسترس قرار گیرند که محتوای وب به وسیله آن ها رصد شود تا بتوان تجربه کاربر را به بهترین شکل درک کرد. بنابراین، چالش ها فراتر از حفظ و نگهداری محتوا رفته، و شامل حفظ ساختار و رویارویی با محتوا نیز می شود.

منابع

Aizen, J., Huttenlocher, D., Kleinberg, J., Novak, A. (2004). Traffic-based feedback on the Web. . 1
Proceedings of the National Academy of Sciences of the United States of America, 101(1), 5254-5260

Billig, M. (2001). Humour and hatred: The racist jokes of the Ku Klux Klan. Discourse .2

- Bloustein, G. (2002). Fans with a lot at stake: Serious play and mimetic excess in Buffy the Vampire Slayer. *European Journal of Cultural Studies*, 5(4), 427–449 .3
- Carswell, A. D. Venkatesh, V. (2002). Learner outcomes in an asynchronous distance education environment. *International Journal of Human–Computer Studies*, 56, 475–494 .4
- Chandler, D. Roberts–Young, D. (1998). The construction of Identity in the Personal Homepages of Adolescents. Retrieved April 15, 2002 from <http://www.aber.ac.uk/media/Documents/short/strasbourg.html> .5
- Creswell, J. W. (1998). *Qualitative inquiry and research design*. Thousand Oaks, CA: Sage .6
- Duffy, M. E. (2003, July). Web of hate: A fantasy theme analysis of the rhetorical vision of hate groups online. *Journal of Communication Inquiry*, 27(3), 291–312 .7
- Ess, C. (2002). Ethical decision making and Internet research: Recommendations from the AoIR working committee. Association of Internet Researchers. Retrieved July 12, 2004, from <http://www.aoir.org/reports/ethics.pdf> .8
- Fair, J. E. Parks, L. (2001). Africa on Camera: Television news coverage and aerial imaging of Rwandan refugees. *Africa Today*, 48(2), 34–57 .9
- Foss, S. K. (1996). *Rhetorical criticism: Exploration and practice*. Prospect Heights, IL: Waveland .10
- Garton, L., Haythornthwaite, C., Wellman, B. (1997). Studying online social networks. *Journal of Computer–Mediated Communication*, 3 (1). Retrieved April 4, 2004 from <http://www.ascusc.org/jcmc/vol3/issue1/garton.htm> .11
- Geertz, C. (1973). *The interpretation of cultures*. New York: Basic Books Graphic, Visualization and Usability (GVU) Center (1998). 10th WWW User Survey. Atlanta, GA: Georgia Institute of Technology. Retrieved March 25, 2004 from <http://www.cc.gatech.edu/gvu/user-surveys/survey-1998-10> .12
- Halavais, A. (2000). National Borders on the World Wide Web. *New Media Society*, 1(3), 7–28 .13
- Haraway, D. (1991). *Simians, cyborgs, and women: The reinvention of nature*. New York: Routledge .14

Henzinger, M. Lawrence, S. (2004). Extracting knowledge from the World Wide Web. Proceedings of . 15
,the National Academy of Sciences of the United States of America

ص: 87

- Hesketh, A. J. Selwyn, N. (1999). Surfing to school: The electronic reconstruction of institutional . 16
identity. *Oxford Review of Education*, 25(4), 501-520
- Howard, P. N. (2002). Network Ethnography and the Hypermedia Organization: New Media, New . 17
Organizations, New Methods. *New Media and Society*, 4(4), 550-574
- Jimroglou, K. M. (1999). A Camera with a view: JenniCam, visual representation, and . 18
cyborgsubjectivity. *Information, Communication and Society*. 2(4), 439-453
- Kroeber, A. (2001). Postmodernism, Resistance, and Cyberspace: Making Rhetorical Spaces for . 19
Feminist Mothers on the Web. *Women's Studies in Communication*, 24(2), 218-240
- Lee, W., Hong, J., Lee, S. (2003). Communicating with American consumers in the post 9/11 climate: .20
An empirical investigation of consumer ethnocentrism in the United States. *International Journal of
Advertising*, 22, 487-510
- McQuail, D. (2000). *McQuail's mass communication theory* (4th ed.). London, UK: Sage .21
- Miller, D. Slater, D. (2000). *The Internet: An ethnographic approach*. Oxford, UK: Berg .22
- Milne, G. R. Culnan, M. J. (2002). Using the content of online privacy notices to inform public policy: .23
A longitudinal analysis of 1998-2001 US Web surveys. *The Information Society*, 18, 345-359
- O'Neill, E. T., Lavoie, B. F., Bennett, R. (2003). Trends in the evolution of the public Web, 1998-2002. .24
.D-Lib Magazine, 9(4). Retrieved March 29, 2004 from <http://wcp.oclc.org>
- Park, H. W. (2003). Hyperlink network analysis: A new method for the study of social structure on the .25
Web. *Convergence*, 25(1), 49-61
- Pitts, M. J. Harms, R. (2003). Radio websites as promotional tools. *Journal of Radio Studies*, 10(2), .26
270-282
- Purcell, D. Kodras, J. E. (2001). Information technologies and representational spaces at the outposts of .27
the global political economy. *Information, ommunication and Society*, 4(3), 341-369
- Rivett, M. (2000). Approaches to analyzing the Web text: A consideration of the Web site as an .28

emergent cultural form. *Convergence*, 6(3), 34–56

?Sarraf, S. (1999). A survey of museums on the Web: Who uses museum Websites .29

ص: 88

- Silver, R. C., Holman, E. A., McIntosh, D. N., Poulin, M., Gil-Rivas, V. (2002). Nationwide . 30 longitudinal study of psychological responses to September 11. *Journal of the American Medical Association*, 288(10), 1235-1244
- Spooner, T. (2002). Internet use by region in the United States. Pew Internet American Life Project. . 31 Washington, DC. Retrieved March 25, 2004 from <http://www.pewinternet.org/reports/pdfs/PIP-Regional-Report-Aug-2003.pdf>
- Tauscher, L. Greenberg, S. (1997). How people revisit Web pages: Empirical findings and implications . 32 for the design of history systems. *International Journal of Human- Computer Studies*, 47, 97-137
- Van Aelst, P. Walgrave, S. (2002, December). New media, New movements? The role of the Internet in . 33 shaping the 'anti-globalization' movement. *Information Communication and Society*, 5(4), 465-493
- Warnick, B. (1998). Appearance or reality? Political parody on the Web in Campaign '96. *Critical . 34 Studies in Mass Communication*, 15, 306-324
- Whelan, D. (2004, 16). Google Me Not. *Forbes*, 174(3), 102-104 . 35

اطلاع از مشخصات وب ملی یکی از نیازمندی های اصلی کشور در رابطه با سیاست گذاری در حوزه فناوری اطلاعات است. از آن جا که متأسفانه تا به حال، تحقیقی در این خصوص، انجام نشده است، نیاز مبرمی برای پژوهش و کسب اطلاع از وضعیت فعلی وب فارسی، احساس می شود. با درک این ضرورت، در این مقاله نتایج حاصل از مطالعات صورت گرفته در مورد خصوصیات و ویژگی های وب فارسی مطرح خواهد شد. برای این کار سامانه نیمه خودکاری طراحی و پیاده سازی شده است. هدف از این سیستم، استخراج شاخص های مختلف از قبیل حجم محتوای فارسی، تنوع محتوا و نیز عمر محتوا می باشد. از این آمار می توان جهت هدفمند نمودن برنامه های آتی در خصوص ساماندهی وب فارسی و نیز ارائه راهکارهایی برای حل معضلات فعلی، بهره گرفت. آنالیزهای انجام شده توسط سامانه فوق الذکر بر روی حدود یازده هزار وبگاه ثبت شده در دامنه IR، مشتمل بر اغلب سازمان ها و ارگان های دولتی، وزارتخانه ها، شرکت ها و دانشگاه ها است که بالغ بر حدود دو میلیون صفحه می باشند.

کلید واژه: دامنه IR، وب فارسی، مشخصات وب

*خصوصیات وب ایران: مریم پیروزمند (1)

1. مقدمه

توسعه و رشد نمایی وب باعث شده است تا با حجم عظیمی از اطلاعات شامل اسناد با فرمت های متفاوت از جمله متن، صوت، تصویر و غیره در مکان ها و سازمان های مختلف مواجه شویم. در عین حال، گسترش روز افزون نیازمندی به وب که زائیده فزونی عرضه خدمات از طریق وب است، باعث شده تا همواره کاربران بیشتری به استفاده از وب، راغب شوند. بنا بر آمارهای موجود در حال حاضر بالغ بر 130 میلیارد صفحه در محیط وب از طریق جویشگر های مطرح دنیا قابل دسترسی و جستجو هستند [1]. به دلایل مختلف از جمله شروع توسعه وب توسط کاربران انگلیسی زبان، غالب سرویس ها و خدمات عرضه شده از طریق وب بویژه سرویس های جستجو، بصورت انگلیسی ارائه می شوند. گر چه در طی سال های اخیر، موتورهای جستجوی عمده ای مانند Yahoo و Google، سرویس جستجو را به زبان های دیگر نیز عرضه کرده اند، اما متأسفانه به دلایل مختلف بالاخص تحریم های سیاسی و اقتصادی، این خدمات هیچگاه به حوزه زبان فارسی گسترش پیدا نکرده است. از سوی دیگر، در داخل کشور نیز به علت عدم وجود شناخت کافی از مختصات وب فارسی موتورهای جستجوی توانمندی بوجود نیامده اند.

بطور کلی، شناخت مختصات وب ملی یک کشور، نمایان گر شاخص های توسعه یافتگی در آن کشور

ص: 91

نیز محسوب می شود. به عنوان مثال، در کشورهای توسعه یافته، حجم وب ملی شامل تعداد وبگاه های مختلف و همچنین تنوع خدمات قابل عرضه در این رسانه، بسیار زیاد است. همچنین کیفیت و به روز بودن اطلاعات عرضه شده نیز در دامنه وب این قبیل کشورها در مقایسه با دیگر کشورها بسیار بهتر است.

از سوی دیگر، با توجه به اینکه زیر ساخت اصلی اشاعه دولت الکترونیک، محیط وب می باشد، لذا مطالعه و شناسایی ویژگی های این محیط، کمک شایانی به انجام برنامه ریزی مناسب جهت تدوین سیاست های اجرایی جهت تحقق اهداف دولت الکترونیک طی مراحل مختلف، خواهد کرد. با درک این مطلب، در این مقاله سعی می شود تا به ویژگی های اصلی وب ایران پرداخته شود. شاخص های مورد مطالعه، به گونه ای انتخاب شده است تا بتوان از آن ها در جهت بهبود عملکرد سیستم های جستجوی وب، بهره گرفت.

با توجه به مسائل فوق در این مقاله سعی می شود تا حد ممکن، محتوای وب فارسی کشور را به صورت خودکار، تحلیل و ارزیابی کرد (ابزار نظارت خودکار) تا در مقاطع زمانی مختلف، بتوان عمل ارزیابی را با کمترین هزینه، تکرار نمود. لازم به ذکر است که در این مقاله حدود 11 هزار وبگاه با پسوند IR، که شامل تقریباً دو میلیون صفحه است پردازش و تحلیل شده اند. البته تعداد وبگاه های فارسی بیش از این مقدار است اما به دلیل دشواری یافتن وبگاه های فارسی زبان، در این پژوهش به وبگاه های ثبت شده در دامنه IR بسنده شده است. علت این دشواری این است که همانطور که در قسمت نتایج، به تفصیل ذکر خواهد شد. متأسفانه دامنه IR به عنوان تنها دامنه فارسی زبان نیست و بسیاری از وبگاه های فارسی زبان، در دامنه های مختلف و حتی غیر مرتبطی از قبیل .com, net, org و دامنه های دیگر ثبت شده اند.

همانگونه که پیش از این ذکر شد، هدف اصلی این مطالعه، تعیین وضعیت فعلی وب فارسی است.

بر این اساس، هدف نهایی، استخراج شاخص های زیر از "وب ایران" می باشد:

1. پاسخگویی به سؤالات کلی درباره وب ایران از جمله:

1-1. درصد پیوندهای معتبر چقدر است؟ تعیین این شاخص بطور ضمنی، نرخ تغییرات وب را نیز تعیین خواهد کرد و زمان بندی برای خزشگر وب را دقیق خواهد کرد.

2-1. توزیع صفحات وب فارسی از لحاظ محتوا، چگونه است؟ یعنی چه حجمی از محتوای وب در کلاس های علمی، تجاری، روزنامه، خبر، وبلاگ و غیره، قابل طبقه بندی است؟

3-1. ساختار محتوایی صفحات وبگاه ها چگونه است؟ بطور مشخص می خواهیم بدانیم که:

• درصد کدینگ های زبانی مختلف استفاده شده نظیر Windows- 1252, Windows - 1256 و UTF-8 به چه صورت است؟ بدین ترتیب مشخص خواهد شد که تبدیل کدینگ ها در سیستم های بازایی وب چقدر حائز اهمیت است؟

• چه کسری از صفحات، عنوان (1) مناسب دارند؟ وجود صفحات با عنوان مناسب و گویا، موجب بهبود کیفیت بازایی توسط سیستم و نیز سهولت دسترسی توسط کاربران خواهد شد.

1-4. فایل های غیر متنی مانند Doc PPT PDF و Image چند درصد از صفحات را تشکیل می دهند؟

1-5. نرخ به روزآوری، تغییر، ایجاد (عمر صفحه) چقدر می باشد؟ استخراج این اطلاعات، نرخ تغییرات وب را تعیین خواهد کرد و موجب برنامه ریزی و زمان بندی برای خزش گر وب را فراهم خواهد آورد.

1-6. سرعت دسترسی به وبگاه چقدر است؟ این آگاهی نیز در تنظیم عملکرد خزشگر وب موثر خواهد بود.

1-7. تعداد دسترسی به وبگاه ها و صفحات در مقاطع زمانی مختلف یا عبارت دیگر الگوی دسترسی به وبگاه ها و صفحات چگونه می باشد؟

1-8. تعداد کل صفحات فارسی، میانگین تعداد صفحات هر وبگاه و حجم آن ها چقدر است؟

1-9. تعداد لغاتی که در تمام صفحات محاسبه شده چقدر می باشد؟

1-10. چند درصد صفحات شامل هر دو محتوای فارسی و انگلیسی است؟

لازم به ذکر است در این فاز به دلایلی مانند پیچیدگی کار استخراج شاخص های مربوط به فازهای تراکنش و تبدیل به کارهای آینده موکول شده است. نتایج این مقاله را می توان در تعریف و پیاده سازی پروژه هایی مانند موتور جستجوی ملی و درگاه دولتی استفاده کرد. به علاوه این پروژه در تدوین راهکارهای آینده جهت تحقق سریع دولت الکترونیک در کشور مفید فایده خواهد بود.

2. کارهای مرتبط در داخل و خارج

تاکنون در داخل کشور مکانیزم ارزیابی وب بدین صورت انجام نشده است، اما بعضی از کشورها کاری شبیه به این پروژه را انجام داده اند. برای مثال در تایلند پروژه ای تحت عنوان "ابزار نظارت خودکار بر پروژه دولت الکترونیک تایلند" [2] انجام شده است. هدف این کار استخراج تمام شاخص های دولت الکترونیک از وب تایلند می باشد. بر اساس نتایج این مطالعه که حاصل بررسی حدود 150 وبگاه دولت تایلند در سال 2002 است، حدود 31% این سایت ها، صرفاً به ارائه اطلاعات می پردازند، حدود 57% امکان تعامل محدود کاربران را با وبگاه فراهم می کنند و تنها حدود 11%، بستر لازم برای اجرای تراکنش های مورد نیاز کاربران را در اختیار وی قرار می دهند.

به صورت مشابه، کارهایی برای استخراج مشخصه های وب در کشورهای اسپانیا [3]، کره جنوبی [4]، استرالیا [5]، پرتغال [6]، اروپا [7]، [8] و آرژانتین [9] انجام شده است. در کنار این تحقیقاتی نیز در مورد مطالعه ویژگی های کلی وب صورت گرفته است [10]، [11]، [12] در فعالیت های فوق بیشتر وب کشورها از دیدگاه ساختاری و شکل گراف وب مورد بررسی قرار گرفته است و به علاوه پارامترهایی مانند توزیع اندازه صفحات و وبگاه ها، نرخ بروزآوری آن ها، رتبه آن ها در موتورهای جستجو و عمر صفحات مورد بررسی قرار گرفته ولی از منظر محتوا و سرویس های ارائه شده کاری انجام نگرفته است. در این مقاله علاوه بر استخراج پارامترهای فوق برای وب فارسی، از دید دولت الکترونیک نیز به وب

ایران توجه گردیده است.

در تحقیق صورت گرفته در مرجع [10]، ساختار کلی گراف وب جهانی مورد مطالعه قرار گرفته است. بررسی های صورت گرفته حاکی از وجود ساختاری موسوم به مدل پایونی (1) در مورد گراف وب است. بعنوان مثال، طی بررسی انجام شده در سال 2000 میلادی روی حدود دویست میلیون صفحه وب، همانگونه که شکل شماره یک نشان می دهد، گراف وب، شامل چهار بخش کلی است که عبارتند از:

الف- بخش هسته مرکزی (2) که حدود 25% کل وب را تشکیل می دهد.

ب- بخش ورودی (3) که لینک هایی به بخش هسته مرکزی دارد و حدود 25% کل وب را تشکیل می دهد.

ج- بخش خروجی (4) که لینک هایی از بخش هسته مرکزی دارد و حدود 25% کل وب را تشکیل می دهد.

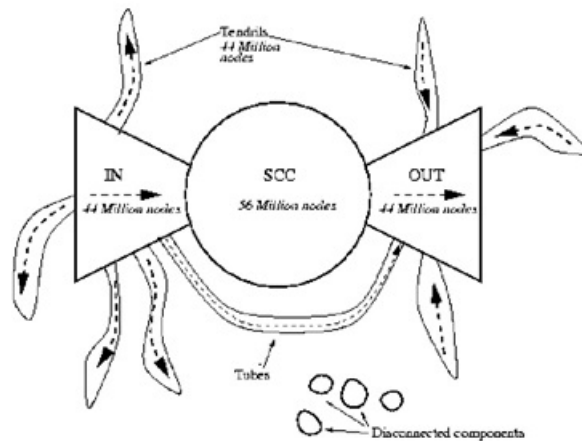
د- بخش های پراکنده که یا به صورت جزایر پراکنده (و بعضاً به سختی قابل دسترس) هستند یا واسطه هایی بین بخش های ورودی و خروجی هستند. این قسمت نیز حدود 25% کل وب را تشکیل می دهد.

عکس

ایران توجه گردیده است.

در تحقیق صورت گرفته در مرجع [۱۰]، ساختار کلی گراف وب جهانی مورد مطالعه قرار گرفته است. بررسی‌های صورت گرفته حاکی از وجود ساختاری موسوم به مدل پایونی^۱ در مورد گراف وب است. بعنوان مثال، طی بررسی انجام شده در سال ۲۰۰۰ میلادی روی حدود دویست میلیون صفحه وب، همانگونه که شکل شماره یک نشان می‌دهد، گراف وب، شامل چهار بخش کلی است که عبارتند از: الف-بخش هسته مرکزی^۲ که حدود ۲۵٪ کل وب را تشکیل می‌دهد.

ب-بخش ورودی^۳ که لینک‌هایی به بخش هسته مرکزی دارد و حدود ۲۵٪ کل وب را تشکیل می‌دهد. ج-بخش خروجی^۴ که لینک‌هایی از بخش هسته مرکزی دارد و حدود ۲۵٪ کل وب را تشکیل می‌دهد. د-بخش‌های پراکنده که یا به صورت جزایر پراکنده (و بعضاً به سختی قابل دسترس) هستند یا واسطه‌هایی بین بخش‌های ورودی و خروجی هستند. این قسمت نیز حدود ۲۵٪ کل وب را تشکیل می‌دهد.



شکل ۱. مدل پایونی گراف وب

۳. سامانه خودکار ارزیابی وب ایران

در این قسمت، ابتدا معماری سیستم ارزیاب خودکار، تشریح می‌شود. این سامانه، نسخه گسترش یافته ابزار استفاده شده در مرجع [۱۳] است که با هدف بررسی مشخصات و بگانه‌های ثبت شده در دامنه IR اعم از دولتی و غیر دولتی می‌باشد. شکل شماره دو شمای کلی این سامانه و تعامل اجزای آن را با

1. Bow-Tie model
2. Central core
3. In
4. Out

شکل ۱. مدل پایونی گراف وب

۳. سامانه خودکار ارزیابی وب ایران

در این قسمت، ابتدا معماری سیستم ارزیاب خودکار، تشریح می‌شود. این سامانه، نسخه گسترش یافته ابزار استفاده شده در مرجع [13] است که با هدف بررسی مشخصات و بگانه‌های ثبت شده در دامنه IR اعم از دولتی و غیر دولتی می‌باشد. شکل شماره دو شمای کلی این سامانه و تعامل اجزای آن را با

Bow-Tie model -1

Central core -2

In -3

Out -4

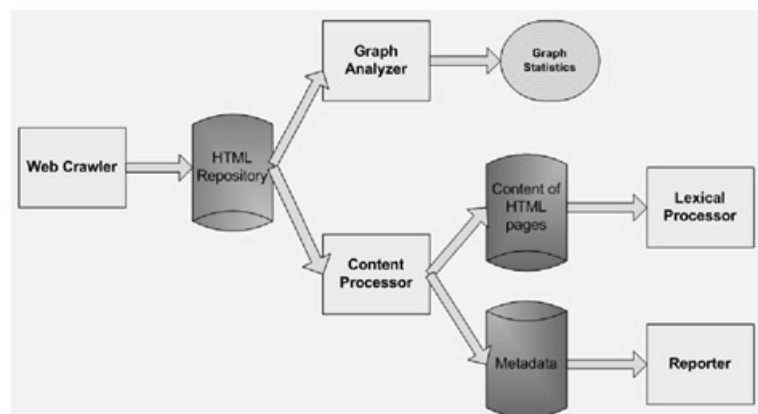
یکدیگر، نشان می دهد.

همان گونه که در شکل شماره دو دیده می شود، ابتدا یک خزشگر وب، بر اساس یک لیست اولیه از وبگاه های ثبت شده در دامنه IR که از قبل تهیه شده است و با توجه به پارامترهای تنظیم عملکرد آن از قبیل عمق خزش، حداکثر صفحات وبگاه و غیره، خزش را با توجه به گراف حاصل از پیوندها، انجام می دهد و صفحات بدست آمده را در یک مخزن موقت، شاخص بندی و ذخیره می کند. در مرحله بعد، این مخزن توسط واحد تحلیل گر گراف وب، مورد بررسی قرار می گیرد و آمارهای مختلفی را مورد ویژگی های گراف متناظر از قبیل قطر گراف، متوسط فاصله بین هر دو گره و غیره، استخراج می کند.

عکس

یکدیگر، نشان می‌دهد.

همانگونه که در شکل شماره دو دیده می‌شود، ابتدا یک خزشگر وب، بر اساس یک لیست اولیه از وبگاه‌های ثبت شده در دامنه IR که از قبل تهیه شده است و با توجه به پارامترهای تنظیم عملکرد آن از قبیل عمق خزش، حداکثر صفحات وبگاه و غیره، خزش را با توجه به گراف حاصل از پیوندها، انجام می‌دهد و صفحات بدست آمده را در یک مخزن موقت، شاخص‌بندی و ذخیره می‌کند. در مرحله بعد، این مخزن توسط واحد تحلیل گر گراف وب، مورد بررسی قرار می‌گیرد و آمارهای مختلفی را مورد ویژگی‌های گراف متناظر از قبیل قطر گراف، متوسط فاصله بین هر دو گره و غیره، استخراج می‌کند.



شکل ۲. معماری کلی سامانه خودکار ارزیابی وب ایران

از سوی دیگر، واحد تحلیل گر محتوا با استفاده از تجزیه کننده HTML، صفحات این مخزن را بررسی نموده و محتوا و داده‌های توصیفی آنها را استخراج می‌کند. این قسمت که توسط زبان جاوا پیاده‌سازی شده است، بصورت مستقل و از طریق پارامترهای قابل تنظیم، اجرا می‌شود. ورودی این بخش، کل صفحات HTML است که شامل برچسب‌های مختلف و نیز نویسه‌های موجود در صفحه می‌باشد. شایان ذکر است که در هنگام جمع‌آوری صفحات از وب، توسط خزشگر وب بصورت خودکار به هر صفحه، یک شناسه عددی نسبت داده می‌شود که آن صفحه را بصورت یکتا مشخص می‌کند. جدول شماره یک، نمونه‌ای از ورودی‌های تجزیه کننده را نشان می‌دهد.

محتوای استخراج شده از صفحات، برای پردازش در اختیار واحد تحلیل گر واژگان، قرار داده می‌شود تا کلمات فارسی را از آن استخراج کند و در یک Lexicon ذخیره نماید. همچنین از اطلاعات بدست آمده می‌توان برای استخراج آمار متنوع از کلمات استفاده کرد.

1. Parser
2. Script

شکل ۲. معماری کلی سامانه خودکار ارزیابی وی ایران

از سوی دیگر، واحد تحلیل گر محتوا با استفاده از تجزیه کننده (1) Parser صفحات این مخزن را بررسی نموده و محتوا و داده‌های توصیفی آن‌ها را استخراج می‌کند. این قسمت که توسط زبان جاوا پیاده‌سازی شده است، بصورت مستقل و از طریق پارامترهای قابل تنظیم اجرا می‌شود و ورودی این بخش کل صفحات HTML است که شامل برچسب‌های مختلف و نیز نویسه (2) های موجود در صفحه می‌باشد. شایان ذکر است که در هنگام جمع‌آوری صفحات از وب توسط خزشگر وب بصورت خودکار به هر صفحه، یک شناسه عددی نسبت داده می‌شود که آن صفحه را بصورت یکتا مشخص می‌کند. جدول شماره یک، نمونه‌ای از ورودی‌های تجزیه کننده را

نشان می دهد.

محتوای استخراج شده از صفحات برای پردازش در اختیار واحد تحلیل گر واژگان، قرار داده می شود تا کلمات فارسی را از آن استخراج کند و در یک Lexicon ذخیره نماید. همچنین از اطلاعات بدست آمده می توان برای استخراج آمار متنوع از کلمات استفاده کرد.

ص: 95

HTML -1

Script -2

بخش مهم دیگر این معماری، واحد تهیه گزارش است. داده‌های توصیفی استخراج شده از صفحات ورودی به عنوان ورودی به این واحد، ارسال می‌شود تا مورد بررسی قرار گیرد. در نتیجه این فرآیند آمارهای ارزشمند مختلفی نظیر درصد استفاده از کدینگ‌های مختلف توزیع صفحات و وبگاه‌ها در طبقه‌بندی‌های مختلف سرویس‌های ارائه شده نظیر جستجو، امنیت و غیره بدست می‌آید.

بخش مهم دیگر این معماری، واحد تهیه گزارش است. داده‌های توصیفی استخراج شده از صفحات ورودی به عنوان ورودی به این واحد، ارسال می‌شود تا مورد بررسی قرار گیرد. در نتیجه این فرآیند، آمارهای ارزشمند مختلفی نظیر درصد استفاده از کدینگ‌های مختلف، توزیع صفحات و وبگاه‌ها در طبقه‌بندی‌های مختلف، سرویس‌های ارائه شده نظیر جستجو، امنیت و غیره، بدست می‌آید.

جدول ۱. نمونه‌ای از ورودی‌های تجزیه‌کننده به ازای وبگاه <http://www.whc.ir>

```

<!-- DOCID: 166182 URL: www.whc.ir -->
<head>
<meta http-equiv="Content-Language" content="fa">
<meta http-equiv="Content-Type" content="text/html; charset=utf-8">
<meta name="author" content="xx">
<meta name="copyright" content="2004-2005 by xxx">
<title>Send to friend</title>
<style type="text/css"></style>
<script type="text/javascript"></script>
</head>
<body>
<div id="container">
<form method="POST" action="">
<table>
<tr>
<td>.....</td>
</tr>
<tr>
<td><input type="submit" name="submit">
</td><td>&nbsp;</td></tr>
</table>
<input type="hidden" name="action" value="1">
<input type="hidden" name="newsid" value="1117">
</form>
</div>
</body>

```

۴. نتایج بدست‌آمده

ورودی سامانه تحلیل‌گر محتوا، بدین صورت انتخاب شد که به‌ازای هر وبگاه، حداکثر بیست‌هزار صفحه مورد بررسی قرار گرفت. علت این امر این است که وبگاه‌های با تعداد صفحات بالاتر، معمولاً فقط وبگاه‌های خبری نظیر IRNA، ISNA، IRIB و غیره هستند. در این خصوص، آمار حاکی از آن است که از حدود ۸ میلیون صفحه موجود در دامنه IR، حدود ۶ میلیون آن، مربوط به وبگاه‌های خبری است. لذا محتوای مناسب به جز خبر، تنها حدود ۲ میلیون صفحه است که این رقم در مقابل اغلب کشورهای که بیش از ۱۰۰ میلیون صفحه دارند، رقم بسیار ناچیزی می‌باشد (به‌عنوان مثال، دولت الکترونیکی کره جنوبی، شامل بیش از ۱۰۸ میلیون صفحه است).

بواسطه حجم محاسباتی بالای به‌منظور انجام ارزیابی‌های مختلف، یک جامعه آماری شامل حدود ۶۰۰ هزار صفحه، در نظر گرفته شد و بررسی‌های مختلف، روی این مجموعه، انجام شد. در زیر آمار صفحات و وبگاه‌ها از نظر حجم، عمر و نوع محتوا برای تمام صفحات غیر از اخبار شامل دو میلیون صفحه ارائه شده‌است.

ورودی سامانه تحلیل گر محتوا، بدین صورت انتخاب شد که به ازای هر وبگاه، حداکثر بیست هزار صفحه مورد بررسی قرار گرفت. علت این امر این است که وب گاه های با تعداد صفحات بالاتر، معمولاً فقط وبگاه های خبری نظیر IRIB ISNA IRNA و غیره هستند. در این خصوص، آمار حاکی از آن است که از حدود 8 میلیون صفحه موجود در دامنه IR. حدود 6 میلیون آن، مربوط به وبگاه های خبری است. لذا محتوای مناسب به جز خبر، تنها حدود 2 میلیون صفحه است که این رقم در مقابل اغلب کشورها که بیش از 100 میلیون صفحه دارند، رقم بسیار ناچیزی می باشد (به عنوان مثال دولت الکترونیکی کره جنوبی شامل بیش از 108 میلیون صفحه است).

بواسطه حجم محاسباتی بالای به منظور انجام ارزیابی های مختلف، یک جامعه آماری شامل حدود 600 هزار صفحه، در نظر گرفته شد و بررسی های مختلف، روی این مجموعه، انجام شد.

در زیر آمار صفحات و وبگاه ها از نظر حجم، عمر و نوع محتوا برای تمام صفحات غیر از اخبار شامل دو میلیون صفحه ارائه شده است.

جدول شماره دو، آمار میانگین حجم و سن (1) و عمق و تعداد صفحات وبگاه ها را نشان می دهد.

عکس

خصوصیات وب ایران ۹۷

۴,۱. آمار وبگاهها

جدول شماره دو، آمار میانگین حجم و سن^۱ و عمق و تعداد صفحات وبگاهها را نشان می دهد.

جدول ۲. خلاصه آمار سایت

تعداد وبگاههای معتبر	۱۰۰۰۰
در داخل کشور ۱۲. درصد وبگاههای	٪۲۷
در خارج کشور ۱۲. درصد وبگاههای	٪۷۳
میانگین تعداد صفحات در هر سایت	۲۳۰
میانگین تعداد صفحات ایستا در هر سایت	۱۲۰
میانگین تعداد صفحات پویا در هر سایت	۱۱۰
میانگین سن مسن ترین صفحه (بر حسب ماه)	۱۳/۶
میانگین سن جوان ترین صفحه (بر حسب ماه)	۶/۸۳
میانگین سن متوسط ترین صفحه (بر حسب ماه)	۸/۵
میانگین عمق ماکزیمم سایت	۴/۵۲
میانگین حجم سایت (بر حسب مگابایت)	۳/۰۶

شکل شماره سه تعداد صفحات وبگاهها را بر حسب درصد سایت نشان می دهد. همانطور که نشان داده شده است بیشتر وبگاهها دارای تعداد صفحات بین ۱۰۰ و ۱۰۰۰ (میانگین ۲۳۰) صفحه می باشند. نمودار دارای توزیع تقریباً خطی می باشد. شکل شماره چهار، توزیع تجمعی شکل شماره سه را نشان می دهد.

شکل شماره پنج، توزیع تجمعی حجم محتوای وبگاهها را نشان می دهد که از این آمار می توان برای مدل سازی محتوای وبگاهها استفاده کرد.

۱. سن صفحه عبارت است از مدت زمان میان ایجاد یا تغییر محتوای صفحه و زمان فعلی

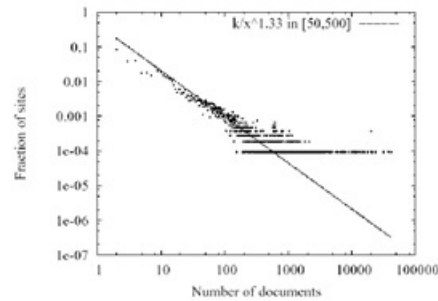
شکل شماره سه تعداد صفحات وبگاهها را بر حسب درصد سایت نشان می دهد. همان طور که نشان داده شده است بیشتر وبگاهها

دارای تعداد صفحات بین 100 و 1000 (میانگین 230) صفحه می باشند. نمودار دارای توزیع تقریباً خطی می باشد. شکل شماره چهار، توزیع تجمعی شکل شماره سه را نشان می دهد.

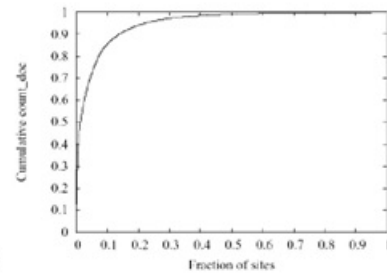
شکل شماره پنج، توزیع تجمعی حجم محتوای وبگاه ها را نشان می دهد که از این آمار می توان برای مدل سازی محتوای وبگاه ها استفاده کرد.

ص: 97

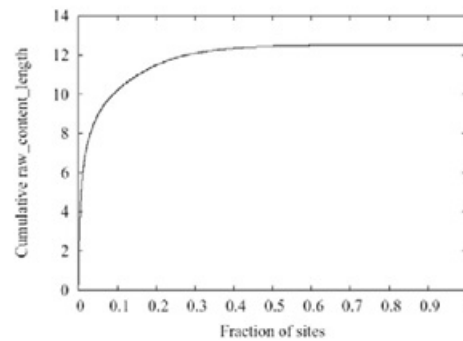
1- سن صفحه عبارت است از مدت زمان میان ایجاد یا تغییر محتوای صفحه و زمان فعلی



شکل ۳. تعداد اسناد در وبگاهها

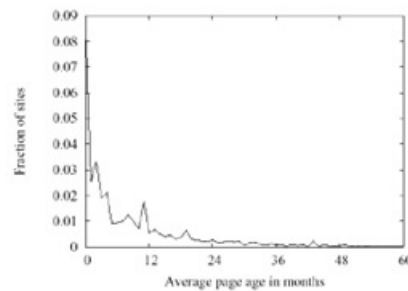


شکل ۴. توزیع تجمعی اسناد در وبگاهها

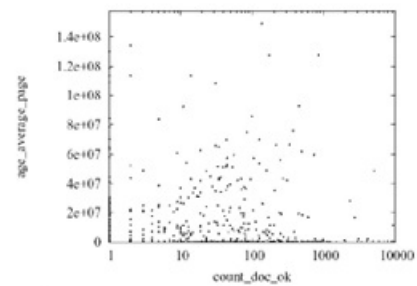


شکل ۵. توزیع تجمعی حجم محتوای وبگاهها

شکل شماره شش نیز میانگین سن صفحات موجود در وبگاهها را نشان می‌دهد. همچنین رابطه بین تعداد اسناد و سن یک سایت در شکل شماره هفت، نشان داده شده است.



شکل ۶. میانگین سن صفحات وبگاهها بر حسب ماه



شکل ۷. تعداد اسناد بر حسب میانگین سن آنها

شکل ۳. تعداد اسناد در وبگاهها

شکل ۴. توزیع تجمعی اسناد در وبگاهها

شکل ۵. توزیع تجمعی حجم محتوای وبگاهها

شکل شماره شش نیز میانگین سن صفحات موجود در وبگاه ها را نشان می دهد. همچنین رابطه بین تعداد اسناد و سن یک سایت در شکل شماره هفت، نشان داده شده است.

شکل 6. میانگین سن صفحات وب گاه ها بر حسب ماه

شکل 7. تعداد اسناد بر حسب میانگین سن آن ها

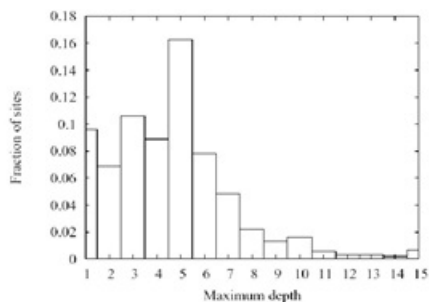
ص: 98

شکل شماره هشت، توزیع عمق وب گاه ها را نشان می دهد. درصد قابل توجهی از وب گاه ها دارای عمق 5 (23%) می باشند. شایان ذکر است که با توجه به شاخص های بدست آمده از وب گاه ها نظیر عمق، تعداد صفحات، حجم محتوا و سن صفحات، به راحتی می توان یک سایت را مدل سازی کرد. با داشتن یک مدل مناسب از وبگاه ها می توان الگوریتم های مختلف را به راحتی تحت آزمون قرار داد.

عکس

۹۹ خصوصیات وب ایران

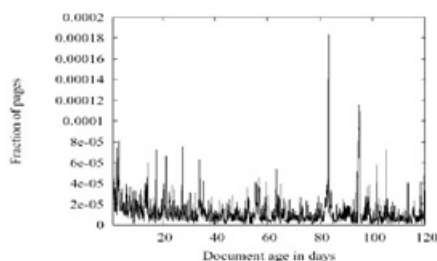
شکل شماره هشت، توزیع عمق وبگاهها را نشان می دهد. درصد قابل توجهی از وبگاهها دارای عمق ۵ (۲۳٪) می باشند. شایان ذکر است که با توجه به شاخص های بدست آمده از وبگاهها نظیر عمق، تعداد صفحات، حجم محتوا و سن صفحات، به راحتی می توان یک سایت را مدل سازی کرد. با داشتن یک مدل مناسب از وبگاهها می توان الگوریتم های مختلف را به راحتی تحت آزمون قرار داد.



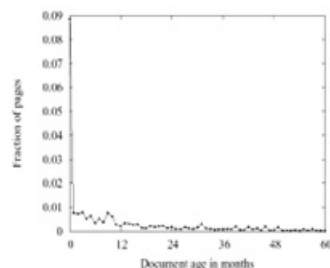
شکل ۸. بیشینه عمق وبگاهها

۴.۲. آمار صفحات

شکل های ۱۰ و ۱۱، به ترتیب، نمایانگر سن درصد صفحات برای حسب روز، ماه و سال است. همانطور که نشان داده شده سن صفحات نسبتاً بالا می باشند. به عبارت دیگر نرخ بروزآوری صفحات، پایین است.



شکل ۹. عمر صفحات بر حسب روز



شکل ۱۰. عمر صفحات بر حسب ماه

شکل 8. بیشینه عمق وب گاه ها

4.2. آمار صفحات

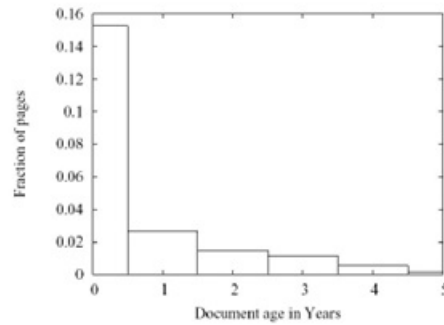
شکل های 9 ، 10 و 11 به ترتیب نمایان گر سن درصد صفحات برای حسب روز، ماه و سال است. همانطور که نشان داده شده سن صفحات نسبتاً بالا می باشند. به عبارت دیگر نرخ بروزآوری صفحات، پایین است.

شکل 9. عمر صفحات بر حسب روز

شکل 10. عمر صفحات بر حسب ماه

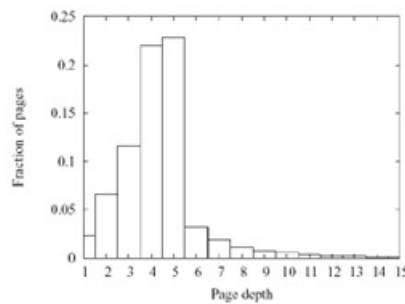
ص: 99

۱۰۰ مدیریت منابع اطلاعاتی وب

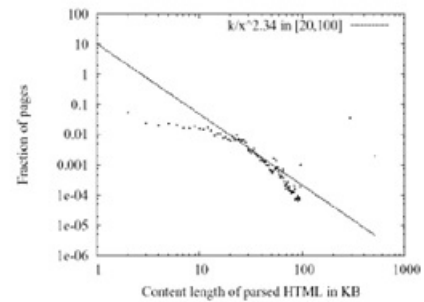


شکل ۱۱. عمر صفحات بر حسب سال

شکل شماره دوازده، عمق صفحات را نشان می‌دهد. با توجه به شکل، اکثر صفحات در عمق ۴ و ۵ قرار دارند. قابل توجه است که عمق صفحات، دارای توزیع پواسون می‌باشد. حجم محتوای صفحات که یک توزیع تقریباً خطی است، در نمودار شماره ۱۳ نشان داده شده است. بیشتر صفحات دارای حجم بین ۱۰ تا ۱۰۰ کیلوبایت را دارا می‌باشند.



شکل ۱۲. عمق صفحات (توزیع پواسون)



شکل ۱۳. حجم محتوای صفحات بر حسب کیلوبایت (تقریباً خطی)

شکل شماره چهارده رابطه بین سن و حجم صفحات را نشان می‌دهد. همچنین رابطه بین عمق با سن و حجم صفحه در شکل‌های ۱۵ و ۱۶ نشان داده شده است. با داشتن شاخص‌های فوق (عمر، حجم، عمق)، به راحتی می‌توان مدل مناسبی از وب ایران بدست آورد.

شکل ۱۱. عمر صفحات بر حسب سال

شکل شماره دوازده، عمق صفحات را نشان می‌دهد. با توجه به شکل، اکثر صفحات در عمق ۴ و ۵ قرار دارند. قابل توجه است که عمق صفحات، دارای توزیع پواسون می‌باشد.

حجم محتوای صفحات که یک توزیع تقریباً خطی است، در نمودار شماره 13 نشان داده شده است. بیشتر صفحات دارای حجم بین 10 تا 100 کیلوبایت را دارا می باشند.

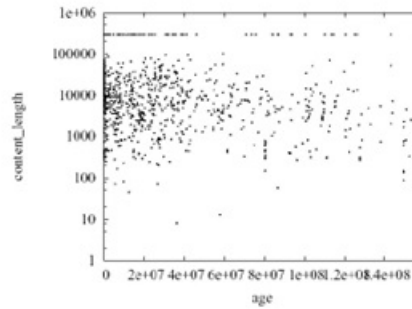
شکل 12. عمق صفحات (توزیع پواسون)

شکل 13. حجم محتوای صفحات بر حسب کیلوبایت (تقریباً خطی)

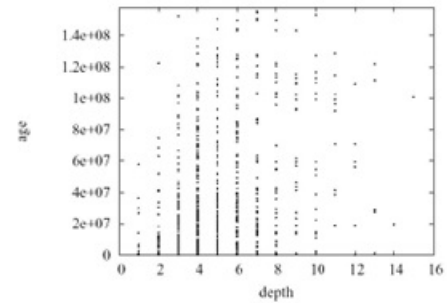
شکل شماره چهارده رابطه بین سن و حجم صفحات را نشان می دهد. همچنین رابطه بین عمق با سن و حجم صفحه در شکل های 15 و 16 نشان داده شده است. با داشتن شاخص های فوق (عمر، حجم، عمق)، به راحتی می توان مدل مناسبی از وب ایران بدست آورد.

ص: 100

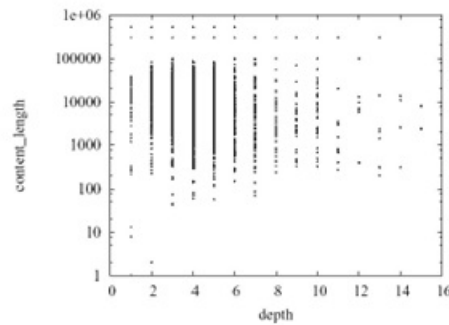
خصوصیات وب ایران ۱۰۱



شکل ۱۴. سن صفحات بر حسب حجم محتوای صفحات



شکل ۱۵. سن صفحات بر حسب عمق



شکل ۱۶. عمق صفحات بر حسب حجم آنها

۴.۳. آمار نوع صفحات

جدول ۳، ۴، ۵، ۶، ۷، ۸ و ۹ به ترتیب، آمارهای مختلفی از صفحات، اعم از ایستا، پویا، صوتی، ویدیویی، تصویری و اجرایی را نشان می‌دهد. طبق جدول ۵، حدود ۵۰ درصد از صفحات، پویا هستند. در مجموع حدود ۸ میلیون آدرس صفحه ایستا و ۱۷ میلیون آدرس صفحه پویا کشف شده است. از بین صفحات پویا، بیشترین درصد را انواع PHP و ASP تشکیل داده‌اند.

جدول ۳. آمار انواع صفحات

تعداد صفحات (بر حسب میلیون)	۲
درصد صفحات ایستا	۵۰/۹۴
درصد صفحات پویا	۴۹/۰۶
درصد صفحات دوتخته‌ای (Duplicate)	۱

شکل ۱۴. سن صفحات بر حسب حجم محتوای صفحات

شکل ۱۵. سن صفحات بر حسب عمق

شکل ۱۶. عمق صفحات بر حسب حجم آنها

4.3. آمار نوع صفحات

جدول 3، 4، 5، 6، 7، 8 و 9 به ترتیب، آمارهای مختلفی از صفحات، اعم از ایستا، پویا، صوتی، ویدیویی، تصویری و اجرایی را نشان می دهد. طبق جدول 5، حدود 50 درصد از صفحات، پویا هستند. در مجموع حدود 8 میلیون آدرس صفحه ایستا و 17 میلیون آدرس صفحه پویا کشف شده است. از بین صفحات پویا، بیشترین درصد را انواع PHP و ASP تشکیل داده اند.

جدول 3. آمار انواع صفحات

ص: 101

جدول ۴. آمار نوع صفحات پویای کشف شده

Filename extension	Number of links found	Percent
php	8,643,428	49.45%
asp	8,429,913	48.22%
jsp	158,477	0.91%
cfm	108,619	0.62%
cgi	81,769	0.47%
shtml	45,847	0.26%
pl	7,507	0.04%
dll	3,130	0.02%
fcgi	1,046	0.01%

جدول ۵. آمار درصد صفحات ایستای کشف شده

Filename extension	Number of links found	Percent
html	7,807,297	95.84%
pdf	220,004	2.70%
xml	63,679	0.78%
doc	15,493	0.19%
rdf	13,951	0.17%
ppt	5,731	0.07%
mso	5,715	0.07%
txt	5,335	0.07%
xls	2,763	0.03%
ps	2,328	0.03%
tex	1,641	0.02%
dvi	1,379	0.02%

جدول ۶. آمار اسناد صوتی کشف شده

Filename extension	Number of links found	Percent
mp3	27,839	44.14%
ram	16,436	26.06%
wma	16,216	25.71%
mid	1,744	2.77%
wav	513	0.81%
pls	173	0.27%
asf	78	0.12%
au	64	0.10%

جدول ۷. آمار ویدیوی کشف شده

Filename extension	Number of links found	Percent
swf	324,678	91.51%
wmv	26,666	7.52%
mpg	2,181	0.61%
avi	1,179	0.33%
mov	99	0.03%
qt	6	0%

جدول ۴. آمار نوع صفحات پویای کشف شده

جدول ۵. آمار درصد صفحات ایستای کشف شده

جدول ۶. آمار اسناد صوتی کشف شده

جدول 7. آمار ویدیوی کشف شده

ص: 102

خصوصیات وب ایران ۱۰۳

جدول ۸. آمار فایل‌های تصویری کشف‌شده

Filename extension	Number of links found	Percent
gif	11,441,330	81.13%
jpg	2,310,009	16.38%
png	300,407	2.13%
ico	38,371	0.27%
bmp	11,536	0.08%
wmf	694	0%
tiff	198	0%
img	2	0%
pbm	1	0%

جدول ۹. آمار نرم‌افزارهای کشف‌شده

Filename extension	Number of links found	Percent
exe	40,344	97.11%
iso	1,066	2.57%
rpm	72	0.17%
patch	42	0.10%
diff	9	0.02%
pdb	9	0.02%
deb	4	0.01%

۵. نتیجه‌گیری و کارهای آینده

هدف این مقاله، بررسی ویژگی‌های وب فارسی است. بر این اساس، سامانه نیمه خودکاری جهت انجام ارزیابی‌های مورد نظر، طراحی و پیاده‌سازی شد. از خصوصیات این سامانه خودکار بودن آن می‌باشد که در زمان‌های مختلف می‌توان آنرا اجرا و آمارهای مورد نظر را استخراج کرد. بررسی‌ها بر روی حدود ۱۱ هزار سایت با پسوند IR و مشتمل بر دو میلیون صفحه انجام شد. این مجموعه، اغلب سازمان‌های دولتی و غیر دولتی را پوشش می‌دهد. هدف اصلی پژوهش، استخراج و تعیین ویژگی‌های وب اعم از خصوصیات ساختاری نظیر ویژگی‌های گراف وب و نیز خصوصیات محتوایی وب فارسی است. شاخص‌های استخراج شده شامل توزیع محتوای وب‌وبگاه‌ها، حجم محتوای فارسی، نوع محتوا، سرویس‌های ارائه شده (جستجو، RSS و غیره) نرخ بروزآوری محتوا، توزیع مکانی وبگاه‌ها در داخل و خارج کشور و غیره می‌باشد. با انجام این پروژه بهتر می‌توان راهبردهای آینده مربوط به ICT را تعیین و تبیین کرد. از نتایج بارز این آمار عبارتند از: کم بودن حجم محتوای فارسی وب در مقایسه با سایر کشورها (نسبت ۱۰٪)، نیاز به یک مرکز داده اینترنتی در کشور (۷۳٪) سایت‌ها در خارج از کشور هستند، نیاز به موتورهای جستجوی بومی (محتوا بیشتر از کدینگ‌های عربی استفاده کرده است)، کم بودن سرویس‌های دولت الکترونیک (جستجو، امنیت و غیره). علاوه بر موارد فوق خروجی‌هایی مانند مجموعه لغات فارسی موجود در وب برای استفاده در خطایاب‌ها و موتورهای جستجو بدست آمده است. در ادامه این کار، در نظر است تا با مطالعه تطبیقی وب فارسی در طی دوره‌های مختلف زمانی،

جدول ۸. آمار فایل‌های تصویری کشف شده

جدول ۹. آمار نرم‌افزارهای کشف شده

۵. نتیجه‌گیری و کارهای آینده

هدف این مقاله، بررسی ویژگی های وب فارسی است. بر این اساس، سامانه نیمه خودکاری جهت انجام ارزیابی های مورد نظر طراحی و پیاده سازی شد. از خصوصیات این سامانه خودکار بودن آن می باشد که در زمان های مختلف می توان آن را اجرا و آمارهای مورد نظر را استخراج کرد. بررسی ها بر روی حدود 11 هزار سایت با پسوند IR و مشتمل بر دو میلیون صفحه انجام شد. این مجموعه اغلب سازمان های دولتی و غیر دولتی را پوشش می دهد. هدف اصلی پژوهش، استخراج و تعیین ویژگی های وب اعم از خصوصیات ساختاری نظیر ویژگی های گراف وب و نیز خصوصیات محتوایی وب فارسی است. شاخص های استخراج شده شامل توزیع محتوای وب و وب گاه ها، حجم محتوای فارسی، نوع محتوا، سرویس های ارائه شده (جستجو، RSS و غیره) نرخ بروزآوری محتوا، توزیع مکانی وبگاه ها در داخل و خارج کشور و غیره می باشد. با انجام این پروژه بهتر می توان راهبردهای آینده مربوط به ICT را تعیین و تبیین کرد. از نتایج بارز این آمار عبارتند از: کم بودن حجم محتوای فارسی وب در مقایسه با سایر کشورها (نسبت 10%)، نیاز به یک مرکز داده اینترنتی در کشور (73%) سایت ها در خارج از کشور هستند، نیاز به موتورهای جستجوی بومی (محتوا بیشتر از کدینگ های عربی استفاده کرده است)، کم بودن سرویس های دولت الکترونیک (جستجو، امنیت و غیره). علاوه بر موارد فوق خروجی هایی مانند مجموعه لغات فارسی موجود در وب برای استفاده در خطایاب ها و موتورهای جستجو بدست آمده است.

در ادامه این کار، در نظر است تا با مطالعه تطبیقی وب فارسی در طی دوره های مختلف زمانی،

نرخ تغییرات شاخص های مختلف تعیین گردد. این اطلاعات، جهت گیری تغییرات وب فارسی را نشان خواهد داد و کمک شایانی به تصمیم گیری های کلان حوزه فناوری اطلاعات خواهد کرد. در این عین حال، در نظر است تا با بررسی های بیشتر در رابطه با محتوای صفحات، اطلاعات بیشتری نظیر متوسط طول حروف در کلمات فارسی، متوسط تعداد کلمات در اسناد، مشخصات گراف (مانند توزیع درجه های ورودی و خروجی) برای زبان فارسی و آمارهای متنوع دیگر، تعیین گردد.

منابع

اشاره

<http://www.worldwidewebsite.com/>, January 2013

Krootkaew C., A. Vongpakaymas, A. Jeawpoung 2002. Services E-readiness Explorer (SEE): Automatic ..monitoring tool for thailand e-government project in proceeding of EurAsia-ICT, Shiraz, Iran, Oct

Baeza-Yates R., C. Castillo and V. LÓpez, 2005. Characteristics of the Web of Spain. Journal of) ..Cybernetics, 9(1

Baeza-Yates, R. and F. Lalanne.2004.Characteristics of the korean web. Technical report, Korea Chile IT ..Cooperation Center ITCC, 2004

Bordino, I. and D. Donato.2009 .Dynamic characterization of a large Web graph, Dynamic characterization ..of a large Web graph. Proceedings of the WebSci'09: Society On-Line, pp. 198-202

Broder A. Z., S.R. Kumar ,F. Maghoul ,P., Raghavan,S. Rajagopalan ,R. Stata , A. Tomkins and J.L.Wiener ..2000.Graph structure in the web. Computer Networks, 33(1): 309-320

Gabriel Tolosa, G., F. Bordignon, R. Baeza-Yates, C. Castillo .2007.Characterization of the argentinian) ..web. International Journal of Scientometrics, Informetrics and Bibliometrics, 7(1

Gomes, D. and M.J. Silva.2003. A characterization of the Portuguese Web.In Proceedings of 3rd ECDL ..Workshop on Web Archives, Trondheim, Norway

Keyhanipour A.H., A.M. Zare Bidoki, M. Mahmoudi and M. Azadnia.2007.Evaluation of Iran's web content from e-government perspective. Proceedings of the 12th International CSI Computer Conference, ..pp. 2081-2086, Tehran, Iran

Rauber, A., O. Aschenbrenner, O. Witvoet, R.M. Bruckner, and M. Kaiser.2002.Uncovering information)

Thelwall, M. and D. Wilkinson.2003.Graph structure in three national-academic webs: Power laws with anomalies. Journal of the American Society for Information Science and

.Technology, 54(8): 706-712

Thelwall, M. and A. Zuccala. 2008. A university-centred European Union link analysis. *Scientometrics*,
.75(3): 407-420

Shestakov, D.2011.Sampling the national deep web, database and expert systems applications. Lecture
.Notes in Computer Science, 6860: 331-340

ص: 105

مقاله حاضر گزارشی است که توسط پژوهشگران مؤسسه اینترنت آکسفورد برای کنفرانس بین المللی حفاظت اینترنت تهیه شده است. هدف از آن ایجاد انگیزه بحث و تبادل نظر بیشتر مابین آرشیویست ها و پژوهشگران وب در مورد روش های آرشیو وب است در بخش اول مقاله نمای کلی 4 سناریو احتمالی از آرشیوهای وب شامل نیروانا، آپو کالیپس انفرادی و آرشیوهای غبار آلوده را معرفی می کند. سپس به توصیف انواع مختلفی از پژوهش هایی که در مورد وب پویا در حال انجام هستند، می پردازد. و در بخش آخر چالش های فعلی و پیش روی آرشیوهای وب را پوشش می دهد و در نهایت به ارائه پیشنهادها برای رفع چالش ها و ارائه راه حل های میان مدت و دراز مدت برای مقابله با تغییرات احتمالی را ارائه می کند.

اشاره

*آینده آرشیو وب (1)

نوشته: اریک تی. مهیر (2)، آرتور توماس (3)، رالف شرودر (4) | مؤسسه اینترنت آکسفورد (5)

ترجمه: رضا خانی پور (6)، محبوبه قربانی (7)

خلاصه اجرایی

این گزارش، توسط پژوهشگران مؤسسه اینترنت آکسفورد برای کنسرسیوم بین‌المللی حفاظت اینترنت (8) به نگارش در آمده است. هدف از آن ایجاد انگیزه بحث و تبادل نظر بیشتر مابین آرشیویست‌ها و پژوهشگران وب در مورد روش‌های آرشیو وب است که می‌تواند مورد استفاده پژوهشگران قرار گیرد.

بخش اول. نمای کلی از چهار سناریوی احتمالی برای آینده:

اشاره

- نیروانا (9): مکانی که آرشیوهای وب توسط بسیاری از گروه‌ها به طور گسترده استفاده، استاندارد سازی، و قابل مرور می‌شود و رابط کاربر قوی برای دسترسی دارند.

ص: 107

Web archives: the future(s) -1

Eric T. Meyer -2

Arthur Thomas -3

Ralph Schroeder -4

Oxford Internet Institute -5

6- عضو هیئت علمی سازمان اسناد و کتابخانه ملی ج. ا. ا.

7- دانشجوی دکترای علم اطلاعات و دانش‌شناسی واحد علوم و تحقیقات تهران دانشگاه آزاد اسلامی

8- (the International Internet Preservation Consortium) IIPC

9- Nirvana

- آپوکالیپس (1): آرشیوها تجزیه می شوند و استاندارد نشده باقی می مانند، به سختی بازیابی شده و در دسترس قرار می گیرند و در نتیجه مفید نیست و به سختی استفاده می شوند.

- انفرادی: در این سناریو، آرشیوها مانند یک هوش هم بسته منفرد که قادر است بین اشیای رقومی و افراد ارتباط برقرار کند، غیر ضروری می شوند.

- آرشیوهای غبارآلود: در این سناریو، جامعه آرشیوی وب هرگز پاسخگوی سؤال (که چی؟) نیست آرشیوهای وب به طور گسترده، مورد استفاده قرار نمی گیرد، بلکه در حال جمع آوری غبارهای و رقومی هستند.

این سناریوها ما را به تفکر درباره تعامل بین آرشیوها پژوهشگران و محققان در روش های گوناگون قادر می سازد.

بخش 2. انواع مختلف پژوهش هایی را توصیف می کند که در مورد وب پویا در حال انجام هستند. فنی که حالا به طور گسترده ای، بیشتر از استفاده از آرشیوهای وب مورد بهره برداری قرار می گیرد.

مقصود این است که استفاده از وب پویا می تواند الهام بخش تفکر درباره کاربردهای بالقوه آرشیوهای وب باشد. این کاربردها عبارت اند از:

مصورسازی: که از طریق آن پیوندها، نه تنها بین وبگاه ها، بلکه میان انواع مختلف اطلاعات نیز

می توانند ایجاد شوند، همانگونه که سازمان ها و دیدگاه های کلی آرشیوها را توانمند می سازد.

سنجش های دگرساز (2): دانشمندان حوزه علم سنجی شروع به کسب داده ها از منابع مجزا، از طریق تحلیل های استنادی نموده اند. به طور مثال: بلاگ های محققان و پیوندهای بین این بلاگ ها.

فنون دیگری، همچون نقشه برداری محتوای ایجاد شده توسط کاربر و تحلیل های شبکه اجتماعی در اینجا ارائه شده است.

بخش 3. چالش های فعلی و پیش رو را پوشش می دهد. نخستین قسمت از این بخش روش های تغییر وب، پیشنهادهایی در این مورد و راه حل های میان مدت و دراز مدت برای مقابله آرشیوها با این تغییرات را توصیف می نماید. این بخش با پیشنهادهایی درباره گام های پیش رو پایان می یابد.

مقدمه

این گزارش، به صورت پیش نویس در می 2012 نگارش یافت و در مجمع عمومی آی. آی. پی. سی. (3) در

ص: 108

2- ALT- metrics : ایجاد و مطالعه سنجش های جدید بر اساس تحلیل های شبکه های اجتماعی و اطلاع رسانی علمی (دسترسی در <http://altmetrics.org/about>)

3- IIPC (آی آپی سی) کنسرسیوم بین المللی حفاظت اینترنت (<http://www.netpreserve.org>) که این کار را پایه گذاری و زیربنایی برای بحث های بیشتر ارائه نمود که با ارائه در نشست مجمع عمومی IIPC در سال 2011 و برگزاری کارگاهی که ویرایشی نهایی از این گزارش خواهد بود آغاز می شود.

این گزارش، در یک نشست جامع، تلخیص و در یک کارگاه نیز ارائه شده است. همچنین، از طریق پست الکترونیک برای جامعه محققان اینترنت و جامعه کتابداری و اطلاع رسانی ارسال شده است. هدف این بود که پیش نویس گزارش منجر به تهییج و ترغیب شود، افکار را برانگیزد. برای اینکه آرشیویست های وب و پژوهشگران را به خروج از سکون وادار کند. افراد مشابه و دیگران را به فعالیت وادار کند. برای ایجاد تغییر، پیش از این برای بحث و تبادل نظر ایجاد انگیزه شده است؛ که فعالیت را به سوی تغییری ملموس ترغیب می کند.

چرا تغییر ضروری است؟ وقتی که IIPC توجه ما را به اجرای این پروژه جلب کرد، این ضرورت که جامعه آرشیوی وب و IPC به طور ویژه، باید روش هایی نوین برای تشویق کاربران و کاربردهای جدید آرشیوهای وب و مدل های جدید آرشیوسازی وب و سبک های جدید تعامل با محققان بررسی نماید، احساس می شد.

این مباحث، پیش تر در قالب دو گزارش به وسیله (2) JISC تدوین شده بود که بر روی شرایط فعلی هنر آرشیوهای وب (دو گرتی و دیگران، 2010) (3) و فرصت های سرمایه گذاری جدید متمرکز شده بود.

بعضی از نتایج این دو مقاله در ادامه آمده است، اما موضوعات عمومی سرتاسر اثر این بود که هنوز فاصله و شکافی بین جامعه پژوهشگران بالقوه ای که دلیل خوبی برای استفاده، تحلیل و اشتراک آرشیو وب دارند و جامعه واقعی (به طور کلی هنوز کوچک) پژوهشگران که در حال حاضر فعالیت می کنند وجود دارد (دو گرتی و دیگران، 2010، ص.5). تجربه کار ما در این گزارش و صحبت با اعضای آی. آی. پی. سی. و جامعه پژوهشگران اینترنتی اندکی باعث تغییر نظر ما در مورد این موضوع شده است، مسلماً ما از اینکه همیشه کاربردهای آرشیوهای وب به خوبی مورد بحث قرار گیرد و جامعه پژوهشی هنوز به صورت معنی داری به آن دست نیافته اند، اطمینان بیشتری یافته ایم. این گزارش، به نوبه خود این [دیدگاه] را کمی تغییر خواهد داد، اما اگر پاره ای از پیشنهاد های ارائه شده در آن به طور جدی توسط جوامع مرتبط، به کار گرفته شود احتمال اینکه در آینده آرشیوهای منابع اینترنتی برای پژوهشگران اهمیت بیشتری یابد، افزایش خواهد یافت.

این گزارش، در ابتدا به منظور اینکه افکار نظری را به آینده احتمالی وب معطوف کند، به صورت یک تمرین که ما را وارد به تفکر درباره آن چه که برای انجام در حال حاضر به آن نیاز داریم و اطمینان از اینکه بتوانیم به صورت مطمئن و پر ثمر از آرشیو وب در آینده استفاده کنیم، شکل یافته است. سپس، ما بر روی روش های ابزاری مورد استفاده برای پژوهش وب پویا تمرکز می کنیم، همانند یک اشاره گر بر روی انواع چیزهایی که می تواند برای کمک در فهم وب آرشیوی، گسترش یابد. در ادامه، ما بر روی یک سری از موضوعات و سؤالاتی که پژوهشگران می خواهند یا ممکن است بخواهند برای استفاده از وب آرشیوی

ص: 109

the Hague - 1

2 - ISC کمیته سیستم های اطلاعاتی مشترک (<http://www.jisc.ac.uk/>) که بخش پژوهش حوزه اطلاعات و ارتباطات و توسعه زیربنایی آموزش و پژوهش در بریتانیا را پایه گذاری نمود.

Dougherty - 3

به آن‌ها بپردازند تمرکز می‌کنیم. در پایان این بخش، برخی از چالش‌های انفرادی، سازمانی و سازمان‌های بین‌المللی که می‌توانند برای افزایش توانایی ما برای توضیح این موضوعات و پاسخ به این سؤالات، مورد توجه قرار گیرند را تبیین می‌نماییم.

ما گزارش مذکور را با نتایج بر اساس آن چه از این تجربه آموخته‌ایم، به پایان می‌بریم.

ساختن آینده

(بهترین راه برای پیش‌بینی آینده، ساختن آن است) (کای، 1995) (1)

برای آغاز بحث از آینده‌شناسی بهره می‌بریم. هدف پیش‌بینی آینده نیست. که اشتباه خواهد بود. در واقع، ما کاملاً در مورد تلاش‌هایی که برای پیش‌بینی آینده انجام شده است، اطمینان نداریم. ساختن سناریوها و تلاش‌های دیگری که برای پیش‌بینی آینده (ادعاهایی درباره آینده) طراحی شده‌اند، معمولاً در دانش، نهفته می‌مانند، به طوری که بیشتر چنین تلاش‌هایی، هرگز جدی گرفته نخواهد شد.

با وجود این، حداقل یک نوع از آینده‌شناسی که از نظر ما مناسب است وجود دارد. [و] آن زمانی است که هدف از تجربه، پیش‌بینی آینده نباشد؛ بلکه الهام بخشی باشد برای افرادی که مسئول ساختن نظام‌هایی هستند که زیربنای آینده خواهند بود. به منظور تفکر درباره نتایج تصمیمات فعلی در دوره‌هایی با تأثیرات احتمالاً طولانی مدت. آی.پی.سی. متشکل از بسیاری از افرادی است که در حال حاضر عهده دار توسعه و گسترش نظام‌ها، ابزارها، استانداردها، و تفاهم‌نامه‌هایی برای حفاظت از محتوای اینترنت با نیم‌نگاهی به سوی مفید ساختن آن برای فهم جامعه‌ای که در آن زندگی می‌کنیم، هستند.

در گسترش نظام‌های رایانه‌ای، «انتخاب‌های معماری» بسیاری وجود دارد. (کلینگ، مک‌کیم و کینگ 2003) (2)، در طول مسیر نقاطی که تصمیمات در آن جا اتخاذ می‌شوند، همان‌هایی هستند که یکی در دوراهی‌های طول مسیر به غیر از انتخاب‌های دیگر گزیده می‌شوند. شواهدی وجود دارد از اینکه جامعه آرشیوی وب اهمیت نقاط انتخاب در حال حاضر و در آینده نزدیک را می‌پذیرد. (پذیرفته است). یک سری از نقاط انتخاب با پیشنهادهایی که نوید تغییر زلزله‌گونه‌ای در آرشیو وب می‌دهند، مرتبط هستند: گذر از دسترسی به وبگاه‌ها و صفحه‌های انفرادی به سوی ساختن و استفاده از یک مجموعه همانند یک مجموعه در مقایسه با ایجاد دسترسی به بخش‌هایی از یک مجموعه. به چه تصمیماتی برای افزایش احتمال اینکه آرشیو - در - جعبه، مفید، قابل استفاده، ادامه‌پذیر و تاثیرگذار باشد، نیاز خواهد بود؟ ما دوباره به ایده آرشیو - در یک جعبه در این مقاله، موازی با چالش‌های دیگری که انتخابها را می‌طلبد، باز خواهیم گشت.

در ادامه، هدف این تمرین این است که جامعه آرشیوی وب کدام راه را برای آن انتخاب‌ها که بر آینده تأثیر خواهند گذاشت و برای پیشنهاد گام‌ها و انتخاب‌هایی که آینده را در یک جهت یا جهت‌گیری هدایت [ترغیب] می‌کند، درخواست می‌کند.

ص: 110

Kay -1

Kling, McKim, King -2

اشاره

ما می توانیم آینده های بسیار مختلفی را برای آرشیوهای وب و کاربردهای شان تصور نماییم، برای فرض بحث، ما چهار سناریوی بالقوه را که می توانند در دهه یا دو دهه آینده اجرا شوند، بررسی الزامات شان و راه های پیشنهادی که جامعه آرشیوی وب باید با آن ها هماهنگ شوند، ارائه داده ایم:

در ادامه این متن، به طور ویژه تعدادی از عناصر که این سناریوها را تکمیل خواهند کرد، نشان می دهیم، چالش هایی را که در سر راه به کارگیری آن ها وجود دارد شناسایی می نمائیم، و مثال هایی از گستره متنوع ابزارهای توسعه وب پویا را نشان می دهیم. اگر آن ها برای داده های تاریخی به کار گرفته شوند، بین بهترین و بدترین طرح می توانند تفاوت ایجاد نمایند.

سناریوی نیروانا

در بهترین شرایط، آرشیوهای وب، مستحکم، استاندارد شده و به صورت ایمن، حفاظت شده خواهند بود. در حالی که، در زمان مشابه به صورت باز، انعطاف پذیر و گسترده، و به عنوان بخشی از ابزار استاندارد پژوهشی در علم اینترنت، علوم سیاسی، اقتصادی، جامعه شناسی، تاریخ معاصر (و در آینده، تاریخ اواخر قرن بیستم و اوایل قرن بیست و یکم) روزنامه نگاری، زبان شناسی، ارتباطات، تجارت، مطالعات رسانه و دیگر رشته ها استفاده می شوند.

علاوه بر دانشگاه، آرشیوهای وب برای عموم، حکومت ها، واحدهای سیاسی، و گروه های مشاوران

(متفکران) و سازمان های غیر دولتی مفید و کاربردی خواهد بود. متأسفانه، از خیلی از جنبه های بسیاری، کم احتمال ترین سناریوست، به این سبب که به نظر می رسد در حال حاضر، برای امکان پذیر ساختن ورود آن به جامعه آرشیوی وب به تلاش گسترده تر و منابع وسیع تری نیاز خواهد بود. اگر چه، ممکن است برای اینکه آن را در ذهن خود به صورت مطلوب نگهداری کنیم، مفید به نظر برسد، همان گونه که ایجاد توازن را بین آن چه که می توانیم و آن چه باید باشد بررسی می کنیم.

به منظور ارائه این سناریو، حتی در طرح کلی، نیاز است که اتفاق هایی رخ دهد (در بخش بعدی، مثال هایی از وب پویا در زمان وقوع، ارائه خواهیم کرد) این موارد عبارت اند از:

- توسعه قوی تر و مؤثرتر ابزارهای جست و جوی متن، تحلیل و استخراج اطلاعات، مصورسازی، وبلاگ نویسی اجتماعی، تحلیل های طولی و تحلیل های نظری.

- توسعه روش های بهتر برای اینکه کاربران، گشتالت مجموعه های چندتایی و انفرادی را درک کنند. در حالی که محتوای متنی می تواند جست و جو شود، نیاز به فراداده غنی برای پشتیبانی خطوط کلی محتوا و یا روش های جدید سازماندهی آن می باشد.

افراد به طور ویژه در تشخیص الگوهای بصری توانمند هستند، و به نظر می رسد ابزارهای گرافیکی یکی از بهترین راه های کسب آن می باشند. می توانیم ایجاد فضای مجازی را تصور کنیم که به [fly through](#) های سه بعدی و سایر روش های فضایی حسی امکان

1 - Fly through: نوعی شبیه سازی کامیوتری است که در فضاهای مجازی به کاربران اجازه مشاهده مدل های مجازی سایت ها را می دهد. (دسترسی در: <http://dictionary.reference.com/browse/fly-through>)

حالت، محیط های مجازی کاملاً همه جانبه از نوع غار (1) (شرودر، 2011) می توانند همکاری گروه های توزیع شده فضایی (مکانی) افراد را پشتیبانی کرده و امکان اشتراک مؤثر و تعامل اجتماعی را ایجاد کنند. از این دیدگاه، تمامیت آرشیوهای وب می تواند به صورت یک کلیت بزرگ در نظر گرفته شود، که افراد می توانند در آن به صورت منفرد یا گروهی سیر نمایند. وب فعلی (و به همین ترتیب آرشیوهای وب) درک ناچیزی از سازماندهی فضایی ارائه می دهند و نشانه های مناسب فضایی یا قابلیت های دیگر کمک به مرور و کاوش را محدود (سلب) می کنند. بر خلاف کتابخانه فیزیکی، اسناد رقومی در محیط جدا سازی شده از یکدیگر نگهداری می شوند، در مکانی که اسناد در دنیایی سه بعدی سازماندهی می شوند تا امکان مرور آن ها به واسطه گردش (جست و جو) در محیط را میسر ساخته و نشانگرهایی همچون همجواری فضایی (مکانی) به کشف آن ها کمک می کنند.

- در حالی که ابزارهای وب نوشت نویسی اجتماعی پدید می آیند، آرشیو های وب، فاقد سایر ابزار های همکاری هستند، مانند موتورهای پیشنهادی (مخزن پیوسته آمازون مثال عالی از این امکان است).

- به طور فزاینده ما نیازمند محتوای آماده شده توسط کاربران (وب 2) در مقیاسی بزرگ (فیس بوک (2)) هستیم. اما ساختاری حیرت انگیز بر اساس چنین محتوای سازمان نیافته نامتجانس و ذاتاً سازمان نیافته ای، در این مقیاس با مشکل مواجه است. انجام آن به وسیله ماشین از نظر فناوری بسیار چالش برانگیز است، بخشیدن غنای معنایی، به عنوان جایگزینی برای پشتیبانی رهیافت انباشت منبع، به کاربران آرشیو اجازه می دهد که محتوا را سازماندهی کنند این شکل نهایی وب نوشت نویسی (حاشیه نگاری) اجتماعی است که کاربران نه فقط داده بلکه فراداده را نیز تولید می کنند (گازان، 2008) (3).

در نیروانا، انتخاب های امروز توسط پژوهشگران آینده ستایش خواهد شد، کسانی که به اتکای اطلاعات و شواهد، تلاش های انسانی انجام شده در اینترنت را، تجسم داده، حفاظت کرده، و افزایش داده اند تا به همه روش های فنون قدرتمند پژوهش دست یابند و توانمند شوند.

سناریوی آپوکالیپس (آخر الزمان)

در بدترین حالت، تغییرات روزافزون اینترنت برای بسط و توسعه فناوری های جدید (HTML5)، محتوای قابل اجرا، ویدئوی جاسازی شده و اشیای تعاملی، وبگاه های مبتنی بر پایگاه داده، نرم افزارهای موبایل غیر وابسته به HTML.TTP و مانند آن) با سرعتی گیج کننده، ادامه خواهند یافت، و ابزارهای آرشیو سازی وب در حفظ همگامی با آن ها شکست خواهند خورد، و بیش از پیش عقب می افتند. حتی اگر فناوری های آرشیو سازی وب بتوانند سرعت را حفظ کنند، تغییرات دائمی (پیوسته) همه اشکال مطرح شده، چالشی حل نشدنی را ایجاد می نماید. در این طرح به طور صادقانه ای، فقط اندکی از محتوای واقعی ذخیره می شود، و حتی اگر ذخیره شود، افزودنی های اختصاصی برای مشاهده آن حفظ نشده یا قابل حفظ

ص: 112

CAVE - 1: فضایی شبیه غار در محیط مجازی (دسترسی در: <http://en.wikipedia.org/wiki/Cave-automatic-virtual> -

(environment

Facebook -2

Gazan -3

شدنی نیستند، و مشاهده محتوا غیر ممکن می شود. بیشتر پیشینه های پیوسته در طول دوره ما سرانجام همچون کارت پانچ های دهه 1960 و نوارهای مغناطیسی ریل به ریل (1)، غیرقابل خواندن (استفاده) خواهند شد. علاوه بر مشکل قالب، مشکل حجم زیاد نیز در حال رشد و فزونی است. همان طور که اینترنت به سمت تکامل استفاده از IPv6 (2) حرکت می کند، اشیای قابل نشانه گذاری (به طور روزافزونی شامل، اشیای فیزیکی در «(3) Web of Things» به راستی، به وسیله دستورات متعدد ذخیره سازی - حتی نشانی ها - بیش از ظرفیت موجودمان، در حال ازدیاد حجم هستند (1038)، که به تنهایی اجازه آماده سازی همه محتوا را می دهد. در نتیجه، ما نمی توانیم موارد بیشتری را جست و جو کنیم، و به این دلیل نمایه سازی و فناوری جست و جوی ما ناامیدانه شکست می خورد.

همزمان با توسعه وب معنایی، همه مفهوم محتوا تغییر می کند. محتوا چیزی بیشتر از متن و تصویر نیست، اما در حال حاضر اقلام داده قراردادی و پیوندهای بین آن ها را در بر می گیرد. حتی حالا در سال 2011، جهان داده های پیوند خورده عمومی (4) (که شامل مجموعه هایی همچون data.gov.uk و data.gov است). ده ها میلیون اقلام داده ای (به طور روز افزون قالب (5) RDF) را در بر می گیرد که به وسیله صدها میلیون پیوند قابل ارجاع مجدد، به هم پیوسته اند. چالش آرشیوسازی این سری داده ها از زمان نشانی دهی - آغاز می شود (برای مثال: به وسیله پروژه پرونوم (6) آزمایشگاه های آرشیوهای ملی بریتانیا (7))، اما احتمالاً جهان داده های پیوند خورده سریع تر از اینکه از نظر مجموعه سازی یا تحلیل، قابل مدیریت باشد، رشد خواهند کرد.

در این سناریو، حتی منابع عظیم شرکتی چون گوگل تحت الشعاع مشکلات قرار می گیرد. بنابراین، پاسخ بدیهی برای آرشیوسازی این حجم («بگذارید گوگل کارش را بکند» (8)) چیزی بیشتر از یک راه حل نیست.

اگر انتخاب ها ما را به این مسیر هدایت کند، پژوهشگران آینده می آموزند که گذشته وب را به صورت غیر قابل دسترسی و بدون اعتماد و داستانی که حکایت می شود و مدرک دست دومی از زمان، تصور کنند. امروزه، حجم بزرگی از اطلاعات به مقیاس جهانی ایجاد می شود که ممکن است به صورت نوشته های روی تکه های کاغذ در میلیون ها جعبه کفش ذخیره شده باشد، که همگان، پیشرفت های صورت گرفته در دنیا را همان گونه که در محتوای اینترنت منعکس می شود به نحو مناسبی درک خواهند کرد.

ص: 113

1- reel to reel

2- IPv6 آخرین نسخه پروتکل ارتباطات در اینترنت (دسترسی در <http://en.wikipedia.org/wiki/IPv6>)

3- Web of Things (http://en.wikipedia.org/wiki/Web_of_Things)

4- Linked Data - Connect Distributed Data across the Web, at www.linkeddata.org

5- (Resource description Framework) RDF (چارچوب توصیف منابع)

6- PRONOM

7- <http://www.nationalarchives.gov.uk/PRONOM/Default.aspx>

8- let Google do it

در دنیایی کاملاً متغیر، افراطی ترین سناریو سناریوی انفرادی است که در آن اینترنت همانگونه که ما می شناسیم به سمت پدیده ای کاملاً جدید و احتمالاً در نوع خود هوشمند (کورزویل 2005) (1) تکامل می یابد. در مسیر منحصر به فرد شدن، به سمت توسعه و تبدیل شدن به موجود مجازی پیچیده ای حرکت می کند، که ما فهم کمی از آن داریم و برای آرشو کردن آن، در حال حاضر روشی بهتر از آن چه که با هوش انسانی خود از آرشو می دانیم، وجود ندارد. حتی اکنون، در سال 2011، برای غلبه بر آن آغاز به تشخیص تمایز بین پردازش انسانی و مصنوعی کرده ایم. خدماتی همچون: ریکاچا (2) (ون آون (3)، مورر (4)، مک میلن (5)، آبراهام (6) و بلام (7) (2008) و آمازونز مکانیکال ترک (8) که از موجودیت های انسانی «در یک چرخه» برای حل مشکلاتی که برای ماشین ها سخت است، استفاده می کنند، به ما راه هایی را به سوی دنیایی که در آن هوش انسان و ماشین به صورت دوقلوهای جدایی ناپذیری هستند، نشان می دهند و مرز بین آن ها نامشخص است. در چنین دنیایی، مشخص نیست که آرشوسازی چه مفهومی می تواند داشته باشد. بنابراین، همزمان با حرکت زمان به سمت جلو، گذشته به ناچار و به صورت جبران ناپذیری از دست می رود. این سناریو، ممکن است شبیه به یک داستان علمی به نظر برسد. اگر چه، ارزشمند است که به خاطر داشته باشیم آینده غیر قابل پیش بینی است، حتی اگر تأثیر بر روی شرایط انتخاب در طول مسیر را برای اینکه در یک مسیر یا مسیر دیگری قرار گیرد - مدیریت نماییم. انتخاب هایی که ما انجام می دهیم ممکن است، برای اجرای وظیفه ای جدید، همچون اینترنت هوشمند ناکافی باشد. این بدین معنا نیست که ما باید انتخاب های صحیح را در نظر بگیریم و تلاشی از خود به خرج ندهیم.

سناریوی غبار آلود

متأسفانه این سناریو، احتمالاً در ابتدای امر [این اندیشه را به ذهن متبادر می کند] که آرشوهای وب معادل آرشو رقومی غبار آلود است، اگر چه اغلب به خوبی حفاظت و نگهداری می گردد. به سختی از آن استفاده می شود.

حتی اگر جامعه آرشوسازی، وب، به توسعه استانداردها و تجربیات برای حفظ بخشهای اینترنت ادامه دهد، کاربردهای مؤثر ناچیزی از جامعه پژوهشی سر می زند و جامعه پژوهشی به طور مؤثر از آن استفاده نمی کند.

ممکن است به صفحه های اینترنت به صورت مجزا با استفاده از ابزارهای پیوسته رجوع شود و بعضی از پژوهشگران به ایجاد آرشوهای کوچک برای موضوعات پژوهشی خاص مبادرت ورزند، اما پژوهش اینترنت تنها با اولویت تمرکز بر روی وب پویا ادامه خواهد داشت و در آینده نزدیک توجه اندکی به

ص: 114

Kurzweil -1

2 - reCAPTCHA: برنامه ای برای رقومی کردن کتب، نشانیات و ... دسترس می در

(<http://www.google.com/recaptcha/learnmore>)

Von Ahn -3

Maurer -4

Mc Millen -5

Abraham -6

Blum -7

(Amazon's Mechanical Turk(Amazon Corp., Mechanical Turk at www.mturk.com -8

استفاده از وب قبلی برای پژوهش جدی، گسترش خواهد یافت.

این سناریو با سناریوی آخر الزمان متفاوت است. در آن سناریو، فناوری آرشیوسازی وب با تغییرات فناورانه روی اینترنت همگام نخواهد بود. در این سناریو سرعت آرشیوسازی وب هماهنگ با فناوری تحویل وب تنظیم می شود.

اگر چه داده حفاظت شده به صورت اسطوره ای حفاظت شده برای استفاده در آینده نامطمئن باقی می ماند.

در مراحل نگارش این گزارش، مشخص شده است که به جای آرشیوهای وب مرجع، به طور روزافزونی کاربران و پژوهشگران با وب پویا همانند آرشیو برخورد می کنند. وب پویا، به رشد خود ادامه می دهد و برای بیشتر بخش ها، داده هایی که ناپدید می شوند، از نظر بسیاری به صورت یک مشکل ساده دیده می شود و مهم تر از آن، برای بخشی بزرگ تر، حجم عظیم دیگری از داده ها بر روی وب در هر زمان باقی می ماند.

تصویر ما از آرشیو، تصویری از اقلام فیزیکی همچون کاغذها و اسناد ذخیره شده در یک مکان فیزیکی است. با وجود این، این ماهیت وب نیست. پژوهشگران می توانند حجم زیادی از مواد را از منابع مختلف وب پویا برای نیل به اهداف پژوهشی خود ذخیره کنند. ما آرشیوها را به صورت چیزهایی که برای آیندگان حبس شده اند، درک می کنیم. در حالی که وب، به تنهایی پدیده ای در حال رشد و منبعی متنوع از انواع مواد است که به صورت بالقوه مورد توجه پژوهشگران است. به طوری که، آن را نه به عنوان یک آرشیو سنتی، بلکه به سادگی به صورت یک منبع داده می بینند.

این یک سناریوی بدبینانه است، اما به نظر می رسد که وزنی از شواهد را در جانب خود داشته باشد. در راینی هایی با محققان برجسته، به یک فقدان علاقه دیرپای بر پرسش از وب قبلی و درک اینترنت به عنوان یک توسعه تاریخی، رسیدیم. البته استثناهایی وجود دارد که بعداً در این گزارش شرح داده خواهد شد، اما قادر بوده ایم تا گرایش نهفته کار با آرشیوهای وب را کشف کنیم که انتظار ساده ای برای [ظهور] فناوری مناسب برای بیدار شدن آن وجود دارد.

ممکن است انتظار تغییر از گوشه و کنار برای آمادگی ایجاد یک گام تغییری در تصورات پژوهشگران، بر اساس جلوه های استفاده جدید یا فناوری تازه، وجود داشته باشد. اگر این امر رخ ندهد، بیم آن می رود که به جمع آوری گردوغبارهای رقومی ادامه دهیم.

اگر از این سناریو اجتناب شود، نیازمند نوعی جدیدی از آرشیویست هستیم که با پژوهشگران و عموم [کاربران]، در استخراج داده هایی که آن ها از وب پویا نیاز دارند در تعامل باشد و زمانی که این داده ها از وب پویا ناپدید شدند، قادر به ذخیره سازی مجدد آن ها به صورتی که قابل مشاهده و استفاده بوسیله ابزارهای وب پویا باشند، هستند. بسیار فراتر، از زمانی که رقومی سازی اسناد تاریخی، از آرشیوها حجم زیادی از مواد تاریخی قابل دسترسی بر روی وب در دهه گذشته ایجاد کرده است (مهیر 2001 (1))،، تانر 2010 (2))، تانر و دیگران 2011 (3))، آرشیوهای وب به انتقال وب به محفظه ها (جعبه ها) نیاز ندارند، بلکه در عوض نیازمند برگشت به وب در زمانی که دارای محتواست، هستند.

ص: 115

Tanner -2

Deegan -3

در حین این که ما مسیر خود را به سوی آینده آرشیو وب مرور می کنیم، پرسش های متعددی در مورد چگونگی دیدگاه مان نسبت به آینده، می توانیم از خود پرسیم.

برای مثال، آیا آرشیو آینده، باغ دیوار کشیده شده ای خواهد بود که از آسیب ها محفوظ و در امان است، اما با دسترسی محدود است؟ آیا فضایی کاملاً باز و قابل دسترسی برای تازه واردان خواهد بود؟ آیا مجموعه ای از مخازن است یا یک عرصه ارتباط متقابل منفرد؟ یا اینکه یک شهر ارواح خالی از سکنه باز و مرتبط خواهد بود؟

بخشی از آرشیو ها برای محققان و دانشمندان و بخشی دیگر برای عموم است (1). آیا این دو می توانند در فضایی سایبرنتیک با هم مرتبط باشند - تا اینکه دانشمندان به صورت دائمی، آن چه که عموم [شامل دانشمندان] به آن دسترسی دارند و از آن استفاده می کنند و آن را ایجاد می کنند، پایش نمایند (بنابراین درک از آگاهی جهانی یا دریافت جمعی را افزایش می دهند) در حالی که در زمان مشابه شکل گیری چنین فضایی با این سبک، برای گسترش، مجموعه سازی و تأثیرگذاری و لذت بخش بودن و کاربرد غنی همه دنیا، غنی سازی می شوند؟

حتی در صورتی که منحصر به فرد بودن، با شکست مواجه شود، اینترنت به طور فزاینده ای اجازه ارتباطات در این سبک را که از نظر استفاده برای هوش جهانی خیلی کم قابل فهم می باشد، می دهد (شرودر و مهیر، 2009). در هوش جهانی - با خوراک ها و پیوندهای پویا - هوش های انفرادی از طریق ابزارهای ورودی و خروجی به یکدیگر مرتبط هستند. چگونه آرشیوهای وب ماهیت ارتباطات درونی هوش جهانی را به صورت غیر از اسناد به ظاهر جدا از هم منعکس می نماید؟

اینها و بسیاری از سؤالات دیگر ما را با حرکتی رو به جلو مواجه می کند. در بخش های آینده این سند، ما نگاهی خواهیم داشت به بعضی از فنون فهم وب پویا که می تواند الهام بخش جامعه آرشیوسازی وب باشد. سپس به چالش های حرکت رو به جلویی که وب آرشیوی را به صورت بالقوه برای تحقیق، ارزشمندتر می کند، اشاره خواهیم نمود.

یادگیری از وب پویا

باید پرسید که چرا در دنیای پژوهشی اینترنت، آرشیوهای وب به صورت شهروندان درجه دوم به نظر می رسند؟ در مقایسه با کسانی که در وب پویا مطالعه می کنند، به مراتب تعداد کمتری از پژوهشگران از آرشیوهای وب استفاده می کنند و تعداد کمی از ابزارهای غیر تخصصی برای آرشیوهای وب در حال ساخت هستند، به ویژه در مقایسه با ابزارهایی که برای مطالعه وب پویا ساخته می شوند.

چالش عمومی برجسته در این بخش این است که، جامعه آرشیوسازی وب به ایجاد ارتباط بین منابعی که تولید می کنند با ابزارهای لبه برش (2) نیاز دارند، که توسط متخصصان رایانه، پژوهشگران، توسعه

1- ما در این جا با کاربرد های تجاری و دولتی آرشیوی وب اسناد که غالباً برای برآوردن نیازهای حقوقی طراحی می شوند، سر و کار نداریم، همچنان که این مورد فراتر از حیطه اختیار و تخصص ماست.

2- (نام نوعی فناوری) cutting edge

دهندگان مستقل، و هکرها برای مطالعه وب پویا گسترش یافته اند.

در حال حاضر، انواع ابزارهای توسعه یافته برای مطالعه وب پویا، روی هم رفته، به آسانی برای مطالعه داده های قابل دسترسی آرشیوهای وب به کار گرفته می شوند. این [امر] بزرگ ترین مانع برای فهم وب، نه تنها به عنوان یک تصویر لحظه ای بلکه به عنوان یک اکوسیستم در حال توسعه است.

مصورسازی

هر آرشیوی که ایجاد شده و مورد استفاده قرار می گیرد به [مرور] وسیع، غیر قابل نظر اجمالی و بدون نقشه یا دسترسی بصری و روش حسی برای مشاهده آرشیوها و چگونگی پیوندهای آن ها، خواهد شد. مصورسازی در اینجا یک راه حل اساسی است، اما انواع متعددی از آن ها در دسترس است.

چالش ها

ابزارهای بسیاری برای بررسی ارتباط بین کاربران رسانه های اجتماعی، ابزارهایی برای مرور پیوستار زمانی، ترکیب نقشه ها با ردپای کاربران، ابزارهای تصویری برای نمایش نحوه پیوند داده ها و نظیر آن وجود دارد؛ اما اینها به منظور کار کردن با آرشیوهای وب نیاز به متمرکز شدن دارند. مصور سازی اطلاعات حوزه پیشرفته پژوهشی است، اما بهترین نحوه مشاهده یک آرشیو (یا جست و جوی یکی از آن ها، یا دیدن ارتباطات درونی و ما بین آن ها) شامل رابط های حسی (بصری)، تغییرات دیداری تصاویر سه بُعدی و پویاگران زمانی، هنوز دست نیافتنی هستند؛ و برای مثال آیا احتمال شناسایی موضوع هایی درون مجموعه به وسیله روش های بازبینی تصویری یا سازماندهی تصویری یک سری از ارتباطات وجود دارد؟ به بیان دیگر، ایجاد ابزارهای مصورسازی برای کمک به پژوهشگران ممکن است؟

مثال ها: [\(2\) Touch graph](#)، [\(1\) Apple Time Machine](#)

عکس

دهندگان مستقل، و هکرها برای مطالعه وب پویا گسترش یافته‌اند.

در حال حاضر، انواع ابزارهای توسعه یافته برای مطالعه وب پویا، روی هم رفته، به آسانی برای مطالعه داده‌های قابل دسترسی آرشیوهای وب به کار گرفته می‌شوند. این [امر] بزرگ‌ترین مانع برای فهم وب، نه تنها به عنوان یک تصویر لحظه‌ای بلکه به عنوان یک اکوسیستم در حال توسعه است.

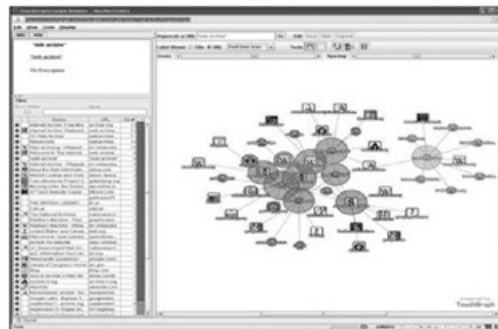
مصورسازی

هر آرشیوی که ایجاد شده و مورد استفاده قرار می‌گیرد به [مرور] وسیع، غیرقابل نظر اجمالی و بدون نقشه یا دسترسی بصری و روش حسی برای مشاهده آرشیوها و چگونگی پیوندهای آنها، خواهد شد. مصورسازی در اینجا یک راه‌حل اساسی است، اما انواع متعددی از آنها در دسترس است.

چالش‌ها

ابزارهای بسیاری برای بررسی ارتباط بین کاربران رسانه‌های اجتماعی، ابزارهایی برای مرور پیوستار زمانی، ترکیب نقشه‌ها با ردپای کاربران، ابزارهای تصویری برای نمایش نحوه پیوند داده‌ها و نظیر آن وجود دارد؛ اما اینها به منظور کار کردن با آرشیوهای وب نیاز به متمرکز شدن دارند. مصورسازی اطلاعات حوزه پیشرفته پژوهشی است، اما بهترین نحوه مشاهده یک آرشیو (یا جست‌وجوی یکی از آنها، یا دیدن ارتباطات درونی و ما بین آنها) شامل رابط‌های حسی (بصری)، تغییرات دیداری، تصاویر سه بُعدی و پوششگران زمانی، هنوز دست نیافتنی هستند؛ و برای مثال، آیا احتمال شناسایی موضوع‌هایی درون مجموعه به وسیله روش‌های بازبینی تصویری یا سازماندهی تصویری یک سری از ارتباطات وجود دارد؟ به بیان دیگر، ایجاد ابزارهای مصورسازی برای کمک به پژوهشگران ممکن است؟

مثال‌ها: 'Touch graph'، 'Apple Time Machine'



تصویر ۱. Touchgraph، در اینجا، توانایی بررسی پیوندهای میان وبگاه‌ها با استفاده از رابط گرافیکی را ارائه می‌کند. داده‌ها از وب پویا ترسیم شده‌اند.

1. Touchgraph(<http://www.touchgraph.com/>)
2. Apple Time Machine (<http://www.apple.com/macosx/what-is-macosx/time-machine.html>)

تصویر ۱. Touchgraph، در اینجا، توانایی بررسی پیوندهای میان وب گاه‌ها با استفاده از رابط گرافیکی را ارائه می‌کند. داده‌ها از وب پویا ترسیم شده‌اند.

برنامه های کاربردی جست و جو همانند شکارچی

همان طوری که اطلاعات اینترنت به تکثیر و افزایش، هم در حجم و هم در تنوع انواع محتوا ادامه می دهند، جست و جوهای به مراتب پیچیده تری برای قابلیت استخراج هر چیز با معنا و استفاده از این مجموعه عظیم مورد نیاز خواهد بود. جست و جو مانند جست و جوی تصویر و ویدئو به سوی تکلیفی پیچیده تر سوق می یابد.

چالش: ایجاد سطح بلند پروازانه تری از برنامه ها با هزینه ای تقریباً ناچیز، به ویژه ابزارهایی که برای مجموعه ها می توانند به کار گرفته شوند. این، امر ممکن است نیازمند طراحانی برای استفاده جسورانه از موتورهای جست و جوی مبتنی بر فناوری ابر و زبان های جست و جوی بهتر با قابلیت انعطاف در پرسش سؤال های پیچیده از داده های موجود در آرشیوهای وب، باشد.

مثال ها: (در حال حاضر آپاچی) یا هو! طرح زیربنایی [1](http://pig.apache.org/) (PIG Latin) برای پشتیبانی موردی تحلیل سری داده های خیلی بزرگ (تصویر 2)

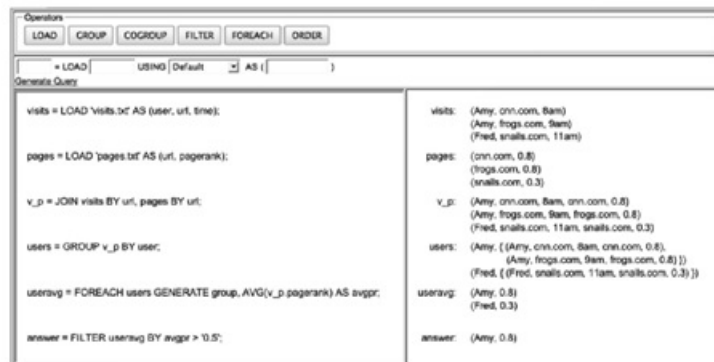
عکس

برنامه‌های کاربردی جست‌وجو همانند شکارچی

همان‌طوری که اطلاعات اینترنت به تکثیر و افزایش، هم در حجم و هم در تنوع انواع محتوا ادامه می‌دهند، جست‌وجوهای به مراتب پیچیده‌تری برای قابلیت استخراج هر چیز با معنا و استفاده از این مجموعه عظیم مورد نیاز خواهد بود. جست‌وجو مانند جست‌وجوی تصویر و ویدئو به سوی تکلیفی پیچیده‌تر سوق می‌یابد.

چالش: ایجاد سطح بلند پروازانه تری از برنامه‌ها با هزینه‌ای تقریباً ناچیز، به‌ویژه ابزارهایی که برای مجموعه‌ها می‌توانند به‌کار گرفته شوند. این امر، ممکن است نیازمند طراحی برای استفاده جسورانه از موتورهای جست‌وجوی مبتنی بر فناوری ابر و زبان‌های جست‌وجوی بهتر با قابلیت انعطاف در پرسش سؤال‌های پیچیده از داده‌های موجود در آرشیوهای وب، باشد.

مثال‌ها: (در حال حاضر آپاچی) یاهو! طرح زیربنایی^۱ PIG Latin برای پشتیبانی موردی تحلیل سری داده‌های خیلی بزرگ (تصویر ۲)



تصویر ۲. نمایی از Pi6 Pen نشان دهنده برنامه‌ای است که افراد متقاضی مشاهده صفحه‌های با رتبه‌بندی بالا را می‌یابد. (اولستون^۲، رید^۳، سریواستاوا^۴، کومر^۵ و تامکینز^۶، ۲۰۰۸)

تحلیل‌های شبکه اجتماعی

تحلیل شبکه اجتماعی (SNA) حوزه قابل توجهی از علاقه و فعالیت پژوهشی در میان پژوهشگران اینترنت، جامعه‌شناسان، فیزیکدانان و بسیاری دیگر است. تنوع موضوع‌های مورد توجه بسیار گسترده

1. PIG Latin(<http://pig.apache.org/>)
2. Olston
3. Reed
4. Srivastava
5. Kummer
6. Tomkins

تصویر ۲. نمایی از Pi6 Pen نشان دهنده برنامه‌ای است که افراد متقاضی مشاهده صفحه‌های با رتبه‌بندی بالا را می‌یابد. (اولستون^۲، رید^۳، سریواستاوا^۴، کومر^۵ و تامکینز^۶، ۲۰۰۸)

تحلیل‌های شبکه اجتماعی

تحلیل شبکه اجتماعی (SNA) حوزه قابل توجهی از علاقه و فعالیت پژوهشی در میان پژوهش‌گران اینترنت، جامعه‌شناسان، فیزیک دان

ها و بسیاری دیگر است. تنوع موضوع های مورد توجه بسیار گسترده

ص: 118

PIG Latin -1

Olston -2

Reed -3

Srivastava -4

Kummer -5

Tomkins -6

است که شامل فهم ارتباطات بین دوستان در شبکه های اجتماعی مانند فیس بوک (هوگان 2010 (1)) ، بررسی وابستگی های سیاسی مشارکت کنندگان در مباحث سیاسی (هیندمن (2)) ، (2007) ، و کشف شبکه های توطئه (طرح) در ادبیات انگلیسی (مورتنی 2011 ، (3) 2005) است.

ابزارهای جست و جوی مبتنی بر تحلیل های شبکه اجتماعی شامل (7) NodeXL (6) Voson ، (5) Pajec ، (4) UCINET ، و بسیاری دیگر، در حال تکثیر و توسعه هستند.

اگر چه تعداد کمی از این ابزارها، هر چه که باشد، برای استفاده به وسیله آرشيوهای وب توانمند یا بهینه شده اند.

چالش: نخست، کار با طراحان اصلی ابزارهای تحلیل شبکه اجتماعی برای توانمند و بهینه سازی آن ها در کار با داده های آرشیوی وب.

همچنین، توسعه روش های جدید ابداعی با احتمال یک بار در بعد زمان، که به داده های شبکه برای ردیابی مواردی همچون تکامل تدریجی شبکه های اجتماعی در طول زمان، افزوده می شود. به وسیله آرشيوسازی، نه تنها وضعیت سایت های شبکه های اجتماعی، بلکه زمانی که افراد پیوندها را ایجاد، نگهداری و حذف نموده، با دیگران ارتباط برقرار کرده، به گروه ها می پیوندند و گروه ها یا وبگاه ها را ترک می کنند.

ما باید به خاطر داشته باشیم که وب شبکه ای از پیوندهاست و تحلیل وب بینشی، کلی به ماهیت آن شبکه برای ما فراهم می کند.

مثال ها:

تحلیل های فیس بوک: ابزارهای بسیاری برای تحلیل تعاملات بین کاربران فیس بوک، جریان نفوذ و منحنی اجتماعی در دسترس هستند.

نمایش مثالی از فیس بوک (8):

عکس

است که شامل فهم ارتباطات بین دوستان در شبکه‌های اجتماعی مانند فیس بوک (هوغان^۱، ۲۰۱۰)، بررسی وابستگی‌های سیاسی مشارکت کنندگان در مباحث سیاسی (هیندمن^۲، ۲۰۰۷)، و کشف شبکه‌های توطئه (طرح) در ادبیات انگلیسی (مورتی^۳، ۲۰۱۱، ۲۰۰۵) است.

ابزارهای جست‌وجوی مبتنی بر تحلیل‌های شبکه اجتماعی شامل^۴ UCINET،^۵ Pajec،^۶ Voson،^۷ و NodeXL^۸، و بسیاری دیگر، در حال تکثیر و توسعه هستند. اگر چه تعداد کمی از این ابزارها، هر چه که باشد، برای استفاده به‌وسیله آرشیوهای وب توانمند یا بهینه شده‌اند.

چالش: نخست، کار با طراحان اصلی ابزارهای تحلیل شبکه اجتماعی برای توانمند و بهینه‌سازی آنها در کار با داده‌های آرشیوی وب.

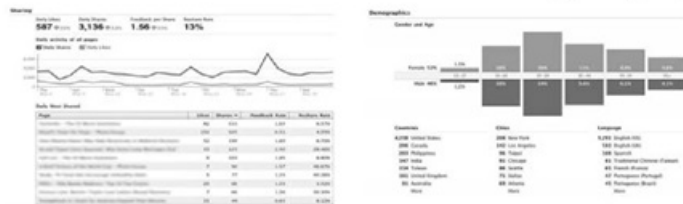
همچنین، توسعه روش‌های جدید ابداعی با احتمال یک‌بار در بعد زمان، که به داده‌های شبکه برای ردیابی مواردی همچون تکامل تدریجی شبکه‌های اجتماعی در طول زمان، افزوده می‌شود، به وسیله آرشیوسازی، نه تنها وضعیت سایت‌های شبکه‌های اجتماعی، بلکه زمانی که افراد پیوندها را ایجاد، نگهداری و حذف نموده، با دیگران ارتباط برقرار کرده، به گروه‌ها می‌پیوندند و گروه‌ها یا وبگاه‌ها را ترک می‌کنند.

ما باید به خاطر داشته باشیم که وب شبکه‌ای از پیوندهاست و تحلیل وب بینشی، کلی به ماهیت آن شبکه برای ما فراهم می‌کند.

مثال‌ها:

تحلیل‌های فیس بوک: ابزارهای بسیاری برای تحلیل تعاملات بین کاربران فیس بوک، جریان نفوذ و منحنی اجتماعی در دسترس هستند.

نمایش مثالی از فیس بوک^۸:



تصاویر نمودارهای اجتماعی فیس بوک:

1. Hogan
2. Hindman
3. Moretti
4. <http://www.analytictech.com/ucinet/>
5. <http://pajec.imfm.si/doku.php>
6. <http://voson.anu.edu.au/>
7. <http://nodexl.codeplex.com/>
8. <http://www.facebook.com/insights/>

تصاویر نمودارهای اجتماعی فیس بوک:

[/http://www.analytictech.com/ucinet](http://www.analytictech.com/ucinet) -4

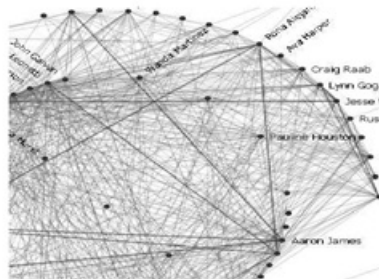
<http://pajek.imfm.si/doku.php> -5

[/http://voson.anu.edu.au](http://voson.anu.edu.au) -6

[/http://nodexl.codeplex.com](http://nodexl.codeplex.com) -7

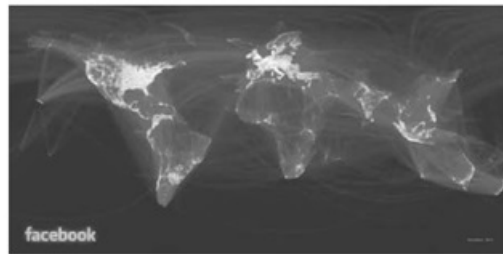
<http://www.facebook.com/insights/> -8

۱۲۰ مدیریت منابع اطلاعاتی وب



تصویر ۳. منبع:

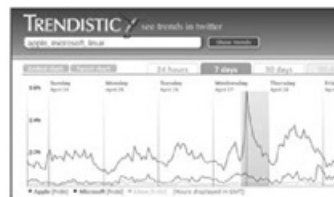
http://infosthetics.com/archives/2008/03/facebook_social_network_graph.html



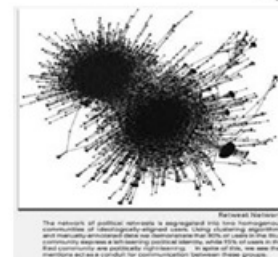
تصویر ۴. منبع:

<http://www.facebook.com/notes/facebook-engineering/visualizing-friendships/469716398919>

به همین ترتیب، طیف گسترده‌ای از تحلیل‌های توییت‌ر مانند Twitalyzer^۱، Trendistic^۲، و دیگران وجود دارد.



تصویر ۵. trndistic.



تصویر ۶. Truthy. (<http://truthy.indiana.edu/>)

1. <http://www.twitalyzer.com/>
2. <http://trendistic.com/>

تصویر ۳. منبع: http://infosthetics.com/archives/2008/03/facebook_social_network_graph.html

تصویر ۴. منبع: <http://www.facebook.com/notes/facebook-engineering/visualizing-friendships/469716398919>

به همین ترتیب، طیف گسترده ای از تحلیل های توییتز مانند [\(2\) Twitalyzer \(1\)](#)، [Trendistic](#)، و دیگران وجود دارد.

تصویر 5. trndistic.

تصویر 6. Truthy (<http://truthy.indiana.edu/>)

ص: 120

1- [/http://www.twitalyzer.com](http://www.twitalyzer.com)

2- [/http://trendistic.com](http://trendistic.com)

آینده آرشیو وب ۱۲۱

Far Left	Moderate Left	Center	Moderate Right	Far Right
#healthcare	#arp #women	#democrats #social	#rangel #vaste	#912project #twisters
#judaism #hollywood	#citizensunited	#seniors #dnc	#saveamerica	#gop2112 #israel
#2010elections	#democratic	#budget #political	#american #gold	#foxnews #mediabias
#capitalism #recession	#banksters #energy	#goproud #christian	#repeal #mexico	#constitution
#security #dceasact	#sarahpalin	#media #nobel	#terrorian #gopleader	#patriots #cednov
#publication	#progressives		#palin12	#abortion
#topprogs	#stopbeck #iraq			

تصویر ۸. برچسب‌های نشانه‌گذاری شده به وسیله کاربران حوزه سیاست (کونور^۱، راتکوویچ^۲، فرانسیسکو^۳، و دیگران: ۲۰۱۱)



تصویر ۷. توییت کردن و سیاست (کونور، راتکوویچ، گونکالوز، فلامینی و مننزر، ۲۰۱۱)

سنجش‌های دگرساز

اصطلاحی در حال پدیدار شدن برای راه‌های جدید اندازه‌گیری تأثیر علمی، فراتر از معیارهای کتابسنجی، وب‌سنجی و علم‌سنجی است.

ارتباطات بین ما و بین دانشمندان و درون گروه‌های علمی به‌طور روزافزونی بر روی وب در حال شکل‌گیری است. جامعه در حال ظهوری از پژوهشگران که مطالعه پژوهشی انجام می‌دهند، از توالی‌ها و پیوندهایی که به وسیله ابزارهایی چون توییت، مندلی^۷، فرند فید و بسیاری دیگر بر جای گذاشته شده است، برای درکی که بسیار سریع‌تر از تأثیرات سنتی می‌توانند توسعه یابند، استفاده می‌شوند.

فراتر از چیزی که می‌توان به‌عنوان سنجش دگرساز تصور کرد: چگونگی ردیابی مشارکت‌های غیرآکادمیک به سمت دانش است.

آیا می‌توان ابزارهای تحقیقاتی پژوهشگران را برای سایر حوزه‌های غیر تخصصی به‌کار برد؟ برای مثال، آیا می‌توانیم تأثیر مشارکت کنندگان انفرادی بر گروه‌های سرگرمی در طول زمان را با استفاده از معیارهای مشابه برای درک چگونگی تحول تأثیر پژوهشگر، ارزیابی کنیم؟

چالش: فعال‌سازی روش‌هایی به مراتب ساده‌تر برای مشخص کردن طیف زمانی مواد رقومی، بنابراین، تحلیل سنجش دگرساز می‌تواند به روش‌های کاملاً مشابه کتابسنجی انجام شود: در الگوی

1. Conover
2. Ratkiewicz
3. Francisco
4. Goncalves
5. Flammim
6. Menczer
7. Mendeley

تصویر ۸. برچسب‌های نشانه‌گذاری شده به وسیله کاربران حوزه سیاست (کونور^(۱)، راتکوویچ^(۲)، فرانسیسکو^(۳)، و دیگران؛

(2011)

اصطلاحی در حال پدیدار شدن برای راه های جدید اندازه گیری تأثیر علمی، فراتر از معیارهای کتاب سنجی، وب سنجی و علم سنجی است.

ارتباطات بین ما و بین دانشمندان و درون گروه های علمی به طور روزافزونی بر روی وب در حال شکل گیری است. جامعه در حال ظهوری از پژوهشگران که مطالعه پژوهشی انجام می دهند، از توالی ها و پیوندهایی که به وسیله ابزارهایی چون توئیتر مندلی (4)، فرند فید و بسیاری دیگر بر جای گذاشته شده است، برای درکی که بسیار سریع تر از تأثیرات سنتی می توانند توسعه یابند، استفاده می شوند.

فراتر از چیزی که می توان به عنوان سنجش دگر ساز تصور کرد: چگونگی ردیابی مشارکت های غیرآکادمیک به سمت دانش است.

آیا می توان ابزارهای تحقیقاتی پژوهشگران را برای سایر حوزه های غیر تخصصی به کار برد؟ برای مثال، آیا می توانیم تأثیر مشارکت کنندگان انفرادی بر گروه های سرگرمی در طول زمان را با استفاده از معیارهای مشابه برای درک چگونگی تحول تأثیر پژوهش گر، ارزیابی کنیم؟

چالش: فعال سازی روش هایی به مراتب ساده تر برای مشخص کردن طیف زمانی مواد رقومی، بنابراین، تحلیل سنجش دگر ساز می تواند به روش های کاملاً مشابه کتاب سنجی انجام شود: در الگوی

ص: 121

Conover -1

Ratkiewicz -2

Francisco -3

Mendeley -4

انتشارات رسمی، هر نشریه ای یک نویسنده و تاریخ نشر دارد و برای ردیابی توالی استنادها به یک اثر انفرادی استفاده می شود. نشریات، نظام آرشیوی دارند که آزمایش شده و به نحو قابل اعتماد و مناسبی پذیرفته شده اند: نشریه تخصص با وجود این، نشریات غیر رسمی وب، روش توسعه یافته مناسب و مشابهی برای آرشیو سازی سهم ها (مشارکت ها) با دانش، به شکلی که قابل استناد باشند و در طی زمان دوباره جایگزین شوند، ندارند. این شکافی است که انتظار می رود رفع شود و آرشیوهای وب زمینه های بارزی برای شروع هستند.

همچنین، برای مشارکت های غیر تخصصی چه ابزارهایی برای تحلیل ویکی ها و موجودیت های مشترک استفاده می شود که می تواند فهم این تغییرها را در دوره های زمانی توسعه دهد؟

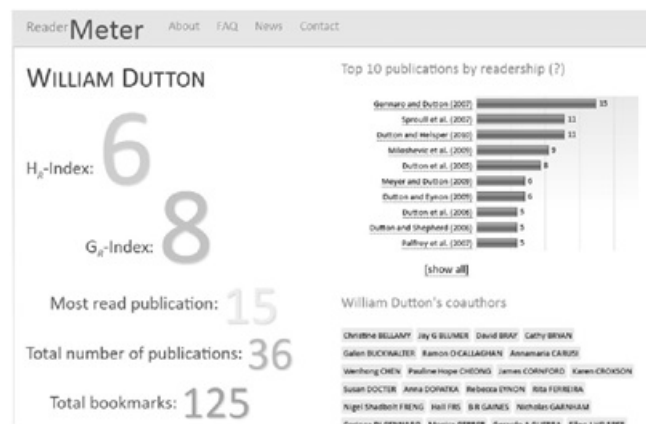
مثال ها: [\(1\) Data Cite](#) و [\(2\) Reader Meter](#)

عکس

انتشارات رسمی، هر نشریه‌ای یک نویسنده و تاریخ نشر دارد و برای ردیابی توالی استنادها به یک اثر انفرادی استفاده می‌شود. نشریات، نظام آرشیوی دارند که آزمایش شده و به نحو قابل اعتماد و مناسبی پذیرفته شده‌اند: نشریه تخصصی با وجود این، نشریات غیر رسمی وب، روش توسعه یافته مناسب و مشابهی برای آرشیو سازی سهم‌ها (مشارکت‌ها) با دانش، به شکلی که قابل استناد باشند و در طی زمان دوباره جایگزین شوند، ندارند. این شکافی است که انتظار می‌رود رفع شود و آرشیوهای وب زمینه‌های بارزی برای شروع هستند.

همچنین، برای مشارکت‌های غیر تخصصی، چه ابزارهایی برای تحلیل ویکی‌ها و موجودیت‌های مشترک استفاده می‌شود که می‌تواند فهم این تغییرها را در دوره‌های زمانی توسعه دهد؟

مثال‌ها: ¹Reader Meter و ²Data Cite



تصویر ۲. reader meter نوعی از ابزار Alt metric برای درک نحوه خواندن یک نویسنده، بر اساس

آمارهای مندی (<http://www.mendeley.com/>) منبع: <http://readermeter.org>

وب‌نوشت (حاشیه نگاری) اجتماعی

کاربران مایل‌اند که بتوانند پیوندها و برجسب‌هایشان^۳ و همچنین نظرها و حاشیه‌نویسی‌هایشان را بر روی منابع نشر دهند. برای پژوهشگران، فهم چگونگی توسعه این جوامع در طول زمان و نگهداری خودشان سؤال مهمی است. برای مثال Reddit بیش از ۸ میلیون خواننده انحصاری و یک میلیون بازدید صفحه در هر ماه دارد (جسرا، ۲۰۱۱).

چالش‌ها: گستره‌ای راکه آرشیوها می‌توانند نه فقط وبگاه‌ها و مجموعه‌هایشان، بلکه پیوندها و

1. ReaderMeter
2. Data Cite
3. bookmark

تصویر 2. reader meter نوعی از ابزار Alt metric برای درک نحوه خواندن یک نویسنده بر اساس آمارهای مندی (<http://www.mendeley.com/>) منبع: <http://readermeter.org>

وب نوشت (حاشیه نگاری) اجتماعی

کاربران مایل‌اند که بتوانند پیوند ها و برجسب های شان (3) و همچنین نظر ها و حاشیه نویسی های شان را بر روی منابع نشر دهند. برای

پژوهشگران فهم چگونگی توسعه این جوامع در طول زمان و نگهداری خودشان سؤال مهمی است. برای مثال Reddit بیش از 8 میلیون خواننده انحصاری و یک میلیون بازدید صفحه در هر ماه دارد (جسرا، 2011).

چالش ها: گستره ای را که آرشیوها می توانند نه فقط وب گاه ها و مجموعه های شان، بلکه پیوند ها و

ص: 122

Data Cite -1

ReaderMeter -2

bookmark -3

حاشیه نویسی های آن صفحه ها و مجموعه ها را ذخیره نمایند، در نظر بگیرید. توانایی پاسخ به این سؤال که «چگونه افراد این منابع را برای یکدیگر نشانه گذاری می نمایند، و چگونه در طول زمان تغییر می کند؟» و بررسی استفاده از فناوری های موجود ترکیب و تطبیق دادن ابزارهای حاشیه نویسی اجتماعی موجود. گامی بیشتر بردارید، آیا می توان اجتماعی از پیوندها و وب نوشت ها برای مجموعه های آرشیوی با استفاده از ابزارهای مشابه مورد تأکید قرار داد تا معلوم شود افراد چگونه موارد موجود روی وب پویا را برای یکدیگر نشانه گذاری می کنند؟

مثال ها: (1) Delicious، ردایت بر مبنای بوک مارکلت (2) مثال: (3) Madcow

معماران جدید

برای استفاده از داده های وب به منظور درک اتکای دنیا بر منابعی چون (4) API، و داده های پیوند شده، فعالیتی تعریف شده است تا داده ها را برای استفاده و هدف گذاری مجدد و ترکیب [آماده سازی]، قابل نماید. دلیل مهمی که چرا وب پویا نسبت به وب آرشیوی بسیار فعالانه تر مورد جست و جو قرار می گیرد، این است که طراحان ابزار به داده های وب پویا یا از طریق خزش گره های وبگاه ها به طور مستقیم و یا از طریق APIها به گوگل / یاهو، توئیتر، فیسبوک و مانند آن می توانند دسترسی یابند. حتی با وجود محدودیت بعضی از این APIها دلیل اصلی، پژوهش هایی هستند که با استفاده از وب پویا رونق گرفته اند.

API روش قدرتمندی برای ساخت نرم افزار های کاربردی جدید است که بر روی داده ها طراحی شده و منابع چندگانه داده را ترکیب می نماید.

پژوهشگری به ما گفت: من در مورد استفاده از برجسب نشانه گذاری اصلی در توئیتر تحقیق می کنم، و محدودیت شان را در استفاده از API در حالتی که بیشترین آشفستگی را در حین تویت کردن دارند می یابم، من می خواهم به آن ها در حالی که هنوز برخط و در دسترس هستند دسترسی داشته باشم، اگر چه استقرار آن ها، یا به عبارت دیگر اجرای گزارش ها یا انباشت آن ها کاملاً مشکل است. برای مثال آن ها Twapper keeper را محدود به خدمت قابل دسترسی کردند که اجازه آماده سازی گزارش هایی برای کار در مورد این برجسب ها که از آن ها شناخت داشتم را می داد.

سؤال بعدی این است که چگونه API ها می توانند آرشیو شوند و چه زمانی محدود یا متوقف می شوند؟ چگونه منابعی که تاکنون آرشیو شده اند از API تأثیر می گیرند؟ آیا روش هایی برای حفاظت محتوا همراه با حفظ و احترام به حقوق مالکان و دوره های مجوز وجود دارد؟

چالش: چگونه داده های وب قدیمی از طریق API ها باز شده و پیوند داده می شوند، تا افراد هوشمند خارج از آن بتوانند برای کاربرد راه های جدید تولید دانش آن ها را ترکیب کنند [؛ البته] به وسیله ارائه ابزارهایی به افراد، برای تعریف روال های کاری به صورتی انعطاف پذیرتر از اینکه وادار به استفاده از

(/Delicious(<http://www.delicious.com> -1

bookmarklet -2

(/MadCow(<http://www.web-notes.com> -3

-4 (رابطہ یا میانجی برنامہ های کاربردی (API: Application Program Interface)

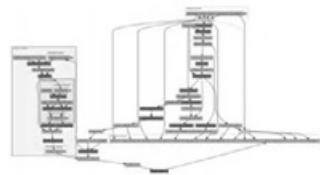
ابزارهایی یک منظوره یکپارچه بشوند. یک رهیافت، پیاده سازی کارکردهای تحلیلی مانند خدمات وب خواهد بود که می تواند با گردش کار قانونی موتور ترکیب شود. این رهیافت به صورت گسترده ای در bioinformatic که با مشکلات مشابهی در خصوص یکپارچه سازی داده ها از مخازن متعدد رویه روست استفاده می شود.

مثال: گردش کار قانون - مصوب موتور [1 Taverna](#) برای ترکیب [تلفیق] خدمات وب پیشنهادی توسط پردازش توزیع شده و نظام های ذخیره سازی استفاده می شود، سپس گردش کارها می توانند به اشتراک در آمده، دوباره استفاده و هدفگذاری شوند. myexperience [2](#) مثالی از یک مخزن قابل اشتراک گردش کارهای علمی:

عکس

ابزارهایی یک منظوره یکپارچه بشوند. یک رهیافت، پیاده‌سازی کارکردهای تحلیلی مانند خدمات وب خواهد بود که می‌تواند با گردش کار قانونی موتور ترکیب شود. این رهیافت به‌صورت گسترده‌ای در bioinformatic که با مشکلات مشابهی درخصوص یکپارچه سازی داده‌ها از مخازن متعدد روبه‌روست، استفاده می‌شود.

مثال: گردش کار قانون - مصوب موتور Taverna^۱ برای ترکیب [تلفیق] خدمات وب پیشنهادی توسط پردازش توزیع شده و نظام‌های ذخیره‌سازی استفاده می‌شود، سپس گردش کارها می‌توانند به اشتراک در آمده، دوباره استفاده و هدفگذاری شوند. ^۲ myexperience مثالی از یک مخزن قابل اشتراک گردش کارهای علمی:



تصویر ۱۱. گردش کارهای تاورنا



تصویر ۱۰. منبع: کارول گلوبال و دیگران^۲

ماشین‌های اجتماعی

تیم برنزی و همکارانش در مورد اینکه وب در حال تبدیل شدن به یک ماشین اجتماعی است، بحث کرده‌اند. به این معنی که فقط یک مخزن اطلاعات نیست، بلکه زیرساختی برای همکاری در رفع مشکل اجرای وظایف موجودیت‌های انسانی است که به‌آسانی نمی‌توانند توسط ماشین انجام شوند. دانشمندان علوم اجتماعی علاقه‌مند به درک تعامل بین فناوری و علوم اجتماعی، می‌خواهند بدانند که چگونه افراد و فناوری برای حل وظایف پیچیده‌ای که هر یک به تنهایی از عهده رفع آن بر نمی‌آیند، با یکدیگر همکاری می‌کنند.

چالش: چگونه می‌توانیم تجربه و تعامل بین کاربران و ماشین‌های اجتماعی بر روی وب را ذخیره و درک کنیم؟ همه امور اجتماعی در مورد تعامل هاست. اگر نتوانیم تعاملات را درک کنیم، هرگز نمی‌توانیم آنچه را که درباره ماشین اجتماعی است درک کنیم.

مثال: ^۳ Amazon's Mechanical Turk مکانیسمی برای توزیع مشکلات بین انسان‌های خیره و جمع‌آوری راه‌حل هاست. گردآوری منابع، می‌تواند برای رفع مشکلات سخت به‌کار رود. به‌عنوان مثال

1. Taverna (<http://www.taverna.org.uk/>)

2. <http://www.myexperiment.org/>

3. Carol Global (http://nar.oxfordjournals.org/content/38/suppl_2/W677.full)

4. Amazon's Mechanical Turk (<https://www.mturk.com/mturk/welcome>)

تصویر 11. گردش کارهای تاورنا

تصویر 10: منبع: کارول گلوبال و دیگران (3)

تیم برنزی و همکارانش در مورد اینکه وب در حال تبدیل شدن به یک ماشین اجتماعی است، بحث کرده اند. به این معنی که فقط یک مخزن اطلاعات نیست، بلکه زیرساختی برای همکاری در رفع مشکل اجرای وظایف موجودیت های انسانی است که به آسانی نمی توانند توسط ماشین انجام شوند. دانشمندان علوم اجتماعی علاقه مند به درک تعامل بین فناوری و علوم اجتماعی، می خواهند بدانند که چگونه افراد و فناوری برای حل وظایف پیچیده ای که هر یک به تنهایی از عهده رفع آن بر نمی آیند، با یکدیگر همکاری می کنند.

چالش: چگونه می توانیم تجربه و تعامل بین کاربران و ماشین های اجتماعی بر روی وب را ذخیره و درک کنیم؟ همه امور اجتماعی در مورد تعامل هاست. اگر نتوانیم تعاملات را درک کنیم، هرگز نمی توانیم آن چه را که درباره ماشین اجتماعی است درک کنیم.

مثال: [Amazon's Mechanical Turk \(4\)](#) مکانیسمی برای توزیع مشکلات بین انسان های خبره و جمع آوری راه حل هاست گردآوری منابع می تواند برای رفع مشکلات سخت به کار رود. به عنوان مثال

ص: 124

1 - <http://www.taverna.org.uk> (/Taverna)

2 - <http://www.myexperiment.org>

3 - http://nar.oxfordjournals.org/content/38/suppl_2/W677.full (Carol Global)

4 - <https://www.mturk.com/mturk/welcome> (Amazon's Mechanical Turk)

1) CAPTCHA برای شناسایی نویسه نوری کمک کننده انسان است. افراد چگونه با این ابزارها و از طریق اینها با یکدیگر تعامل می کنند؟

برای مثال، در مورد «بازهای هدف دار» همچون (2) Foldit، چالش نه در آرشیوسازی وبگاه ها و نه فقط در بازی بلکه در چگونگی بازی افراد است. کاربران چگونه با بازی تعامل دارند؟ از تحقیق در مورد چگونگی بازی بازیکنان در تعامل با فضای پیوسته، دروسی را می توان طراحی نمود (برای مثال: ویلیامز (3)، یی (4) و کاپلان (5)، 2008).

شبکه های نقشه برداری

به طور روزافزونی جغرافی دانان از داده های اینترنت برای دریافت اطلاعات مکان ها، جریان ها، جهات و ثروت، فقر و تغییر شکل محتوا، و تأثیر در طول زمان و فضا استفاده می کنند.

چالش: استخراج خودکار اطلاعات جغرافیایی از پیوندهای درونی و بیرونی درون یک مجموعه که می تواند نقشه برداری شود. این چالش، در حال حاضر در وب پویا وجود دارد و حتی در زمان افزایش پیچیدگی تغییرات در طول زمان، بیشتر هم می شود. در حال حاضر، بخش عمده ای از اطلاعات که می تواند با استفاده از روش های دو بعدی نمایان شود به داده های سه بعدی و چهار بعدی (همانند اسلایدر (6) ها) نیاز خواهند داشت. تا اطلاعات جغرافیایی را که در طول زمان تغییر می یابند، ایجاد کند برای مثال، درک تأثیر جغرافیایی درون و مابین دانشگاه ها، دولت ها، و شرکت ها در طول زمان از نظر تئوری محتمل است، اما به استخراج اطلاعات جغرافیایی از داده های ساختار نیافته وب نیاز دارد.

مثال: (7) Floating Sheep

ص: 125

1 - CAPTCHA: برنامه ای کامپیوتری است برای اجرای آزمون های نرم افزاری که فقط انسان می تواند به آن پاسخ دهد. (دسترسی در <http://www.google.com/recaptcha/captcha>)

2 - Foldit

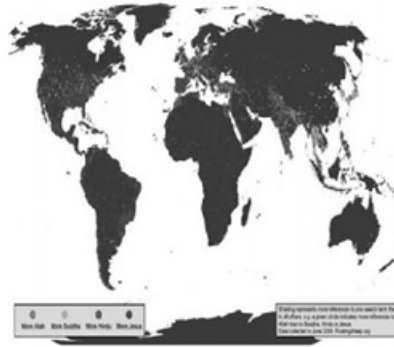
3 - Williams

4 - Yee

5 - Caplan

6 - Slider: رابط های کاربر گرافیکی برای ارتباط افراد با نرم افزار در ابزارهایی نظیر موبایل (دسترسی در <http://www.wisegeek.com/in-computing-what-is-a-slider.htm>)

7 - Floating Sheep (<http://www.floatingsheep.org>)



تصویر ۱۳. نقشه فلوتینگ شیپ «نقشه جغرافیایی مذهبی گوگل»

علم وب

علم وب تلاشی است که توسط محققان برای مطالعه وب به عنوان یک «مصنوع اطلاعاتی»، فهم چگونگی رشد و تکامل آن و کیفیت توسعه جوامع آن، انجام می‌گیرد.

چالش‌ها: نیاز به ابزارهای قدرتمندی برای تحلیل «نمودار وب» به عنوان شیء ریاضی. مکان شناسی آن چیست؟ چگونه گروه‌ها تکامل می‌یابند؟ چه نوع رتبه‌بندی قانونی به کار می‌رود؟ آیا واقعاً وب به وسیله قانونی قوی مدیریت می‌شود؟ چگونه اطلاعات روی وب منتشر می‌شوند؟

وب تنها یک فضای اطلاعاتی نیست، بلکه مجموعه‌ای پیچیده و خانواده‌ای از روابط درونی فضاهای فرعی است که محتوای اطلاعاتی آن در بعضی مواقع توسط جوامع مجزا تعیین می‌شوند. چگونه اطلاعات، بین این صفحه‌ها به اشتراک گذاشته و منتشر می‌شوند؟

برای پاسخ به این سؤال، به توسعه ابزارهایی نیاز داریم که قادر به ردیابی تکامل و انتقال مفاهیم در طول زمان و فضاهای مختلف باشند (برای مثال: بین وبلاگ‌نویسی و رسانه «mainstream»).

مثال: ^۱ Media Cloud و ^۲ Recorded future

1. Media Cloud (<http://cyber.law.harvard.edu/research/mediacloud>)

2. Recorded future (<https://www.recordedfuture.com/>)

تصویر 13. نقشه فلوتینگ شیپ «نقشه جغرافیایی مذهبی گوگل»

علم وب تلاشی است که توسط محققان برای مطالعه وب به عنوان یک «مصنوع اطلاعاتی»، فهم چگونگی رشد و تکامل آن و کیفیت توسعه جوامع آن، انجام می گیرد.

چالش ها: نیاز به ابزارهای قدرتمندی برای تحلیل «نمودار وب» به عنوان شیء ریاضی مکان شناسی آن چیست؟ چگونه گروه ها تکامل می یابند؟ چه نوع رتبه بندی قانونی به کار می رود؟ آیا واقعاً وب به وسیله قانونی قوی مدیریت می شود؟ چگونه اطلاعات روی وب منتشر می شوند؟

وب تنها یک فضای اطلاعاتی نیست، بلکه مجموعه ای پیچیده و خانواده ی از روابط درونی فضاهای فرعی است که محتوای اطلاعاتی آن در بعضی مواقع توسط جوامع مجزا تعیین می شوند. چگونه اطلاعات، بین این صفحه ها به اشتراک گذاشته و منتشر می شوند؟

برای پاسخ به این سؤال، به توسعه ابزارهایی نیاز داریم که قادر به ردیابی تکامل و انتقال مفاهیم در طول زمان و فضاهای مختلف باشند (برای مثال: بین وبلاگ نویسی و رسانه «mainstream»).

مثال: (1) Media Cloud و (2) Recorded future

ص: 126

1- (Media Cloud (<http://cyber.law.harvard.edu/research/mediacloud>

2- (/Recorded future (<https://www.recordedfuture.com>

آینده آرشیو وب ۱۲۷



تصویر ۱۴. Recorded Futur سیر اخبار را در طی زمان را دنبال می کند.

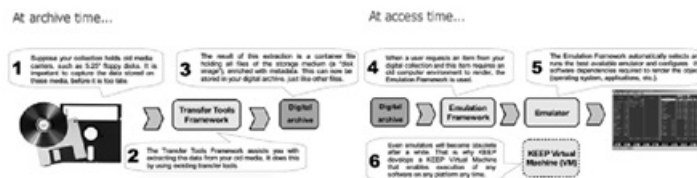
تصویر ۱۵. Media Cloud سیر اخبار جغرافیایی را دنبال می کند.

درک تجربه به جای محتوا

به طور فزاینده، دانشمندان به سمت درک اهمیت چگونگی استفاده افراد از محتوای وب و نه فقط محتوا به خودی خود، سوق پیدا کرده اند. این امر، به حساب وضعیت تجربه وب و محتوای قابل اجرا گذاشته می شود: تجربه درباره اینکه کدام طرح، کدام مرورگر یا افزونه‌ها یا رمز گذارهای دوطرفه مورد استفاده قرار می گیرد و به طور فزاینده به موقعیت مکانی کاربران وابسته است.

چالش‌ها: برای درک تجربه‌ها به توانایی تکرار تجربه نیازمندیم، طرح‌های زیربنایی، سیستم‌های عملیاتی، مرورگرها، و به همین ترتیب تغییر تجربه وب.

مثال‌ها: Browsershots^۲ و KEEP^۳ (نگهداری نمونه سازی محیط‌های قابل انتقال)



تصویر ۱۶. KEEP (نگهداری نمونه سازی محیط‌های قابل انتقال)

تحلیل وب معنایی و مجموعه داده‌های پیوند شده

مجموعه داده‌های پیوند شده به سرعت با حداقل ۲/۷۵ بیلیون پیوند سه تایی در مجموعه‌های شناخته شده در

1. plugins
2. Browsershots (<http://browsershots.org>)
3. KEEP (<http://www.keep-project.eu>)

تصویر ۱۴. Recorded Futur سیر اخبار را در طی زمان را دنبال می کند.

تصویر ۱۵. Media Cloud سیر اخبار جغرافیایی را دنبال می کند.

درک تجربه به جای محتوا

به طور فزاینده دانشمندان به سمت درک اهمیت چگونگی استفاده افراد از محتوای وب و نه فقط محتوا به خودی خود، سوق پیدا کرده اند. این امر، به حساب وضعیت تجربه وب و محتوای قابل اجرا گذاشته می شود: تجربه درباره اینکه کدام طرح کدام مرورگر یا افزونه ها (1) یا رمزگذارهای دوطرفه مورد استفاده قرار می گیرد و به طور فزاینده به موقعیت مکانی کاربران وابسته است.

چالش ها: برای درک تجربه ها به توانایی تکرار تجربه نیازمندیم، طرح های زیربنایی، سیستم های عملیاتی، مرورگرها، و به همین ترتیب تغییر تجربه وب.

مثال ها: (2) Browsershots و (3) KEEP (نگهداری نمونه سازی محیط های قابل انتقال)

تصویر 16. KEEP (نگهداری نمونه سازی محیط های قابل انتقال)

تحلیل وب معنایی و مجموعه داده های پیوند شده

مجموعه داده های پیوند شده به سرعت با حداقل 28/5 بلیون پیوند سه تایی در مجموعه های شناخته شده در

ص: 127

plugins -1

(Browsershots (<http://browsershots.org> -2

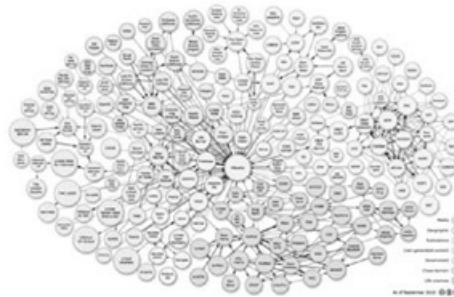
(KEEP (<http://www.keep-project.eu> -3

حال رشد هستند. ابزارها توسط وب معنایی یا جامعه داده های پیوند شده توسعه می یابند، به طوری که می توانند به سرعت مشکل مدیریت فراداده را آسان کرده و به این فراداده ها اجازه دهند که در مقیاس بزرگ قابل جست و جو بوده و همچنین برای یکپارچه سازی داده ها در مجموعه به روشی بسیار پیچیده تر استفاده شوند.

عکس

۱۲۸ مدیریت منابع اطلاعاتی وب

حال رشد هستند. ابزارها توسط وب معنایی یا جامعه داده های پیوند شده توسعه می یابند، به طوری که می توانند به سرعت مشکل مدیریت فراداده را آسان کرده و به این فراداده ها اجازه دهند که در مقیاس بزرگ قابل جست و جو بوده و همچنین برای یکپارچه سازی داده ها در مجموعه به روشی بسیار پیچیده تر استفاده شوند.



تصویر ۱۷. منبع: نمودار ابر گونه داده های باز پیوند خورده

توسط ریچارد کی گانیاک^۱ و آنجا ینتس^۲ <http://lod-cloud.net/>

مثالها: ^۳ [Sig.ma / Sindic](http://sig.ma/) جست و جوی مبتنی بر آر دی اف^۴



تصویر ۱۸. [Sig.ma / Sindic](http://sig.ma/) جست و جوی مبتنی بر RDF.

منبع: <http://sig.ma/search?q=Tim/۲۰Berners/۲۰Lee>

1. Richard Cyganiak
2. Anja Jentzsch
3. <http://sig.ma/>
4. Resource Description Framework (RDF)

تصویر 17. منبع: نمودار ابر گونه داده های باز پیوند خورده توسط ریچارد کی گانیاک (1) و آنجا ینتس (2) <http://lod-cloud.net>

مثال‌ها: (3) [Sig.ma / Sindic](http://sig.ma) جست و جوی مبتنی بر آر دی اف (4)

تصویر 18. [Sig.ma / Sindic](http://sig.ma) جست و جوی مبتنی بر RDF

منبع: <http://sig.ma/search?q=Tim%20Berners%20Lee>

ص: 128

Richard Cyganiak -1

Anja Jentzsch -2

/http://sig.ma -3

Resource Description Framework (RDF) -4

اما گام های پیش رو برای جامعه آرشیوی وب چیست؟

بسیاری از مواردی که محققان [به آن ها] نیاز دارند، به وضوح قابل مشاهده هستند، ولی این به معنی در دسترس بودن آن ها نیست. در گزارش های قبلی GISC که با همکاران دیگر نوشتیم و در قسمت های پیشین توضیح دادیم (دوگرتی و دیگران 2010؛ توماس و دیگران، 2010) توجه های متعددی پیرامون موضوع اصلی ارائه دادیم:

ساخت جامعه، ساخت ابزار، و منابع و طراحی تجربه ها (دوگرتی و دیگران، 2010، ص 27-29). قصد نداریم که فهرستی کامل از پیشنهاد های گزارش های قبلی را در اینجا مطرح کنیم. 22 پیشنهاد در گزارش دوگرتی و همکارانش و بیش از 20 پیشنهاد در گزارش توماس و همکارانش وجود دارد. بنابراین، خواننده را به مرور این گزارش ها همانند این گزارش توصیه می کنیم با این حال می توانیم برای نیل به هدفمان بعضی از موضوع های مهمی را که محققان به آن ها علاقمند هستند، برجسته کرده و چالش هایی را شناسایی کنیم که آرشیوها برای کمک به آرشیوهای وب در تبدیل آن ها به بخشی از ابزارهای استاندارد برای محققان رشته های متنوع، با آن ها مواجه هستند.

این بخش از گزارش، موضوع ها و سؤال هایی را بر می شمرد که گروه های مختلف محققان می خواهند پرسند یا از آرشیو وب - آرشیو - در - جعبه - می خواهند پرسند. ما آن دسته را که چالش ها و راه حل های بالقوه ای دارند شناسایی کرده ایم بعضی از این راه حل ها، به طور ویژه راه حل هایی کوتاه مدت هستند که می توانند در سطح مؤسسه ها انجام شوند. بسیاری از رهیافت های نه چندان گسترده به نگاه وسیع تری در سطح ملی منطقه ای یا بین المللی توسط سازمان هایی چون IIPC نیاز دارند.

ما از آرشیو - در - جعبه یاد کردیم، زیرا در میان بعضی افراد این تفکر وجود دارد که وب آرشیوی فراتر از روزهای اولیه خودش در ایجاد صفحه های قابل دسترسی برای تحلیل های آینده، (تفکر تعامل سنتی (1) برای پایگاه Wayback Machin که به کاربر اجازه می داد به طور عمده به صفحه های منفرد آرشیو دسترسی داشته و آن ها را ببیند) در حال حرکت به سمت ایجاد مجموعه های قابل دسترسی به عنوان ابزارهای تحقیقی، است. برای مثال، زمانی که با دامنه «the UK government. Gov.uk» از 2011 - 2020 مواجه هستیم، یک پژوهش گر چه تصویری می تواند از توانایی انجام آن داشته باشد؟ چه سؤال هایی می تواند از یک مجموعه پرسیده شود؟ برای نمونه، محتوای کامل وب در مورد بانکدارهای اصلی وال استریت و سایت هایی که به طور مستقیم با آن پیوند دارند، در صورتی که آرشیو، دوره ای را پوشش دهد که طی آن بحران های بانکی توسعه یافته اند. به عبارت دیگر، به جای تحلیل یک وبگاه منفرد در سطحی خرد یا تحلیل همه وب در سطح کلان، با مجموعه های متمرکز وب در سطح میانی چه کار می توانیم بکنیم؟ قسمت عمده ای از پژوهش علوم اجتماعی در فضای برخط، در تعاملات سطوح میانی مشاهده می شود. آیا می توانیم برای فهم این که چگونه تغییر وب، واقعیت اجتماعی را منعکس و تقویت می کند یا تغییر می دهد، به صورت یکسان عمل کنیم؟

پاره ای از اقدامات برای پشتیبانی آرشیوهای وب موجود غیر ممکن خواهد بود - ممکن است برای انجام آن داده یا محتوا گردآوری نشده و تقریباً از بین رفته باشد. با این حال، رو به سوی تا آینده چه تغییراتی را می توانیم امروز و در سال های آینده برای آرشیوهای وب، ایجاد کنیم محققان در سالهای 2015، 2020 یا 2050 قادر به طراحی منابعی باشند که در حال حاضر برای پاسخ به این سؤال ها گردآوری می کنیم. محققان آینده، چه چیزی را از ما طلب می کنند که حالا در سال 2011، و در آینده، انجام نمی دهیم؟ مؤسسه های خاص چه می توانند انجام دهند؟ اگر IIPC به صورت جمعی در بهره داری از قدرت آرشیوهای چندگانه، عمل نماید چه اتفاق بهتر و مؤثرتری ممکن است رخ دهد؟

وب مجتمعه: زندگی آرشیو وب

سؤال: چرا آرشیوهای وب به آرشیو شدن نیاز دارند؟ چرا آن ها نمی توانند با وب پویا یکپارچه گردند و به طور شفاف برای عموم و محققان قابل دسترسی باشند؟ امکان تصور وبی که با وب پویای فعلی بر روی صفحه ای قابل دسترسی به صورت منبعی پیش فرض از داده ها و اطلاعات لایه بندی شده باشد، وجود دارد. با این حال، این سطح می تواند بر روی لایه های زیرین وب قدیمی ایجاد شود که به آسانی برای علاقه مندان با حرکت نزولی یک لایه یا تعداد بیشتری از لایه ها به سمت پایین، قابل دسترسی است. اگر ده ها هزار ابزار قابل دسترسی برای پژوهش وب پویا بتواند برای این لایه های زیرین با استفاده از مکانیسم های ساده به کار گرفته شود احتمال کشف کاربردهایی برای اطلاعات و داده های موجود در وب گذشته توسط محققان افزایش می یابد.

احتمالاً- این بزرگترین و بلند پروازانه ترین چالش در این گزارش است، زیرا نیاز به تغییر زیادی در زیر ساخت وب دارد. در عین حال، به معنی احتمال ضعیف وقوع آن است، اما نتایج مورد انتظار آن بسیار خواهد بود فراتر از ارزش پژوهشی که وب قدیمی به عنوان لایه های زیرین وب فعلی، دارد، به صورت بالقوه نیز باعث تغییر ساختاری در نحوه نگرش کاربران به وب خواهد بود. وب فعلی از نظر بسیاری برای اشاعه پیوندهای از دست رفته، از دست دادن اطلاعات، گم شدن صفحه ها، تغییر نشانی ها و تغییر اطلاعاتی که روی ویرایش های قبلی بدون هیچ امکانی برای مشاهده یا برگرداندن به ویرایش های قبلی، بازنویسی می شوند، غیر قابل اطمینان است. اگر ساختار اینترنت در برگشت زمان به عقب، به یکی از لایه های چندگانه تبدیل شود، به نحوی که حفره های موجود در لایه بالایی باعث ایجاد حفره در وب نشود، و به جای آن لایه پایین تر آشکار شود امکان اینکه وب به صورت منبعی با ثبات قابل اعتماد و مقاوم به از دست رفتن اطلاعات باشد، وجود خواهد داشت.

مشکل پیوندهای از دست رفته که تحلیل پیوند (1) نامیده می شود، مشکل دیرپایی برای کاربران وب پویاست. مشکل، زمانی پیچیده می شود که توجه بر وب آرشیوی متمرکز شود که روی هم رفته شناسه- گره های ثابتی برای ویرایش های آرشیوی صفحه ها وب ندارد. تلاش های متعددی صورت گرفته است،

ص: 130

به عنوان مثال 1) / (Web Cite (<http://www.webcitation.org>)) به نویسندگان امکان آرشیو نسخه ای از صفحه وب و ایجاد پیوند محافظت شده یا تجزیه کننده DOI را می دهد. 2) (Dead URL (<http://deadurl.com>)) رهیافت متفاوتی با اتکا بر آرشیو اینترنت و مخزن گوگل، در میان منابع دیگر برگزیده است تا سعی کند نسخه های پیوندهای از دست رفته را بیابد. اگر چه، تلاش های این چینی بوسیله اکثریت وسیعی از محققان که عمدتاً از وجود آن آگاه نیستند، بدون استفاده باقی می ماند. فراتر از عامل مشکل ساز، نتیجه ناخواسته و غیر منتظره دیگری نیز وجود دارد. عادت محققان به اینکه تا حد امکان نشانی های متعددی را در آثار علمی شان درج کنند. این امر اثرات متعددی دارد:

نخست: زمانی که سعی در ارزیابی تأثیر منابع برخط با استفاده از فنونی چون وب سنجی دارند، فقدان پیوندها باعث کاهش تأثیر منابع شان می شود.

دوم: خوانندگان را وادار به تلاشی جدی تر برای پیگیری منابع اطلاعاتی می کند، به طوری که سعی می کنند نه تنها منبع صحیح استناد، بلکه منبع استناد شده ای را که ممکن است ویرایش صفحه هایش می کنند به طور قابل توجهی تغییر کرده باشد نیز کشف کنند. اگر آرشیوهای وب تبدیل به منبعی قابل اتکا برای استناد پیوسته اطلاعات بشوند، این امر باعث افزایش پژوهش می شود و وجهه (نیمرخ) آرشیوهای وب را به صورت عمومی تر رشد خواهد داد.

در این وب مجتمع، دانش هویدا در وب، رشد و تکامل می یابد، اما به همان روش مشابه وب فعلی، به دور انداخته نمی شود. وب مجتمع، با استفاده از موتورهای جست و جویی نظیر گوگل، قابل خزش، زدودن پیوند شدن و تحلیل و جست و جو خواهد شد. پیوندها از بین نخواهند رفت، بلکه به سوی تبدیل شدن به و موادی هدایت می شوند. که مدت زیادی از وجودشان به صورت فعال بر روی وب فعلی نمی گذرد.

چالش بلند مدت: دو چالش در این سؤال نهفته است، که هر دوی آن ها مستلزم مشارکت کنندگان و فعالان زیادی خواهد بود. نخست: ما باید دوباره در مورد اینکه چگونه اینترنت را ببینیم و مهندسی کنیم فکر کنیم، در حال حرکت از موجودیتی تک لایه با پیوندهای جانبی بسیار به سوی موجودیتی چند لایه با پیوندهای جانبی فعلی، بلکه با پیوندهای فعلی به مواد قدیمی و پیوندهای قدیمی به مواد قدیمی یا فعلی [باشیم]. این چالش کم اهمیتی نیست و متقاعد کردن بسیاری از نقش آفرینان که در ساختن حال آینده اینترنت، سرمایه گذاری کرده اند، برای انتخاب، مشکل خواهد بود. با این حال، نتیجه آن منجر به زیر ساختی خواهد شد که وب گذشته را بسیار قابل دسترس تر برای ارجاع و جست و جو خواهد ساخت. چالش بزرگ بعدی این خواهد بود که آرشیویست های وب نیاز به تعریف مجدد نقش شان خواهند داشت، در حقیقت نه برای آرشیویست بودن به شکل کاملاً سنتی اش، بلکه به صورت متخصصانی که می توانند به پژوهشگران در درک روندها و منابع اینترنت در طول زمان کمک کند و در ابزارهایی که برای دسترسی و دستکاری لایه های این اینترنت چند لایه به آن ها نیاز دارند، برای راهنمایی پژوهشگران و عموم در پاسخ به پرسش های پیش بینی نشده در مورد رشد وب همانگونه که به صورت جدید توسعه یافته اند خبره هستند.

ص: 131

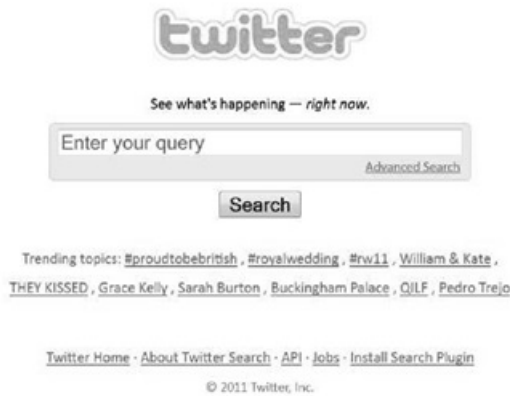
سؤال: محققان چگونه به تغییر رویدادها در جهان پاسخ می دهند و یا حوادث جاری را پیش می کنند؟ به طور فزاینده، رویدادهای محلی و جهانی خارج از وب اتفاق می افتد. این حوادث ممکن است ---ت رویدادهای مهم بین المللی و جالب توجه باشند. ممکن است مانند رویدادهای سیاسی اخیر در آفریقای شمالی و خاور میانه یا زلزله هایی در هائیتی ژاپن و سایر مکان ها رویدادهای محلی یا منطقه ای مهمی باشند. ممکن است کوچک اما گسترده، یا رویدادهای جاری باشند که در اولویت توجه گروهی کوچک یا حتی یک محقق باشند. این امر، به طور بالقوه انواع بینش ها را در مورد ماهیت و چیرستی اطلاعاتی که مردم به اشتراک می گذارند، موضوع ها و حوادثی که به طور برجسته توسعه می یابند، و چگونگی پاسخ افراد و سازمان ها و دولت ها به بحران ها و در طی زمان، چگونگی استحاله و افول حوادث در گفت و گوهای عامه مردم، به بار می آورد. این مورد ساده ترین چالش ولی واضح ترین آن هاست.

علاوه بر این، تعداد زیادی از محققان در این حوزه با هم کار می کنند، پدیده برداشت (1) موضوع چندین گفتگو در نشست IIPC در سال 2011 بود، توسط سخنرانانی که مثال های زیادی از پدیده برداشت در ارتباط با انقلابهای 2011 از سراسر دنیای عرب تا نشت نفت در آب های عمیق، تا المپیک سال 2012 لندن را بیان می داشتند.

یکی از موضوع های محوری کار با آرشیوهای وب، این است که ما نیاز به حرکت از فهم وب به صورت مجموعه ای از داده های انتخاب شده داریم و به جای آن باید آن را به صورت یک شبکه در حال تغییر و تحول ببینیم که به دوره های زمانی و رهیافت های طولانی نیاز دارد در این زمینه، تلاش هایی انجام گرفته است. برای نمونه، پروژه اروپایی تحلیل های طولی آرشیوهای وب (2) در حال ایجاد یک رصد خانه مجازی وب برای انجام تحلیل های طولی است.

چالش فوری: ایجاد سازوکارهایی برای محققان به منظور پیشنهاد سریع گرانولیتته (3) افزایشی و دامنه مناسبی برای استفاده در سایت های آرشیوی و موضوع های در حال تغییر، در حال حاضر، ممکن است یک پژوهشگر ماهر، ابزارهایی را برای تکرار اجرای خزش، راه اندازی کند، اما پژوهشگری با مهارت کمتر از نظر فنی بدون پشتیبانی قوی سازمانی، یادگیری نزولی خواهد داشت. زمانی که پدیده ای به سرعت متغیر، در حال توسعه است، پژوهشگر متخصص علاقه مند به آن پدیده که هیچ تجربه ای در وب آرشیوی ندارد، باید راهی برای گردآوری داده برای تحلیل آن، قبل از دست رفتن داشته باشد. سازمان های ماهر در راه اندازی ابزارهایی برای گردآوری داده های موجود بر روی این موقعیت های توسعه یافته، در سطح مناسبی از گرانولیتته، می توانند راه هایی برای محققان یا دیگران فراهم کنند که وبگاه ها، موضوع ها، کلید واژه ها، مانند آن را برای پاسخگویی سریع به رویدادهای وب برگزینند.

آینده آرشیو وب ۱۳۳



تصویر ۱۹. روند موضوعات تویتر ، ۲۹ آوریل ۲۰۱۱

چالش در حال توسعه: استفاده از ابزارهایی نظیر خوراک‌های R.S.S. برای به‌کارگیری نشانه‌ها، بیانگر این است که تغییرات صفحه‌های وب نیاز به آرشیو شدن دارند.

در ارتباط با توسعه رویدادها، امکان پایش مواردی نظیر خوراک‌های R.S.S. یا نرم‌افزارهای کاربردی اخیراً توسعه یافته، برای داشتن سیستم‌هایی که به فعالیت تکثیر به وسیله افزایش تکرار ذخیره صفحه وب، پاسخ داده یا بوسیله آگاه‌سازی متصدیان انسانی حوزه‌های در حال توسعه بالقوه مورد علاقه، وجود دارد. چالش طولانی مدت: ساخت الگوریتم‌هایی که از روندهای فعالیت برخط (نظیر روندهای گوگل، یا روند موضوع‌های تویتر) به‌منظور راه‌اندازی گرانولیت‌آرشیوسازی افزایشی برای صفحه‌های وب مرتبط با آن موضوع‌های استفاده می‌کنند.

این امر مستلزم مهارت و پختگی بیشتری برای آرشیو‌هایی است که به‌طور مناسبی برای استفاده مجدد و به اشتراک گذاری و استانداردسازی ایجاد شده و استمرار دارند. پژوهشگران علاقه‌مند به استفاده از چنین آرشیو‌های گردآوری شده الگوریتم گونه‌ای (دارای الگوریتمی)، می‌خواهند منطبق حذف یا افزودن را بدانند و قادر به فهم ماهیت و محتوای مجموعه باشند. سؤال محوری که پرسیده خواهد شد این است که چگونه مجموعه‌ها با اکوسیستمی از منابع پژوهشی قابل استفاده و قابل دسترسی، متناسب شوند؟

کاربردهای آرشیوها و وبگاه‌ها

سؤال: چگونه افراد آرشیو‌های وب را بسیار مهم‌تر از وبگاه‌ها به‌کار می‌گیرند؟ در حال حاضر، احتمال زیادی وجود دارد که وبگاه‌های بر روی وب را در مقاطع مشخص زمانی، با استفاده از آرشیو‌های وب و زیر ساخت آرشیوی وب، مشاهده نمود. برای پژوهشگران دانشگاهی و صنعت تحلیل‌های

تصویر 19 روند موضوعات تویتر ، 29 آوریل 2011

چالش در حال توسعه: استفاده از ابزارهایی نظیر خوراک‌های R.S.S. برای به‌کارگیری نشانه‌ها بیان گر این است که تغییرات صفحه‌های وب نیاز به آرشیو شدن دارند.

در ارتباط با توسعه رویدادها، امکان پایش مواردی نظیر خوراک های R.S.S. یا نرم افزارهای کاربردی اخیراً توسعه یافته برای داشتن سیستم هایی که به فعالیت تکثیر به وسیله افزایش تکرار ذخیره صفحه، وب پاسخ داده یا بوسیله آگاه سازی متصدیان انسانی حوزه های در حال توسعه بالقوه مورد علاقه، وجود دارد.

چالش طولانی مدت: ساخت الگوریتم هایی که از روندهای فعالیت برخط (نظیر روندهای گوگل یا روند موضوع های توییتر) به منظور راه اندازی گرانولیتة آرشیو سازی افزایشی برای صفحه های وب مرتبط با آن موضوع های استفاده می کنند.

این امر مستلزم مهارت و پختگی بیش تری برای آرشیو هایی است که به طور مناسبی برای استفاده مجدد و به اشتراک گذاری و استاندارد سازی ایجاد شده و استمرار دارند پژوهش گران علاقه مند به استفاده از چنین آرشیوهای گردآوری شده الگوریتم گونه ای دارای الگوریتمی ، می خواهند منطق حذف یا افزودن را بدانند و قادر به فهم ماهیت و محتوای مجموعه باشند. سؤال محوری که پرسیده خواهد شد این است که چگونه مجموعه ها با اکوسیستمی از منابع پژوهشی قابل استفاده و قابل دسترسی متناسب شوند؟

کاربرد های آرشیو ها و وب گاه ها

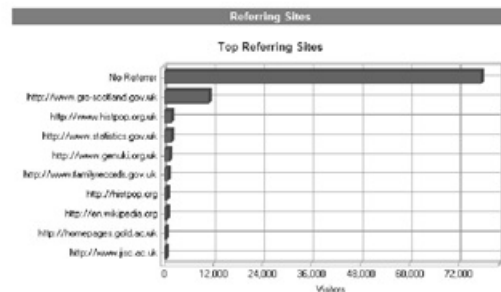
سؤال چگونه افراد آرشیوهای وب را بسیار مهم تر از وب گاه ها به کار می گیرند؟ در حال حاضر، احتمال زیادی وجود دارد که وب گاه های بر روی وب را در مقاطع مشخص، زمانی با استفاده از آرشیو های وب و زیر ساخت آرشیوی، وب مشاهده نمود برای پژوهش گران دانشگاهی و صنعت تحلیل های

سرویس دهنده ثبت وقایع و تحلیل های تحلیل گران تکنیکی پر تکرار برای ارزیابی کاربردها تأثیر و الگوهای ترافیک در وب پویاست با وجود، این فنون برای وب آرشیوی امکان پذیر نیست زیرا داده ها برای درک استفاده های وب قدیمی و آرشیو های وب به سادگی در دسترس نیستند.

عکس

۱۳۴ مدیریت منابع اطلاعاتی وب

سرویس دهنده ثبت وقایع و تحلیل های تحلیل گران، تکنیکی پر تکرار برای ارزیابی کاربردها، تأثیر و الگوهای ترافیک در وب پویاست. با وجود این، این فنون برای وب آرشیوی امکان پذیر نیست. زیرا داده ها برای درک استفاده های وب قدیمی و آرشیو های وب به سادگی در دسترس نیستند.



تصویر. ۲۰ نمونه ای از فایل داده ثبت وقایع برای سایت histpop.org منبع: مه پر و دیگران، ۲۰۰۹

چالش فوری: آرشیو سرویس دهنده ثبت وقایع سایت های وب آرشیوی، که محققان بتوانند نحوه استفاده از آرشیو های وب را مطالعه نمایند. این راه حلی ساده و سراسر است و در دسترس برای مؤسسه هایی است که آرشیو های وب را می سازند. سرویس دهنده ثبت وقایع آرشیو های وب می تواند برای پژوهشگران علاقمند به فهم چگونگی مرور آرشیو های وب توسط کاربران و چگونگی دسترسی آنها به منابع و اینکه چه بخش هایی از آرشیو بیشترین تکرار استفاده را دارند، ذخیره، نگهداری و قابل دسترسی شود. این اطلاعات به طور عمده مورد توجه جامعه آرشیوی وب خواهد بود. اما این گام نخست است.

چالش بلند مدت: تلاش بلند پروازانه تری است، اما علاقه بالقوه بسیار گسترده تری خواهد بود: راه اندازی زیر ساختی برای امکان آرشیوسازی سرویس دهنده های ثبت وقایع و تحلیل های مرتبط و پیوند با وبگاه ها، به طوری که پژوهشگران نه فقط آنچه که بر روی وب موجود بوده، بلکه نحوه استفاده از آن را نیز بتوانند ببینند. این هدفی بسیار بلند پروازانه تر است، زیرا سرویس دهنده های ثبت وقایع و اعتبار های تحلیل ها فقط به صورت داخلی برای سرور و مدیران اعتبار در وضعیتی حفاظت شده قابل مشاهده هستند. تجربیات عمومی مدیران سرور برای ذخیره سرویس دهنده ثبت وقایع در بلند مدت ضروری نیستند، به طوری که آنها به طور عادی برای ذخیره سازی فضا و دوری از انباشتگی و آشفستگی سرور، حذف و بازنویسی می شوند. اگر چه این داده ها، به طور بالقوه برای پژوهشگرانی ارزشمند است که نه فقط می خواهند وضعیت یک سایت آگاه شوند، بلکه می خواهند از چگونگی و مقدار استفاده از آن، ترافیک منابع و دیگر حقایقی که از ثبت وقایع و داده های تحلیلی می تواند گردآوری شود، آگاهی یابند. راه حل های اجتماعی شامل ایجاد مکانیسم هایی برای مدیران سرور برای مشارکت در ثبت وقایع

تصویر. 20 نمونه ای از فایل داده ثبت وقایع برای سایت histpop.org منبع: مه پر و دیگران، 2009

چالش فوری: آرشیو سرویس دهنده ثبت وقایع سایت های وب آرشیوی، که محققان بتوانند نحوه استفاده از آرشیوهای وب را مطالعه نمایند این راه حلی ساده و سر راست و در دسترس برای مؤسسه هایی است که آرشیوهای وب را می سازند. سرویس دهنده ثبت وقایع آرشیوهای وب می تواند رای پژوهش گران علاقمند به فهم چگونگی مرور آرشیوهای وب توسط کاربران و چگونگی دسترسی آن ها به منابع و اینکه چه بخش هایی از آرشیو بیش ترین تکرار استفاده را دارند، ذخیره، نگهداری و قابل دسترسی شود. این اطلاعات به طور عمده مورد توجه جامعه آرشیوی وب خواهد بود. اما این گام نخست است.

چالش بلند مدت: تلاش بلند پروازانه تری، است اما علاقه بالقوه بسیار گسترده تری خواهد بود راه اندازی زیر ساختی برای امکان آرشیو سازی سرویس دهنده های ثبت وقایع و تحلیل های مرتبط و پیوند با وب گاه ها به طوری که پژوهش گران نه فقط آن چه که بر روی وب موجود بوده، بلکه نحوه استفاده از آن را نیز بتوانند ببینند این هدفی بسیار بلند پروازانه تر است، زیرا سرویس دهنده های ثبت وقایع و اعتبارهای تحلیل ها فقط به صورت داخلی برای سرور و مدیران اعتبار در وضعیتی حفاظت شده قابل مشاهده هستند. تجربیات عمومی مدیران سرور برای ذخیره سرویس دهنده ثبت وقایع در بلند مدت ضروری نیستند به طوری که آن ها به طور عادی برای ذخیره سازی فضا و دوری از انباشتگی و آشفستگی سرور، حذف و بازنویسی می شوند. اگر چه این داده ها به طور بالقوه برای پژوهش گرانی ارزشمند است که نه فقط می خواهند از وضعیت یک سایت آگاه شوند، بلکه می خواهند از چگونگی و مقدار استفاده از آن، ترافیک منابع و دیگر حقایقی که از ثبت وقایع و داده های تحلیلی می تواند گردآوری شود، آگاهی یابند. راه حل های اجتماعی شامل ایجاد مکانیسم هایی برای مدیران سرور برای مشارکت در ثبت وقایع

است که می تواند با وبگاه های آرشیو ارتباط داشته باشد و برای ارائه کنندگان تحلیل هایی همچون گوگل به منظور ارائه گزینه ای به منظور مشارکت سایت های تحلیلی با آرشیوهای وب و تا حد امکان تعیین دوره محدودیت قبل از این که داده ها انتشار یابند.

چالش بلند پروازانه طراحی نظام هایی برای آرشیو نه تنها سرویس دهنده های ثبت وقایع، بلکه خود ترافیک وب در یک سبک ناشناخته حفاظت شده این راه حل حتی بسیار بلند پروازانه تر است؛ زیرا سازو کارهای تحلیل ترافیک، وب روی هم رفته عموماً در دسترس نیستند. نگرانی های نهفته ای از مواردی نظیر بازرسی بسته های عمیق که به تحلیلگران درباره ترافیک روان، وب می گوید وجود دارد.

پرسش محوری برای سؤال کردن این است که چه زمان مزایای فهم نحوه رفتار افراد در وب بر خطرات غلبه مینماید، بنابراین توجه به این امر که روش هایی وجود دارند که این داده ها می توانند برای تحلیل های بعدی آرشیو و ذخیره شوند در زمانی که خطرات برای افراد و سازمان ها به حد کافی به واسطه گذر زمان و گمنام کردن داده ها کاهش می یابد ارزشمند است.

متخصصی وب

سؤال: آیا امکان شناسایی مجموعه هایی که اندازه و شکل وب تاریخی را در رابطه با حوزه های تخصصی مورد توجه می سنجند، وجود دارد؟ فرض می کنیم که بسیاری از گروه ها و تنالگان ها وب گاه هایی را در طول دوره های زمانی ایجاد خواهند کرد چگونه امکان دارد که شناسایی تداوم حضور آن ها در وب و جمع آوری بدنه مناسبی از سایت های پژوهشی، برای آن ها کوچک و بی اهمیت باشد؟

مثال های متعددی را می توان بر شمرد گروه های علاقه مند به سرگرمی موضوع های تخصصی دانشگاهی میراث صنایع دستی نظیر دیداری و شنیداری سایت های گروه های سیاسی، و مانند آن. مطالعه این موارد می تواند موجب ساخت ابر مجموعه ها شود مجموعه مجموعه ها، آرشیو - در - جعبه چه راهنمایی در زمان ساختن مجموعه های آرشیوی از منابع متنوع وب برای افزایش ارزش این مجموعه ها مورد نیاز است؟ یک ابر مجموعه تا چه اندازه می تواند با دیگری مقایسه و پیوند داده شود و حتی به بخش هایی از ابر مجموعه های بزرگ تر تبدیل شوند؟

متخصص وب، ضرورت مقیاس کوچک تری از داده های انتخاب شده بر اساس رویداد ها یا موضوع هایی هم چنان، انتقادی و موازی با تجربیات و انتظارات موجود پژوهشی را تشخیص می دهد. مجموعه قوانینی از این دست به صورتی قابل مشاهده و جست و جو بوسیله افراد یا تیمی باقی میماند که این رشته های داده را از یک چشم اندازه ویژه رشته یا موضوع پژوهشی مهمی در دنیای آکادمیک، ایجاد و تحلیل می کنند به ویژه در زمینه هایی که باوری دیرینه وجود دارد که پژوهش گران، خالقان مدیران و تحلیل گران مجموعه قوانین خود شان باشند با این حال حتی این مجموعه های تخصصی به طور بالقوه در زمانی که به طور استاندارد برای پاسخ به سؤال های محققان دوباره ترکیب و دسته بندی می شوند، می توانند ارزش آفرین باشند

این، امر مسئله مقیاس پذیری را در مسیر پرسش های مرتبط با آرشیو های وب متخصص پرورش

می دهد. آیا اندازه بحرانی در طول دوره های پوشش و دامنه مجموعه وب آرشیوی برای اینکه در طول زمان مورد استفاده دیگر پژوهش گران باشد وجود دارد؟ چگونه راهبرد ها و سیاست های خزندگی وب مؤسسه های مختلف (همچون برداشت های عمده و حوزه رویدادی یا انتخابی) انتظارات پژوهش گران را برآورده کرده، و چه راهبرد های متنوعی به مناسب ترین وضع برای پشتیبانی، استفاده می شوند؟

بعضی از مثال های ویژه حوزه مجموعه های تخصصی در ابزارهای علوم اجتماعی برای ترکیب بانک جهانی داده با تحلیل های موجود در ولف رام آلفا (1) در ارتباط با سلامتی یا دیگر اطلاعات جمعیتی می توانند توسعه یابند. این مستلزم استفاده از ابزاری (همچون ولف رام آلفا) دارد که به طور مداوم با داده های جدید روزآمد می شود به همان صورت الگوریتم ها و ابزار های مصور سازی آن، بستر تحلیلی داده های پویا را ارائه می کنند در علوم، طبیعی پژوهش گران به گرفتن اطلاعات آب و هوا (دما) علاقه مند خواهند بود و این اطلاعات جمع آوری شده را با طبیعت شناسان غیر حرفه ای (مانند مشاهده کنندگان پرندگان) در گروه های مختلف (با استفاده از ابزارهایی چون گوگل پلاس یا فیس بوک) در سراسر دنیا، به اشتراک گذاشته و به گروه هایی می پیوندند که چگونگی تغییر الگوهای مهاجرت پرندگان با آب و هوا یا آن چه مهاجرت پرندگان درباره تغییرات آب و هوا می تواند بگوید را تحلیل می نمایند. این امر نیاز به آشنایی با علم، شهری جامعه ای برخط و محیط گرایی دارد در علوم انسانی به عنوان عقاید تاریخی ممکن است مصاحبه های آلن مک فارلن (2) با متفکران معاصر اصلی که در سخنرانی های (3) iTunesU آکسفورد نمایش داده می شود مقایسه شود تا الگوهای چگونگی تفکر اجتماعی متفکران مهم در مقابل آن هایی که در iTunesU مشهور هستند شفاف شود. (به عبارت دیگر مقایسه افکار اصلی در مقابل با افکار مردمی).

چالش فوری: ارائه راهنمایی برای انواع متفاوتی از ابر مجموعه ها به طوری که عناصر استاندارد در آرشیو وب وارد شوند و در نتیجه دامنه و انسجام آن ها تضمین و تعیین شود.

چالش توسعه یافته: ارائه ابزارها و زیر ساخت های سازمانی برای این که پژوهش گران بتوانند بر تردیدها غلبه یابند تا ابر مجموعه هایی را بسازند که حداقل به صورت بالقوه توسط یک پژوهش گر منفرد مورد استفاده قرار گیرد. استانداردهایی تعریف شود تا ابر مجموعه های آرشیوی توسط یکدیگر قابل استفاده باشند.

چالش بلند مدت: تشویق به پدید آمدن سازمان هایی که ابر مجموعه ها را به طور مفید و گسترده تشویق و حمایت می کنند.

وب دیداری

سؤال: اگر بخواهم از تصاویر برای فهم چگونگی تغییرات دنیا استفاده کنم آیا می توانم برای فهم این فرآیندها به صورت دیداری تصاویر را از آرشیوهای وب استخراج کنم؟ برای مثال: آیا این امکان وجود دارد که با استخراج تصاویر تغییر یافته در طول زمان از مکان های یکسان بر روی صفحه هایی نظیر

ص: 136

Wolfram Alpha -1

Alan MacFarlane -2

iTunes U -3

(1) [Flicker](#) تحلیل های تصویری انجام گیرد؟ تصویر برداری مجدد، تجربه ای است که طی آن مکانی که قبلاً از آن عکس برداری شده بازدید کرده و برای مستند کردن تداوم و تغییرات در طول زمان، عکس جدیدی از آن گرفته می شود. یکی از نخستین پروژه ها برای انجام این کار بازدید مجدد از 1200 سایت در آمریکای غربی و تصویر برداری از صد سال قبل به وسیله پژوهشگران دولتی بوده است. (کلت، (2) منچستر (3) و وریورگ 1984. (4)).

حالا، تصور کنید 100 سال بعد از حال را که قادر باشیم نه تنها عکس های یکسانی را از مکانی مشخص مقایسه کنیم بلکه بتوانیم از آرشیو های وب سری کامل عکس های مستند شده ای را از دنیای در حال تغییر و ثابت در طول زمان استخراج کنیم.

عکس

Flicker¹، تحلیل های تصویری انجام گیرد؟ تصویر برداری مجدد، تجربه ای است که طی آن مکانی که قبلاً از آن عکس برداری شده، بازدید کرده و برای مستند کردن تداوم و تغییرات در طول زمان، عکس جدیدی از آن گرفته می شود. یکی از نخستین پروژه ها برای انجام این کار، بازدید مجدد از ۱۲۰۰ سایت در آمریکای غربی و تصویر برداری از صد سال قبل به وسیله پژوهشگران دولتی بوده است. (کلت، منچستر و وربورگ، ۱۹۸۴)

حالا، تصور کنید ۱۰۰ سال بعد از حال را که قادر باشیم نه تنها عکس های یکسانی را از مکانی مشخص مقایسه کنیم، بلکه بتوانیم از آرشیوهای وب سری کامل عکس های مستند شده ای را از دنیای در حال تغییر و ثابت در طول زمان استخراج کنیم.



تصویر ۲۱. تصویر قصر باکینگهام که از ۲۲۰ تصویر مختلف گردآوری شده است.

منبع: <http://photosynth.net/view.aspx?cid=34e49d3e-2d1e-4118-bbad-d2f5d74ce340>

چالش فوری: تضمین اینکه، تصاویری که بیشتر در صفحه های وب آرشیو شده، از دست می روند، در اولویت حفاظت قرار گیرند.

چالش توسعه یافته: ساخت فناوری هایی همانند photosynth^۲، که قادر به اتصال تعداد زیادی از تصاویر برای نمایش پانورامای یک مکان یا شی برای کار با اطلاعات زمانی به منظور جمع آوری مناظر مشابه در طول زمان باشند.

چالش بلند مدت: ایجاد یک آرشیو از تصاویر جهان، شامل اطلاعات بسیار از زمان و مکان که تا حد امکان برگرفته از EXIF^۳ و داده های صفحه های وب باشند، به طوری که تصاویر بتوانند برای پژوهش استفاده شوند. ابزارهایی برای جای دادن، استخراج، ترکیب و تصاویر مورد نیاز است.

1. Flickr
2. Klett
3. Manchester
4. Verburg
5. <http://photosynth.net/Background.aspx>
6. EXIF

تصویر 21. تصویر قصر باکینگهام که از 220 تصویر مختلف گردآوری شده است.

منبع: <http://photosynth.net/view.aspx?cid=34e49d3e-2d1e-4118-bbad-d2f5d74ce>

چالش فوری: تضمین، این که تصاویری که بیش تر در صفحه های وب آرشیو شده، از دست می روند، در اولویت حفاظت قرار گیرند.

چالش توسعه یافته ساخت فناوری هایی همانند (5) <http://photosynth.net/Background.aspx>. که قادر به اتصال تعداد زیادی از

تصاویر برای نمایش پانورامای یک مکان یا شی برای کار با اطلاعات زمانی به منظور جمع آوری مناظر مشابه در طول زمان باشند.

چالش بلند مدت: ایجاد یک آرشیو از تصاویر جهان شامل اطلاعات بسیار از زمان و مکان که تا حد امکان برگرفته از [EXIF 6](#) و داده های صفحه های وب باشند به طوری که تصاویر بتوانند برای پژوهش استفاده شوند. ابزارهایی برای جای دادن استخراج ترکیب و تصاویر مورد نیاز است.

ص: 137

Flicker -1

Klett -2

Manchester -3

Verburg -4

photosynth -5

EXIF -6

سؤال: چگونه من می توانم وب را همان گونه که بود ببینم؟ اگر من بخواهم وب را همان گونه که، در یکم ژانویه 2011 بود مرور کنم و قادر باشم تا روی صفحه ها، تصاویر پیوندها و دیگر محتوای آن همان گونه که در آن روز ظاهر شده بود کلیک کنم چگونه این کار را می توانم انجام دهم؟

عکس

۱۲۸ مدیریت منابع اطلاعاتی وب

وب همان گونه که بود

سؤال: چگونه من می توانم وب را همانگونه که بود ببینم؟ اگر من بخواهم وب را همانگونه که، در یکم ژانویه ۲۰۱۱، بود مرور کنم و قادر باشم تا روی صفحه ها، تصاویر، پیوندها و دیگر محتوای آن همانگونه که در آن روز ظاهر شده بود، کلیک کنم، چگونه این کار را می توانم انجام دهم؟



تصویر ۲۲- ویرایش بنای پاسخ WayBack Machine

(<http://replay.web.archive.org/20041010185532/http://netpreserve.org/about/index.php>)

ویرایش بنای فعلی نسخه پاسخ [پایگاه] Wayback Machin چنین کارکردی را نوید می دهد (گشت و گذار در وب همانگونه که بود، نسخه بتا، که در حال حاضر در سایت تبلیغ می شود)، اما IIPC یا آرشیوهای انفرادی چه تلاش های دیگری می توانند برای امکان تکرارپذیری وب، انجام دهند؟ چالش: گسترش و افزایش تلاش هایی برای ایجاد وب فعلی به صورت تکرارپذیر در آینده که نیاز به گردآوری، ذخیره سازی و اشاعه مجدد وب فعلی همانگونه که در گذشته بود، خواهد داشت. سؤال محوری برای پرسش در اینجا این است که [وب] چگونه می تواند فراتر از این که منبعی ارجاعی صرفاً برای رفع کنجکاوی یا مراجعه گاه به گاه باشد، مورد استفاده قرار گیرد. چه نوع نیاز بکر یا سؤال پژوهشی غیر متصور بر جست و جوی های دستی و گشت و گذار در وب گذشته (قدیمی) اتکا دارد؟ آیا مورخان آینده، همانگونه که در دنیای امروز خواندن اخبار و انتشارات و دوران گذشته فانی وب مورد علاقه بوده است، به آن علاقه مند هستند؟ آیا آنها می خواهند از نرم افزارهای کاربردی استفاده کنند یا صرفاً از منابع مرجع؟ به عبارت دیگر، بزرگترین سؤال، ساختن موارد استفاده برای وب تکرارپذیر و سپس ساختن رابط هایی که این موارد استفاده را پشتیبانی می کند، می باشد. انجام این امر مستلزم مشورت متخصصان آرشیو سازی وب با متخصصان حوزه هایی شامل مورخان و دیگر افرادی است که به بازسازی مجدد گذشته علاقه مند هستند.

(منبع <http://replay.web.archive.org/20041010185532/http://netpreserve.org/about/index.php>)

ویرایش بتای فعلی نسخه پاسخ [پایگاه] Wayback Machin چنین کارکردی را نوید می دهد (گشت و گذار در وب همان گونه که ، بود نسخه بتا که در حال حاضر در سایت تبلیغ می شود) اما IIPC یا آرشیو های انفرادی چه تلاش های دیگری می توانند برای امکان تکرار پذیری وب، انجام دهند؟

چالش: گسترش و افزایش تلاش هایی برای ایجاد وب فعلی به صورت تکرار پذیر در آینده که نیاز به گردآوری، ذخیره سازی و اشاعه مجدد وب فعلی همان گونه که در گذشته بود، خواهد داشت. سؤال محوری برای پرسش در این جا این است که [وب] چگونه می تواند فراتر از این که منبعی ارجاعی صرفاً برای رفع کنجکاوی یا مراجعه گاه به گاه باشد مورد استفاده قرار گیرد چه نوع نیاز بکر یا سؤال پژوهشی غیر متصور بر جست و جوی دستی و گشت و گذار در وب گذشته (قدیمی) اتکا دارد؟ آیا مورخان آینده ، همان گونه که در دنیای امروز خواندن اخبار و انتشارات و دوران گذشته فانی وب مورد علاقه بوده است ، به آن علاقه مند هستند؟ آیا آن ها می خواهند از نرم افزارهای کاربردی استفاده یا صرفاً از منابع مرجع؟ به عبارت دیگر بزرگ ترین سؤال ساختن موارد استفاده برای وب تکرار پذیر و سپس ساختن رابطه ای که این موارد استفاده را پشتیبانی می کند می باشد. انجام این امر مستلزم مشورت متخصصان آرشیو سازی وب با متخصصان حوزه هایی شامل مورخان و دیگر افرادی است که به بازسازی مجدد گذشته علاقه مند هستند

سؤال: وب چگونه مقایسه شده و در طول زمان تغییر می کند؟

افزایش تلاش هایی برای فهم وب به صورت یک سیستم مستلزم قابلیت های افزایش به مقیاس بزرگ است که قادر خواهد بود الگوها و روندها را در طی زمان دست بی اندازد (تغییر دهد). برای انجام این امر، ما نیاز داریم که پرسیم چه رهیافت هایی برای توسعه روش های تحلیلی معتبر در دسترس هستند؟ چگونه ما می توانیم فرضیات ساخته شده درباره داده های وب را به صورت یک سری داده اعتبار سنجی کنیم؟ چه ابزارهای آماری برای مجموعه های وب آرشیوی می تواند به کار گرفته شود و چه ابزارهایی نیاز به توسعه دارند؟

حتی آمارهای ساده نیز برای استخراج درباره وب بی اهمیت نیستند برای نمونه، چه تعداد از وب گاه ها به طور سالانه (سراسر جهان در یک کشور مشخص روی موضوع مشخص) برای X تعداد سال گذشته بوده اند؟

در مجموعه آرشیوی (آرشیو در جعبه) تاریخ ایجاد صفحه ها چیست؟ صفحه های آن به چه زبان هایی هستند؟ آیا در زمان ایجاد صفحه ها روندهایی وجود داشته است؟ آیا خوشه ای است؟ آیا یک فرایند ساختمانی ثابت بوده است؟ آیا موضوع های مشخصی بیش تر از دیگران پیوند داده شده اند؟ آیا برخی از انواع مجموعه ها احتمال کم تر یا بیش تری برای پیوند به منابع خارجی دارند؟ آیا وب گاه ها می توانند به رده هایی که می توانیم با استفاده از تحلیل های خوشه ای کشف کنیم تقسیم بندی شوند؟

آیا ما می توانیم سایت ها را به وسیله آمار هایی همچون اندازه میانگین وب گاه ها در رده های مختلف، میانگین تعداد پیوند ها میزان داده های غیر متنی (عکس ها، تصاویر و مانند آن) سن محتوا در سن محتوا در فاصله بین روز آمدسازی ها دفعات روز آمدسازی، نوع رابط (ثابت در مقابل دینامیک به طور مثال)، مقایسه کنیم؟

چگونه این آمارها به ما در فهم ساختار مجموعه ها و وب کمک می کنند؟

چالش: ایجاد ابزارها و روش هایی برای استفاده از وب به عنوان یک رشته داده عظیم به جای مجموعه ای از اسناد در حال حاضر در صورت وجود داده ها اگر شخص بخواهد بداند چه چیزی باید به طور بسیار اساسی درباره اندازه و ساختار وب فعلی یا وب به صورت گذشته اش مورد سؤال باشد داده ها برای پژوهش گران قابل دسترسی نیستند. بنابراین ابزارهایی باید ایجاد شوند که آمارگیری روی وب یا روی صفحه های یک مجموعه آرشیوی را انجام دهند.

ایده ها چگونه تکثیر می شوند

سؤال: ایده ها چگونه روی اینترنت کشش یافته و تکثیر می شوند؟ یکی از جنبه های قابل توجه اینترنت توانایی شگفت انگیز آن برای پشتیبانی انتقال از الگوی رفتاری ایده هایی که رشد می یابند و گسترش فرهنگی است اگر علاقه مند هستیم به چگونگی ویروسی شدن ویدئو ها یا چگونگی گسترش شوخی ها یا چگونگی وارد شدن یک بیت از اطلاعات یا اطلاعات غلط به آگاهی عمومی، چه ابزارهایی به ما کمک خواهد نمود؟ چگونه ابزارهایی را با توانایی امکان این که یک آرشیو ساخته شود نه بر اساس

جغرافیایی فیزیکی یا مجازی بلکه بر اساس حرکت از یک ایده می‌سازیم؟ هر کس می‌تواند تصور کند که قادر به تعیین یک ایده باشد و برای دنبال نمودن آن ایده به صورتی که در طول زمان توسعه می‌یابد، حرکت کند

هم چنین چه زمینه گسترده‌ای در اطراف محتوایی که ما در آرشیو می‌بینیم، وجود دارد؟ برای مثال: مردم روی وب در زمان پیدایش آن چه چیزی را جست و جو می‌کردند؟ گوگل زیتگیست (1) و گوگل ترندز (2) در مورد چیزهایی که مردم درباره آن جست و جو می‌کردند به ما می‌گویند. چه چیزهایی دیگری را می‌توانیم برای فهم زمینه جمع‌آوری کنیم؟ برای مثال- در تویتر اشاره شده است - به فهم زمینه محتوا بوسیله مشاهده چیزهایی که با یکدیگر آمده‌اند، کمک می‌کند.

برای نمونه واتسون آی بی ام (3)، نظام اختصاصی است که از مواد آرشیوی زیادی برای ساخت موتور دیپ کیو ای (4) خودش استفاده می‌کند که به آن برای برنده شدن در بحران 2011 کمک می‌کند. چگونه این نوع از ابزار به طور گسترده تری برای جست و جوی پیشرفته در دسترس خواهند شد؟

چالش: بعد زمان وب همان گونه که ایجاد شده نیاز دارد که حفاظت شده قابل استخراج و قابل تحلیل، باشد گرانولیت بهتری نیاز است تا بیند ایده‌ها از کجا آغاز شده، چگونه گسترش می‌یابند و این که چه فعالیت‌هایی سرعت تکثیر آن‌ها را کاهش یا افزایش می‌دهد؟ ایده‌ها به صورت موجودی زنده در دنیا پدید می‌آیند و تنها پس از این که نگهداری، شوند برای برگشت به اصلیت شان مورد توجه قرار خواهند گرفت با این حال بدون گرانولیت و عمق مناسب آرشیوی ایده اصلی ممکن است در زمانی که شخص به جست و جوی آن فکر می‌کند، از دست رفته باشد.

وب غیر قانونی

سؤال: چگونه وب برای پشتیبانی و توانایی فعالیت‌های غیر قانونی استفاده می‌شود و چگونه در طول زمان تغییر می‌کند؟ نوعی محتوا که به صورت برخط تکثیر می‌شود و کم‌تر مورد توجه دانشمندان است، مواد غیر قانونی وب است این طیف گسترده‌ای از محتوای جنسی تا اطلاعاتی درباره استفاده از مواد، مخدر قمار منابع گروه‌های تندرو محتوای مرتبط با تروریسم و دیگر منابعی است که یا غیرقانونی یا دارای مشکلات اجتماعی هستند سؤال این جاست که چه کسی باید محتواهای غیر قانونی یا قانونی ولی کم‌تر از نظر اجتماعی پذیرفته شده وب را آرشیو کند؟ چگونه بدون شکستن قانون می‌توانیم این کار را انجام دهیم؟ و چگونه می‌توان برای محققان بدون در خطر افتادن هم پژوهش‌گر و هم مؤسسه‌هایی که دسترسی را ایجاد ساخته‌اند قابل دسترسی شود؟

دانستن این که چه فعالیت‌های غیر قانونی‌ای از نظر حجم و عمومیت یافتن در حال رشد هستند کدام یک در طول زمان رنگ می‌بازند و چه فعالیت‌های غیرقانونی غیر منتظره‌ای پدید می‌آیند نه فقط

ص: 140

Google Zeitgeist (<http://www.google.com/press/zeitgeist> 2010/) -1

Google Trends (<http://www.google.com/trends>) -2

IBM's Watson -3

Deep QA (<http://www.research.ibm.com/deepqa/deepqa.shtml>) -4

برای پژوهش گران، بلکه برای سیاست گذاران عمومی، متخصصان حوزه سلامت که احتیاج به پیگیری نتایج رفتارهای خطر آفرین دارند خبرگان سلامت عمومی، متخصصان و مؤسسه های حمایت اجتماعی و کسانی که مسئول حفظ رفاه قشر آسیب پذیر هستند مفید می باشد

چالش: بزرگ ترین چالش در این جا این است که حتی اگر منابع غیر قانونی بر روی اینترنت رواج داشته باشند سازمان های اندکی هستند که خارج از اعمال فشار قانونی مایل به پذیرش خطر مشارکت در جمع آوری داده های مربوط به این منابع هستند. هستند تابوهای فرهنگی و خطرات قانونی در ارتباط با دسترسی و ذخیره سازی منابع غیر قانونی حتی برای اهداف مثبت همانند پژوهش فهم جنبه های جامعه مدرن بیشتر پژوهش گران وب و آرشیویست های وب را از چنین منابعی بر حذر می دارد.

به نظر می رسد مهم ترین سازوکاری که می تواند در این جا وجود داشته باشد نظامی با حفاظت قانونی خواهد بود که به طور مناسب و تا حد امکان افراد و سازمان ها را برای آرشیو و جست و جوی داده های غیرقانونی قابل دسترسی اینترنت بدون در معرض خطر قرار گرفتن سازمان ها یا پژوهش گران استفاده کننده از این مجموعه ها تأیید کند

رد پای رقومی

سؤال: چگونه می توانیم (و باید) رد پای رقومی یک شخص را آرشیو کنیم؟ فعالیت ها و اقدامات یک شخص به طور پیوسته مورد توجه بالقوه است به ویژه اگر شخص (مشهور باشد یا بشود) قبلاً در این مورد بحث شده است که (گارفینکل (1) و کاکس 2009 (2) فهرست آثار زندگی یک شخص برای آرشیویست ها یک وظیفه خواهد بود آرشیو وب یک شخص می تواند شامل صفحه های وب پروفایل های شبکه اجتماعی و پست هایش، ارتباطاتش، انتشارات، و منابع دیگر در مورد زندگی رقومی او باشد.

چالش: چالش محوری، فهمیدن چگونگی ساختن ابزارهایی است که به افراد اجازه می دهد تا به صورت دستی مشخص کنند که چگونه به صورت خودکار آثار رقومی خود را جمع آوری کنند. آیا ابزارها به طور خودکار ردپای رقومی را تا حد امکان بر مبنای انتخاب، جمع آوری می کنند؟ چگونه می توانیم نظام هایی را ایجاد کنیم که نه فقط یادآوری بلکه امکان فراموشی به وسیله حذف های بعدی توسط افراد (و فراموشی) بخش ها یا همه رد پاها را داشته باشد به صورتی که برخی دانشمندان از آن به عنوان یک حق اساسی ذکر کرده اند مهیر (3) شونبرگر 2009 (4)؟ پیشرفت ها به ویژه در این حوزه مشکل است به دلیل این که مسائل بسیار زیاد حقوقی و خصوصی آن نیاز به بررسی خواهد داشت.

ص: 141

Garfinkel -1

Cox -2

Meyer -3

Schonberger -4

سؤال: چگونه داده‌ها از آرشیوهای وب دوباره استخراج می‌شوند؟ در سال‌های اخیر، رشد چشم‌گیری در وب مبتنی بر داده (در مقابل وب اسنادی) صورت گرفته است مسائل متعددی مطرح شده است، مانند این که چگونه داده‌ها در کنار اسناد آرشیو بشوند؟ چه نوع ابزار پاک‌کننده داده‌ها برای کار با داده‌ها مورد نیاز خواهد بود؟

برای مثال به فرآیندهای بین حوزه‌های علمی یا صنعتی برای فراهم‌آوری دانش مانند طراحی، هوافضا، کشف مواد مخدر و غیره، فکر کنید زمانی که یک هواپیما دچار سانحه می‌شود و بازرسان می‌خواهند محاسبات اولیه مهندسی را دوباره ارزیابی کنند چگونه می‌توانیم توانایی درک داده‌ها از یک زنجیره تأمین طراحی مهندسی هفتاد ساله را حفظ کنیم؟ طراحی رقومی، بود دانش توسط 100 نفر همکار تدارک می‌شد که بعضی از آن‌ها در فواصل زمانی از حوزه کار خارج شده‌اند و همه آن توسط مجموعه کاملاً متفاوتی از افراد به کار گرفته می‌شد در این مورد چرخه زندگی دانش بسیار طولانی‌تر از چرخه زندگی کسب و کار است و آرشیوها نقش اساسی در حفظ این اطلاعات ایفا می‌نمایند.

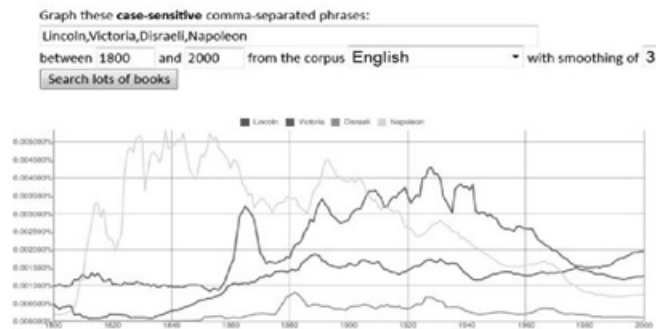
سؤال‌های مرتبط این است که چگونه می‌توانیم اطلاعات اختصاصی را آرشیو و تحلیل کنیم؟ بعضی از داده‌ها اختصاصی هستند و لازم است برای تعهد به موافقت نامه‌های حقوقی نگهداری شوند. با وجود، این اگر نهادهای قابل، اعتماد داده‌های خام را ذخیره کنند و فقط به ابزارهای تحلیلی اجازه دسترسی به آن‌ها را بدهند تحلیل‌های بسیاری بدون دسترسی به داده‌های خام می‌توانند انجام گیرند. سپس نتایج گردآوری و تلخیص شده می‌توانند بدون افشای داده‌های حفاظت شده در دسترس پژوهش‌گر قرار گیرند. ابزارهایی برای سرویس‌دهی به عنوان الگویی که امکان دسترسی به بخش‌هایی از داده‌های خام را می‌دهند به جای اینکه داده‌های خام را بارگذاری کنید وجود دارند همانند الیکسر، (1) برنامه‌ای برای دسترسی به اطلاعات نجوم مثال، دیگر مشاهده گر n-gram کتاب‌های گوگل (2) که به کاربران امکان تحلیل داده بدون دسترسی به داده خام را می‌دهد.

اگر یکبار به داده‌هایی از مخازن داده‌های قابل، تحلیل دسترسی داشته باشید، امکانات افزوده فزونی خواهند یافت. همانند پیوند داده‌ها از طریق، ابزارها ایجاد کتابخانه‌هایی متشکل از اجزایی که به پژوهش‌گران اجازه می‌دهند که داده‌ها را از راه‌های گوناگون تحلیل کنند و امکاناتی برای ترکیب و تطابق آن‌ها را به روش‌های جدید فراهم می‌کنند. اگر در ایجاد این سری داده‌ها و ابزارها، دقت شود، الگوها و همبستگی‌هایی که قبلاً کشف نشده‌اند، افزایش می‌یابد.

ص: 142

آینده آرشیو وب ۱۴۳

Google labs Books Ngram Viewer



تصویر ۲۳- مشاهدهگر n-gram کتابهای گوگل

چالش: برای بخش‌هایی از اینترنت که حاوی داده‌ها به جای اسناد هستند، تغییر نگرش از اینکه آرشیوهای وب برای نگهداری اسناد هستند به سمت مشاهده آنها به صورت روشی برای آرشیوسازی داده‌هایی است که این اسناد در بردارند. این امر مستلزم مدل‌های جدیدی برای ذخیره‌سازی و استخراج داده‌هایی است که از مدل‌های متمایل به تحلیل‌های داده‌های داده‌های ساختار یافته، به جای تحلیل‌های داده‌های ساختار نیافته، پیروی می‌کنند.

وب‌های ملی: چه ارزشی در ایجاد یک وب ملی زمانی که وب پدیده‌ای فراتر از مرزهاست، وجود دارد؟ تلاش‌هایی در حوزه آرشیو- در جعبه- برای اجرا در سطوح ملی، با توجه به محدودیت‌های بودجه‌ای و قانونی، اساسی است. در حال حاضر، در بریتانیا، کتابخانه بریتانیا، خودش را برای پیش‌بینی قوانینی تأثیرگذار بر روی آرشیوسازی فضای وب بریتانیا به صورت کتابخانه واسپاری برای همه انتشارات بریتانیا آماده می‌کند. در می ۲۰۱۱، وزارت پژوهش و اختراع دانمارک، تصمیم به ایجاد زیرساخت‌های پژوهشی ملی در حوزه‌های پژوهشی مختلف (علوم طبیعی، علوم انسانی) که متمرکز بر استفاده از ابزارهای تحلیلی منابع وب آرشیوی خواهد بود.

چالش فوری: عدم شفافیت نحوه استفاده پژوهشگران از آرشیوهای ملی، بسیاری از آنها هنوز در مراحل طراحی هستند. ما در مورد یکی از مهم‌ترین مواردی که باید انجام شود، یعنی تعامل دامنه پژوهشگران با تخصص‌هایی نه فقط در پژوهش اینترنت، بلکه در زمینه‌هایی نظیر جامعه‌شناسی، علوم سیاسی، و دیگر حوزه‌های علوم اجتماعی، فیزیک و سایر علوم، هنر و علوم انسانی، بحث خواهیم کرد. همان‌طور که این زیرساخت‌ها به منظور انعکاس نیازهای پژوهشگران ملی در مجموعه‌های ایجاد شده، طراحی می‌شوند. این فرآیندی زمان‌بر است و تعامل خبرگان هر حوزه می‌تواند مشکل باشد. با وجود

تصویر 23- مشاهده گر n-gram کتاب های گوگل

چالش: برای بخش‌هایی از اینترنت که حاوی داده‌ها به جای اسناد هستند، تغییر نگرش از این که آرشیوهای وب برای نگهداری اسناد هستند به سمت مشاهده آن‌ها به صورت روشی برای آرشیوسازی داده‌هایی است که این اسناد در بردارند. این امر مستلزم مدل‌های جدیدی برای ذخیره‌سازی و استخراج داده‌هایی است که از مدل‌های متمایل به تحلیل‌های داده‌های ساختار یافته، به جای

وب های ملی: چه ارزشی در ایجاد یک وب ملی زمانی که وب پدیده ای فراتر از مرزهاست وجود دارد؟ تلاش هایی در حوزه آرشیو - در جعبه-برای اجرا در سطوح، ملی با توجه به محدودیت های بودجه ای و، قانونی اساسی است در حال حاضر، در بریتانیا، کتابخانه بریتانیا خودش را برای پیش بینی قوانینی تأثیر گذار بر روی آرشیوسازی فضای وب بریتانیا به صورت کتابخانه واسپاری برای همه انتشارات بریتانیا آماده می. کند در می، 2011 وزارت پژوهش و اختراع دانمارک، تصمیم به ایجاد زیر ساخت های پژوهشی ملی در حوزه های پژوهشی مختلف (علوم طبیعی علوم انسانی) که متمرکز بر استفاده از ابزارهای تحلیلی منابع وب آرشیوی خواهد بود.

چالش فوری: عدم شفافیت نحوه استفاده پژوهش گران از آرشیو های ملی. بسیاری از آن ها هنوز در مراحل طراحی هستند ما در مورد یکی از مهم ترین مواردی که باید انجام شود، یعنی تعامل دامنه پژوهشگران با تخصص هایی نه فقط در پژوهش، اینترنت بلکه در زمینه هایی نظیر جامعه شناسی، علوم سیاسی و دیگر حوزه های علوم اجتماعی فیزیک و سایر علوم هنر و علوم انسانی، بحث خواهیم کرد. همان طور که این زیرساخت ها به منظور انعکاس نیازهای پژوهشگران ملی در مجموعه های ایجاد شده، طراحی می شوند. این فرآیندی زمانبر است و تعامل خبرگان هر حوزه می تواند مشکل باشد. با وجود

این شکست اجرای آن احتمال کسب کاربرد گسترده زیر ساخت های جدید را کاهش می دهد.

نتایج: مسیر پیش رو

آن چه گفته شد، فقط تعداد معدودی از مواردی است که گروه کوچک ما توانست به آن فکر کند و موارد بسیاری هنوز وجود دارد. بعضی از مباحث عمومی است زیرا فنون گام به گام ویژه ای برای مشخص نمودن مجموعه وب یا روندهای تحلیلی از اینکه چگونه محتوای وب در طول زمان تغییر می کند نیاز به منابعی از پروژه تحقیقاتی خاص، یک تیم از متخصصان حوزه و مجموعه های مرتبط برای آزمون این روش ها خواهد داشت، بنابراین ما گوی جادویی نخواهیم داشت با این حال تعدادی از چالش های عمومی را نشان داده ایم که پژوهشگر علاقه مند به کار در حوزه آرشیوهای وب با آن مواجه است یکی از موارد اصلی که در طول زمان حاصل می شود و دوباره در مصاحبه ها و بحث ها وجود دارد، فقدان پایداری و رابطه ای کاربر پسند برای ایجاد آرشیو های وب و یکبار ساخته شده به منظور دسترسی و تحلیل داده های موجود در آن هاست.

در دراز مدت، ما امیدوار به نتایج دیگری هستیم که ممکن است از این تلاش ها حاصل شود. برای مثال، می توانیم گروه کاری پست هاگ (1) را تصور کنیم که برای توسعه کارگاه بحث و ایده های مورد استفاده احتمالی در آینده و تمرکز بر توسعه ابزاری ایجاد شد. این گروه نه تنها برای اعضای IIPC، بلکه برای انواع جوامع پژوهش گرانی که با جامعه آرشیوی وب در تعامل نبوده و بیش تر تمایل به استفاده از وب آرشیوی را دارند آگاهی رسانی خواهد کرد. مثال ها شامل پژوهش گران اینترنت (مانند اعضای 2) AOIR) دانشمندان اطلاعات (مانند 3) IFIP و 4) ASIS T) و طیف فهرست ها و انجمن های علاقه مند به علوم انسانی رقومی هستند. ما قویاً توصیه می کنیم که IPC نمایندگانی را به نشست های سالانه سازمان هایی از این نوع و پانل ها رهنمون می سازد و کارگاه هایی را برای تعامل پژوهشگران با امکانات آرشیوهای وب سازماندهی کند. ما چند ایده را به صورت برجسته مطرح ساخته ایم اما آن جوامع موارد بیش تری را رانند آماده نمایند. برای آمدن آن ها به IIPC منتظر نمایند. IPC باید به سمت آن ها برود.

ایده دیگر برای آینده فعالیت، نوعی برنامه مشترک کدنویسی نرم افزاری (5) است که برنامه نویسان رایانه ای هکرها با همدیگر و با پژوهشگران برای 2 الی 3 روز گردهم آمده و به ایجاد دسترسی به داده های آرشیوی وب توجه و همکاری می کنند آن ها باید بتوانند به گروه ها و وظایف مرتبط با کشف رهیافت های خلاقانه و ابتکاری توجه کنند تا بتوانند داده ها و ابزارهای موجود را برای ایجاد سریع ابزارها و رابطه ای جدید، به کار گیرند پژوهش گران انتخاب خواهند شد به دلیل این که آن ها سؤالاتی دارند

ص: 144

1- (گروهی کاری در لاهه) post-Hague

2- AoIR)

3- IFIP

4- ISAST

5- Hackathon: برنامه ای است که طی آن تعدادی از برنامه نویسان به صورت مشترک و فشرده در یک دوره کوتاه زمانی اقدام به کدنویسی برای برنامه نرم افزاری خاص می نمایند منبع قابل دسترسی در:

<http://www.techopedia.com/definition/23193/hackathon>

که تمایل به پاسخ‌گویی به آن‌ها دارند و برنامه نویسان رایانه‌ای همه مهارت‌شان را برای کمک به آن‌ها در رسیدن (نیل یا نزدیک شدن) به اهداف پژوهشی‌شان به کار می‌برند دوباره برنامه نویسان رایانه‌ای که با داده‌های وب پویا کار می‌کنند، مهارت‌هایی برای اجرای طرح‌های خلاقانه با ابزارها دارند که نتایج بسیار بیشتری از انتظار برای درخواست از آرشیوهای وب به سمت آن‌ها گسیل خواهد کرد.

آیا ما به نیروانا خواهیم رسید یا محکوم به آخر الزمان خواهیم شد، شد یا با منحصر به فردی (فنی) جایگزین خواهیم شد یا اینکه با آرشیوهای غبار آلوده مواجه خواهیم شد؟ هیچ راهی برای دانستن آن نداریم با این حال ما در نقطه‌ای هستیم که این سؤال را ایجاب می‌کند امروز برای اطمینان از این که در آینده به چه چیزی می‌توانیم دسترسی داشته باشیم چه گام‌هایی می‌توانیم برداریم؟ تجمیع تصمیمات به ندرت قابل توجه، ساده، نبود اما بخشی از تلاشی است برای تضمین این که آرشیوهای وب قوی پایدار قابل دسترسی، ارزشمند و بالاتر از همه توسط پژوهش‌گران آینده قابل استفاده باشند.

منابع

Conover, M. D., Ratkiewicz, J., Francisco, M., Gonçalves, B., Flammini, A., Menczer, F. (2011). .1 Political Polarization on Twitter. Paper presented at the ICWSM: International Conference on Weblogs and Social Media 2011, Barcelona

Conover, M. D., Ratkiewicz, J., Gonçalves, B., Flammini, A., Menczer, F. (2011). The Echo Chamber. .2 Paper presented at the Journal of Information Technology Politics Conference 2011: The Future of Computational Social Science, Seattle

Dougherty, M., Meyer, E. T., Madsen, C., Van den Heuvel, C., Thomas, A., Wyatt, S. (2010). Researcher .3 Engagement with Web Archives: State of the Art. Report. London: JISC. Retrieved from <http://ssrn.com/abstract=1714997> and <http://ie-repository.jisc.ac.uk/544>

Garfinkel, S. Cox, D. (2009, 9–11 February). Finding and Archiving the Internet Footprint. Paper .4 presented at the First Digital Lives Research Conference: Personal Digital Archives for the 21st Century, London

(Gazan, R. (2008). Social annotations in digital library collections. D-Lib Magazine, 14(11/12) .5

Hindman, M. (2007). "Open-source politics" Reconsidered: Emerging Patterns in Online Political .6 Participation. In V. Mayer-Schönberger D. Lazer (Eds.), Governance and information technology: From electronic government to information government (pp. 183–207). Cambridge: The MIT Press

- Hogan, B. (2010). Analyzing Facebook Networks. In D. Hansen, M. Smith B. Schneiderman (Eds.), . 7
Analyzing Social Media Networks with NodeXL. New York, NY: Morgan Kaufman
- Jasra, M. (2011, 3 February). Reddit Surpasses 1 Billion Monthly Page Views Retrieved 30 April, 2011, .8
from <http://www.webanalyticsworld.net/2011/02/reddit-surpasses-1-billion-monthly-page.html>.
(Archived by WebCite® at <http://www.webcitation.org/5yKdMBKNC>)
- Kay, A. (1995). The Best Way to Predict the Future is to Invent it. *Mathematical Social Sciences*, 30, .9
.326-326
- Klett, M., Manchester, E., Verburg, J. (1984). *Second View: The Rephotographic Survey Project*. . 10
.Albuquerque: University of New Mexico Press
- Kling, R., McKim, G., King, A. (2003). A Bit More to IT: Scholarly Communication Forums as Socio- .11
Technical Interaction Networks. *Journal of the American Society for Information Science and Technology*,
.54(1), 46-67
- .Kurzweil, R. (2005). *The Singularity is Near: When Humans Transcend Biology*. New York: Viking .12
- Mayer-SchÖnberger, V. (2009). *Delete: the virtue of forgetting in the digital age*. Princeton, NJ: . 13
.Princeton Univ Press
- Meyer, E. T. (2006). Socio-technical Interaction Networks: A discussion of the strengths, weaknesses .14
and future of Kling's STIN model. In J. Berleur, M. I. Numinem J. Impagliazzo (Eds.), *IFIP International
Federation for Information Processing, Volume 223, Social Informatics: An Information Society for All?* In
.Remembrance of Rob Kling (pp. 37-48). Boston: Springer
- Meyer, E. T. (2011). *Splashes and Ripples: Synthesizing the Evidence on the Impact of Digital* . 15
.Resources. Report. London: JISC. Retrieved from <http://ssrn.com/abstract=1846535>
- Meyer, E. T., Eccles, K., Thelwall, M., Madsen, C. (2009). *Final Report to JISC on the Usage and* . 16
*Impact Study of JISC-funded Phase 1 Digitisation Projects the Toolkit for the Impact of Digitised Scholarly
Resources (TIDSR)*. Retrieved from [http://microsites.oi.ox.ac.uk/tidsr/system/files/TIDSR-Final-Report-
20July2009.pdf](http://microsites.oi.ox.ac.uk/tidsr/system/files/TIDSR-Final-Report-20July2009.pdf)
- .Moretti, F. (2005). *Graphs, Maps, Trees: Abstract models for a literary history*. London: Verso Books .17

- .Moretti, F. (2011). Network Theory, Plot Analysis. *New Left Review*, 68, 80–102 .18
- Olston, C., Reed, B., Srivastava, U., Kumar, R., Tomkins, A. (2008). Pig Latin: A Not- So-Foreign .19
Language for Data Processing. Paper presented at the ACM SIGMOD'08 Conference, Vancouver, BC,
.Canada
- Schroeder, R. (2011). *Being There Together: Social Interaction in Shared Virtual Environments*. New .20
.York, NY: Oxford University Press USA
- Schroeder, R. Meyer, E. T. (2009). An Emerging Global Brain: How the Internet is Revolutionising .21
Scientific Research. *Britain in 2009 (Economic Social Research Council Annual Magazine)*, 113.27
- Tanner, S. (2010). *Inspiring Research, Inspiring Scholarship*. Report. London: JISC. Retrieved from .22
<http://www.jisc.ac.uk/media/documents/programmes/digitisation/12pag>
.efinaldocumentbenefitssynthesis.pdf
- Tanner, S. Deegan, M. (2011). *Inspiring Research, Inspiring Scholarship: The value and benefits of .23*
digitised resources for learning, teaching, research and enjoyment. Report. London: JISC. Retrieved from
[http://www.kdcs.kcl.ac.uk/fileadmin/documents/ Inspiring-Research-Inspiring-Scholarship-2011-Simon](http://www.kdcs.kcl.ac.uk/fileadmin/documents/Inspiring-Research-Inspiring-Scholarship-2011-Simon)
.Tanner.pdf
- Thomas, A., Meyer, E. T., Dougherty, M., Van den Heuvel, C., Madsen, C., Wyatt, S. (2010).. 24
Researcher Engagement with Web Archives: Challenges and Opportunities for Investment. Report. London:
.JISC. Retrieved from <http://ssrn.com/abstract=1715000> and <http://ie-repository.jisc.ac.uk/543>
- van den Heuvel, C. (2009). MAPS: Manuscript Map Annotation and Presentation System: Linking .25
formal ontologies with social tagging to (re-) construct relationships between manuscript maps and
contextual documents. *Digital Humanities 2009 (University of Maryland, Maryland Institute for*
.Technology in the Humanities (MITH) Abstracts), 138– 141
- Von Ahn, L., Maurer, B., McMillen, C., Abraham, D., Blum, M. (2008). reCAPTCHA: Human-Based .26
.Character Recognition via Web Security Measures. *Science*, 321(5895), 1465–1468
- Williams, D., Yee, N., Caplan, S. E. (2008). Who plays, how much, and why? Debunking the .27
stereotypical gamer profile. *Journal of Computer-Mediated Communication*, 13(4), 993–1018. doi:
10.1111/j.1083-6101.2008.00428.x

فصل دوم: تجارب جهانی و مسائل بومی در آرشيو سازی وب

اشاره

ص: 149

کتابخانه ملی، استرالیا مؤسسه راهبر آرشیو و حفاظت رقومی در استرالیاست. آرشیو پاندورا (1)، که بیش از 10 سال مخزن آرشیو منابع وب استرالیا بوده است، سیستمی کامل و در حال پیشرفت و توسعه مداوم است. پانداس PANDAS، سیستم مدیریت، آرشیوی که آرشیو را پشتیبانی می کند از سال 2007، در حال سومین بازبینی خود. می باشد فعالیت های دیگر آرشیو وب از جمله گردآوری سالانه دامنه استرالیا و استفاده از آرشیو - آی تی (2) با همکاری آرشیو اینترنت اداره می شوند.

این مقاله به بررسی وضعیت کنونی آرشیو وب در استرالیا می پردازد و نشان می دهد که چگونه کتابخانه ها خدمات خود را با ورود فزاینده مواد برخط در مجموعه های شان وفق می دهند سالیان متمادی تصور می شد که با آرشیو کردن فقط می توانیم نمونه ای کوچک اما روایت گر از اینترنت را ضبط کنیم. امروزه شکاف بین آن چه موجود است و آن چه می تواند آرشیو شود در حال کاهش است. درعین حال هر چه آرشیو ها بیشتر می شوند و توانایی ما در آرشیو کردن افزایش می یابد بیش تر با فناوری های جدید و برنامه های کاربردی وب 2/0 درگیر خواهیم بود به عنوان مثال در انتخابات فدرال سال 2007 که تعداد وسیعی از سایت های متعامل نظیر MySpace Kevin و YouTube آرا را آرشیو کردند، معلوم شد که آرشیوکننده های استرالیایی وب به سازگاری ادامه می دهند و با چالش های جدید روبه رو می شوند.

ص: 150

*آرشیو وب در دنیای وب 2/0 شعبه آرشیو وب و حفاظت رقومی کتابخانه ملی استرالیا (1)

ادگار کروک (2) | ترجمه: مرجان هادی زاده (3)

مقدمه

مقاله حاضر آرشیو وب در استرالیا را مورد بحث و بررسی قرار می دهد. از آن جا که کتابخانه ملی استرالیا نقش اصلی را در این زمینه دارد و همچنین با توجه به این که نویسنده خود در آن جا شاغل است این مقاله می تواند تأثیر زیادی بر مسائل این مؤسسه داشته باشد. لازم به ذکر است که در استرالیا پروژه های آرشیو وب دیگری نیز وجود دارند نظیر پروژه Our Digital Island در منطقه تاسمانی (<http://odi.statelibrary.tas.gov.au/>) و پروژه خدمات Territory Stories (<http://www.territorystories.nt.gov.au/>) در منطقه شمالی.

اخیراً کتابخانه ملی، استرالیا برای اجرای آرشیو وب از سه روش شناسی مختلف استفاده می کند. آرشیو انتخابی در پاندورا (آرشیو وب استرالیا)، طی قراردادی با آرشیو، اینترنت گردآوری کل دامنه و بهره برداری از خدمات آرشیو - آی تی را با همکاری یکدیگر انجام داده اند. با این روش به سمت ایجاد مجموعه کامل و جامع نشریات برخط استرالیا حرکت می کنیم البته با افزایش چالش های نوین فناوری کتابخانه مجبور است برای تداوم این امر مهم روش هایی برای تعدیل طرح های آینده اتخاذ کند؛ نظیر محدوده مجموعه ها و گسترش همکاری های جدید

ص: 151

Web Archiving in a WEB 2/0 World -1

Edgar Gruk -2

3- کارشناس ارشد سازمان اسناد و کتابخانه ملی ایران

آرشیو پاندورا (<http://pandora.nla.gov.au>)، انتشارات وب استرالیا را از سال 1996 بایگانی کرده است؛ و در همان زمان به عنوان مخزن مورد قبول در میان تمام ملت ها معرفی شد. اصلی ترین موفقیت پاندورا ساخت شبکه ای مشتمل بر 9 گروه آرشیو استرالیایی بوده است مانند تمام کتابخانه های دولتی قاره استرالیا و AIATSIS (مؤسسه استرالیایی مطالعات بومیان و جزیره نشینان تنگه تارس)، آرشیو ملی فیلم و صدا و یادمان جنگ استرالیا، بنیانگذاری پانداس - سیستم گردش کار آرشیو که سیستمی است برای شناسایی مداوم مطالب آرشیو شده و ارتباط با طیف وسیعی از سازمان های نمایه سازی و چکیده نویسی - و همچنین تعداد بی شماری از ناشران استرالیایی؛ به گونه ای که در اول جولای 2008، آرشیو، در برگیرنده 307،19 عنوان شامل 531,140,080 فایل بالغ بر 2/2 ترابایت داده بوده است.

آرشیو پاندورا موضوع های زیر را آرشیو کرده است:

- انتشارات برخط منتخب استرالیا در سطح جهانی نظیر مجله های الکترونیکی، انتشارات دولتی و وب گاه های مهم پژوهشی و یا فرهنگی؛

- خط مشی ها، شیوه ها و دستورالعمل های منتخب برای مجموعه و تأمین دسترسی بلند مدت به مواد موجود در آرشیو؛

- رویکرد ملی مشترک برای آرشیو و حفاظت بلند مدت انتشارات پیوسته استرالیا، شامل مشارکت کتابخانه های دولتی و دیگر مؤسسه های فرهنگی؛

- سیستم آرشیو رقومی (پانداس) برای جمع آوری و بارگذاری ساده و مؤثر انتشارات در آرشیو، ذخیره اطلاعات در مورد آن ها و مدیریت دسترسی عمومی به آن ها؛

- طرحی برای نام گذاری مستمر تمام اشیای موجود در آرشیو و ارائه خدمات با دقت به آن ها؛

- تنظیم مقررات همکاری با مؤسسه های نمایه سازی و چکیده نویسی در آرشیو پاندورا، به منظور استناد و دسترسی دائمی به مواد بایگانی شده و شناسه های ماندگار برای انتشارات نمایه سازی و چکیده نویسی در آن ها؛ و

- محتوا که در اول جولای 2008 بالغ بر 307،19 عنوان مشتمل بر 531,120,080، فایل به میزان 2/2 ترابایت داده می باشد

گستره وسیعی از انتشارات وجود دارد نیمی از آن ها از وب گاه های دولت فدرال و ایالتی آرشیو شده اند و نیمی دیگر بازتابی از تنوع کامل فرهنگ و پژوهش استرالیاست نوع انتشارات درون آرشیو می تواند یک تک مدرک PDF یا یک وبگاه کامل سازمانی شامل هزاران فایل باشد در آرشیو، و بنوشت، پادکست، و فیلم نیز وجود دارند.

جمع آوری گزینشی هدایت شده انتشارات، وب بر اساس ارزش بلند مدت فرهنگی و پژوهشی، آن ها بدان معناست که تنها بخش کوچکی از دامنه استرالیا آرشیو شده است کتابخانه، از سال 2005، این مطلب را دریافت و از آن زمان با آرشیو اینترنت (<http://www.archive.org>) قرارداد بست که خزش های هدایت شده سالانه را به منظور جمع آوری هر ماده ممکن از دامنه استرالیا،

انجام دهد. این

خزش ها هر سال حدود یک ماه صورت می گیرد و مقدار داده هایی که جمع آوری می کند به اندازه ای است که محتوای پاندورا را کم جلوه می دهند؛ به عنوان مثال در سال 2007 از کل گردآوری دامنه 18 ترابایت داده در یک ماه برداشت شد در حالی که پاندورا در 11 سال 2 ترابایت داده جمع آوری کرده است. انتظار می رود جمع آوری سال 2008 یک بیلیون فایل باشد.

عکس

آرشیو وب در دنیای وب ۱۵۳

خزش ها، هر سال حدود یک ماه صورت می گیرد و مقدار داده هایی که جمع آوری می کند به اندازه ای است که محتوای پاندورا را کم جلوه می دهند؛ به عنوان مثال، در سال ۲۰۰۷، از کل گردآوری دامنه ۱۸ ترابایت داده در یک ماه برداشت شد، در حالیکه، پاندورا در ۱۱ سال ۲ ترابایت داده جمع آوری کرده است. انتظار می رود جمع آوری سال ۲۰۰۸ یک بیلیون فایل باشد.

گردآوری دامنه وب استرالیا: تجزیه و تحلیل کمی مقدماتی آرشیو داده، NLA، ۲۰۰۸ -

تاریخ گردآوری دامنه	۲۰۰۵	۲۰۰۶	۲۰۰۷
مدارک (فایل های) تکی خزش شده	۱۸۵۵۴۹۶۶۲	۵۹۶۲۳۸۹۹۰	۵۱۶۰۶۴۸۲۰
کل مدارک (فایل های) خزش شده	۱۸۹۸۲۴۱۱۹	۶۲۱۶۶۴۸۷۶	۵۲۳۵۱۰۹۴۵
میزان ها	۸۱۱۵۲۳	۱۲۶۰۵۵۳	۱۲۴۷۶۱۴
اندازه داده های خام	ترابایت ۶,۶۹	ترابایت ۱۹,۰۴	ترابایت ۱۸,۴۷
اندازه فایل های ARC فشرده شده	ترابایت ۴,۵۲	ترابایت ۱۰,۴۸	ترابایت ۱۰,۱۸

HTTP://PANDORA.NLA.GOV.AU/DOCUMENTS/AUSCRAWLS.PDFKOERBIN, P.

گردآوری هایی که با استفاده از HERITRIX هدایت شده اند (<http://crawler.archive.org>) بزرگ است، ولی کاملاً جامع نیستند، زیرا فقط یک ماه در هر سال گردآوری می شوند (و مقادیر زیادی می توانند در این فاصله زمانی در اینترنت بیابند و بروند)، آنها از قوانین robots.txt پیروی می کنند، و اگرچه HERITRIX عملکردی قوی دارد، و نگاه هایی وجود دارد که از لحاظ فنی خزش به آنها مشکل است. نقص ها هر چه باشند، گردآوری در اندازه قابل توجهی انجام می گیرد و بنابراین مقادیر زیاد داده گردآوری شده باعث می شود هرگونه تلاشی در جهت ارزشیابی کیفی و نگاه های شخصی مشکل گردد. بنابراین، به طور معکوس، پاندورا، جایی است که می توان مشکلات را درون هر عنوان شناسایی کرد و تحلیل کیفی را انجام داد - در صورت لزوم صفحه ها را گردآوری کرد و تثبیت کرد - که در اینجا ممکن نیست. یکی دیگر از اشکال ها این است که برخلاف پاندورا - که در آن اجازه ناشر برای آرشیو کسب می شود - در اینجا چنین امکانی وجود ندارد. بنابراین، با توجه به قانون کپی رایت (حق تألیف) فعلی استرالیا که طبق آن نشریات برخط شامل واسپاری قانونی نیستند، در حال حاضر، قادر نیستیم آنچه را که حفظ شده به عموم نمایش دهیم. البته، این عبارت به این معنی نیست که هیچ کاری در این آرشیو انجام نمی شود، چرا که دانشگاهیان در حال کار بر روی این داده ها هستند و به نظر می رسد تحقیقات آنها

گردآوری دامنه وب استرالیا: تجزیه و تحلیل کمی مقدماتی آرشیو داده، NLA، 2008 -

گردآوری هایی که با استفاده از HERITRIX هدایت شده اند (<http://crawler.archive.org>) بزرگ است، ولی کاملاً جامع نیستند زیرا فقط یک ماه در هر سال گردآوری می شوند (و مقادیر زیادی را می توانند در این فاصله زمانی در اینترنت بیابند و بروند آن ها از قوانین robots.txt پیروی می کنند، و اگر چه HERITRIX عملکردی قوی دارد وب گاه هایی وجود دارد که از لحاظ فنی خزش به آن ها مشکل است. نقص ها هر چه باشند گردآوری در اندازه قابل توجهی انجام می گیرد و بنابراین مقادیر زیاد داده گردآوری شده باعث می شود هر گونه تلاشی در جهت ارزشیابی کیفی وب گاه های شخصی مشکل گردد.، بنابراین به طور معکوس، پاندورا جایی است که می توان مشکلات را درون هر عنوان شناسایی کرد و تحلیل کیفی را انجام داد- در صورت لزوم صفحه ها را گردآوری کرد و تثبیت کرد - که در این جا ممکن نیست یکی دیگر از اشکال ها این است که بر خلاف پاندورا- که در آن اجازه ناشر برای آرشیو کسب می شود- در اینجا چنین امکانی وجود ندارد، بنابراین با توجه به قانون کپی رایت (حق تألیف) فعلی استرالیا که طبق آن نشریات برخط شامل واسپاری قانونی نیستند در حال حاضر، قادر نیستیم آن چه را که حفظ شده به عموم نمایش دهیم البته این عبارت به این معنی نیست که هیچ کاری در این آرشیو انجام نمی شود، چرا که دانشگاهیان در حال کار بر روی این داده ها هستند و به نظر می رسد تحقیقات آن ها

ارزشمند باشد آرشیو-آی تی سرویس آرشیو وب های میزبانی شده توسط آرشیو اینترنت است. اولین (و تا کنون تنها) سازمانی که در استرالیا از این سرویس استفاده کرده و می کند بخش مجموعه های آسیایی کتابخانه ملی (<http://www.nla.gov.au/asian/asianwebarchive.html>) است.

از آرشیو-آی تی برای جمع آوری مجموعه ای از وب گاه های خارج از کشور استفاده می شود که رویدادهای خاص اجتماعی و سیاسی را ضبط می کنند؛ زیرا انتظار نمی رود که هیچ سازمان دیگری در منطقه این نقش را به انجام برساند گزینه میزبانی نیز از آن رو انتخاب شد که به نظر می رسد می تواند راهی سریع و آسان برای جمع آوری و نگهداری مجموعه ها باشد و به مهارت های فنی نیاز ندارد و میزان زیادی از وقت کارکنان را نمی گیرد. در حالی که خیلی زود معلوم شد فقط قسمتی از این واقعیت دارد چرا که برای ساخت موفقیت آمیز مجموعه ها زمانی بسیار زیادتر از آن چه تصور می شد صرف شد گزینش وب گاه ها برای ، خزش فعالیتی است که اغلب با سو تفاهم همراه است و می تواند به طور شگفت آوری زمان زیادی بطلبد.

در حالی که استفاده از آرشیو آی تی دارای این مزیت است که دیگر در مورد میزبانی و حفظ محتوای جمع آوری شده نگرانی وجود نخواهد داشت. البته برخی مشکلات عمده در مورد عدم کنترل مواد آرشیوی و نمایش توابع وجود دارد. آرشیو آی تی اجازه می دهد تا برخی URL های هسته انتخابی جمع آوری شوند؛ با وجود این، اگر فایل های جمع آوری شده از بین بروند چه جمع آوری بشوند چه جمع آوری نشوند، هیچ راهی برای اصلاح خرابی یا از دست دادن محتوا وجود ندارد که بتوانید با سیستم خود انجام دهید به همین ترتیب، کنترل واقعی یا مالکیت فرآیند نمایش نیز وجود ندارد؛ به طوری که به عنوان مثال یک پیوند به URL هسته که جمع آوری نشده ، هم چنان درون یک مجموعه ظاهر می شود یکی دیگر از اشکال ها این است که اگر شما اشتراک سالیان خدمات خود را قطع کنید مجموعه های شما به مخزن محتوای آرشیو اینترنت عمومی باز می گردد. با وجود این مسائل کتابخانه در نظر دارد به آرشیو با استفاده از این روش ادامه دهد.

سایت های آرشیو شده با استفاده از آرشیو آی تی برای مجموعه های آسیا - کتابخانه ملی استرالیا

Papua New Guinea Government and Research Websites

(<http://archive-it.org/collections/1039>)

وب گاه های انتخابی مؤسسه های دولتی و پژوهشی مهم پاپوا گینه نو که از سال 2008 آرشیو شده اند. برخی سایت های آرشیوی که در زمان ضبط کردن جاری نبودند.

(<http://archive-it.org/collections/918>)

(<http://archive-it.org/collections/1040>)

وب گاه انتخابی بین المللی بین الدولی مرتبط با انتخاب عمومی سال 2007 تایلند

وب گاه های منتخب دولتی، احزاب سیاسی و رسانه های مرتبط با انتخابات ملی کامبوج سال 2008

(<http://archive-it.org/collections/937>)

(<http://archive-it.org/collections/1054>)

وب گاه های منتخب بین المللی مرتبط با شورش راهبان برمه سپتامبر - اکتبر 2007

وب گاه های انتخابی دولتی و غیر دولتی جمهوری دموکراتیک خلق لائوس که از سال 2008 آرشیو شده اند.

ص: 154

(<http://archive-it.org/collections/920>)

وب گاه های انتخابی بین المللی بین دولتی و پاپوآگینه نو مرتبط با انتخابات پارلمانی پاپوآگینه نو در سال 2007. شامل تصاویری از وبگاه کمیسیون انتخابات پاپوآگینه نو، دادخواست انتخابات پاپوآگینه، نو گزارشی از گروه ارزشیابی انتخابات انجمن جزایر مشترک المنافع اقیانوس آرام و بیانیه رسانه ای شفافیت بین المللی (پاپوآگینه نو).

(<http://archive-it.org/collections/912>)

وب گاه های منتخب بین المللی بین دولتی مرتبط با انتخابات ریاست جمهوری و پارلمانی تیمور شرقی در سال 2007 شامل آموزش رأی دهندگان و مواد تبلیغاتی سیاسی

وب گاه های منتخب اندونزیایی

جمع آوری فایل ها

(<http://archive-it.org/collections/664>)

وقتی پاندورا پا به عرصه وجود گذاشت به علت ناتوانی اولیه نرم افزارهای جمع آوری وب، بسیاری از وب گاه ها نمی توانستند جمع آوری شوند وب گاه های با پایه HTML می توانستند جمع آوری شوند، ولی سایت هایی با چنین قالب های ساده ای در ابتدا بسیار مشکل به وجود آوردند تا یک مشکل حل می گردید مشکل دیگری پیدا می شد بنابراین برای آرشویست های وب پیشرفت شگفت انگیز اینترنت جاوا، اسکریپت اپلت ها، شیوه نامه های آبخاری Shockwave flash، و ده هزار فایل دیگر و انواع قالب ها فقط مشکلاتی پس از دیگری بود در حالی که مشکلات اغلب انواع فایل ها بر طرف شده اند، محتوای چندرسانه ای یک مسئله باقی مانده است پیش از این در مورد فایل های پخش (Real Player) و در حال حاضر، درباره پادکست ها مجبوریم نه با پیچیدگی خود فایل ها که با پیچیدگی سیستم های تحویل آن ها مقابله کنیم در حال حاضر این امر در مورد فیلم ها این مسئله لاینحل مانده است.

مجموعه انتخابات فدرالی سال 2007 بزرگ ترین تلاش تاکنون بوده است کتابخانه ملی مسئول آرشو تمام منابع ملی مرتبط با انتخابات بود از جمله وب گاه های احزاب، گروه های لابی وب گاه های برخی نامزدها، وب نوشت ها، فیلم ها و وب گاه های رسانه ای کتابخانه های ایالتی مسئول جمع آوری وب گاه های نامزدها، احزاب، و رسانه های محلی در ایالت خود بودند در مجموع بیش از 350 وبگاه توسط کتابخانه ملی و همکارانش آرشو شده است، بسیاری از این سایت ها چندین بار جمع آوری شدند تا محتوای در حال تغییر را ضبط کنند بزرگ ترین چالش در آرشو کردن این انتخابات تعداد زیاد فیلم ها بود؛ البته مشکلی در خود آن ها نبود، بلکه مشکل در مکانیزم تحویل و فناوری های تعبیه شده در آن ها به منظور استفاده مفید کاربران بود.

رویکردهای متفاوتی برای آرشو کردن فیلم ها بسته به ماهیت وب گاه ها به کار گرفته شده است. برای وب گاه های عمومی که در آن فیلم های مجزا در یک صفحه وب موجود است فایل های فیلم را جداگانه با استفاده از ابزارهای متفاوت و رایگان موجود بارگیری کردیم (گردآورندگان وب به طور کلی نمی توانند فیلم ها را به طور خودکار جمع آوری کنند)؛ و سپس با استفاده از مبدل های فایل فایل ها را از flv. به چیزی کاربر پسند تر مانند قالب Mpeg تغییر دادیم جایی که تعدادی فیلم به یک صفحه وب پیوند داشت، فیلم ها در قالب اصلی

خودشان باقی ماندند و یک پخش کننده flv. در فایل های جمع آوری شده نصب گردید، به طوری که فایل ها می توانند به راحتی وب گاه زنده ارائه شوند وقتی که وب گاه انتخابات یوتیوب

ص: 155

(<http://nla.gov.au/nla.arc-76644>) شامل بیش از 700 فیلم) را جمع آوری کردیم، با مراجعه به مهارت های فنی در بخش فناوری اطلاعات توانستیم برای فیلم ها URL ها را از سایت جاری استخراج کنیم آن ها را بارگیری کنیم و تغییرات ضروری را برای گردآوری صفحات وب انجام دهیم.

هیچ یک از این فرآیند ها سریع نبود و همه به مقدار نسبتاً خوبی از مهارت های فنی نیاز داشتند، از جمله کد گذاری مجدد صفحات آرشیو شده با تغییراتی که باید انجام می دادیم تا فیلم ها درون آرشیو قابل پخش باشند. همچنین به خاطر انتخابات حفاظت از اسناد برخط در دولت پیشین نیز به ما سپرده شد. با پیش بینی این امر طرح هایی (همان گونه که برای انتخابات پیشین داشتیم) ساخته بودیم و تمام وب گاه های وزارتی دولت را آرشیو کرده بودیم؛ در حالی که آن ها درست قبل از تاریخ انتخابات در حالت مستحفظ بودند این کار صورت گرفت زیرا گمان می رفت با یک تغییر احتمالی در دولت حذف کلی محتوا از اینترنت وجود داشته باشد همان گونه که در حوزه های قضایی دیگر اتفاق افتاده بود. از آن جا که انتشارات وب برخی بخش های دولتی به ویژه بخش هایی که دارای اسناد تغییر یافته بودند، از دید عموم حذف شده، بود این دور اندیشی اجر نهاده شد.

دستورالعمل های جمع آوری

با توجه به گردآوری دامنه استرالیا توسط کتابخانه ملی و آرشیو انتخابی پاندورا، می توان گفت که حجم قابل توجهی از نشریات استرالیا و یا وب گاه ها آرشیو شده اند اگر چه نمی توانیم بگوییم تا چه حد این مجموعه جامع و کامل است با این حال می دانیم که شکاف های بزرگی باقی مانده است.

ما، وب گاه های تولیدی توسط استرالیایی های خارج از دامنه استرالیا را به طور جامع آرشیو نمی کنیم (اگر چه امیدواریم آرشیو اینترنت بسیاری از آن ها را به دست خواهد آورد). خارج از نشریات مرسوم، اما تا اندازه ای دارای اهمیت بیشتر، ما هم چنین میزان بسیار زیادی از محتوای خلاقانه ای را که توسط افراد تولید شده و در وب گاه ها، وب نوشت ها، دنیای مجازی و سایت های شبکه های اجتماعی فیلم، عکس و هنر میزبانی شده اند، جمع آوری نکرده ایم.

تلاش هایی برای جمع آوری برخی از این محتواها وجود دارد اما پروژه های هدایت شده کوچکی هستند. یکی از آن ها مجموعه رقص استرالیایی است که به دنبال گردآوری نمونه های از رقص استرالیایی به محض قرار گرفتن بر روی وب گاه های مختلف و سایت های میزبان فیلم می باشد.

اگر چه کتابخانه ملی با Flickr قرارداد بسته است و از MySpace و YouTube اجازه آرشیو کردن دریافت کرده ایم کتابخانه فقط بخش کوچکی از این منابع را آرشیو یا گردآوری کرده است. هم چنین، دیگر منابع برخط هر جایی که فعالیت های استرالیایی وجود دارد نظیر فضاهای مجازی (Second Life غیره) و شبکه های اجتماعی (Facebook, Bebo و غیره) نیز آرشیو نشده اند. نخستین دلیل این که معمولاً محتوا دارای حق نشر و حفظ حریم خصوصی است که آرشیو کردن آن را مجاز نمی کند و یا به دلیل ماهیت منابع که آن را خارج از اینترنت عمومی قرار می دهد.

ما به طور جداگانه به عنوان کتابدار نیز می توانیم تفاوتی نسبت به حفظ میراث برخط خودمان به وجود آوریم؛

البته با پذیرش این مسئولیت که می‌توانیم درون سازمان خود یا سازمان مادر از حفظ و پشتیبانی آن چه در وب‌گاه سازمان قرار گرفته اطمینان حاصل نماییم این، جنبش به ویژه به طور قابل توجهی در دولت و دانشگاه‌ها در مورد حرکت نشریات از چاپی به برخط ادامه می‌یابد آن چه ما پیدا کرده ایم این است که نباید مطمئن باشیم که نشریات برخط در وب‌گاه ناشران در بلندمدت (و یا حتی کوتاه مدت) در دسترس باقی بمانند دانشگاه‌ها در حال حاضر، مجبور به ایجاد مخازن رقومی برای خروجی فکری خود هستند و در این راه نشریات آن‌ها دسترس پذیر باقی می‌مانند. در حالی که ممکن است انتظار داشته باشیم که وب‌گاه‌های دولتی دسترسی عمومی خود را حفظ نمایند؛ تجربه نشان می‌دهد که این موضوع همیشه صادق نیست بنابراین اگر نشریه‌ای برای مجموعه و یا کاربران شما ارزشمند است، عاقلانه این است که تلاشی در جهت حفظ دسترسی بلندمدت آن انجام شود.

دستور عمل‌های آینده

تا زمانی که چشم انداز در حال رشد اینترنت وجود دارد فناوری‌های جدید و شکاف‌هایی که تازه شناسایی شده‌اند در گردآوری ما وجود دارند که پرداختن به آن‌ها ضروری است بنابراین، به نظر نمی‌رسد آرشیو وب به منطقه‌ای تبدیل شود که در آن شیوه‌ها و یا پروتکل‌های توسعه مجموعه کاملاً تأسیس و یا برقرار گردد. ما همیشه به شناسایی و جمع‌آوری مطالب نیاز داریم نه این که منتظر بمانیم تا آن‌ها به سوی ما بیایند. وب بسیار پویاست و فناوری بسیار متغیر تعداد ناشران آن قدر وسیع است که بعید به نظر می‌رسد ما قادر به ایجاد سبکی مقیاس پذیر - نظیر سبک سیستم سپرده فیزیکی که نیاکان ما برای مواد چاپی انجام داده‌اند باشیم. هنگامی که کتابخانه ملی برای اولین بار شروع به آرشیو کردن وب کرد، ابزار و نهادهای بسیار کمی در خارج وجود داشت که ما می‌توانستیم برای یادگیری با آن‌ها کار کنیم. در نتیجه، باید مکانیسم‌ها و ابزارها را خودمان اختراع می‌کردیم برای این، منظور کتابخانه ملی در درون خودش سیستم مدیریت آرشیوی را - PANDAS - ساخت این سیستم در حال حاضر، در سومین و طبق تصور ما آخرین مرحله خود است چرا که کتابخانه دیگر نمی‌تواند به طور مستقل چنین سرمایه‌گذاری را برای توسعه متحمل شود. با این حال در حال حاضر که آرشیو وب در طیف وسیعی از نهاد‌های بین‌المللی به طور عملی ایجاد شده، است شرکای زیادی وجود دارند که می‌توان این بار را با آن‌ها تقسیم کرد. کنسرسیوم بین‌المللی حفاظت از اینترنت (1) که اعضای آن متشکل از تمام کتابخانه‌های ملی راهبر آرشیو وب از جمله کتابخانه ملی استرالیا و دیگر نهاد‌های مرتبط می‌باشد این توسعه را هدایت و راهبری می‌کند با این، روش کتابخانه می‌تواند به نقش راهبری خود در آرشیو، وب از طریق تطابق با ابزارهای در حال توسعه و سازگاری با کتابخانه‌ها و مؤسسه‌های همکار ادامه دهد.

شرح حال مختصری از پدید آورنده

ادگار کروک از 1999، در کتابخانه ملی استرالیا کار می‌کند. او از سال 2000 بر روی پاندورا: آرشیو وب استرالیا کار می‌کند. قبل از آن در کند. قبل از آن در کتابخانه عمومی ATC کار می‌کرده است.

همان گونه که در بسیاری از متون نیز مطرح است زبان فارسی به واسطه ویژگی های خاص خود چه از نظر رسم الخط و چه از نظر صرف و معنا با چالش های منحصر به فردی در زمینه ذخیره و بازیابی اطلاعات رو به روست در مقالات مختلف به این مشکلات در قالب ها و سطوح مختلف بدون توجه به نوع شناسی این چالش ها اشاره شده است. به عبارت بهتر کمتر مقاله ای در زبان فارسی در حوزه اطلاع رسانی از دیدی مبتنی بر علم زبان شناسی مدرن به بحث دسته بندی این مشکلات و پاسخ دهی به آن ها پرداخته است. نوشتار حاضر بر آن است تا عمده ترین چالش های مطرح این زمینه را در سه گروه رسم الخط مسائل صرفی و مسائل معنایی مورد اشاره قرار داده و پس از بیان نمونه هایی در پیوند با هر یک از این چالش ها و ارائه ساختواره ای درختی از انواع مسائل مطرح در این سه گروه با نگاهی به پژوهش های صورت گرفته در زمینه زبان فارسی و عربی به ارائه راهکار هایی جهت حل این مسائل بپردازد.

کلیدواژه ها: چالش های صرفی زبان فارسی؛ چالش های رسم الخط زبان فارسی؛ چالش های معنایی زبان فارسی؛ بازیابی اطلاعات

شعله ارسطو پور (1) | فاطمه احمدی نسب (2)

درآمدی بر مشخصه‌های زبان فارسی

از دیدگاه دستور زبان زایشی (3)، زبان یا دانش زبانی قوه مستقلی از قوای ذهنی به شمار می‌رود که از دیگر قوای شناختی انسان مستقل بوده و خود نیز دارای بخش‌های مستقل نحو، معنا، واژگان و بخش واجی است به عبارت دیگر این رویکرد یک رویکرد حوزه‌ای (4) به زبان است که بین دانش زبانی و کاربرد آن تمایز قائل است (دبیر مقدم 1383، 18-19). در همین راستا در این نوشتار نیز منظور از زبان فارسی دانش زبانی فارسی‌زبانان ایرانی و هم‌چنین تبلور این دانش زبانی در قالب گونه‌نوشتاری فارسی رایج در ایران یعنی خط فارسی است. زبان فارسی یکی از زبان‌های هندو اروپایی و از شاخه زبان‌های ایرانی جنوب غربی است که زبان رسمی ایران و تاجیکستان و یکی از دو زبان رسمی افغانستان است (کامری 1990، 13-16) این در حالی است که تفاوت‌هایی میان این انشقاق‌ها وجود دارد مثلاً فارسی تاجیکی و فارسی رسمی ایرانی از لحاظ دستوری یکسان هستند

ص: 159

1- دکترای علم اطلاعات و دانش‌شناسی مرکز منطقه‌ای اطلاع‌رسانی علوم و فناوری 2 arastoopoor@ricest.ac.ir

2- دکترای زبان‌شناسی همگانی مرکز منطقه‌ای اطلاع‌رسانی علوم و فناوری ahmadinasabricest.ac.ir

3- Generative grammar

4- Modular

و فقط مقداری تفاوت های واژگانی دارند، اما از لحاظ نظام نوشتاری کاملاً متفاوت هستند چرا که خط فارسی ایرانی از الفبای عربی اقتباس شده است و خط فارسی تاجیکی الفبای سیریلیکی دارد. از دیدگاه رده شناسی (1)، زبان فارسی از برخی ویژگی های سه رده زبانی زبان های تحلیلی (2)، زبان های تصریفی (3) و زبان های پیوندی (4) برخوردار است. این تمایز سه گانه بر اساس ساختار واژگانی زبان تعریف می شود و از آن جا که ذخیره و بازایی اطلاعات و نمایه سازی غالباً بر اساس واژه صورت می گیرد به نظر می رسد که توجه به این تمایز از اهمیت ویژه ای برخوردار باشد. زبان تحلیلی زبانی است که واژه های آن بر اساس نقش (فاعلی، مفعولی و...) هیچ گونه تغییری نکند و صرفاً جایگاه آن در جمله نشان دهنده نقش واژه باشد کریستال (2008، 256) در زبان های، تصریفی واژه غالباً از چند تکواژ (5) تشکیل می شود که نشان دهنده روابط دستوری است در این زبان ها جایگاه واژه در جمله نشان دهنده نقش واژه نیست بلکه تصریف و حالت (6) واژه نقش واژه را به تصویر می کشد (همان، 244) و بالاخره زبان پیوندی زبانی است که واژه های آن غالباً از بیش از یک تکواژ تشکیل شده اند و هر کدام از تکواژها نشان دهنده یک نقش یا رابطه دستوری است (همان، 16). البته لازم به ذکر است که تمایز زبان تصریفی و پیوندی در آن است که تقطیع و تعیین تکواژهای واژه در آن امکان پذیر است و به راحتی می توان مرز آن ها را تعیین نمود اما در زبان های تصریفی واژه غالباً صرف شده و تکواژها نه به صورت منظم و به ترتیب بلکه به شکل تجمعی در یک صورت کلمه (7) ظاهر می شوند صورت های مفرد و جمع واژگان روایت / روایات عاقبت / عواقب مدرسه / مدارس و ده ها مثال دیگر که البته وام واژه هایی (8) از زبان عربی هستند نشان دهنده وجود تصریف در زبان فارسی است. است از دیگر سو صد ها مثال مانند واژه دانشگاهی (دان + ش + گاه + ی) نشان گر وجود ویژگی زبان پیوندی در فارسی است اگر چه تحقیق جامعی در تمایل بیش تر زبان فارسی به هر یک از این سه رده صورت نگرفته اما به نظر می رسد که زبان فارسی به زبان تحلیلی نزدیک تر بوده و از ترتیب واژه برای نمایش نقش های دستوری بیش از شیوه های دیگر بهره می برد. از بحث های مرتبط با بحث توالی اجزای واژگانی می توان به بحث هسته (9) واژه مرکب اشاره نمود در زبان فارسی واژه مرکب از نظر معنایی و بر اساس هسته واژگانی خود به چهار دسته برون، مرکز درون مرکز (هسته آغازی، هسته پایانی) دو سویه و متوازن تقسیم می شود. در واژه مرکب برون مرکز هیچ کدام از اجزای واژه هسته را تشکیل نمی دهند برای مثال «آب سیاه» نام یک بیماری است و نه آب است و نه سیاه رنگ در واژه مرکب درون مرکز یکی از اجزای واژه هسته است و به دو نوع هسته - آغازی و هسته پایانی

ص: 160

Typology -1

Isolating/analytic -2

Inflectional -3

Agglutinating -4

Morpheme -5

Case -6

Word form -7

Loan word -8

Head -9

تقسیم می شود. برای مثال «موش خرما» هسته آغازی و «مویرگ» هسته پایانی است. در واژه مرکب دوسویه هر دو جزء به یک مرجع واحد اشاره می نماید مانند «سرباز - معلم» و بالاخره اینکه در واژه مرکب متوازن هر دو جزء به یک اندازه در ساخت کلمه مرکب نقش دارند؛ مانند «دخل و خرج» شقاقی (123، 1386-124).

نکات پیش گفته تنها بخشی از مهم ترین ویژگی های زبان فارسی است که به صورت بالقوه قابلیت ایجاد بزرگ ترین چالش ها را در بحث ذخیره و بازیابی اطلاعات دارد عمده ترین دلیل چنین امری آن است که ماشین قابلیت تشخیص ترتیب و نقش دستوری را به خودی خود ندارد و همین امر منجر به بازیابی منابعی می شود که بعضاً مورد نیاز کاربر نبوده و عملاً پاسخی به نیاز وی نخواهند بود. بنابراین بی شک بعد زبان شناختی فرایند ذخیره و بازیابی اطلاعات یکی از مهم ترین و در عین حال برانگیز ترین مسئله طراحی و کاربرد پذیری پایگاه ها است بازیابی اطلاعات با مسائل زبانی گره خورده و نمی توان بدون توجه به زبان مبدأ و ویژگی های آن نظام های بازیابی اطلاعات را در جهت افزایش جامعیت و مانعیت از یک سو و بهبود ربط از دیگر سو ارتقا بخشید. این نوشتار بر آن است تا از سه بعد رسم الخط صرف و معنا به بررسی مشکلات و چالش های زبان فارسی پرداخته و به صورت همزمان به نمونه هایی عینی از این مشکلات و نتیجه عدم توجه به آن ها در سطح وب و برخی از پایگاه های مقالات فارسی اشاره کرده و به آسیب شناسی زبان و خط فارسی در رابطه با این حوزه بپردازد.

پیشینه پژوهش

تلاش در جهت اصلاح و تقویت زبان و خط فارسی را می توان به تأسیس فرهنگستان اول نسبت داد. در این نوشتار به منظور پرهیز از طولانی شدن بحث صرفاً به برخی از تلاش های اخیر اشاره شده است. به عنوان نمونه آشوری (1375) برای اصلاح خط فارسی پیشنهاد می دهد که صورت های صرفی زمان حال فعل بودن یعنی «آم»، «ای»، «است»، «ایم»، «اید»، «اند» و ضمائر متصل «آم»، «ات»، «آش»، «مان»، «تان»، «شان» جدا از کلمه های قبل از خود نوشته شوند معصومی همدانی (1381) در مقاله ای تحت عنوان خط فارسی و رایانه یکسان سازی رسم الخط فارسی را برای استفاده از خطایاب امکان جستجوی واژه در متن و تهیه نمایه ضروری می داند وی راه حل را در تهیه یک دستورالعمل واحد برای خط فارسی دانسته و با اشاره به دستورالعمل فرهنگستان می نویسد: «متأسفانه میزان آزادی که این دستور خط به استفاده کنندگان داده به قدری است که می توان گفت کاربرد آن در مواردی بر تشنگت موجود خواهد افزود». اسلامی (1381) اعمال اصلاحاتی را در خط فارسی ضروری می داند؛ اصلاحاتی از قبیل وضع یک نشانه اصلی برای نمایش کسره اضافه استفاده از نشانه های متفاوت برای نمایشی نکره و ی اسم-ساز و صفت ساز سرهم نویسی کلمات غیر بسیط، قرار دادن علامت های جداگانه برای واژه بست (1) های ربطی فعل بودن و ضمائر ملکی، و بالاخره نشان

ص: 161

دادن اسامی خاص محقق زاده و زارعیان (1383) با الگوگیری از زبان های لاتین چپ نویسی کاهش حروف از طریق جدانویسی قائل شدن به دو شکل کوچک و بزرگ حروف را در مورد زبان فارسی پیشنهاد می نمایند. علاوه بر این نگارش حروفی که خوانده شده ولی نوشته نمی شوند، ایجاد علامتی برای کسره اضافه و عدم تمایز بین «آ» و «ا» را نیز در جهت سازگار ساختن زبان فارسی با محیط های رایانه ای ضروری می دانند اسلامی (1386) نقطه دار بودن برخی از نویسه های فارسی را یکی از ایرادهای خط فارسی می داند؛ چرا که علاوه بر دشوار نمودن املا، کلمات این ایراد اساسی را دارد که استفاده از برنامه نشانه خوان نوری (1) را برای متون فارسی دشوار می نماید. صفار مقدم (1386) فاصله گذاری را در ترکیبات فارسی ضروری دانسته و مفصلاً به مبحث فاصله گذاری درون کلمه ای و برون کلمه ای در انواع کلمات فارسی پرداخته است. گل تاجی و بذرگر (1389) با استفاده از یک سیاهه 17 کلید واژه ای مشکلات ریخت شناسی خط فارسی را در سه پایگاه اطلاعاتی مرکز منطقه ای اطلاع رسانی علوم و فناوری پژوهشگاه اطلاعات و مدارک علمی ایران و جهاد دانشگاهی مقایسه کرده و نتیجه می گیرند که هیچ کدام از پایگاه های مذکور به ور جامع به این مسائل توجه نداشته اند.

همان طور که از مرور آثار پیشین پیداست برخی از پژوهش ها به ضرورت اصلاح خط فارسی با توجه به محیط های رایانه ای و برخی به آموزش زبان فارسی توجه داشته اند که برخی از پیشنهادات مطرح شده در این راستا افراطی بوده و در صورت اجرا چهره خط فارسی را کاملاً دگرگون می نمایند برای مثال چپ نویسی استفاده از حروف بزرگ و کوچک و یا نگارش صداهایی که خوانده می شود. خط دارای ماهیتی محافظه کارانه بوده و به راحتی تغییر را نمی پذیرد. علاوه بر این تغییر آن به راحتی امکان پذیر نیست و به بسیج امکانات مادی و معنوی عظیمی نیاز دارد به نظر می رسد که در راستای حفظ گنجینه زبان فارسی باید خط فارسی را به همین صورت حفظ کرد و تنها با اعمال برخی اصلاحات و پیروی از دستور العملی واحد در جهت یکسان سازی خط فارسی از یک سو و کسب توانمندی های بیش تر در حوزه پردازش زبان طبیعی از دیگر سو به سمت رفع مشکلات بازیابی اطلاعات در سطح وب حرکت کرد.

رسم الخط فارسی و بازیابی اطلاعات

در مورد مسائل مرتبط به خط فارسی به طور مطلق و همچنین ارتباط آن با بازیابی اطلاعات آثار متعددی به نگارش درآمده است که هر کدام از آن ها به گوشه ای از این مشکل عمده توجه داشته است در این نوشتار تنها به دو مسئله فاصله گذاری و نگارش «الف» پرداخته می شود یکی از راه حل های ارائه شده برای یکسان سازی خط فارسی و بهبود بازیابی اطلاعات فارسی در وب پیشنهاد عدم تمایز صورت های مختلف الف آغازی است یعنی «آ» و «ا» به یک صورت یعنی «ا» نوشته شود (اسلامی، 1381). اما این راهکار نمی تواند مفید باشد و به ابهام معنایی منجر می شود. برای مثال در پایگاه RICEST بین «آ» و «ا» تمایزی وجود ندارد و باعث می شود که جستجوی کلید واژه های آسم

ص: 162

(بیماری تنفسی) و اسم (مقوله دستوری) در موتور جستجوی جامع به نتایج کاملاً یکسان منجر شود و حال آن که منطقاً این دو کلیدواژه باید به بازیابی نتایج کاملاً متفاوتی منجر شود. به عبارت بهتر، قائل شدن این تمایز سطح ربط نتایج بازیابی شده به میزان چشم‌گیری افزایش می‌یابد. مثلاً در پایگاه مگ ایران این تمایز در نظر گرفته شده است و جستجوی کلیدواژه‌ها به نتایج متفاوت و مرتبط به کلیدواژه مورد نظر می‌انجامد. (1)

چالش عمده دیگر در پیوند با رسم الخط فارسی متوجه نحوه فاصله‌گذاری است. همان طور که می‌دانیم اکثر نویسه (2) های فارسی با توجه به جایگاه نویسه در واژه به صورت منفصل یا متصل نگاشته می‌شوند و این خود امر آموزش و همچنین ذخیره و بازیابی اطلاعات را دشوار می‌سازد. فرهنگستان زبان و ادب فارسی دستورالعملی را تحت عنوان دستور خط فارسی تدوین کرده است تا به یکسان‌سازی و بهبود نگارش فارسی منجر شود یکی از مباحث این دستورالعمل به فاصله‌گذاری مربوط است. طبق این دستورالعمل در خط فارسی از دو فاصله یعنی نیم فاصله (درون کلمه) و فاصله کامل (برون کلمه) استفاده می‌شود رعایت فاصله کامل و نیم فاصله در پرهیز از بدخوانی بسیار راه‌گشا است (فرهنگستان زبان و ادب فارسی 1389: 10). برای مثال در مورد نام و نام خانوادگی ایرانیان تنها فاصله‌گذاری دقیق است که باعث خوانایی و رفع ابهام می‌شود. زنجیره «علی رضا خانی» دو خوانش محتمل «علیرضا خانی» و «علی رضاخانی» دارد که تنها فاصله‌گذاری صحیح می‌تواند خوانش درست را تعیین نماید صفا مقدم 1386، 125) در بازیابی اطلاعات فاصله‌گذاری مخصوصاً در مورد واژه‌های مشتق و مرکب نقش تعیین‌کننده‌ای دارد چرا که پیوسته نویسی جدا نویسی کلید واژه‌های جستجو نتایج متفاوتی را بدست می‌دهد به عنوان نمونه پایگاه مگ ایران قابلیت جستجوی واژگان دو قسمتی را از طریق استفاده از نقطه فراهم آورده است اما کاربر با جستجوی واژه «خاکبرداری» در سه حالت «خاک برداری»، «خاکبرداری» و «خاک برداری» به نتایج متفاوتی دست خواهد یافت با جستجوی واژه «خاک برداری» 5 یافته بازیابی می‌شود در حالی که با جستجوی «خاک برداری» نتایج جستجو به 16 صفحه نیز می‌رسد. در میان این 16 صفحه یک نتیجه در برگیرنده واژه «خاک برداری» است در حالی که در میان نتایج 5 یافته اول نبوده و از 5 یافته اول حاصل از جستجوی «خاک برداری» مورد در نتایج حاصل از جستجوی «خاک برداری» مشاهده نشد. این در حالی است که نتایج حاصل از جستجوی «خاکبرداری» 7 مورد بوده که کاملاً متفاوت از دو جستجوی قبلی است. در واقع با جستجوی «خاک برداری» آن دسته از نتایج مرتبط در سایر جستجو‌ها به عنوان مرتبط ترین شناسایی نشده و در صفحات 2 و 4 و 8 پراکنده اند. در پایگاه RCeST نیز جستجوی کلیدواژه «خاک + برداری» 496، «خاکبرداری» 18 و عبارت «خاک برداری» 9 نتیجه را بازیابی می‌نماید بررسی نتایج جستجوهای «خاکبرداری» و «خاک برداری» و مقایسه آن‌ها با یکدیگر نشان داد که تنها یکی از عنوان‌ها مشترک بوده است. این بدان معنا است که در هر دو پایگاه

ص: 163

1- البته در نتایج بازیابی شده برای این عبارت، جستجو هم چنان مسائلی همچون عدم رعایت درست تقطیع به چشم می‌خورد

Letter -2

جستجو با نگارش های متفاوت به ریزش کاذب منتهی می شود. بدیهی است تعداد زیاد بازیافت ها و کاهش دقت در 10 نتیجه اول، یافتن نتایج مرتبط را برای کاربر زمان بر و دشوار خواهد کرد. باید توجه نمود که واژه هایی مانند خاک برداری کتاب، شناسی خودکشی و غیره از طریق فرایند ترکیب و اشتقاق ناپایگانی ساخته شده اند (شقایق 1386، 99) و اجزای دوم آن ها یعنی «-برداری»، «-شناسی» به طور مستقل بکار نرفته و استقلال واژگانی و معنایی ندارند بنابراین نظام های بازیابی اطلاعات باید راهکاری را برای جستجو در نظر گیرند تا این اجزا بصورت مجزا جستجو و بازیابی نشده و در ترکیبات مورد نظر به صورت یک واحد جستجو شوند.

مسائل صرفی و بازیابی اطلاعات

صرف (1) یکی از شاخه های زبان شناسی است که به مطالعه واژه و ساختمان درونی آن ها می پردازد. در این نوشتار به دو موضوع صرفی مرتبط به بازیابی اطلاعات یعنی واژه مرکب و وام گیری واژگانی پرداخته می شود. همان طور که پیش تر گفته شد واژه مرکب در فارسی چهار نوع مختلف برون مرکز، درون مرکز، دوسویه و متوازن دارد ویژگی های خط فارسی در اغلب موارد باعث می شود که نتوان از لحاظ نوشتاری تمایزی بین واژه مرکب و گروه نحوی قائل شد. در نتیجه نظام های بازیابی اطلاعات نیز نمی توانند واژه مرکب را از گروه نحوی تشخیص دهند. برای مثال واژه «زیست شیمی» (2) یک واژه مرکب درون مرکز هسته پایانی است که «شیمی» هسته آن را تشکیل می دهد. در حالی که نظام بازیابی اطلاعات این واژه مرکب را هم گروه نحوی «زیست شیمی» و هم «زیست+ شیمی» قلمداد می کند. جستجوی این واژه مرکب در گوگل 11800000 بازیافت در پی دارد که بررسی 10 نتیجه اول نشان می دهد که نظام بازیابی هر سه حالت را مد نظر داشته است. این در حالی است که تنها 50000 بازیافت مرتبط محسوب می شود. لازم به ذکر است که در واژه مرکب «زیست شیمی»، «زیست» پیشوند است و نه واژه مستقل و «شیمی» هسته است در حالی که در واژه «محیط زیست» «زیست» واژه ای مستقل به شمار می رود و «محیط» هسته است. در واقع «زیست» در دو مثال فوق الذکر از دو مقوله و نقش متفاوت برخوردار است که اگر نظام های بازیابی مجهز به امکاناتی برای تشخیص و تعیین نوع آن ها بودند، مسلماً بازیابی بهبود چشم گیری می یافت. از دیگر چالش های صرفی بازیابی اطلاعات وام گیری واژگانی است. همان طور که می دانیم زبان ها از شیوه های مختلفی برای افزایش و غنی سازی واژگان خود بهره می برند که یکی از آن ها وام گیری است. در وام گیری، زبانی زبان مقصد عناصر زبانی را از زبان مبدأ به صورت های مختلفی مانند وام گیری، مستقیم غیر مستقیم ترجمه قرضی تعبیر، قرضی تغییر قرضی، ترجمه و تعبیر، قرضی آمیزش قرضی و تبادل قرضی وام می گیرد یکی از معضلات بازیابی اطلاعات خصوصاً در حوزه های تخصصی و علمی حاصل از وام گیری مستقیم است. در وام گیری مستقیم واژه ای از زبان

ص: 164

الف مستقیماً وارد زبان ب می شود این واژه پس از ورود به زبان ب بر طبق قواعد آوایی زبان مقصد تغییراتی پیدا می کند تا تلفظ آن برای گویشوران آسان شود (شقاقی 1386، 127-131). با توجه به این که در حوزه های تخصصی سرعت استفاده و جذب علم و فناوری بالاتر از سرعت واژه گزینی در زبان مقصد می باشد، لذا بسیاری از واژگان به همان صورت زبان مبدأ در زبان مقصد آوانگاری می شوند. بدیهی است چنین فرایندی عمیقاً تحت تأثیر نحوه تلفظ متخصص زمینه موضوعی بوده به طوری که امکان دارد یک واژه با املاهای متفاوتی نوشته شود؛ برای مثال سندروم / سندرم کلسیم / کالسیم، نیدروژن / هیدروژن این چندگانه نویسی باعث می شود که نه تنها ذخیره بلکه بازیابی اطلاعات نیز با چالش های فراوانی رو به رو شود به عنوان نمونه جستجوی معادل های آوانگاری شده واژه Psychology در پایگاه تخصصی مجلات نور ادعای پیش گفته را تأیید می کند برای واژه «پسیکولوژی» 95 یافته «پسیکولوژی» 3 یافته و «سایکولوژی» 8 یافته به دست آمد. لازم به ذکر است این نگارش های متفاوت به ترتیب بر اساس تلفظ، فرانسه آلمانی و انگلیسی وارد زبان فارسی شده است.

مسائل معنایی و بازیابی اطلاعات

زبان طبیعی از بُعد معنایی آکنده از ابهام است با وجود این کاربران زبان با استفاده از دانش زبانی و زمینه گفتمان در برخورد با این مسائل به ابهام زدایی می پردازند فلاحتی (1385) به طور کلی انواع ابهام های واژگانی را که می توانند به ابهام معنایی بی انجامند به 5 گروه عمده تقسیم می کند. ابهام مقوله ای ابهام حاصل از هم آوایی و هم نگاشتی ابهام چند معنایی و ابهام انتقالی ابهام مقوله ای در پیوند با معانی متفاوت یک واژه در بافت ها و نقش های مختلف همچون اسم، فعل، صفت و قید ایجاد می شود. در زبان فارسی این گونه از ابهام ها معمولاً به واسطه گذاشتن علائم واکه های کوتاه تا حد زیادی برطرف می شوند؛ اما مشکل اساسی این جاست که در متون مخصوص بزرگسالان و افراد باسواد (از جمله متون تخصصی و نوشتارهای علمی) استفاده از این علائم و مشخص کردن آن ها مرسوم نیست بنابراین در بسیاری از مواقع این خواننده است که با مراجعه به متن و بستر واژه، کار ابهام زدایی را انجام می دهد لازم به ذکر است در پاره ای از موارد امکان دارد تفاوتی میان تلفظ دو واژه وجود نداشته باشد اما همچنان به واسطه متفاوت بودن نقش دستوری هر واژه در جمله معنای متفاوتی از واژه مورد نظر دریافت گردد. به عنوان نمونه می توان به واژه «بردار» اشاره کرد.

بَردار: وزنه بردار دیابتی! قوی ترین مرد جهان (ورزش)

بُرदार: پیشرفت هایی در کنترل بردار رانش (ریاضیات)

بَردار: داستان بردار کردن حسنگ وزیر (ادبیات)

چنان چه این واژه در بستر های مختلف جستجو شود با توجه به این که در متن اصلی استفاده از علائم نشان دهنده واکه های کوتاه مرسوم نیست و حتی در دو مورد اول و سوم نمایش واکه نیز تفاوتی در ابهام معنایی این تک واژه ایجاد نمی کند لذا در صورت جستجو، حداقل در سه گروه

متفاوت مقالات و مطالبی بازیابی می شود این در حالی است که فعل امر «بَرَدار» از این گروه حذف شده است و حال آن که در جستجو و بازیابی اطلاعات احتمال بازیابی گروهی از مطالب که دارای این واژه هستند نیز وجود دارد تصاویر 1 (الف و ب) شمایی از نتایج جستجو در دو پایگاه اطلاعاتی فارسی را به نمایش می گذارد همان گونه که در این تصویر مشاهده می شود الگوریتم های بازیابی اطلاعات در نظام های جستجو و بازیابی اطلاعات در شرایط معمول با استفاده از میزان حضور واژه در متن و این که واژه در کجای متن باشد به رتبه بندی نتایج جستجو می پردازند. چنین روشی در پیوند با حالت های معمول زبانی مشکلی ایجاد نمی کند اما در مواجهه با مشکلات زبانی می تواند منجر به بروز ریزش کاذب در نتایج شود. لازم به توضیح است دو نمونه ارائه شده در تصویر الف هر دو برگرفته از صفحه اول نمایش نتایج هستند که انتظار می رود مرتبط ترین نتایج را در اختیار کاربر قرار دهند.

ابهام های حاصل از هم آوایی هم نگاشتی و چند معنایی گونه های دیگری از ابهام های زبان فارسی به شمار می آیند در این گروه از ابهام های واژگانی واژه ها یا تلفظ یکسانی دارند و یا به گونه یکسانی نوشته می شوند در نظام های بازیابی مبتنی بر جستجوی متنی ابهام های حاصل از هم نگاشتی مشکلاتی را برای جستجوگران اطلاعات به همراه دارد این مشکلات در پاره ای از اوقات با ابهام مقوله ای هم پوشانی پیدا می کنند یکی از نمونه های بارز این نوع از ابهام حداقل در پیوند با حوزه ای فیزیک واژه «گرم» و «گِرم» است (تصویر 2)

1 **پیشرفت هایی در کنترل بردار دانش**

رشدی، علی
 شماره مدرک: 1920583
 نوع مدرک: مقاله فارسی
 عنوان: روشی جدید - 20 - مسائل، 242
 سال: 1390
 تعداد صفحه: 5
 رتبه: 0.35

[ممن کمال](#) [اطلاعات بیشتری](#)

2 **وزنه بردار دبستانی آقوینوس مرد جهان**

شماره مدرک: 1238580
 نوع مدرک: مقاله فارسی
 عنوان: دبایت جسمانی، 47 - دوره 12 - مسائل، 46
 سال: 1389
 تعداد صفحه: 2
 رتبه: 0.35

[ممن کمال](#) [اطلاعات بیشتری](#)

« تصویربرداری ارتکی پلازما سولوی از ضایعات رنگدانه ای پوست: مقایسه، نوآوری و چالش »

آرزو ناگزی، محمدحسین صراری بگی، پوریا منصوری
 فصلنامه: انوار پزشکی، سال هشتم، شماره 2، بهار 1390، صص 13-6
[مشاهده متن](#) [PDF 5395B]

« برآورد بیرونی تعمیر یافته سیمپلکس میله‌گین توزیع نرمال چندمتغیره با مارتینس کوپولایم سه‌جمله »

شکوفه رین الهی، احمد بارسیان
 مجله علوم انسانی، شماره 4، بهار و تابستان 1389، ص 183
[مشاهده متن](#) [PDF 4775B]

« بچه دار، بچه رو بردار و بیار »

محمد آید
 ماهنامه شعبدار، شماره 51، شهریور 1391، ص 67

« جنگ شوکی بردار نیست »

هنرمندانه، شماره 142، زمستان و بهار 1391، ص 55
[مشاهده متن](#) [PDF 1644B]

تصویر 1 (الف و ب). نتایج بازیابی شده از دو پایگاه مقالات فارسی مرکز منطقه ای اطلاع رسانی علوم و فناوری نشریات کشور (مگ ایران). در هر دو مورد بازیابی بر مبنای الگوریتم‌های معمول محاسبه ربط با توجه به حضور واژه در جمله انجام گرفته است.

تصویر 1 (الف و ب). نتایج بازیابی شده از دو پایگاه مقالات فارسی مرکز منطقه ای اطلاع رسانی علوم و فناوری و بانک اطلاعات نشریات کشور (مگ ایران). در هر دو مورد بازیابی بر مبنای الگوریتم‌های معمول محاسبه ربط با توجه به حضور واژه در جمله انجام گرفته است



تصویر ۲. نمایی از نتایج جستجو در پایگاه مجلات تخصصی نور

در همین راستا در حوزه معنایی^۱ چالش دیگری تحت عنوان هم‌معنایی مطرح می‌شود. هم‌معنایی به وضعیتی گفته می‌شود که دو یا چند واژه به یک مفهوم واحد اشاره داشته باشند (صفوی ۱۳۸۷، ۱۶۰). بایستی به این نکته توجه داشت که بسیاری از علوم در ایران بومی نبوده و پژوهشگران و دانشجویان غالباً از طریق ترجمه آثار غربی با آنها آشنا می‌شوند. بنابراین، اصطلاح‌گزینی برای واژگان علمی توسط یک مرجع واحد صورت نگرفته و گاه مترجمان و پژوهشگران با سلیقه شخصی دست به واژه‌گزینی می‌زنند. به همین دلیل با وجود اینکه یکی از مهم‌ترین عوامل موفقیت ارتباط علمی شفافیت واژگانی است، در بسیاری از زمینه‌های موضوعی برای بیان یک مفهوم واحد از چندین واژه متفاوت فارسی استفاده می‌شود. برای مثال در حوزه زبان‌شناسی برای هر مفهوم به طور متوسط سه واژه علمی فارسی وجود دارد (احمدی نسب ۱۳۹۰، ۲۷۶). به همین دلیل، همانگونه که پیشتر نیز بیان شد، بازبایی اطلاعات در بسیاری از نظام‌ها برپایه واژگان موجود در متن صورت می‌پذیرد، و از آنجا که در متون مختلف مفاهیم به صورت‌های مختلف ثبت می‌شوند، لذا با جستجوی یک صورت از مفهوم، کاربر امکان دسترسی به صورت‌های دیگر همان مفهوم که منطقیاً مرتبط با نیاز اطلاعاتی وی خواهد بود را از دست می‌دهد. برای مثال برای مفهوم *Ontology* در فارسی سه معادل هستی‌شناسی، هستان‌شناسی و آنتولوژی بکار می‌رود که جستجوی هرکدام از این کلیدواژه‌ها تنها مدارکی را بازبایی می‌کند که صرفاً همان صورت واژگانی در آنها بکار رفته است در حالی که هر سه واژه به یک مفهوم اشاره دارد (جدول ۱).

1. Synonymy

تصویر ۲. نمایی از نتایج جستجو در پایگاه مجلات تخصصی نور

در همین راستا در حوزه معنایی^(۱) چالش دیگری تحت عنوان هم‌معنایی مطرح می‌شود. هم‌معنایی به وضعیتی گفته می‌شود که دو یا چند واژه به یک مفهوم واحد اشاره داشته باشند صفوی (۱۳۸۷، ۱۶۰) بایستی به این نکته توجه داشت که بسیاری از علوم در ایران بومی نبوده و پژوهش‌گران و دانشجویان غالباً از طریق ترجمه آثار غربی با آن‌ها آشنا می‌شوند بنابراین اصطلاح‌گزینی برای واژگان علمی توسط

یک مرجع واحد صورت نگرفته و گاه مترجمان و پژوهش‌گران با سلیقه شخصی دست به واژه‌گزینی می‌زنند به همین دلیل با وجود این که یکی از مهم‌ترین عوامل موفقیت ارتباط علمی شفافیت واژگانی است در بسیاری از زمینه‌های موضوعی برای بیان یک مفهوم واحد از چندین واژه متفاوت فارسی استفاده می‌شود. برای مثال در حوزه زبان‌شناسی برای هر مفهوم به طور متوسط سه واژه علمی فارسی وجود دارد (احمدی نسب 1390، 276). به همین دلیل همان‌گونه که پیش‌تر نیز بیان شد بازایی اطلاعات در بسیاری از نظام‌ها بر پایه واژگان موجود در متن صورت می‌پذیرد، و از آن‌جا که در متون مختلف مفاهیم به صورت‌های مختلف ثبت می‌شوند، لذا با جستجوی یک صورت از مفهوم کاربر امکان دسترسی به صورت‌های دیگر همان مفهوم که منطقاً مرتبط با نیاز اطلاعاتی وی خواهد بود را از دست می‌دهد. برای مثال برای مفهوم *Ontology* در فارسی سه معادل هستی‌شناسی، هستان‌شناسی و آنتولوژی بکار می‌رود که جستجوی هرکدام از این کلیدواژه‌ها تنها مدارکی را بازایی می‌کند که صرفاً همان صورت واژگانی در آن‌ها بکار رفته است در حالی که هر سه واژه به یک مفهوم اشاره دارد (جدول 1)

ص: 168

آسیب شناسی زبان و خط فارسی ... 169

جدول ۱. تفاوت بازیافت‌ها در پیوند با سه صورت واژگانی متفاوت از یک مفهوم واحد

بازگاه مقالات فارسی مرکز منطقه ای (Ricest)	بانک اطلاعات نشریات کنسور (مگ ایران)	بازگاه مجلات تخصصی نور
276	420	3877
8	7	69
8	5	100

ابهام انتقالی در بحث ترجمه از یک زبان به زبان دیگر مطرح می‌شود. زمانی که واژه از زبان مبدأ به زبان مقصد آمده اما در زبان مقصد چندین برابر نهاده پیدا می‌کند. برابرنهاده‌هایی که می‌توانند معانی متفاوت داشته و یا معانی مشابه و نگاهت‌های متفاوتی داشته باشند. این دست از واژگان ابهام‌برانگیز زمانی ایجاد مشکل می‌کنند که ابزار جستجو و بازیابی اطلاعات در قالب یک نظام بازیابی اطلاعات بین‌زبانی نیز فعالیت کند. به عنوان نمونه ۴ واژه *Transpiration* (تعرق)، *Sweating* (تعرق)، *Evapotranspiration* (تعرق و تبخیر)، و *Guttation* (تعریق) هر کدام در زبان انگلیسی به یک معنا به کار برده می‌شوند مثلاً در زبان مبدأ *Transpiration* برای گیاهان و در حوزه کشاورزی، و *Sweating* در پیوند با انسان به کار می‌رود. این درحالی است که واژه‌های پیش گفته در زبان فارسی معادل نسبتاً یکسانی دارند. هم اکنون گوگل تاحدی قابلیت جستجوی بین‌زبانی در سطح وب را فراهم آورده است. اما زمانی که معادل واژه *Transpiration* در صفحات فارسی جستجو شود نتایج مشابه با تصویر ۳ بدست خواهد آمد. این درحالی است که نتایج زمانی که واژه *تعرق* به فارسی برای بازیابی صفحات انگلیسی به کار رود به مراتب از پراکندگی بیشتری برخوردار خواهد بود؛ چرا که واژه *تعرق* در زبان فارسی برای بیان سه مفهوم متفاوت در زبان انگلیسی انتخاب شده است.



تصویر ۳. نمایش از صفحه اول نتایج بازیابی شده بر اساس اولین برابر نهاده پیشنهادی گوگل برای واژه

Transpiration

جدول 1 تفاوت بازیافت‌ها در پیوند با سه صورت واژگانی متفاوت از یک مفهوم واحد

ابهام انتقالی در بحث ترجمه از یک زبان به زبان دیگر مطرح می‌شود. زمانی که واژه از زبان مبدأ به زبان مقصد آمده اما در زبان مقصد چندین برابر نهاده پیدا می‌کند برابر نهاده‌هایی که می‌توانند معانی متفاوت داشته و یا معانی مشابه و نگاهت‌های متفاوتی داشته باشند. این دست از واژگان ابهام‌برانگیز زمانی ایجاد مشکل می‌کنند که ابزار جستجو و بازیابی اطلاعات در قالب یک نظام بازیابی اطلاعات بین

زبانی نیز فعالیت . کند به عنوان نمونه 4 واژه Transpiration (تعرق)، Sweating (تعرق)، Evapotranspiration (تعرق و تبخیر) و Guttation (تعریق) هر کدام در زبان انگلیسی به یک معنا به کار برده می شوند مثلاً در زبان مبدأ Transpiration برای گیاهان و در حوزه کشاورزی، و Sweating در پیوند با انسان به کار می رود این در حالی است که واژه های پیش گفته در زبان فارسی معادل نسبتاً یکسانی دارند هم اکنون گوگل تا حدی قابلیت جستجوی بین زبانی در سطح وب را فراهم آورده است. اما زمانی که معادل واژه Transpiration در صفحات فارسی جستجو شود نتایجی مشابه با تصویر 3 بدست خواهد آمد. این در حالی است که نتایج زمانی که واژه تعرق به فارسی برای بازیابی صفحات انگلیسی به کار رود به مراتب از پراکندگی بیش تری برخوردار خواهد بود؛ چرا که واژه تعرق در زبان فارسی برای بیان سه مفهوم متفاوت در زبان انگلیسی انتخاب شده است.

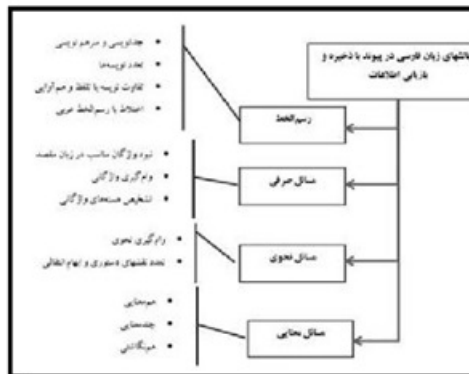
تصویر 3 نمایی از صفحه اول نتایج بازیابی شده بر اساس اولین برابر نهاده پیشنهادی گوگل برای واژه Transpiration

ص: 169

مسئله ایجاد یک نظام برای به هم نزدیک سازی ذهنیت پدیدآورنده و استفاده کننده از مسائل مهم در امر طراحی پایگاه های اطلاعاتی است. دسترسی مناسب به معنای ایجاد شرایط لازم از نظر زبان و واژگان اختصاص یافته به مدارک است؛ واژگانی که به کاربر در کاوش منابع یاری می رساند (مک ایوان 1379، 4) حال مشکل از این جا آغاز می شود که در پاره ای از پایگاه های اطلاعاتی، شخص سومی تحت عنوان نمایه ساز وارد شده و پاره ای از کلیدواژه ها را به مدارک با هدف بازیابی بهتر اختصاص می دهد در این فرایند با توجه به این که حلقه سومی در زنجیره ارتباطی میان پدیدآور و مخاطب ایجاد می شود اگر چه هدف بهبود بازیابی است اما گاه احتمال افزایش ضریب خطا اگر چه در ورود اطلاعات و چه در اختصاص واژگان کلیدی نیز وجود دارد (شاپوری 1379؛ Gross Taylor 2005) مسلماً دانش حوزه ای از طریق مطالعه حوزه ای حاصل می شود تحقیقات متعددی در تأیید این مطلب صورت گرفته است که وایلدرموت (Wildermuth، 2004) ضمن برشمردن برخی از آن ها تأکید می کند دانش حوزه ای و واژگانی افراد که عموماً برگرفته از متون تخصصی، است نقش عمده ای در انتخاب واژگان جستجو دارد در چنین شرایطی وظیفه نمایه ساز به عنوان حلقه واسط نزدیک کردن دو زبان پدیدآور و مخاطب است خصوصاً در هنگامی که نظام بازیابی اطلاعات قابلیت برقراری این ارتباط را بصورت مؤثر ندارد با وجود، این حجم زیاد متون ذخیره شده و کمبود نیروی انسانی توانمند متخصصان حوزه بازیابی اطلاعات را به سمت هرچه هوشمند ساختن بازیابی اطلاعات سوق داده است. تصویر 4 به منظور جمع بندی ساختواره ای درختی از اهم چالش های متصور در زمینه ذخیره و بازیابی اطلاعات را به نمایش می گذارد

عکس

مسئله ایجاد یک نظام برای به هم نزدیک‌سازی ذهنیت پدیدآورنده و استفاده‌کننده از مسائل مهم در امر طراحی پایگاه‌های اطلاعاتی است. دسترسی مناسب، به معنای ایجاد شرایط لازم از نظر زبان و واژگان اختصاص یافته به مدارک است؛ واژگانی که به کاربر در کاوش منابع یاری می‌رساند (سکایوان ۱۳۷۹، ۴۹). حال، مشکل از این جا آغاز می‌شود که در پاره‌ای از پایگاه‌های اطلاعاتی، شخص سومی تحت عنوان نمایه‌ساز وارد شده و پاره‌ای از کلیدواژه‌ها را به مدارک با هدف بازیابی بهتر اختصاص می‌دهد. در این فرایند با توجه به اینکه حلقه سومی در زنجیره ارتباطی میان پدیدآور و مخاطب ایجاد می‌شود، اگرچه هدف بهبود بازیابی است، اما گاه احتمال افزایش ضریب خطا چه در ورود اطلاعات و چه در اختصاص واژگان کلیدی نیز وجود دارد (شاپوری ۱۳۷۹؛ Gross & Taylor ۲۰۰۵). مسلماً دانش حوزه‌ای از طریق مطالعه حوزه‌ای حاصل می‌شود. تحقیقات متعددی در تأیید این مطلب صورت گرفته است که وایلدروموت (Wildermuth, ۲۰۰۴) ضمن برشمردن برخی از آنها، تأکید می‌کند دانش حوزه‌ای و واژگانی افراد که عموماً برگرفته از متون تخصصی است، نقش عمده‌ای در انتخاب واژگان جستجو دارد. در چنین شرایطی وظیفه نمایه‌ساز به عنوان حلقه واسط، نزدیک کردن دو زبان پدیدآور و مخاطب است. خصوصاً در هنگامی که نظام بازیابی اطلاعات قابلیت برقراری این ارتباط را بصورت مؤثر ندارد. با وجود این، حجم زیاد متون ذخیره شده و کمبود نیروی انسانی توانمند، متخصصان حوزه بازیابی اطلاعات را به سمت هرچه هوشمند ساختن بازیابی اطلاعات سوق داده است. تصویر ۴ به منظور جمع‌بندی ساختارهای درختی از اهم چالش‌های متصور در زمینه ذخیره و بازیابی اطلاعات را به نمایش می‌گذارد.



تصویر ۴. ساختار درختی چالش‌های عمده زبانی در پیوند با ذخیره و بازیابی اطلاعات^۱

۱. لازم به توضیح است این ساختار کامل‌تر از بحث‌های مطرح در این نوشتار است. دلیل این امر لزوم توجه به برخی از مهمترین موارد به منظور جلوگیری از طولانی شدن مبحث بوده است.

تصویر ۴. ساختار درختی چالش‌های عمده زبانی در پیوند با ذخیره و بازیابی اطلاعات (۱)

ص: 170

1- لازم به توضیح است این ساختار کامل‌تر از بحث‌های مطرح در این نوشتار است. دلیل این امر لزوم توجه به برخی از مهم‌ترین موارد به منظور جلوگیری از طولانی شدن مبحث بوده است

راهبردهای متفاوتی در بحث رفع این مسائل مطرح شده است که می توان آن ها را به گروه های مختلف دسته بندی نمود از میان آن ها عمده ترین موارد شامل تدوین پیکره زبان فارسی یادگیری ماشینی ریشه یاب ها، نرمال سازی، نویسه ها استفاده از الگوریتم های احتمالاتی همچون N-gram ها و مهار واژگانی می شود برای رفع چالش های معنایی بازیابی اطلاعات به نظر می رسد که باید بر اساس یک پیکره زبانی و استخراج قواعد واژه سازی زبان فارسی و تعیین اجزای واژگانی همچون وندها و ریشه ها عمل کرد این کار از طریق یادگیری ماشینی از یکسو و ایجاد پیکره های زبان فارسی به صورت الکترونیک امکان پذیر است حرکت هایی در این زمینه در پیوند با زبان فارسی صورت گرفته است اما همچنان در مراحل ابتدایی است در این راستا یکی از کارهای ارزشمند انجام شده پارس مورف متعلق به موجی، اسلامی و وزیر نژاد (1390) است. این نرم افزار قابلیت تحلیل صرفی زبان فارسی را داراست. در همین رابطه می توان به ریشه یاب تهیه شده توسط احسان و فیلی (1390) نیز اشاره کرد که میزان دقت در نظام های بازیابی را 5 درصد افزایش می دهد. از دیگر تلاش ها می توان به ریشه یاب برنجیان (1390) برای افعال ماضی و مضارع ناگذر اشاره کرد. لازم به یادآوری است که حرکت های مشابه برای زبان عربی مدت هاست آغاز شده و به مرحله ارزیابی و بهسازی رسیده است. به عنوان نمونه می توان به پژوهش های محمود ماجد و خلدون (Mahmoud, Majed Kaldoun 2011) و یا انیس و دیگران (Anis et al. In Press) اشاره کرد.

از راهبرد های پیشنهادی در پیوند با برطرف سازی مشکلات مرتبط با تعدد نویسه هایی همچون الف و یا همزه در زبان عربی نرمال سازی است. در این راهبرد با استفاده از یک الگوریتم ساده، کلیه صورت های مختلف یک نویسه به طور خودکار به یک شکل یکسان تبدیل می شود به عنوان مثال کلیه نگارش های متفاوت «آ» و «ا» به صورت «1» در نظر گرفته می شود (Mahmoud, Majed Kaldoun 2011). اما با بررسی های صورت گرفته و همان گونه که پیش تر در این مقاله نیز مطرح گردید زبان فارس به واسطه ویژگی های معنایی خاص خود بر تابنده چنین تغییری، حداقل در پیوند با این مثال نمی باشد لذا نرمال سازی در پیوند با نویسه های زبان فارسی بایستی با دقت و مطالعه بیش تری صورت گرفته و نمی توان به یافته های مطالعات مرتبط با زبان عربی در این زمینه استناد کرد راهبرد دیگر در این راستا استفاده از ریشه یاب ها است استفاده از ریشه یاب امکان بازیابی صورت های مختلف واژگانی که به یک حوزه مفهومی واحد اشاره دارند را فراهم می کند. این در حالی است که این راهبرد می تواند به صورت بالقوه پاسخی به مشکلات نرم افزارهای بازیابی اطلاعات در پیوند با صورت های مفرد و جمع واژگان خصوصاً جمع های مکسر و یا واژه های هم خانواده باشد (Mahmoud, Majed Kaldoun 2011) با توجه به این که زبان فارسی از نظر رسم الخط مشابه با زبان عربی است پیشنهاد می شود در این زمینه از مطالعات صورت گرفته در پیوند با خط عربی و بازیابی اطلاعات استفاده شود در حالی که در زمینه ویژگی های صرفی و معنایی زبان فارسی به زبان انگلیسی نیز نزدیک بوده و مطالعات صورت گرفته در این زمینه تا حدود زیادی

اگر چه در بسیاری از متون بحث رعایت فاصله گذاری مناسب به پدیدآور نسبت داده می شود اما هم چنان نیازمند نظام هایی با قابلیت تشخیص مناسب فاصله ها در واژگان مرکب (1) هستیم. در زبان فارسی یکی از شیوه های پر کاربرد واژه سازی استفاده از فرایند ترکیب است لذا تعداد واژگان ترکیبی در زبان فارسی بسیار زیاد است این در حالی است که در زبان فارسی به واسطه ماهیت زبانی و دستوری و نیز سهولت، خواندن واژگان مرکب معمولاً به صورت ناپیوسته همراه با یک فاصله یک حرفی ثبت می شوند. لازم به یادآوری نیست که در زبان انگلیسی این مشکل کمتر دیده می شود. به عنوان نمونه کلیه واژگانی که با «شناسی» همراه هستند. معمولاً به صورت دو جزء جداگانه ثبت می شوند در حالی که در زبان انگلیسی "logy" بخشی از واژه به شمار می آید. بدیهی است رفع این مشکل تنها از طریق استفاده از عملگر And در یک فیلد جستجو امکان پذیر نمی باشد، چرا که منجر به افزایش تعداد بازیافت های غیر مرتبط می شود این در حالی است که چنان چه کاربری کلمات وارد شده در یک فیلد به شکل عبارتی تغییر یابد امکان بروز این خطا کمتر خواهد شد. اما با نگاهی به متون موجود در پیوند با رفع این دسته از چالش ها، حداقل در پیوند با زبان عربی استفاده از الگوریتم های N-gram پیشنهاد می شود. دلیل این امر قابلیت ارائه بافت واژه مورد جستجو همراه با واژه در قالب پیشنهادهایی برای بهسازی جستجوست (Mahmoud, Majed Kaldoun 2011).

مهار واژگانی یکی دیگر از راهبرد های مفید در زمینه بهسازی جستجو از طریق کم اثر سازی تعدد صورت های واژگانی و یا معنایی است بدین معنا که با ارائه صورت های مرجح، نامرجح و واژگان شامل و زیر شمول کلیه صورت های متفاوت یک مفهوم را ذیل یک مفهوم واژه پذیرفته شده گرد می آورد. واژگان مهار شده معمولاً نمایان گر ساختواره ای از روابط معنایی و سلسله مراتبی بوده و استفاده از آن خصوصاً در برطرف ساختن چالش های هم معنایی، چندمعنایی، املاهای متعدد و وام گیری واژگانی کاراست (Svenonius 2003). این در حالی است که استفاده همزمان از اصطلاحنامه ها در کنار کلید واژه های زبان طبیعی یکی از کاراترین روش ها در این راستا به شمار می آید که از جمله پایگاه های موفق در این زمینه می توان به INSPEC، Eric و LISA اشاره کرد.

در پایان خاطر نشان می سازد که تاکنون پیکره های فارسی مختلفی همچون همشهری محک بی جن خان و دادگان زبان فارسی تدوین شده است از این میان درودی و دیگران (1387) مجموعه محک وب را تهیه کرده اند که از آن می توان برای انجام مطالعات در حوزه بازیابی اطلاعات فارسی در وب استفاده نمود. این مجموعه dotIR نام دارد و استفاده از آن برای امور غیر تجاری رایگان است. مجموعه پیش گفته از یک پیکره استاندارد 50 پرس و جوی استاندارد 18 هزار دآوری تعیین ربط پرس و جویها به اسناد پیکره و 50 هزار بردار ویژگی استخراج شده از اسناد، تشکیل شده و بیش تر مناسب پژوهش در رابطه با بازیابی اطلاعات است همشهری (2) مجموعه دیگری است که

1- در این جا برای سهولت بحث واژه مرکب به واژه غیر بسیط اطلاق شده و معنای دقیق زبان شناختی آن مدنظر نیست 2.

<http://ece.ut.ac.ir/dbrg/hamshahri/faindex.html>

2- <http://ece.ut.ac.ir/dbrg/hamshahri/faindex.html>

توسط رهگذر و دیگران بر اساس اخبار روزنامه همشهری تهیه شده و دارای دو نسخه 1 و 2 است که نسخه دوم از امکانات بیش تری مانند پیوند به تصاویر و به اصل صفحات وب برخوردار است. از این مجموعه نیز می توان برای انجام پژوهش در حوزه بازیابی اطلاعات، فارسی، پردازش زبان طبیعی و تدوین الگوریتم های ریشه یابی استفاده نمود. بر اساس این مجموعه پژوهش های متعددی صورت گرفته که از جمله می توان به آل احمد و دیگران (2007) اشاره نمود که با استفاده از این مجموعه مدل فضای برداری بر اساس n-gram و کلمه را ارزیابی کرده و نشان داده اند که بازیابی متون فارسی بر اساس مدل فضای برداری 4-gram و طرح وزن دهی Inu ltu نتایج قابل قبول و تحلیل محلی متن (1) بهترین نتایج را در بر خواهد داشت. از پیکره های زبانی دیگر می توان به پیکره بی جن خان (2) اشاره کرد که از متون خبری و متون عمومی شامل 4300 موضوع تشکیل شده و بر اساس 40 مقوله دستوری فارسی بر چسب دهی شده است و مناسب پژوهش های مرتبط به پردازش زبان طبیعی است. نمونه دیگر پایگاه دادگان فارسی (3) است که توسط عاصی و دیگران در پژوهش گاه علوم انسانی و مطالعات فرهنگی تهیه شده است این پایگاه گونه های مختلف زبانی نوشتاری و گفتاری را در بر می گیرد و در آن امکان جستجوی واژه ها ترکیب ها و بررسی بسامد آن ها به همراه گزارش های آماری از متون وجود دارد. در این جا ذکر این نکته خالی از لطف نیست که متون موجود و پژوهش های انجام شده نشان دهنده این امر است که علی رغم وجود این پیکره های زبانی فارسی، پژوهش های چندانی بر اساس این پیکره ها صورت نگرفته است شاید بتوان این امر را به عدم اطلاع اکثر پژوهش گران از اهمیت و نقش پیکره های زبانی در تحقیقات در حوزه پردازش زبان، طبیعی بازیابی اطلاعات و تهیه ریشه یاب ها و همچنین عدم آموزش کافی جامعه علمی در این زمینه نسبت داد. در این نوشتار نیز با توجه به چارچوب زمانی در نظر گرفته شده برای بررسی، حاضر متأسفانه محققان قادر به بررسی دقیق این پیکره ها و استفاده از آن ها در این بررسی نشدند. بنابراین لازم است تا پژوهشی جامع، میزان موفقیت راهبردهای پیشنهادی این بررسی بر پایه پیکره های پیش گفته را بسنجد. هم چنین توصیه می شود که در طراحی و به سازی پایگاه های اطلاعاتی زبان فارسی کارگروه هایی متشکل از متخصصان علم اطلاعات و دانش شناسی زبان شناسی و علوم رایانه نقش اصلی را ایفا نمایند.

منابع

آشوری، داریوش 1375 چند پیشنهاد دیگر برای اصلاح خط فارسی سرهم نویسی و جدا نویسی. نگاه نو. 28:101-117.

احمدی نسب، فاطمه . 1390 تهیه و تدوین اصطلاحنامه تک زبانه فارسی زبان شناسی رساله دکترا دانشگاه علامه طباطبایی.

اسلامی محرم 1381 دشواری های پردازش رایانه ای خط فارسی نشر دانش دوره 19(3) 28-32

ص: 173

Local context analysis -1

<http://ece.ut.ac.ir/DBRG/Bijankhan> -2

<http://ece.ut.ac.ir/DBRG/Bijankhan> -3

اسلامی، محرم 1386 خط فارسی و رسانه های گروهی دو فصل نامه پردازش علائم و داده ها. 8(2):93-98.

برنجیان شاپور رضا 1390. ریشه یاب ماضی و مضارع از مصدر افعال ناگذر در زبان فارسی. شیراز نوید شیراز مرکز منطقه ای اطلاع رسانی علوم و فناوری

دبیر مقدم، محمد. 1383. زبان شناسی نظری پیدایش و تکوین دستور زایشی (ویراست دوم) تهران انتشارات سمت.

درودی احسان و دیگران 1387 پیکره محک استاندارد برای تحقیقات بازیابی اطلاعات وب فارسی گزارش، فنی گروه تحقیقاتی پایگاه داده ها دانشگاه تهران شماره: <http://ccc.ut.ac.ir/DBRG/webir/files/Papers/WebIR.pdf>. DBRG-TR-138702 (دسترسی در 1391/10/19).

شاپوری، سودابه (1379) بررسی مشکلات جستجوی موضوعی استفاده کنندگان رایانه ای کتابخانه مرکزی دانشگاه فردوسی مشهد پایان نامه کارشناسی ارشد مشهد دانشکده علوم تربیتی و روانشناسی

شقایقی ویدا. 1386. مبانی صرف تهران: انتشارات سمت.

صفر مقدم، احمد. 1386. فاصله گذاری در خط فارسی نامه فرهنگستان 9/4: 123-137.

صفوی، کورش 1387. درآمدی بر معنی شناسی. تهران: انتشارات سوره مهر (حوزه هنری سازمان تبلیغات اسلامی)

فرهنگستان زبان و ادب فارسی 1389. دستور خط فارسی چاپ نهم تهران: فرهنگستان ادب و زبان فارسی (نشر آثار).

فلاحی فومنی، محمدرضا 1385 ابهام در ماشین ترجمه کتابداری و اطلاع رسانی. 9(3): 21 - 38

کامری، برنارد جان ماونتفورد ویوین، لا و د. آ. کرو ز 1990. زبان های دنیا چهار مقاله در زبان شناسی ترجمه کورش صفوی 1384 تهران: انتشارات سعاد.

گل تاجی، مرضیه و سعیده برزگر. 1389 بررسی مشکلات ریخت شناسی زبان فارسی در سه پایگاه اطلاعاتی مرکز منطقه ای اطلاع رسانی علوم و فناوری پژوهشگاه اطلاعات و مدارک علمی ایران و جهاد دانشگاهی کتابداری و اطلاع رسانی 13(2): 191-214

مک ایوان اندرو. (1379) استفاده از سرعنوان های موضوعی کتابخانه کنگره: هزینه همکاری برای رسیدن به دسترس پذیری ترجمه مجتبی اسدی گزیده مقالات ایفلا 98 (آمستردام: 21 - 16 اوت 1998) (ص. 59 - 47) تهران کتابخانه ملی جمهوری اسلامی ایران.

محقق زاده محمد صادق و کاظم زارعیان 1383 ارائه راه حل برای برخی مسائل اتوماسیون و نگارش فارسی. فصلنامه اطلاع-رسانی 19(3 و 4): 1-10.

معصومی همدانی حسین 1381 خط فارسی و رایانه نشر دانش. 19(2) 2-6

مواجی وحید محرم اسلامی و بهرام وزیر نژاد 1390. پارس مورف: تحلیل گراف صرفی زبان

- Aleahmad, Abolfazl. et al. 2007. N-Gram and Local Context Analysis for Persian Text Retrieval, International Symposium on Signal Processing and Its Applications, Sharjah U.A.E. Retrieved 2013-01-8
From <http://ece.ut.ac.ir/dbrg/hamshahri/files/Papers/isspa.pdf>
- Anis, Z. et al. In Press. Contribution to Semantic Analysis of Arabic language. Advances in Artificial Intelligence. Retrieved 2012-11-10 From <http://www.hindawi.com/journals/aai/aip/620461>
- .Crystal, David. 2008. A dictionary of linguistics and phonetics. Six Edition. Blackwell Publishing
- Gross, Tina, Arlene G. Taylor. 2000. What Have We Got to Lose? The Effect of Controlled Vocabulary on Keyword Searching Results. College Research Libraries. 66(3): 212 -230. Retrieved August 2, 2006, From www.ala.org/ala/acrl/acrlpubs/crljournal/backissues2005a/crlmay05/Gross.pdf
- Mahmoud, R., S. Majed, and Z. Kaldoun. 2011. Improving Arabic information retrieval system using N-gram method. WSEAS Transactions on Computers 10 (4): 125-133
- Svenonius, E. 2003. Design of Controlled Vocabularies. Encyclopedia of Library and Information Science) (822-838). New York: Marcel Dekker (reprinted from the first edition, 1989
- Wildermuth, Barbara M. 2004. The Effects of Domain Knowledge on Search Tactic Formulation. Journal of the American Society for Information Science and Technology. 55(3): 246-258

هدف: هدف اصلی از انجام این پژوهش ارزیابی کاربرد پذیری وب سایت نهاد کتابخانه های عمومی کشور است.

روش شناسی: پژوهش حاضر که از نوع کاربردی است، با استفاده از شیوه ارزیابانه به بررسی و ارزیابی معیارها و مولفه های مطرح در کاربرد پذیری وب سایت نهاد کتابخانه های عمومی کشور می پردازد. ابزار گردآوری اطلاعات در این پژوهش سیاهه ارزیابی محقق ساخته مشتمل بر 11 معیار و 160 مؤلفه است که وب سایت مورد نظر با آن سنجیده شد. در سیاهه مزبور از روش روایی صوری و محتوایی استفاده شده منظور تجزیه و تحلیل یافته های پژوهش از آمار توصیفی فراوانی درصد و میانگین برای به توصیف کشیدن وضعیت موجود وب سایت مورد بررسی استفاده شد.

یافته ها: در مجموع نتایج به دست آمده نشان داد در وب سایت نهاد کتابخانه های عمومی کشور نیمی از استانداردها 346 امتیاز (52/2 درصد) از 663 امتیاز وب سایت شاخص رعایت گردیده است.

کلید واژه ها: ارزیابی، کاربرد پذیری وب سایت، نهاد کتابخانه های عمومی کشور

صدیقه محمد اسماعیل (1) | ماهرخ ناصحی اسکویی (2)

مقدمه

امروزه از وب به عنوان یکی از مهم ترین و اصلی ترین ابزارها برای دسترسی به اطلاعات استفاده می کنند لذا ضروری است، کاربران وب میزان کاربرد پذیری منابع موجود بر روی آن (به ویژه وب سایت های آن) را مورد بررسی و ارزیابی مستمر قرار دهند در این راستا پژوهش حاضر بر آن است تا وب سایت نهاد کتابخانه های عمومی کشور را مورد ارزیابی قرار دهد تا از این طریق مشخص گردد وب سایت نهاد به لحاظ کاربردپذیر بودن در چه شرایطی قرار دارد و تا چه حد در طراحی آن اصول و ضوابط طراحی وب سایت به لحاظ رعایت مؤلفه های مربوط به این امر (اعم از: اعتبار، صحت، روزآمد بودن، سطح پوشش و مخاطبان خاص، وجود نماهای تعاملی و تبادلی عینیت، اطلاعات قابلیت ناویری، نماهای غیر متنی، دسترس پذیری، کارآمدی، ویژگی های ظاهری رعایت گردیده است. بدیهی است، پرداختن به وب سایت نهاد کتابخانه های عمومی کشور با عنایت به جایگاه محوری آن در بحث اطلاعات و اطلاع رسانی امری حائز ارزش و توجه است.

ص: 177

1- استادیار گروه کتابداری و اطلاع رسانی دانشگاه آزاد اسلامی واحد علوم تحقیقات تهران M.esmaeili2@gmail.com

2- دانشجوی کارشناسی ارشد کتابداری و اطلاع رسانی دانشگاه آزاد اسلامی واحد علوم و تحقیقات تهران

mahrokh.nasehi@gmail.com

2- روش پژوهش و توجیه روایی آن

روش پژوهش حاضر، روش کتابخانه ای (سندی) و پیمایشی از نوع ارزیابانه است. بدین معنا که، با استفاده از سیاهه کنترل یا به ارزیابی کاربردپذیری وب سایت نهاد کتابخانه های عمومی کشور (بر مبنای معیارهای متعدد موجود در چک لیست) پرداخته شده است. در توجیه روایی استفاده از چنین روشی جهت انجام این پژوهش ذکر این نکته ضروریست که از آن جایی که شاخص های بیرونی متعددی در قالب سیاهه واریسی بوده که چگونگی کاربردپذیری صفحات وب نهاد کتابخانه های عمومی کشور با آن ها سنجیده شده است لذا مناسب ترین راه برای آگاهی از وضعیت این وب سایت ها (بر اساس معیارها و مؤلفه های ذکر شده) روش ارزیابانه فوق بوده است.

3- شیوه گردآوری اطلاعات و تجزیه و تحلیل آن ها

در این پژوهش گردآوری اطلاعات با استفاده از روش مشاهده مستقیم وب سایت و بر اساس سیاهه ای با 11 معیار و 160 مؤلفه صورت گرفت و بر این مبنای وب سایت نهاد کتابخانه های عمومی کشور از نقطه نظر کاربردپذیری و معیارهای مطرح در آن ارزیابی گردید. برای این منظور، تلاش شد تا ابتدا نسبت به تهیه سیاهه مربوطه تحت عنوان سیاهه ارزیابی کاربردپذیری وب سایت (1) اقدام گردد. از این رو، متون و منابع و پیشینه های موجود در زمینه کاربردپذیری وب سایت ها در داخل و خارج کشور مورد مطالعه و تفحص قرار گرفت و به منظور اطمینان از روایی محتوایی هر چه بیش تر در اختیار بررسی 5 تن از مدیران وب سایت های کتابخانه ای قرار گرفت سپس نظرات آن ها، بررسی و نسبت به تهیه نسخه نهایی سیاهه (مشتمل بر 11 معیار کلی و 160 مؤلفه) اقدام شد و با استفاده از آن وضعیت وب سایت نهاد کتابخانه های کشور در طول مدت زمان انجام این پژوهش سه ماهه اول (1391)، مورد بررسی و ارزیابی قرار گرفت. به منظور تجزیه و تحلیل یافته های پژوهش از آمار توصیفی (فراوانی درصد و میانگین) برای به توصیف کشیدن وضعیت موجود استفاده شد. در سیاهه یاد شده از دو مقیاس وجود و نبود بلی (ii) و خیر (-) استفاده شده که امتیازات در نظر گرفته شده برای آن ها به ترتیب عبارتند از بلی = 1 (یک) و خیر = 0 (صفر). (نکته: اعداد مندرج در قسمت تواتر در جدول به شماره منابع ارجاع خورده است).

4- تجزیه و تحلیل داده ها و ارائه یافته ها

همان گونه که قبلاً نیز گفته شد در این پژوهش از یک وب سایت مفروض استفاده شده با این فرض که در طراحی آن هر 160 مؤلفه یا ویژگی به گونه ای رعایت شده است که می توان آن را به مثابه شاخصی برای ارزیابی و سنجش دیگر وب سایت جامعه مورد مطالعه به کار برد. یافته های پژوهش بیان می کند که در وب سایت شاخص امتیاز کل رعایت همه معیارهای مطرح در کاربردپذیری سایت برابر با 663 امتیاز است که از این امتیاز 107 امتیاز مربوط به رعایت مؤلفه های اعتبار اطلاعات (16 درصد)

ص: 178

32 امتیاز مربوط به صحت اطلاعات (5 درصد)، 26 امتیاز مربوط به روزآمد بودن (4 درصد)، 21 امتیاز مربوط به سطح پوشش و مخاطبان خاص (3 درصد) 31 امتیاز مربوط به نماهای تعاملی و تبادلی (5 درصد)، 34 امتیاز مربوط به عینیت اطلاعات (5 درصد)، 166 امتیاز مربوط به رعایت معیار کلی ناوبری در سایت (جمعا: 25 درصد)، به تفکیک شامل: 9 امتیاز مرتبط با ویژگیهای عنوان مرورگر (1 درصد) 22 امتیاز مرتبط با ویژگیهای عنوان صفحه (3 درصد) 92 امتیاز مرتبط با پیوندهای متنی و فرامتنی (14 درصد)، 6 امتیاز مرتبط با نشانه اینترنتی (1 درصد)، 25 امتیاز مرتبط با نقشه سایت یا نمایه (4 درصد) 12 امتیاز مرتبط با موتور جستجوی داخلی (2 درصد)، 64 امتیاز مربوط به نماهای غیرمتنی (10 درصد)، 77 امتیاز مربوط به قابلیت دسترس پذیری (12 درصد)، 99 امتیاز مربوط به کارآمد پذیری (15 درصد) و 6 امتیاز مربوط به ویژگی های ظاهری (1 درصد) است

همان گونه که از جدول 1 مستفاد می شود یافته های تحقیق نشان می دهد که وب سایت نهاد کتابخانه های عمومی کشور از لحاظ رعایت معیارهای «اعتبار اطلاعات» 45 امتیاز (42 درصد)، «صحت اطلاعات» امتیاز (53 درصد)، «روزآمد بودن» 7 امتیاز (26 درصد)، «سطح پوشش و مخاطبین خاص» 10 امتیاز (47 درصد) «نماهای تعاملی و تبادلی» دارای 12 امتیاز (38 درصد)، «عینیت اطلاعات» 15 امتیاز (44 درصد) از لحاظ معیارهای کلی «ناوبری» 103 امتیاز (62 درصد) «نماهای غیر متنی» 18 امتیاز (28 درصد) «دسترس پذیری» 48 امتیاز (62 درصد)، «کارآمدی» 51 امتیاز (51 درصد) و در نهایت در رابطه با آخرین معیار مورد بررسی که ویژگی های «ظاهری» 6 امتیاز (1000 درصد) را به طور کامل به خود اختصاص داده است. یافته های این پژوهش با عصاره که به بررسی وب سایت های 58 کتابخانه ملی جهان با هدف شناسایی ویژگی های شاخص در طراحی وب سایت یک کتابخانه ملی به روش تحلیل محتوا و ارزیابی بر مبنای سیاهه واری (وب سایت کتابخانه ملی آمریکا 80 درصد، وب سایت کتابخانه بیمارستانی انگلستان 28 درصد و وب سایت کتابخانه ملی کشورهای مالزی و دانمارک 72 درصد با سیاهه واری و وب سایت کتابخانه ملی ایران 51 درصد با سیاهه همخوانی داشتند) و نیز با یافته های تحقیق حاجی زین العابدینی که به بررسی وضعیت وب سایت های کتابخانه های ملی جهان از نظر میزان مورد استفاده قرار گرفتن از طریق پیوندهای آن ها می پردازد همسو است. و در بعد دسترس پذیری نیز با یافته های محمد اسماعیل و کاظمی در تحقیقی با عنوان دسترس پذیری وب سایت کتابخانه های ملی، خاورمیانه با استفاده از سیاهه واری ای مشتمل بر 13 مؤلفه که در آن به بررسی تطبیقی وب سایت سازمان اسناد و کتابخانه ملی جمهوری اسلامی ایران با وب سایت های کتابخانه های ملی کشورهای اسلامی در منطقه خاورمیانه می پردازد کاملا مطابقت و همسویی دارد.

۱۸۰ مدیریت منابع اطلاعاتی وب

جدول ۱. سیاهه ارزیابی کاربردپذیری وب سایت نهاد کتابخانه‌های عمومی کشور از نظر میزان رعایت مؤلفه‌ها (به تفکیک)

نهاد کتابخانه‌های عمومی کشور		اعتبار اطلاعات	
ردیف	تواتر	مؤلفه‌ها	
۱	۱,۴۵,۱۰,۱۵,۱۷,۱۸,۱۹,۲۳,۲۵	بیان نام مولف	-
۲	۱,۱۰,۳۵,۱۵,۱۶,۱۷,۱۸,۱۹,۲۳	بیان ویژگی‌های مولف (صلاحیت, شهرت, اعتبار و...)	-
۳	۶,۸,۱۰,۴۵,۱۹,۴۲, ۱	چگونگی تماس با سازمان یا شخص مسئول محتویات سایت	✓
۴	۱,۶,۱۵,۱۸,۲۱,۲۳,۲۶	دارا بودن تاییده رسمی از سوی سازمان یا فرد مسئول محتوی سایت	-
۵	۶,۲۵,۸,۱۵,۱۷,۴۲ و ۱	ذکر مشخصات سازمان یا شخص مسئول محتویات سایت	✓
۶	۱,۸,۱۵,۲۵,۳۷,۴۲	ارائه فهرستی از کارگزاران اصلی و مشخصات آن‌ها	-
۷	۲۴,۲۵,۲۹,۴۰, ۲۷,۱,۸,۱۰, ۱۵, ۱۹	ذکر راهی برای تماس با مولف صفحه (نشانی پستی, شماره تلفن و پست الکترونیک و ...)	✓
۸	۱,۱۵,۱۶,۱۸,۱۹,۲۰	بیان راهی برای اثبات ویژگی‌های مولف (تجارب وی در زمینه موضوعی خاص, عضویت وی در سازمان‌های حرفه‌ای و...)	-
۹	۱,۸,۲۷,۱۵,۱۷,۱۹	ذکر نام سازمان مسئول محتویات سایت	✓
۱۰	۶,۱۸,۱۹, ۱	بیان ماهیت حامی	-
۱۱	۷,۴۲, ۱,۸	ذکر مدت زمان تاسیس سازمان	-
۱۲	۱,۸,۱۵,۳۵,۲۲	بیان حق مالکیت معنوی منابع اطلاعاتی عرضه شده در سایت (ذکر نام)	✓
۱۳	۱,۸,۱۵,۴۲	بیان نام سازمان پشت پرده سایت	-
۱۴	۱,۸,۱۰, ۱۵	ارائه سیاهه‌ای از نام سازمان‌ها یا سایت‌هایی که این سایت را توصیه می‌کنند	-
۱۵	۱,۱۵,۲۹,۱۶,۲۵	مشخص بودن هدف از طراحی و نشر صفحه	-
۱۶	۱,۸,۱۰, ۱۵	ارائه فهرستی از اسامی و مشخصات افراد مسئول نظارت بر سازمان	-
۱۷	۱, ۱۵	ارائه فهرستی از منابع چاپی منتشر شده توسط سازمان	-
۱۸	۲۷,۲۲,۴۰	قرار دادن نام و آرم سازمان سایت در هر صفحه	✓
۱۹	۳۰	نشان دادن آرم سازمان بصورت واضح و برجسته	✓
۲۰	۲۲	قرار دادن نام سازمان سایت در بالای صفحه	✓
۲۱	۱, ۱۵	ذکر نام دارنده حق مالکیت معنوی	✓
۲۲	۹	قراردادن پرچم کشور	✓
۴۵	جمع امتیازها	۱۰۷	
نهاد کتابخانه‌های عمومی کشور		صحت اطلاعات	
ردیف	تواتر	مؤلفه‌ها	
۲۳	۸,۲۵,۳۶,۱۶,۱۹, ۲۷,۲۲,۴۰, ۱,۶, ۱۰	نبود خطاهای املائی, گرامری و تایی	✓
۲۴	۱,۱۰,۲۲,۱۵,۱۷,۲۴,۱۸,۲۳	ذکر نام و مشخصات کتابشناختی منبع اصلی	-
۲۵	۱,۱۰,۱۵,۱۹,۲۸,۱۷,۱۹	وجود شاخصی دال بر بررسی صحت و سقم اطلاعات توسط ویراستار یا مصحح در طول فرآیند بازنگری مجدد منابع	-
۲۶	۸,۳۵, ۲۷,۳۹, ۱, ۱۰	عنوان بندی روشن و واضح گراف, نمودار یا جداول موجود در صفحه	✓
۱۷	جمع امتیازها	۳۲	

جدول 1 سیاهه ارزیابی کاربردپذیری وب سایت نهاد کتابخانه‌های عمومی کشور از نظر میزان رعایت مؤلفه‌ها (به تفکیک)

ارزیابی کاربردپذیری وبگاه نهاد کتابخانه‌های عمومی کشور ۱۸۱

نهاد کتابخانه‌های عمومی کشور		روزآمد بودن اطلاعات	
ردیف	تواتر	مؤلفه‌ها	
۲۷	۷,۱۰,۱۵,۱۷,۱۸,۱۹,۲۳,۲۵,۱	ذکر تاریخ آخرین تجدید نظر در محتویات صفحه	-
۲۸	۸,۱۶,۱۷,۱۸,۱۹,۲۳,۱,۱۰,۳۵	ذکر نخستین تاریخ قرار گرفتن منبع اطلاعاتی (با هر فرمتی) بر روی صفحه وب	-
۲۹	۶,۸,۱۰,۱۵,۱۹,۴۲,۱	ذکر فواصل زمانی به روز کردن اطلاعات دارای حساسیت زمانی	-
۳۰	۱,۶,۱۵,۱۸,۲۱,۲۳,۲۶	وجود اطلاعات آماری در صفحه	✓
۳۱	۱,۶,۲۵,۸,۱۵,۱۷,۴۲	بیان تاریخ گردآوری آمار	-
۳۲	۱,۸,۱۵,۲۵,۳۷,۴۲	ارائه تاریخ‌ها در یک فرمت بین‌المللی	-
۷	۲۶	جمع امتیازها	
نهاد کتابخانه‌های عمومی کشور		سطح پوشش مخاطبان خاص	
ردیف	تواتر	مؤلفه‌ها	
۳۳	۴۴,۴۰,۱,۶,۱۰,۸,۱۵,۱۶,۱۸,۲۸	مشخص بودن نام و نوع منابع موجود در صفحه	✓
۳۴	۱,۶,۱۰,۸,۱۵,۱۸,۲۸	تعیین مخاطبان خاص صفحه	-
۳۵	۱,۶,۱۰,۱۵	درج زمان تخمینی تکمیل صفحه در دست ساخت	-
۱۰	۲۱	جمع امتیازها	
نمایهای تعاملی و تبادل		نهاد کتابخانه‌های عمومی کشور	
ردیف	تواتر	مؤلفه‌ها	
۳۶	۲۷,۳۰,۱,۸,۴,۱۵,۱۱,۲۶,۳۸	وجود نظام مشخصی جهت بازخورد کاربران (feedback)	✓
۳۷	۱,۸,۴,۱۵,۲۶,۳۸	وجود نظام مشخصی برای کاربران برای درخواست اطلاعات بیشتر از سازمان	-
۳۸	۱,۲۷,۱۵,۲۶,۳۸	ذکر مشخصه زمانی لازم برای دریافت پاسخ از سوی سازمان	-
۳۹	۱,۲۶,۳۸	وجود امکان عضویت در سایت	-
۴۰	۱,۲۶,۳۸	وجود نظام مشخصی جهت عضویت کاربران در سایت	✓
۴۱	۱,۱۵	آگاهی کاربر از وجود مکانیزم کوکیز در سایت	-
۴۲	۲۷	پاسخگویی سریع و دقیق سیستم FAQ به پرسشهای کاربر	-
۴۳	۲۷	قابلیت پاسخگویی به سئوالات کاربران	-
۴۴	۲۷	آگاهی کاربران از زمان دریافت پاسخ	-
۱۲	۳۱	جمع امتیازها	
عینیت اطلاعات		نهاد کتابخانه‌های عمومی کشور	
ردیف	تواتر	مؤلفه‌ها	
۴۵	۱,۱۰,۱۵,۱۶,۲۳,۴	مشخص بودن ارتباط میان شخص مولف یا سازمان و فرد مسئول محتویات سایت	-
۴۶	۱,۱۵,۱۶,۴,۲۳	مشخص بودن دیدگاه مولف	-
۴۷	۱,۸,۱۵,۴,۳۸	مشخص بودن دیدگاه شخص یا سازمان مسئول تهیه اطلاعات	✓
۴۸	۱,۱۰,۱۵,۱۸	نبود در تبلیغات در صفحه	✓
۴۹	۱,۶,۱۷,۱۵,۲۷,۳۰	ارائه مطالبی در خصوص اهداف شخص یا سازمان گردآورنده اطلاعات	-
۵۰	۱,۱۵	توجه اطلاعات غیر مرتبط با خدمات سازمان در صفحه	-

✓	وجود وجه تمایزی میان محتوی اطلاعاتی و محتوی تقریباتی	۱۵,۳۹,۲۷	۵۱
✓	شبه نبودن اطلاعات به تبلیغات از لحاظ نوع طراحی	۲۷	۵۲
✓	اجتناب از بکاربردن مفاهیم و اطلاعات نا مرتبط	۴۰,۳۹	۵۳
۱۵	۳۴	جمع امتیازها	
نهاد کتابخانه‌های عمومی کشور		ناوبری	
	الف: ویژگیهای عنوان مرورگر	تواتر	ردیف
✓	مبین نام سازمان یا فرد مسئول محتویات سایت	۱,۱,۲۱	۵۴
✓	مبین صفحه اصلی بودن صفحه	۱,۱۵,۲۶	۵۵
✓	کوتاه بودن عنوان مرورگر	۱,۱۵	۵۶
✓	منحصر بودن عنوان مرورگر برای سایت	۱۵	۵۷
۹	۹	جمع امتیازها	
	ب: ویژگیهای عنوان صفحه	تواتر	ردیف
✓	مبین متعلق بودن صفحه اصلی بودن صفحه	۱,۱۰,۳۴,۱۵,۱۶,۲۱	۵۸
-	مبین متعلق بودن صفحه به یک سایت مشخص (حداقل توسط یک آرم)	۱,۱۰,۲۷,۱۵,۲۶	۵۹
✓	کوتاه بودن عنوان صفحه	۷۰	۶۰
✓	منحصر بودن عنوان صفحه برای سایت	۳۴,۴۰,۴۴,۱۵	۶۱
✓	بکاربردن کلمات کلیدی و مهم برای عناوین	۲۲,۲۷	۶۲
✓	بکارگیری عنوان مناسب جهت بیان دقیق محتوی متن	۲۷,۲۹,۲۲	۶۳
۱۷	۲۲	جمع امتیازها	
	ج: پیوندهای متنی و فرا متنی	تواتر	ردیف
✓	پیکنواختی ظاهر پیوندهای تکراری	۱,۶,۱۰,۴,۱۵,۳۸,۲۵,۲۶	۶۴
✓	امکان ناوبری آسان میان صفحات	۲۷,۳۰,۲۴,۱۶,۱۰,۴,۱۵,۲۱,۲۵	۶۵
✓	قرارگیری هماهنگ پیوندهای مستقیم داخلی در صفحه	۱,۶,۳۰,۱۰,۴,۱۵,۲۱,۲۶	۶۶
✓	وجود تناسب میان عنوان پیوند با آنچه پیوندبان ختم می شود	۱,۱۵,۱۷,۲۱,۲۵,۲۶,۳۵	۶۷
-	ادامه بکار پیوندها متناسب با اهداف از پیش تعریف شده	۱,۱۵,۱۶,۱۷,۱۸,۲۶	۶۸
✓	قابل درک بودن ساختار سایت / صفحه برای کاربران	۲۷,۴۰,۳۳,۱۶,۴,۱۵,۲۱,۲۶	۶۹
✓	تناسب چیدمان صفحات مختلف سایت با یکدیگر	۲۷,۱۰,۱۵,۲۱,۲۵	۷۰
-	وجود مابتهایی در صفحه اصلی جهت دسترسی هرچه آسانتر و موثر تر کاربران به صفحات پرمخاطب سایت	۱,۱۵,۲۱,۲۶,۳۸	۷۱
✓	وجود امکان دسترسی از صفحه اصلی سایت به بخش های اصلی آن	۱,۱۰,۳۹,۱۵,۲۱,۲۶	۷۲
✓	اجتناب از بکاربردن متن های زیر خط دار در کنار پیوندها	۱,۱۰,۱۵,۲۱	۷۳
-	عنوان دهی مناسب بوک مارک ها	۱۵,۲۱,۲۵	۷۴
-	امکان انتخاب اقلام اطلاعاتی مورد نظر از روی فهرست به جای تایپ آنها	۱,۱۵,۱۹	۷۵
✓	قرار دادن پیوند آرم سایت در صفحه اصلی	۱۵,۲۶	۷۶
-	شناساندن پیوندها از طریق زیر خط دار کردن آنها یا کاربرد نوعی رنگ خاص	۳۰	۷۷
✓	وجود تناسب بین عنوان صفحه و پیوند	۲۷	۷۸
-	نمایان بودن پیوندهای بازدید شده و بازدید نشده با استفاده از تغییر رنگ پیوند	۴۰	۷۹

ارزیابی کاربردپذیری وبگاه نهاد کتابخانه‌های عمومی کشور ۱۸۳

✓	قرار دادن پیوند صفحه اصلی جهت مشخص بودن صفحه اصلی	۲۷,۴۴	۸۰
-	منطقی بودن تعداد پیوندها	۴۴	۸۱
-	پیوند دادن تصویر به صفحه مربوطه	۳۰	۸۲
✓	امکان شناسایی آسان پیوندها	۲۴,۲۷	۸۳
✓	استفاده از عبارات مناسب برای پیوندها (عدم استفاده از عبارت , klik here more)	۲۴	۸۴
-	قرار دادن پیوند متنی در ابتدای پاراگراف	۲۲,۳۰	۸۵
✓	وجود تناسب بین عنوان صفحه و پیوند	۲۴	۸۶
-	استفاده از علائم دیداری مانند رنگ، اندازه برای نشان دادن ارتباط میان پیوندها	۲۲	۸۷
-	بکارگیری متن کافی برای توضیح پیوند	۴۰	۸۸
۶۸	۹۲	جمع امتیازها	
	۵: نشانه اینترنتی صفحه	تواتر	ردیف
-	درج نشانه اینترنتی صفحه در بدنه اصلی آن	۱,۱۵,۲۱	۸۹
✓	مختصر بودن نشانه اینترنتی صفحه	۲۷	۹۰
✓	عدم تغییر در نشانه اینترنتی سایت	۲۷	۹۱
✓	کاربر پسند بودن نشانه اینترنتی صفحه	۳۰	۹۲
۳	۶	جمع امتیازها	
	۵: نقشه سایت یا نمایه	تواتر	ردیف
-	وجود نقشه سایت یا نمایه ها بر روی صفحه اصلی و صفحات پیوندی	۲۷,۳۹,۱,۶,۱۵,۲۱,۲۶,۳۷	۹۳
-	دارا بودن موضوعات اصلی سایت	۱,۶,۱۵,۲۱,۲۶,۳۷	۹۴
-	امکان خواندن آسان نمایه ها یا نقشه سایت	۱,۶,۱۵,۲۱,۳۷	۹۵
-	سازماندهی منطقی نمایه ها یا نقشه سایت	۱,۲۷,۶,۱۵,۲۱,۳۷	۹۶
۰	۲۵	جمع امتیازها	
	ی: موتور جستجوی داخلی	تواتر	ردیف
✓	وجود بک موتور جستجوی داخلی	۱,۴,۱۵,۲۱,۲۵,۲۸	۹۷
-	مرتبط بودن اطلاعات بازایی شده توسط موتور جستجوی داخلی	۱,۴,۱۵,۲۱,۲۵,۲۸	۹۸
۶	۱۲	جمع امتیازها	
نهاد کتابخانه های عمومی		نماهای غیرمتنی	
	مؤلفه ها	تواتر	ردیف
✓	عدم استفاده از تصاویر متحرک (انیمیشن) بی مورد	۲۷,۳۹,۱,۶,۱۵,۲۱,۲۶,۳۵,۴۲,۳۷	۹۹
-	بکارگیری تصاویر گرافیکی، فایل های صوتی و تصویری به منظور افزایش کارایی سایت	۲۲,۴۰,۱,۱۰,۱۵,۲۱,۲۶,۳۷,۳۵	۱۰۰
-	بیان نام نرم افزار خاص مورد نیاز و چگونگی دسترسی به آن	۱,۱۰,۱۵,۲۱,۲۶,۳۷,۳۵	۱۰۱
-	وجود جایگزینی برای فایل‌های نیازمند نرم افزار خاص جهت دسترس پذیری اطلاعات برای کلیه کاربران	۱,۱۰,۱۵,۲۱,۲۶,۳۷,۳۵	۱۰۲
-	بیان نام نرم افزار مرورگر مورد نیاز، یا ویرایش خاصی از آن، جهت دسترسی به	۱,۱۰,۱۵,۲۶,۳۵,۳۷	۱۰۳

۱۸۴ مدیریت منابع اطلاعاتی وب

صفحات وب (در صورت لزوم)			
-	وجود جایگزین متنی برای تصاویر و گرافیک های موجود در صفحه برای کلیه کاربران	۱,۱۰,۱۵,۲۷,۲۱,۲۶,۳۷,۳۵	۱۰۴
✓	عدم استفاده از تکنیک فلش (علامت چشمک زن)	۱,۱۰,۳۰,۲۱,۳۵,۳۷,۴۲	۱۰۵
-	آگاهی کاربر از بار شدن فایلی حجیم در صورت پیروی از یک پیوند	۱,۱۵,۲۱,۲۶,۳۵,۳۷	۱۰۶
✓	مرتبط بودن تصاویر گرافیکی با محتوی متن	۲۷	۱۰۷
-	اضافه کردن متون به تصاویر برای درک بیشتر	۲۷	۱۰۸
-	استفاده از تگ ALT مناسب برای تصاویر	۳۰	۱۰۹
-	هدفمند بودن تصاویر گرافیکی	۳۲	۱۱۰
۱۸	۶۴	جمع امتیازها	
تعداد کتابخانه های عمومی کشور		دسترس پذیری	
	مؤلفه ها	تواتر	ردیف
✓	دسترس پذیر بودن با استفاده از مرورگر اینترنت اکسپلورر (Internet explorer 6.0)	۱,۶,۱۰,۱۵,۱۶,۲۵,۴۱,۳۷,۴۳,۳۵	۱۱۱
✓	دسترس پذیر بودن با استفاده از مرورگر نت اسکپ نیویگیتور (Netscape Navigator 6.5)	۱,۶,۱۰,۱۵,۱۶,۲۶,۴۱,۳۷,۴۳,۳۵	۱۱۲
-	دسترس پذیر بودن با استفاده از مرورگر موزیلا برد (Mozilla Bird 0.7)	۳۷,۴۳,۳۵, ۱,۶,۱۰,۱۵,۱۶,۲۶,۴۱	۱۱۳
-	دسترس پذیر بودن با استفاده از مرورگر اپرا (Opera 7.2)	۳۷,۴۳,۳۵, ۱,۶,۱۰,۱۵,۱۶,۲۶,۴۱	۱۱۴
✓	دسترس پذیری سایت از طریق موتورهای جستجوی عمومی	۱,۱۵,۱۶,۱۷,۴۳,۲۵,۲۶	۱۱۵
-	اندازه صفحه کمتر از ۵۰ کیلو بایت	۱,۴,۳۹,۲۶,۳۵,۳۸,۴۱	۱۱۶
✓	امکان پشتیبانی از سکوی عملیاتی کاربر (وبندوز) me,۹۸,۲۰۰۰, لینوکس (REDHAT)	۱۰,۶,۱۵,۲۶	۱۱۷
✓	استفاده از قلم (فونت) های استاندارد	۱۲, ۴, ۱۵, ۳۵, ۳۸, ۲۷, ۳۰, ۲۴, ۴۴, ۳۲, ۳۴	۱۱۸
✓	قابل رویت بودن تمامی اجزای سایت	۱۵, ۳۸	۱۱۹
-	دسترسی کاربر به منابع مورد نیاز در کمتر از ۳ بار کلیک کردن	۲۷	۱۲۰
-	شبه بودن اطلاعات به هرم واژگون (دستیابی سریع به اطلاعات مهمتر)	۲۷	۱۲۱
✓	سهولت استفاده از سایت برای همه کاربران (مبتدی و متخصص)	۲۷ و ۲۴	۱۲۲
✓	عدم استفاده از تصاویر متحرک (انیمیشن) بی مورد	۳۹ و ۳۰	۱۲۳
۴۸	۷۷	جمع امتیازها	
تعداد کتابخانه های عمومی		کارآمدی	
	مؤلفه ها	تواتر	ردیف
✓	استفاده از رنگ های استاندارد در طراحی صفحات	۴۲, ۳۷, ۳۵ و ۱۶ و ۲۱ و ۱۰ و ۳۴ و ۳۲ و ۳۹ و ۲۴ و ۲۷	۱۲۴
✓	عنوان بندی مناسب قابلیت های سایت / صفحه	۲۵, ۲۶, ۱, ۱۵, ۲۱, ۱۷	۱۲۵
✓	امکان تشخیص سایت بر اساس عنوان دامنه صفحه (Domin)	۱۵, ۲۸, ۲۵, ۱۹, ۳۵	۱۲۶
✓	تناسب زبان سایت با فرهنگ و روحیات کاربر	۱, ۱۵, ۳۹, ۲۶, ۲۸, ۳۸	۱۲۷
-	وجود تعاریف حاشیه ای برای تشریح اقلام اطلاعاتی موجود بر روی صفحه	۱, ۱۵, ۲۶, ۳۸	۱۲۸

ارزیابی کاربردپذیری وبگاه نهاد کتابخانه‌های عمومی کشور ۱۸۵

-	اطلاع کاربر از عملیات در حال انجام	۳۰,۲۲,۱,۱۵,۲۶,۳۸	۱۲۹
-	تناسب اطلاعات موجود در صفحات با مأموریت سایت	۱,۱۵,۲۶,۳۸	۱۳۰
✓	امکان انجام تمامی قابلیت های سایت در درون سایت	۱۵,۲۶,۳۷	۱۳۱
-	وجود قابلیت لغو عملیات صورت گرفته در سایت	۱۰,۱۵,۲۵	۱۳۲
-	وضوح موقعیت کاربر در سایت	۱,۱۵,۲۶	۱۳۳
✓	در نظر گرفتن پیشینه ذهنی کاربر	۱,۳۷,۳۹,۲۶	۱۳۴
-	جذابیت سایت	۱۷,۲۶	۱۳۵
-	در نظر گرفتن ابزارهایی جهت کمک به کاربر	۱,۱۵,۲۶,۳۸	۱۳۶
-	ارائه عنوان کامل تمام سرواژه های مهم بکاررفته در صفحات	۱۵,۳۴,۲۵,۳۰	۱۳۷
✓	طومارنوردی آسان صفحه خانگی	۶,۴۲	۱۳۸
✓	امکان چاپ اطلاعات بدون اعمال تغییرات در تنظیم سیستم رایانه	۲۵,۱۵	۱۳۹
✓	بیان تعداد مراجعان سایت در یک بازه زمانی	۱۵	۱۴۰
-	استفاده از عبارات کوتاه برای تشریح موارد موجود در صفحه	۲۷,۴۴,۲۶	۱۴۱
-	ایجاد جلب توجه کاربران از طریق نوع طراحی (نوع رنگ و فونت و ...)	۲۷,۳۴	۱۴۲
✓	بهم فشرده نبودن مطالب در یک صفحه و ایجاد فضای خالی بین آنها	۲۷,۳۹,۴۰	۱۴۳
-	دسترسی آسان به بخش Help	۳۹	۱۴۴
-	ارائه قیمت پیشنهادی به کاربر در صورت انتفاعی بودن سایت	۳	۱۴۵
-	بالا بودن سرعت بارگذاری	۳۰,۲۷	۱۴۶
✓	قابلیت استفاده از عملکرد جستجو	۲۷,۳۹,۳۴,۲۲	۱۴۷
-	امکان ناوبری توسط توروک	۲۹	۱۴۸
-	امکان ناوبری توسط جستجو	۲۹,۲۷	۱۴۹
-	دارا بودن خطی متنی اختصاصی	۲۷,۴۴	۱۵۰
-	قرار دادن مطالب در کمتر از ۳ صفحه	۲۴	۱۵۱
✓	رعایت سلسله مراتب مطالب در صفحات (بر حسب تاریخ و عنوان و ...)	۲۹	۱۵۲
✓	بکار بردن عدد برای نشان دادن ارقام (۲ به جای دو)	۴۴	۱۵۳
✓	عدم استفاده از فایل های pdf بدلیل کاهش جذابیت (به جز برای اسناد)	۴۴	۱۵۴
-	قابل رویت بودن پیام های اخطار	۳۹	۱۵۵
-	پایین بودن زمان داتلود	۳۴,۴۴	۱۵۶
۵۱	۹۹	جمع امتیازها	
نهاد کتابخانه‌های عمومی کشور		ویژگی‌های ظاهری	
	مؤلفه ها	تواتر	ردیف
✓	صفحه آرای مناسب	۲۲,۳۴,۴۰	۱۵۷
✓	وجود هماهنگی بین طرح ها و صفحات	۲۷	۱۵۸
✓	وجود هماهنگی بین رنگ ها و سبک نگارش	۳	۱۵۹
✓	بکارگیری حاشیه در اطراف متن	۲۴	۱۶۰
۶	۶	جمع امتیازها	

در مجموع نتایج به دست آمده از ارزیابی کاربرد پذیری وب سایت نهاد کتابخانه های عمومی کشور به تفکیک معیارهای 11 گانه ی موجود در این پژوهش در مقایسه با وب سایت شاخص فرضی با 663 امتیاز، در سطحی معادل 52 درصد (346 امتیاز)، قرار دارد و این امر گر چه بیان گر آن است که تنها نیمی از استانداردها در وب سایت نهاد کتابخانه های عمومی کشور رعایت گردیده است و این امر با وب سایت شاخص فاصله زیادی دارد لذا لازم است در طراحی این وب سایت دقت و توجه بیش تری به عمل آید افزون بر این با عنایت به کاربردی بودن این پژوهش پژوهش گران در صدند تا با انعکاس نتایج پژوهش حاضر، کتابداران و طراحان وب سایت ها را از نقاط ضعف و قوت وب سایت مطلع ساخته تا با به کارگیری عناصر لازم و ویژگی های مناسب در طراحی وب سایت فضایی را ایجاد نمایند که کاربر به محض ورود به وب سایت امکان بهره گیری از تمامی اطلاعات و خدمات مرتبط با این، محمل بدون وجود هر گونه محدودیتی را داشته باشد.

6- پیشنهاد های پژوهش

با توجه به یافته های پژوهش در باب رعایت معیارهای مطرح در کاربرد پذیری هر چه بیشتر و بهتر وب سایت نهاد کتابخانه های عمومی کشور توجه به رعایت معیارها و مولفه های: بیان نام مولف و ویژگی های او (صلاحیت شهرت و اعتبار و...)، ذکر فواصل زمانی به روز کردن اطلاعات دارای حساسیت زمانی درج زمان تخمینی تکمیل صفحه در دست ساخت نماهای تعاملی و تبدلی، وجود امکان عضویت در سایت، مشخص بودن ارتباط میان شخص مولف با سازمان و فرد مسئول محتویات سایت منحصر بودن عنوان صفحه برای سایت نمایان بودن پیوندهای بازدید شده و بازدید نشده با استفاده از تغییر رنگ پیوند درج نشانه اینترنتی صفحه در بدنه اصلی آن بیان نام نرم افزار خاص مورد نیاز و چگونگی دسترسی به آن بیان تعداد مراجعان سایت در یک بازه زمانی توصیه می گردد

منابع

1. اصغری، پوده احمدرضا 1380. بررسی عناصر و ویژگی های مطرح در طراحی وب سایت کتابخانه های دانشگاهی. پایان نامه کارشناسی ارشد کتابداری و اطلاع رسانی دانشکده علوم تربیتی و روانشناسی، دانشگاه فردوسی مشهد
2. الکساندر ژانت 1383. شناخت وب: چگونه اطلاعات موجود بر روی وب را ارزیابی نموده و صفحاتی اینگونه را پدید آوریم، ترجمه صدیقه محمد اسماعیل، تهران: دبیزش.
3. جانسون، استیو 1382. صفحات سفارش و دریافت مواد منابع کتابخانه ای بر روی وب. ترجمه صدیقه محمد اسماعیل اطلاع شناسی، 1 (پاییز): 187-199.
4. حاجی زین العابدینی، محسن مکتبی، فرد، لیلا عصاره، فریده 1384 تحلیل پیوندهای وب سایت های کتابخانه های ملی جهان مجله مطالعات تربیتی و روانشناسی دوره 7، شماره 1، ص 193-173

5. خوانساری، جیران ساختن سایت های وب موفق برای کتابخانه های کوچک دانشگاهی. صنعت برق 54 (آبان) ص. 24-28.
6. دایره المعارف کتابداری و اطلاع رسانی ج. 1 و 2. 1381 تهران: نهاد کتابخانه های عمومی کشور، 1381.
7. رضایی شریف آبادی، سعید، فرودی، نوشین 1381. ارزیابی صفحات وب کتابخانه های دانشگاهی ایران و ارائه الگوی پیشنهادی فصلنامه کتاب دوره سیزدهم 4 (زمستان) ص. 12-19.
8. عصاره، فریده، مرادمند، علی 1384. شناسایی ویژگی های عمده در طراحی وب سایت های کتابخانه های ملی جهان به منظور ارائه الگویی مناسب جهت ارتقاء کیفی وب سایت نهاد کتابخانه های عمومی کشور. فصلنامه اطلاع شناسی پاییز و زمستان شماره 1 و 2. ص. 170-190
9. عصاره فریده 1381 معیارهای ارزیابی منابع اینترنتی فصلنامه کتاب، دوره سیزدهم 2 (تابستان). ص. 61-73
10. محمد اسماعیل صدیقه 1383. بررسی کاربرد پذیری وب سایتهای دانشگاه های صنعتی کشور. پایان نامه دکتری کتابداری و اطلاع رسانی تهران دانشگاه آزاد اسلامی واحد علوم و تحقیقات
11. محمد اسماعیل صدیقه 1384 بررسی کاربرد پذیری وب سایت های دانشگاه های صنعتی کشور فصل نامه کتاب دوره 16 شماره 1.
12. نویدی، فاطمه 1386. ارزیابی دسترس پذیری وب سایت وزارتخانه های دولت جمهوری اسلامی. ایران پایان نامه کارشناسی ارشد علوم کتابداری و اطلاع رسانی دانشگاه تربیت مدرس
13. Alexander, J. E. ,Tate, M. A.1999. Web Wisdom: who to evaluate and create information quality on the . (web. London, Mahwah, NewJersey.LEA: (Lawrence Erlbaum Associates
14. Back, S. E.1997. Evaluation criteria: the good, The Bad and The Ugly: Or, Why It's a Good Idea to . "Evaluate Web Sources". [on-line] Available: <http://lib.nmsu.edu/instruction/evalcrit.html>
15. Barker. J. 2004 Finding Information on the Internet: A Tutorial University of California. [on-line]. Available: <http://www.lib.berkeley.edu/TeachingLib/Guides/Internet/Evaluate.html>
16. Barker. J. 2003. Evaluating Web Pages: Techniques to Apply Questions to Ask." VC Berkeley - Teaching Library Internet Workshops, 12sep.2003.[on-line]. Available: <http://www.Lib.berkeley.edu/Teaching Lib/Guides/Internet/Evaluated.html>
18. Berger. P. 1999 Web Evaluation Guide: Tramline, Incorporated. [on-line]. Available: .

<http://www.infosearcher.com/cybertours/tours/touro4/-tourlaunch1.html>

Bertot J.C. et al. 2006 Functionality, usability and accessibility: Interactive user-centered evaluation . 19 strategies for digital libraries”, Performance Management and Metrics, Vol. 7

ص: 187

- Braynik. G. 2003. Atomic Web Usability Evaluation: What Need to be done?. (29jan.2003): 1–16. .20
[on–line]. Available: <http://Usable.binghamton.edu.Atomic Thesis. html>
- Clyde. L. A. 1996. The Library as information providers, The Homepage The Electronic Library. Vol. .21
.14 no.6 pp:549–558
- Engle. M. 1996. Evaluating Websites: Criteria and Tools. New York Library Association Conference, .22
Saratoga Springs, Ny. (October 1996): 1–3.[on–line]. Available: <http://www.Library.Cornel.edu/Okuref/research/Webeval.html>
- .The Essential Web Site Usability Checklist. 2007. [on–line]. Available: www.dailybits.com .23
- Hupp.J.2008. Test Your Web Site: A 57– Point Checklist for Maximum Usability .[on– online]. . 24
.Available: www.virtualhosting.com
- .Jafari. A. Optimizing Campus Web Sites. EDUCASE QUARTERLY, No.2 2000:56–58 .25
- .Leggett. D. Quick Usability Checklist.2009, [on–line]. Available: www.uxbooth.com .26
- Meyers. P. J. 25–Point Web Site Usability Checklist, 2008. [on–line]. Available: www.usereffect.com .27
- Nadler.D.M. and Furman.V.M. (2001). Access board issues final standards for disabled access under .28
.Section 508 of Rehabilitation Act, Government Contract Litigation Reporter, Vol. 14 No. 19, p. 14
- National Science foundation (NSF).Universal Design of College Algebra, 2008. [on– line]. Available: .29
www.usablealgebra.landmark.edu
- Nielsen.J.Coyne, K, Tahir, Marie, (2001). Make It Usable–Web Site Usability Magazine, (6 feb.2001): .30
.1–5.[on–line]. Available: <http://www.Pcmag.com/article2/0,4149,33821,00. asp>, 23jan.2003
- Osareh.F. (2003). "A notice on content of library information science (LIS) schools websites, Libri (3): .31
262–265
- .The SEO File: Usability Checklist 2009. [on–line]. Available: www.theseofiles.co.uk .32
- Sloan, John. Rampo College Web Design Standards. Rampo College of New jersey Statements and .33
policies.2001.[on–line].Available: <http://www.guide.rampo.edu/1/content/Webstandards.html>

,Stueart. R.D. Library and Information Center management fifth edition. Endewood .34

ص: 188

- Sullivan, T. Matson R. Barriers to use: Usability and Content Accessibility on the Web's Most Popular .35
.Sites. Interdisciplinary PhD. program in information Science, University of North Texas, Denton, 2000
- Tungar.M. N. Heuristic Evaluation.2002.[on-line]. Available: <http://www.cc.gatech.edu/~manas/cs8803a/Heuristic.pdf> . 36
- .Usability Guidelines Information Science Technology, 2004. [on-line]. Available: www.web.mit.edu .37
- .Web Usability Checklist. The College of New Jersey. 2008. [on-line]. Available: www.teng.edu" .38
- Weibel .S.L. The World Wide Web and Emerging Internet Resource Discovery Standards for Scholarly .39
.Litrature. Library Trends.vol.43, No.4 1995:627-664
- .Wilson, S. 1995 World Wide Web Design Guide. Indianapolis. IN: Hayden Books, 1995 .40
- Zaphiris. P. Darin, Ellis, R. Website Usability and Content Accessibility of The Top USA Universities .41
2001. Dertoit, MI: Institute of Geronotology and Dept of Industrial and Manufacturing Engineering Wayne
.State University, 2001
- Web Usability Tips to Attract and Retain Web Visitors 2009,:[on-line].Available: 50 . 42
www.doshdosh.com

هدف: هدف کلی پژوهش دستیابی به فرایند پردازش و سازماندهی وبگاه‌ها با استفاده از قواعد و استانداردهای مورد استفاده در سازمان اسناد و کتابخانه ملی ایران است.

جامعه آماری: جامعه آماری این پژوهش را 50 وبگاه تشکیل می‌دهند که در 20 رده موضوعی از وبگاه پارس ایندکس و به صورت تصادفی انتخاب شده‌اند.

روش شناسی: این پژوهش به دو روش کتابخانه‌ای و پیمایش توصیفی انجام شده است - جهت گردآوری اطلاعات از سیاهه واریسی و برای تجزیه و تحلیل داده‌ها از آمار توصیفی استفاده شده است.

نتایج: نتایج نشان می‌دهند که استاندارد سازی موضوع سازماندهی و پردازش وبگاه‌ها به عنوان نوعی از منابع الکترونیکی و بومی سازی آن در سازمان اسناد و کتابخانه ملی ایران امکان‌پذیر است.

کلیدواژه‌ها: وبگاه‌ها سازماندهی امکان‌سنجی سازمان اسناد و کتابخانه ملی جمهوری اسلامی ایران

اشاره

دکتر رضا خانی پور (1) | محبوبه قربانی (2) | سهیلا فعال (3)

مقدمه

امروزه شبکه جهانی اینترنت منابع اطلاعاتی متنوعی را در دسترس قرار داده است. وب هر روزه پهنه بیش تری از وسعت دنیا را در بر گرفته و در دور افتاده ترین نقاط جهان نیز رسوخ نموده است. شبکه جهانی وب با دسترسی پذیری خود امکانات زیادی را برای مخاطبان فراهم نموده است. دسترسی به آخرین اخبار دنیا و اطلاعات کشفیات و نوآوری ها آشنایی با فرهنگ و تمدن ملل مختلف، دریافت اطلاعات مورد نیاز روزانه مانند آب و هوا و اخبار وقایع مهم از امتیازات دسترسی به وب است. این گستردگی و دسترسی پذیری در کنار محاسنی که دارد کاربران را با مشکلاتی روبه رو ساخته است. اطلاعات در لابلای صفحات وب نهفته و نیاز به جستجوهای دقیقی برای کشف آن ها وجود دارد.

عمده ترین مشکل کاربران، وب دست یابی به بهترین و با کیفیت ترین اطلاعات مورد نیاز و حصول جامعیت و مانعیت در بازیابی اطلاعات در زمینه های تخصصی است (علی محمدی 1381)

ص: 191

1- دکترای کتابداری و اطلاع رسانی عضو هیات علمی و مدیر کل پردازش و سازماندهی سازمان اسناد و کتابخانه ملی ایران -R
Khanipour@nlai.ir

2- دانشجوی دکترای کتابداری و اطلاع رسانی معاون اداره کل پردازش و سازماندهی سازمان اسناد و کتابخانه ملی ایران -m
ghorbani@nlai.ir

3- کارشناس ارشد کتابداری و اطلاع رسانی رئیس گروه سازماندهی منابع غیر کتابی سازمان اسناد و کتابخانه ملی ایران -s
faal@nlai.ir

کتابخانه ها و مراکز اطلاع رسانی به عنوان متولیان گردآوری سازماندهی و اشاعه اطلاعات، خود را موظف به گردآوری سازماندهی اطلاعات موجود بر روی وب می دانند یکی از چالش های فراروی کتابداران نحوه سازماندهی منابع وب بوده است تلاش های گسترده ای در خصوص سازماندهی اطلاعات بر روی وب صورت گرفته است تهیه انواع استانداردها و دست نامه ها و ابداع فراداده های متنوع سازماندهی منابع وب از جمله تلاش های صورت گرفته در جهت سازماندهی منابع وب است. (نشاط 1382)

وبگاه ها: شناخته شده ترین منابع وب

منابع وب آثاری هستند که توسط افراد مختلف در قالب های متنوع تولید و انتشار یافته اند. این منابع وبگاه ها، وبلاگ ها کتاب های الکترونیکی مقالات تمام نشریات الکترونیکی و... را در بر می گیرد. یکی از مهم ترین منابع تحت وب، وبگاه ها (1) هستند وبگاه مجموعه ای از صفحات وب است که دارای یک دامنه اینترنتی و به صورت مجموعه ای از صفحات مرتبط که داده هایی نظیر، متن صدا، تصویر و فیلم روی آن ها ارائه می شود، روی شبکه اینترنت قرار می گیرد صفحه وب به صورت سندی است که در قالب اچ تی ام ال (2) نوشته می شود و همواره با استفاده از پروتکل اچ تی تی پی (3) می توان به آن دسترسی پیدا کرد. مهم ترین قسمت یک وبگاه در واقع صفحه اصلی یا صفحه خانگی آن است وبگاه یک مؤسسه مرکز تجاری یا سازمان چهره او به سوی جهان و نقطه شروع بیش تر کاربران تلقی می شود نیلسن (2002) (4)

ارزیابی وبگاه ها

عواملی که مجموع قابلیت به رهوری از یک وبگاه را تعیین می کنند عبارتند از: محتوا، زبان، ساختار، طراحی جهت یابی عبور و قابلیت دسترسی روش های گوناگونی جهت ارزیابی قابلیت های بهره برداری از یک وبگاه وجود دارد که عبارتند از ارزیابی با مشارکت کاربر و ارزیابی بدون مشارکت کاربر روش ارزیابی در صورتی مفید خواهد بود که هم زبان و هم ساختار وبگاه مورد ارزیابی به سهولت توسط کاربران درک شود. (پل 2007) (5)

استانداردها و طرحهای ابر داده ای سازماندهی منابع وب

استانداردها و طرح های ابر دادهای مختلفی برای توصیف منابع اینترنتی تهیه شده اند که از انواع آن ها می توان به موارد ذیل اشاره نمود:

1. طرح ابر داده ای دابلین کور (6)

ص: 192

1- WebSites

2- HTML

3- HTTP

4- 4J, Nielsen

5- Roswitha Poll

6- Dublin Core (DC)

2. مارک 21 (1)

3. آر دی اف (2)

4. طرح کدگذاری توصیف آرشیوی (3)

5. قالب ابر داده‌های خدمات مکان یاب اطلاعات دولتی (4)

6. قالب ابر داده ای طرح کدگذاری متن (5)

پردازش منابع الکترونیکی در سازمان اسناد و کتابخانه ملی ایران

در سازمان اسناد و کتابخانه ملی ایران برای پردازش و سازماندهی انواع منابع کتابی و غیر کتابی از ویرایش دوم قواعد فهرست نویسی انگلو امریکن (6)، استاندارد بین المللی توصیف کتاب شناختی (7) و استاندارد یونی مارک (8) استفاده می شود. کاربرد مارک (9) در سازماندهی منابع به طور جدی از سال 1385 آغاز شد. در این راستا از مارک ایران که شکل بومی سازی شده استاندارد یونی مارک است استفاده شده است. بر اساس این استاندارد منابع الکترونیکی در گروهی جداگانه و با کد 1 مشخص می شوند. در کتابخانه ملی ایران انواع کار برگه ها با توجه به این تقسیم بندی طراحی شده اند و هر کار برگه از 10 بلوک و فیلدهای اصلی و فرعی مربوط به توصیف هر منبع تشکیل شده است.

بر اساس طرح مقدماتی برای فهرست نویسی منابع الکترونیکی (2003) (10)، منابع الکترونیکی از لحاظ نوع دسترسی به دو دسته منابع قابل «دسترسی مستقیم» (11) و منابع قابل «دسترسی از راه دور» (12) تقسیم می شوند (عبداللهی 1383). بنابراین وب گاه ها دسته ای از منابع الکترونیکی و از نوع «دسترسی از راه دور» به حساب می آیند منابع الکترونیکی در واحد سازماندهی منابع غیر کتابی اداره کل پردازش، سازماندهی می شوند. تاکنون بیش از 3600 منبع الکترونیکی از نوع «دسترسی مستقیم» در این واحد سازماندهی شده اند.

با توجه به مطالب پیش گفته به نظر می رسد که استانداردها و قواعد توصیف منابع الکترونیکی برای توصیف وبگاه ها نیز کاربردی باشد.

بیان مساله

وبگاه ها به عنوان نوعی از منابع اطلاعاتی وب در چرخه اطلاعات و دانش محسوب می شوند. با توجه به این که پردازش و سازماندهی منابع اطلاعاتی هسته مرکزی و فنی این چرخه به حساب می آید، لازم

ص: 193

Machin Readabel cataloging (MARC) -1

Encoded Archival description(EAD) -2

Resource Description Framework(RDF) -3

Government information locator sources (GILS)	-4
Text Encoding Initiative(TEI)	-5
Anglo American Cataloging Rules, Second Edition (AACR2)	-6
International Standard Bibliographic Description (ISBD)	-7
Universal MARC	-8
Machine-Readable Cataloging	-9
Draft internet guidelines for cataloging electronic resources	-10
Direct access	-11
Remote access	-12

می نماید برای استفاده هر چه کامل تر از محتوای اطلاعاتی وبگاه ها سازمان اسناد و کتابخانه ملی ایران که مسئولیت استاندارد سازی و نظارت بر سازماندهی را بر عهده دارد از پردازش و سازماندهی وبگاه ها نیز غافل نشود با عنایت به اینکه کتابخانه ملی ایران در عرصه پردازش انواع منابع کتابخانه ای اعم از چاپی و الکترونیکی وارد شده است لازم است برای روزآمد سازی کنترل مدیریت و سازماندهی وبگاه ها نیز طرح و برنامه ای داشته باشد. این پژوهش بر آن است تا امکان پردازش وبگاه ها را به عنوان نوعی از منابع الکترونیکی مهم در عرصه های ملی و جهانی توسط سازمان اسناد و کتابخانه ملی ایران مورد پژوهش قرار دهد.

پیشینه پژوهش

پیشینه پژوهش در ایران:

حاجی زین العابدینی (1381) در پژوهشی به بررسی مشکلات اینترنت در زمینه سازماندهی و بازیابی اطلاعات پرداخته و دست نامه فهرست نویسی منابع اینترنتی را ارائه کرده است در تهیه دست نامه از قواعد فهرست نویسی انگلو امریکن قواعد به کار گرفته شده در طرح های فهرست نویسی منابع اینترنتی و الکترونیکی، اصول و قواعد فهرست نویسی منابع فارسی و نمونه های منابع اینترنتی فهرست نویسی شده، استفاده شده است. با استفاده از این دستنامه امکان فهرست نویسی منابع اینترنتی و سازماندهی اطلاعات دل خواه در اینترنت فراهم می شود هم چنین کاربرد های ایجاد شده است که بر اساس فیلدهای اطلاعاتی موجود در منابع، اینترنتی دستورالعمل های دست نامه نمونه های موجود در طرح هایی چون اینترکت، رهنمود های ایجاد مارک و فیلد های اطلاعاتی در راهنمای نرم افزار مارکایت طراحی شده است.

فتاحی حسن زاده (1385) به منظور مطالعه و ارزیابی شیوه های سازماندهی اطلاعات در وبگاه های کتابخانه های دانشگاهی پژوهشی انجام داده اند نتایج کلی پژوهش بیان گر آن است که در سازماندهی، اطلاعات صفحه اول وبگاه ها، شیوه دسته بندی بر اساس نوع خدمات بیش از سایر انواع رایج در حالی که برای سازماندهی سایر صفحات از شیوه های نسبتاً متنوعی استفاده می شود همچنین شیوه الفبایی عنوان و موضوع بیش ترین کاربرد و شیوه دایرکتوری (موضوعی سلسله مراتبی) کم ترین کاربرد را برای سازماندهی انواع منابع اطلاعاتی در وب سایت کتابخانه های دانشگاهی دارند.

پیشینه پژوهش در خارج از ایران:

کوچ (1) و همکارانش (1997) نقش طرح های رده بندی را در توصیف و بازیابی منابع اینترنتی مورد بررسی قرار داده اند. آن ها استفاده از این طرح ها را در سازماندهی محتویات وب گاه ها توصیه کرده اند، ولی به کاربرد شیوه های مختلف سازماندهی و نظریات کتابداران و کاربران در مورد چگونگی آن ها نپرداخته اند.

ویلیامسون (1997) (2) با تأکید بر ساختار دانش موجود در منابع اینترنتی در راستای سازماندهی دانش و

بازیابی اطلاعات، لزوم سازماندهی آن‌ها را بیان نموده است. در تحقیق وی، بررسی امکان کاربرد شیوه‌های متعارف سازماندهی برای سازماندهی محتوای اطلاعاتی وبگاه‌ها توصیه شده است.

وارد (2001) (1) در مقاله‌ای به تشریح اهمیت فهرست نویسی منابع اینترنتی پرداخته است. در ادامه فهرستی از فعالیت‌های انجام شده در کتابخانه‌های ایالات متحده در خصوص سازماندهی منابع اینترنتی را ارائه کرده است.

وایلر (2) و همکارانش (2008) در پژوهشی به ارزیابی هزینه‌های پردازش و سازماندهی منابع وب در کتابخانه و دانشگاه ملی کرواسی پرداخته‌اند نتایج این پژوهش نشان داده است که زمان پردازش منابع وبی با پردازش منابع چاپی یکسان است.

یانگه‌ی (3) (2011) در پژوهشی به بررسی فاصله بین ایجاد کنندگان منابع وب و ایجاد کنندگان فراداده‌های آن‌ها به منظور بهبود عناصر فراداده پرداخته است. یافته‌های پژوهش نشان‌گر این بوده است که رضایت کاربر تا حد بالایی به مفید بودن سهولت استفاده و دریافت اطلاعات و اثر بخشی عناصر اطلاعاتی بستگی دارد.

اهمیت پژوهش

سازماندهی و پردازش وبگاه‌ها در سازمان اسناد و کتابخانه ملی ایران می‌تواند از جنبه‌های زیر حائز اهمیت باشد:

1. ایجاد بانک اطلاعاتی وبگاه‌ها و غنی‌سازی کتاب‌شناسی ملی ایران؛
2. امکان جستجوی وبگاه‌ها با استفاده از قابلیت‌های جستجوی نرم‌افزار کتابخانه‌ای (رسا)؛
3. ایجاد نقاط دسترسی توصیفی و تحلیلی برای برآوردن نیازهای اطلاعاتی کاربران نهایی؛
4. ایجاد جامعیت و مانعیت در بازیابی پیشینه‌های کتاب‌شناختی وبگاه‌ها؛
5. استاندارد سازی و یکپارچه سازی پردازش توصیفی و تحلیلی وبگاه‌ها.

اهداف پژوهش

هدف کلی پژوهش دستیابی به فرایند پردازش و سازماندهی وبگاه‌ها با استفاده از قواعد و استانداردهای مورد استفاده در سازمان اسناد و کتابخانه ملی ایران است به این منظور اهداف فرعی زیر دنبال می‌شوند:

1. بررسی چگونگی دستیابی و دریافت اطلاعات وبگاه‌ها؛
2. توصیف کتاب‌شناختی وبگاه‌ها بر اساس استانداردها و قواعد مورد استفاده در سازمان اسناد و کتابخانه ملی ایران؛
3. دسترسی جامع و مانع به وبگاه‌ها با ایجاد نقاط بازیابی و تحلیل موضوعی مورد استفاده در سازمان اسناد و کتابخانه ملی ایران؛

Ward -1
Willer -2
Younghee -3

4. کاربرد نظام های موضوعی مورد استفاده در سازمان اسناد و کتابخانه ملی ایران برای نمایه سازی وبگاه ها

پرسش های اساسی

1. دستیابی و دریافت اطلاعات وبگاه ها از چه راه هایی امکان پذیر است؟

2. توصیف کتاب شناختی وبگاه ها بر اساس استانداردها و قواعد مورد استفاده در سازمان اسناد و کتابخانه ملی ایران چگونه است؟

3. دسترسی جامع و مانع به وبگاه ها با ایجاد نقاط بازیابی و تحلیل موضوعی مورد استفاده در سازمان اسناد و کتابخانه ملی ایران چگونه است؟

4. کاربرد نظام های موضوعی مورد استفاده در سازمان اسناد و کتابخانه ملی ایران برای نمایه سازی وبگاه ها چگونه است؟

تعاریف عملیاتی

امکان سنجی: منظور از امکان سنجی در این پژوهش بررسی طرح اولیه جهت ارائه و پیاده سازی فرایند پردازش و سازماندهی وبگاه ها در اداره کل پردازش و سازماندهی سازمان اسناد و کتابخانه ملی ایران است.

پردازش: منظور از پردازش در این پژوهش فرایند سازماندهی، منبع شامل دریافت، توصیف و تحلیل است.

جامعه آماری

جامعه آماری این پژوهش را 50 وبگاه تشکیل می دهند این جامعه مورد مطالعه در 20 رده موضوعی از وبگاه پارس ایندکس (1) و به صورت تصادفی انتخاب شده اند.

روش پژوهش و ابزار گردآوری داده ها

در این پژوهش دو روش کتابخانه ای و پیمایش توصیفی به کار گرفته شده است. جهت گردآوری اطلاعات از سیاهه واری و برای تجزیه و تحلیل داده ها از آمار توصیفی استفاده شده است. روایی سیاهه واری بر اساس تجربه پدید آوران مقاله در مدت بیش از یک دهه فعالیت تخصصی در حوزه سازماندهی منابع الکترونیکی و همچنین از طریق مشاوره با استادان به دست آمده است

یافته های پژوهش و پاسخ به پرسش های اساسی

پرسش اول: دست یابی و دریافت اطلاعات وبگاه ها از چه راه هایی امکان پذیر است؟

ص: 196

در این پژوهش جهت دستیابی و دریافت اطلاعات وبگاه‌ها دوروش مورد بررسی قرار می‌گیرد:

1. استفاده از فهرست وب گاه‌ها

برخی وب گاه‌ها و پایگاه‌های اطلاعاتی فهرست وب گاه‌های مختلف را بر اساس تقسیم‌بندی‌های موضوعی ارائه می‌دهند به این ترتیب یکی از راه‌های دست‌یابی به اطلاعات وب گاه‌ها، استفاده از فهرست وب گاه‌ها است. نمونه فهرست‌های وب گاه‌ها را می‌توان در وب گاه ایران (1) پارس ایندکس و الکسا (2) و غیره جستجو کرد.

در این پژوهش برای دسترسی به حجم نمونه از جامعه مورد مطالعه وبگاه پارس ایندکس انتخاب شد و امکان دست‌رسی به بسیاری از وب گاه‌ها امکان‌پذیر گردید.

2. ورود اطلاعات وبگاه‌ها در وبگاه سازمان اسناد و کتابخانه ملی ایران

منابع اینترنتی گروهی از منابع اطلاعاتی به حساب می‌آیند و کتابخانه ملی به لحاظ رسالت خود، لازم است مانند کتاب‌ها و سایر منابع اطلاعاتی اقداماتی در خصوص فراهم‌آوری منابع اینترنتی نیز به انجام رساند. در این راستا وب گاه سازمان اسناد و کتابخانه ملی ایران می‌تواند درگاهی جهت دست‌یابی به اطلاعات وب گاه‌ها باشد. بنابراین با طراحی «فرم الکترونیکی ورود اطلاعات وب گاه‌ها» و قرار دادن آن در وبگاه سازمان می‌توان به نشانی وب گاه‌ها و اطلاعاتی که برای پردازش آن‌ها مورد نیاز است، دست پیدا کرد.

جدول 1 فیلدهای پیشنهادی برای طراحی فرم الکترونیکی ورود وب گاه‌ها را نشان می‌دهد.

عکس

در این پژوهش جهت دستیابی و دریافت اطلاعات وبگاهها دو روش مورد بررسی قرار می‌گیرد:

۱. استفاده از فهرست وبگاهها

برخی وبگاهها و پایگاههای اطلاعاتی، فهرست وبگاههای مختلف را بر اساس تقسیم بندیهای موضوعی ارائه می‌دهند، به این ترتیب یکی از راههای دستیابی به اطلاعات وبگاهها، استفاده از فهرست وبگاهها است. نمونه فهرستهای وبگاهها را می‌توان در وبگاه ایران^۱، پارس ایندکس و الکسا^۲ و غیره جستجو کرد.

در این پژوهش برای دسترسی به حجم نمونه از جامعه مورد مطالعه، وبگاه پارس ایندکس انتخاب شد و امکان دسترسی به بسیاری از وبگاهها امکان پذیر گردید.

۲. ورود اطلاعات وبگاهها در وبگاه سازمان اسناد و کتابخانه ملی ایران

منابع اینترنتی گروهی از منابع اطلاعاتی به حساب می‌آیند و کتابخانه ملی به لحاظ رسالت خود، لازم است مانند کتابها و سایر منابع اطلاعاتی اقداماتی در خصوص فراهم‌آوری منابع اینترنتی نیز به انجام رساند. در این راستا وبگاه سازمان اسناد و کتابخانه ملی ایران می‌تواند در گامی جهت دستیابی به اطلاعات وبگاهها باشد. بنابراین با طراحی «فرم الکترونیکی ورود اطلاعات وبگاهها» و قرار دادن آن در وبگاه سازمان میتوان به نشانی وبگاهها و اطلاعاتی که برای پردازش آنها مورد نیاز است، دست پیدا کرد.

جدول ۱ فیلدهای پیشنهادی برای طراحی فرم الکترونیکی ورود وبگاهها را نشان می‌دهد.

جدول ۱. فیلدهای پیشنهادی برای طراحی فرم الکترونیکی ورود اطلاعات وبگاهها

نام فیلد	عنوان فارسی	عنوان به زبان دیگر	صاحب امتیاز	طراح رایانه ای	سال تولید	زمان آخرین ویرایش	زبان	حوزه موضوعی	کلید واژهها	نشانی اینترنتی	توضیحات
فیلد نوع	ضروری	عادی	ضروری	عادی	ضروری	ضروری	ضروری	ضروری	ضروری	ضروری	عادی

پرسش دوم: توصیف کتاب‌شناختی وبگاهها بر اساس استانداردها و قواعد مورد استفاده در سازمان اسناد و کتابخانه ملی ایران چگونه است؟

جدول ۲ عناصر پیشنهادی توصیف وبگاهها را با استفاده از قواعد فهرست‌نویسی انگلومریکن در بخش منابع الکترونیکی و در بستر فرادادهای یونی مارک نشان می‌دهد.

1. <http://www.iran.ir/directory>

2. <http://www.alexa.com>

جدول 1 فیلدهای پیشنهادی برای طراحی فرم الکترونیکی ورود اطلاعات وبگاهها

پرسش دوم: توصیف کتاب‌شناختی وبگاهها بر اساس استانداردها و قواعد مورد استفاده در سازمان اسناد و کتابخانه ملی ایران چگونه است؟

جدول 2 عناصر پیشنهادی توصیف وبگاهها را با استفاده از قواعد فهرست‌نویسی انگلومریکن در بخش منابع الکترونیکی و در بستر فرادادهای یونی مارک نشان می‌دهد.

<http://www.iran.ir/directory> -1

<http://www.alex.com> -2

۱۹۸ مدیریت منابع اطلاعاتی وب

جدول ۲. فراوانی کاربرد قواعد انگلومریکن و استاندارد یونی مارک در توصیف کتابشناختی وبگاهها

درصد فراوانی	یونی مارک (بلوک)	یونی مارک (فیلد)	انگلوامریکن (قواعد)	عناصر توصیف	ناحیه
%۷۶	اطلاعات توصیفی (۲)	Sar۰۰	B۹/۱	عنوان کامل	عنوان و پدیدآور
		Sbr۰۰	C۹/۱	وجه تسمیه عام: [منابع الکترونیکی]	
		Sdr۰۰	D۹/۱	عنوان به زبان دیگر	
		Ser۰۰	E۹/۱	دیگر اطلاعات عنوان	
		Sfr۰۰	F۹/۱	شرح پدیدآور	
%۶۸		Sar۰۵	B۹/۲	وضعیت ویراست	
بر اساس آخرین ویرایش قواعد انگلومریکن ناحیه ۳ (نوع انتشار) برای منابع الکترونیکی ذکر نمی‌شود					
%۷۰/۶	اطلاعات توصیفی (۲)	Sar۱۰	C۹/۴	محل تولید	چاپ و پدیدآور
		Ser۱۰	D۹/۴	ناشر و صاحب امتیاز	
		Sdr۱۰	F۹/۴	تاریخ تولید	
بر اساس قواعد انگلومریکن، ناحیه مشخصات ظاهری برای منابع الکترونیکی دسترسی از راه دور ذکر نمی‌شود					
%۵۱/۴	یادداشت‌ها (۳)	۳۰۰	۱B۹/۷	ماهیت و دامنه	یادداشت
		۳۰۰	۲B۹/۷	زبان	
		۳۰۴	۳B۹/۷	منبع عنوان کامل	
		۳۰۴	۶B۹/۷	شرح های پدیدآور	
		۳۰۵	۷B۹/۷	ویراست و تاریخچه کتابشناختی اثر	
		۳۰۶	۹B۹/۷	نشر، پخش و غیره	
		۳۰۷	۱۰B۹/۷	مشخصات ظاهری	
		۳۱۲	۴B۹/۷	عنوان های مرتبط	
		۳۱۴	۶B۹/۷	مسئولیت معنوی اثر	
		۳۲۷	۱۸B۹/۷	مندرجات	
		۳۳۰	۱۷B۹/۷	چکیده	
		۳۳۳	۱۴B۹/۷	مخاطبان	
		۳۳۶	۸B۹/۷	نوع منابع الکترونیکی	
۳۳۷	۸B۹/۷	روش دسترسی			
%۶۰	درصد فراوانی کل				

جدول 2 فراوانی کاربرد قواعد انگلو امریکن و استاندارد یونی مارک در توصیف کتاب شناختی وب گاه ها

همان گونه که در جدول 2 مشاهده می شود بیش ترین کاربرد قواعد و استانداردهای مورد مطالعه برای توصیف وب گاه ها در ناحیه عنوان و پدیدآور با 76 درصد و کم ترین آن متعلق به ناحیه یادداشت با 51/4 درصد است. بنابراین در پاسخ به پرسش دوم می توان گفت امکان توصیف کتاب شناختی وب گاه ها اساس استانداردها و قواعد مورد استفاده در سازمان اسناد و کتابخانه ملی، ایران به میزان 60 درصد است.

پرسش سوم: دسترسی جامع و مانع به وب گاه ها با ایجاد نقاط بازیابی و تحلیل موضوعی مورد استفاده در سازمان اسناد و کتابخانه ملی ایران چگونه است؟

بخش دوم قواعد فهرست نویسی انگلوامریکن قواعد مربوط به نقاط بازیابی و تحلیل موضوعی را ارائه می دهد جدول 3 فیلد های بکار رفته در یونی مارک را بدین منظور نشان می دهد.

عکس

همان گونه که در جدول ۲ مشاهده می شود بیشترین کاربرد قواعد و استانداردهای مورد مطالعه برای توصیف وبگاهها در ناحیه عنوان و پدیدآور با ۷۶ درصد و کمترین آن متعلق به ناحیه یادداشت با ۵۱/۴ درصد است. بنابراین در پاسخ به پرسش دوم می توان گفت امکان توصیف کتاب شناختی وبگاهها بر اساس استانداردها و قواعد مورد استفاده در سازمان اسناد و کتابخانه ملی ایران، به میزان ۶۰ درصد است.

پرسش سوم: دسترسی جامع و مانع به وبگاهها با ایجاد نقاط بازیابی و تحلیل موضوعی مورد استفاده در سازمان اسناد و کتابخانه ملی ایران چگونه است؟
بخش دوم قواعد فهرست نویسی انگلوا امریکن قواعد مربوط به نقاط بازیابی و تحلیل موضوعی را ارائه می دهد. جدول ۳ فیلهای بکاررفته در یونی مارک را بدین منظور نشان می دهد.

جدول ۳. فراوانی کاربرد استاندارد یونی مارک در نقاط بازیابی و تحلیل موضوعی وبگاهها

درصد فراوانی	شماره بلوک	شماره فیلد	نام فیلد
%۶۱/۵	۵) آن های مرتبط	۵۱۰	عنوان اصلی به زبان دیگر
		۵۱۷	عنوان های گونه گون دیگر
		۵۳۲	عنوان گسترده
		۵۴۰	عنوان ترجمه شده
%۸۲	۶) موضوعی و تحلیلی	۶۰۰	نام شخص به منزله موضوع
		۶۰۱	نام تنالگان به منزله موضوع
		۶۰۵	عنوان به منزله موضوع
		۶۰۶	موضوع (اسم عام یا عبارت اسمی عام)
		۶۰۷	نام جغرافیایی به منزله موضوع
%۵۳/۳	۷) مقوری اثر	۷۰۲	نام شخص به منزله شناسه افزوده
		۷۱۰	نام تنالگان به منزله سرشناسه
		۷۱۲	نام تنالگان به منزله شناسه افزوده
%۱۰۰	۸) کاربردی	۸۵۶	نشانی اینترنتی
%۷۰/۴	درصد فراوانی کل		

با توجه به جدول ۳، بیشترین کاربرد استاندارد یونی مارک مربوط به بلوک ۸ با ۱۰۰ درصد است، بعد از آن بلوک ۶ با ۸۲ درصد و بلوک ۵ با ۶۱/۵ درصد قرار دارد. کمترین کاربرد نیز در بلوک ۷ با ۵۳/۳ درصد مشاهده می شود. بنابراین در پاسخ به پرسش سوم میتوان گفت امکان دسترسی جامع و مانع

جدول 3 فراوانی کاربرد استاندارد یونی مارک در نقاط بازیابی و تحلیل موضوعی وب گاه ها

با توجه به جدول، بیش ترین کاربرد استاندارد یونی مارک مربوط به بلوک 8 با 100 درصد است بعد از آن بلوک 6 با 82 درصد و بلوک 5 با 61/5 درصد قرار دارد کم ترین کاربرد نیز در بلوک 7 با 53/3 درصد مشاهده می شود بنابراین در پاسخ به پرسش سوم می توان گفت امکان دسترسی جامع و مانع

به وب گاه ها با ایجاد نقاط بازیابی و تحلیل موضوعی مورد استفاده در سازمان اسناد و کتابخانه ملی ایران حدود 70 درصد است.

پرسش چهارم: کاربرد نظام های موضوعی مورد استفاده در سازمان اسناد و کتابخانه ملی ایران برای نمایه سازی وب گاه ها چگونه است؟

در کتابخانه ملی، ایران منابع غیر کتابی با استفاده از زبان کنترل شده نمایه سازی می شوند. جدول 4 نظام های موضوعی مورد استفاده را نشان می دهد لازم به ذکر است علاوه بر نظام های موضوعی زیر، از سر عنوان های موضوعی کتابخانه کنگره (1)، سایر اصطلاحنامه ها و واژه نامه های موضوعی، دایره المعارف ها و بانک های اطلاعاتی عمومی و موضوعی جهت مستندسازی توصیف گر ها و پیشنهاد آن ها به نظام های موضوعی پیش گفته استفاده می شود

عکس

به وبگاه‌ها با ایجاد نقاط بازیابی و تحلیل موضوعی مورد استفاده در سازمان اسناد و کتابخانه ملی ایران حدود ۷۰ درصد است.

پرسش چهارم: کاربرد نظام‌های موضوعی مورد استفاده در سازمان اسناد و کتابخانه ملی ایران برای نمایه‌سازی وبگاه‌ها چگونه است؟
در کتابخانه ملی ایران، منابع غیرکتابی با استفاده از زبان کنترل‌شده نمایه‌سازی می‌شوند. جدول ۴ نظام‌های موضوعی مورد استفاده را نشان می‌دهد. لازم به ذکر است علاوه بر نظام‌های موضوعی زیر، از سرعنوان‌های موضوعی کتابخانه کنگره^۱، سایر اصطلاحنامه‌ها و واژه‌نامه‌های موضوعی، دایره‌المعارف‌ها و بانک‌های اطلاعاتی عمومی و موضوعی جهت مستندسازی توصیف‌گرها و پیشنهاد آنها به نظام‌های موضوعی پیشگفته استفاده می‌شود.

جدول ۴. فراوانی کاربرد نظام‌های موضوعی در نمایه‌سازی وبگاه‌ها

نظام موضوعی	حوزه موضوعی	درصد فراوانی
اصطلاحنامه فرهنگی فارسی (اصفا)	علوم انسانی	٪۶۶
اصطلاحنامه پزشکی فارسی	پزشکی	٪۱۰۰
اصطلاحنامه های علوم	فنی - مهندسی	٪۱۰۰
سرعنوان های موضوعی فارسی	همه علوم	٪۷۲
درصد فراوانی کل		٪۸۴/۵

بر اساس جدول ۴، هر کدام از نظام‌ها در حوزه موضوعی خود در جامعه مورد مطالعه، بررسی شده‌اند. بنابراین در پاسخ به پرسش چهارم می‌توان نتیجه گرفت که کاربرد نظام‌های موضوعی مورد استفاده در سازمان اسناد و کتابخانه ملی ایران برای نمایه‌سازی وبگاه‌ها بالای ۸۴ درصد امکان‌پذیر است.

بحث و نتیجه‌گیری

یافته‌های پژوهش نشان می‌دهند پردازش وبگاه‌ها در سازمان اسناد و کتابخانه ملی ایران، بیش از ۶۵ درصد امکان‌پذیر است. بر اساس نتایج این پژوهش برای طراحی فرایند پردازش وبگاه‌ها می‌توان فرم الکترونیکی ثبت وبگاه‌ها را طراحی کرد و آن را در وبگاه سازمان دسترسی‌پذیر نمود. در برخی موارد نیز با استفاده از فهرست وبگاه‌ها می‌توان نوعی فراهم‌آوری را انجام داد. در مرحله بعدی با استفاده از قواعد فهرست‌نویسی انگلوا امریکن در بستر فراداده‌های یونی‌مارک و بومی‌سازی آن و با استفاده از کاربرگه

1. (LCSH) Library of Congress Subject Headings

جدول 4 فراوانی کاربرد نظام‌های موضوعی در نمایه‌سازی وب‌گاه‌ها

بر اساس جدول 4، هر کدام از نظام‌ها در حوزه موضوعی خود در جامعه مورد مطالعه بررسی شده‌اند. بنابراین در پاسخ به پرسش چهارم می‌توان نتیجه گرفت که کاربرد نظام‌های موضوعی مورد استفاده در سازمان اسناد و کتابخانه ملی ایران برای نمایه‌سازی وب‌گاه‌ها بالای 84 درصد امکان‌پذیر است.

یافته های پژوهش نشان می دهند پردازش وب گاه ها در سازمان اسناد و کتابخانه ملی ایران، بیش از 65 درصد امکان پذیر است. بر اساس نتایج این پژوهش برای طراحی فرایند پردازش وب گاه ها می توان فرم الکترونیکی ثبت وب گاه ها را طراحی کرد و آن را در وب گاه سازمان دسترس پذیر نمود. در برخی موارد نیز با استفاده از فهرست وب گاه ها می توان نوعی فراهم آوری را انجام داد در مرحله بعدی با استفاده از قواعد فهرست نویسی انگلوامریکن در بستر فرادادهای یونی مارک و بومی سازی آن و با استفاده از کاربرگه

ص: 200

منابع الکترونیکی اقدام به توصیف و تحلیل وبگاه‌ها نمود سپس آن‌ها را بر اساس نقاط بازیابی مختلف دسترس پذیر ساخت. با توجه به نتایج پژوهش، نمایه سازی و تحلیل موضوعی وب گاه‌ها نیز با استفاده از نظام‌های موضوعی پیش گفته امکان پذیر خواهد بود

در مقایسه نتایج این پژوهش با پژوهش‌های پیشین می‌توان خاطر نشان کرد نمونه سازماندهی وب گاه‌ها در کتابخانه‌ها در پژوهش «وایلر» و همکارانش و همچنین «وارد» مورد بررسی قرار گرفته است پردازش وب گاه‌ها از لحاظ هزینه و رضایت کاربران همان گونه که «وارد»، «وایلر» و همکارانش و هم چنین «یانگهی» در پژوهش خود اشاره کرده اند از اهمیت و ارزش والایی برخوردار است. فتاحی» و «حسن زاده» و همچنین «کوچ» و همکارانش نیز بر اعمال رده بندی و فهرست موضوعی، جهت دسترسی به وب گاه‌ها تاکید کرده اند. همچنین استفاده از قواعد و استانداردهای مرسوم سازماندهی از جمله قواعد فهرست نویسی انگلواامریکن و مارک برای سازماندهی وب گاه‌ها در پژوهش «حاجی زین العابدینی»، «ویلیامسون» و «کوچ» و همکارانش مورد بررسی و تایید قرار گرفته است

به طور کلی از این پژوهش می‌توان نتیجه گرفت که استاندارد سازی موضوع سازماندهی و پردازش وب گاه‌ها به عنوان نوعی از منابع الکترونیکی و بومی سازی آن در سازمان اسناد و کتابخانه ملی ایران امکان پذیر است

پیشنهاد های برخاسته از پژوهش

پس از انجام این پژوهش و با عنایت به یافته‌های آن پیشنهاد می‌شود سازمان اسناد و کتابخانه ملی ایران نسبت به انجام موارد زیر که جنبه عملیاتی دارند اقدام نماید.

1. طراحی و ایجاد سامانه ورود اطلاعات وب گاه‌ها در وب گاه سازمان اسناد و کتابخانه ملی ایران جهت فراهم آوری و پردازش آن‌ها؛
2. ارائه آموزش‌های تخصصی در حوزه پردازش و سازماندهی وب گاه‌ها به نیروی انسانی و کارشناسان ذیربط؛
3. برنامه ریزی در جهت اشاعه اطلاعات وب گاه‌ها به صورت جامع و مانع.

منابع

پارس ایندکس <http://parsindex.com/default.aspx> (دسترسی در 1391/7/20)

پژوهشگاه علوم و فناوری اطلاعات ایران 1390 اصطلاح نامه‌های علوم <http://thesauri.irandoc.ac.ir> (دسترسی در 1391/7/20)

پل روزیتا. 1390. ارزیابی وب گاه کتابخانه‌ای پژوهش‌های آماری و معیارهای کیفیت ترجمه رضا خانی پور گزیده مقالات ایفلا 2007 دوربان آفریقای جنوبی، 23-19 اوت 2007 تهران نشر کتابدار

حاجی زین العابدینی، محسن . 1381. بررسی مسائل فهرست نویسی منابع اینترنتی و ارائه دست نامه پیشنهادی برای کتابخانه های ایران پایان نامه کارشناسی ارشد دانشگاه علوم پزشکی و خدمات بهداشتی درمانی ایران دانشکده مدیریت و اطلاع رسانی پزشکی.

سازمان اسناد و کتابخانه ملی جمهوری اسلامی ایران بانک مستندات <http://opac.nlai.ir> (دسترسی در 1391/7/20).

عبداللهی، مریم 1383 ارزیابی وضعیت فهرست نویسی منابع الکترونیکی بر اساس قواعد انگلوماریکن در کتابخانه های دانشگاهی و آرشیوهای موجود در شهر تهران پایان نامه کارشناسی ارشد. دانشگاه آزاد، اسلامی واحد تهران شمال

علی محمدی، داریوش 1381. ابر پایگاه های اطلاعاتی در شبکه جهانی. وب فصلنامه کتاب. 49(1): 93-100

فتاحی، رحمت الله. حسن زاده محمد. 1385 نظر سنجی از کتابداران متخصص پیرامون شیوه های سازماندهی اطلاعات در وب سایت کتابخانه های دانشگاهی گزارشی از مرحله دوم یک طرح پژوهشی فصلنامه کتابداری و اطلاع رسانی 36 (4): 6-30.

کمیته ملی مارک ایران 1381 مارک ایران تهران: کتابخانه ملی جمهوری اسلامی ایران

نشاط، نرگس. 1382 چالش های سازماندهی موضوعی منابع وب اطلاع شناسی. 1: 37-54.

IFLA. 2008. UNIMARC concise bibliographic format. Available at: <http://archive.ifla.org/VI/8/unimarc-concise-bibliographic-format-2008.pdf> [September. 2012]

Koch, Traugott et al. 1997. The role of classification schemes in Internet resource description and discovery. [Available at: www.ub.lu.se/desire/radar/reports/D3,2,3/class-v10.html] [September. 2012]

Nielsen, J. 2002. Top 10 guidelines for homepage usability. Available at: <http://www.useit.com/alertbox/20020512.htm> [September. 2012]

Steering Committee of American Library Association, et al.. 2002. Anglo American Cataloging Rules. Joint Chicago: American Library Association

Ward, Diane. 2001. Internet resource cataloging: the SUNY Buffalo Libraries' response. OCLC Systems Services, 17 (1): 19-26

Willer, Mirna. Buzina, Tanja. Holub, Karolina. Zajec, Jasenka. Milinovic, Miroslav. Topolscak, Nebojsa. 2008. Selective Archiving of Web Resources: A Study of Processing Costs. Program: Electronic Library and Information Systems. v42. n4: p341-364

Williamson N. J.1997. Knowledge structures and the Internet. Knowledge organization for information retrieval: proceedings of the sixth international study conference on classification research. University College London 16-18 June: pp23-27

Younghee Noh. 2011. A study on metadata elements for web-based reference resources system developed through usability testing. Library Hi Tech, 29(2):242-265

ص: 203

محاسبه هزینه‌ها دانسته‌های کنونی ما را در خصوص ایجاد آرشیو تحت وب و اینکه این نوع آرشیو در مقایسه با شیوه فراهم‌آوری منابع چاپی، فعالیتی پرهزینه بوده را تأیید می‌کند. مقاله حاضر اطلاعات سودمندی را در خصوص ایجاد یک آرشیو گزینشی منابع تحت وب و چگونگی ارتقای برنامه‌های نرم‌افزاری ارزیابی مجموعه‌های فراهم‌آوری شده را در اختیار خواننده قرار می‌دهد. و در مقاله راهکارهایی ارائه می‌شود که منجر به افزایش دانش در زمینه بهره‌گیری از فناوری‌های مربوط به ایجاد آرشیو تحت وب، و برای افزایش کارآمدی تمامی پروژه‌های مربوط به ایجاد آرشیو تحت وب می‌شود. یکی از این پروژه‌ها که در مقاله بررسی شده است، آرشیو کتابخانه ملی استرالیاست که اخیراً با استفاده از رویکرد گزینشی، موفق شده که هزینه‌های سرانه خود را در یک بازه زمانی کاهش دهد.

طی دهه، گذشته تعدادی اندکی - اما - اما رو به رشد - از کتابخانه های ملی اقدام به تدوین برنامه هایی در زمینه ایجاد آرشیو تحت وب کرده اند. برنامه هایی که یک یا بیش از چهار رویکرد اصلی زیر را در برداشته اند:

1. ایجاد آرشیو گزینشی، برای مثال، آرشیوهای کتابخانه ملی، کانادا، ژاپن و استرالیا؛

2. تهیه و تدارک دامنه وب در کل کشور به طور دوره ای، برای مثال نمونه های انجام گرفته توسط آرشیو های کشور های اسکانديناوی، از جمله سوئد؛

3. گردآوری موضوعی، برگرفته از کتابخانه کنگره مجموعه مینروا از منتخب های سال 2000، 2002، و 11 سپتامبر 2001؛ و

ص: 205

مجموعه های واسپاری نظیر STORS کتابخانه ایالتی تاسمانی و واسپاری الکترونیکی (1) کتابخانه ملی هلند

کتابخانه ملی، استرالیا به دلیل مزایای موجود در رویکرد گزینشی این روش را انتخاب کرد این دلایل عبارت اند از:

• هر یک از اقلام موجود در آرشیو کتابخانه از لحاظ کیفیت مورد ارزیابی قرار گرفته و تا سر حد امکان با بهره گیری از قابلیت های فنی موجود در این زمینه کاربرد پذیر می شوند.

• هر یک از اقلام موجود در آرشیو را می توان به طور کامل فهرست نویسی، و به بخشی از کتاب شناسی ملی تبدیل کرد و داده های کتاب شناختی را به اشتراک گذاشت هم چنین در فهرست کتابخانه، اطلاعات کتاب شناختی منابع تحت وب با منابع دیگر ادغام شده و کاربران می توانند به صورت یکجا به تمامی منابع دسترسی پیدا کنند.

• هر کدام از اقلام موجود در وب می تواند به سرعت از طریق محیط وب برای مخاطبان خود دسترس پذیر شود؛ زیرا پیش تر تلاش شده اجازه انجام این کار از ناشران اخذ شود.

• ویژگی های هر یک از منابع موجود در آرشیو و شیوه طبقه بندی آن ها برای مدیران مجموعه شناخته شده است. این امر در نخستین گام توانایی ما را در توسعه روش ها و ابزار های گردآوری منابع ذخیره سازی و دسترس پذیری به آن ها افزایش می دهد افزون بر این شناخت و آگاهی، مدیران مجموعه را بر آن می دارد که سیاست ها و راهبرد های حفاظت و نگهداری منابع را به گونه ای اتخاذ کنند که امکان دسترس پذیر سازی منابع برای مدت زمان طولانی فراهم شود.

• سایت هایی را که برای روایات های جست و جوگر موجود دسترس پذیر نیستند نیز می توان شناسایی و گردآوری نمود و با استفاده از روش های دیگر - برای مثال بر حسب نام ناشر - مرتب کرد. این امر، نام های تجاری - که نیازمند رمز عبور ناشر هستند - و پایگاه داده ها را نیز در بر می گیرد.

با وجود مزایای فوق هر یک از رویکرد های موجود در زمینه ایجاد آرشیو تحت وب، معایبی نیز دارد. رویکرد گزینشی نیز از این قاعده مستثنا نیست این، امر به کارمندان کتابخانه بستگی دارد که در محیطی مملو از اطلاعات کاملاً جدید به ایفای وظیفه می پردازند، عملکرد آن ها در این زمینه، باید به گونه ای باشد که در نهایت به قضاوت در خصوص آن چه که در آینده باید مورد تحقیق و تفحص قرار گیرد، بی انجامد و عرصه تحقیقات آینده را روشن سازد. رویکرد گزینشی همچنین به استخراج منابع، جدا از محتوا و پیوند میان آن ها با دیگر منابع خارجی می پردازد. از نظر توسعه و مدیریت آرشیو، بزرگ ترین نقطه ضعف رویکرد آرشیوی این است که روند کار فشرده و پر زحمت است و هزینه انجام کار به ازای هر منبع که قرار است آرشیو شود، بالاست.

ایجاد آرشیو گزینشی منابع تحت وب

اشاره

ایجاد آرشیو گزینشی منابع تحت وب (2)

بررسی هزینه های مربوط به فراهم آوری منابع تحت و در کتابخانه ملی استرالیا

تاریخچه

کتابخانه ملی استرالیا در 1996، به صورت گزینشی و آزمایشی اقدام به ایجاد آرشیوهای تحت وب کرد.

ص: 206

e-Depot -1

Selective Archiving of Web Resources: A Study of Acquisition Costs at the National Library of Australia -2

Margaret Philips, Director, Digital Archiving, National Library of Australia -3

4- استادیار کتابداری و اطلاع رسانی و دانشیار گروه کتابداری و اطلاع رسانی دانشگاه آزاد اسلامی واحد علوم و تحقیقات تهران

اقدامی که در تاریخ خود حرکتی بسیار نو محسوب می شد و کمتر نوشته ای بدان پرداخته بود. حتی کسی نبود که بتوان از او آموخت. ما در مورد ایجاد آرشیوگان تحت وب اطلاعات کافی نداشتیم تا بتوانیم در این خصوص اقدام به تدوین طرح یا راهبردهای مناسبی جهت برآورد هزینه های مربوط به این کار کنیم. تنها در سایه برداشتن گام های کوچک عملی و با عنوان خودآموزی توانستیم به راهمان ادامه داده و رشد کنیم هیچ گونه کمک مالی خاصی برای انجام این کار جدید دریافت نمی کردیم، فقط مجبور بودیم آن دسته از کارکنان مجموعه سازی را که جدا از وظایف کاملاً سنتی، خود، به کار جدید ایجاد آرشیوگان تحت وب علاقه مند بوده و از خود استعدادهایی در این زمینه بروز می دادند، برای انجام این کار جدید تربیت کنیم برای مثال شخصی در بخش فناوری اطلاعات باید بخشی از وقت خود را صرف این می کرد که چگونه می تواند منابع را از وب گاه های ناشران دریافت و اطلاعات آن ها را بر روی سرور کتابخانه نگهداری کند ما برای در اختیار گرفتن و مدیریت فایل های دریافت شده، تنها مجبور به استفاده از نرم افزارهایی بودیم که به طور رایگان موجود بود با اعمال چنین تدابیری نه تنها هزینه های اولیه ایجاد آرشیو تحت وب کاملاً اندک، بلکه عمدتاً در بودجه مربوط به کارکنان پنهان می شد.

با گذشت 9 سال عملیات مربوط به مراتب پیچیده تر شدند کتابخانه استرالیا نرم افزار پاندورا (1) یعنی آرشیو تحت وب استرالیا (2) را راه اندازی کرد که به مثابه آرشیوی کارآمد، وارد طرح شراکت با 9 کتابخانه دیگر استرالیا و نیز دیگر سازمان های مسئول گردآوری منابع فرهنگی، یکی پس از دیگری شد. در مجموع می توان گفت که این امر روندی رو به رشد و توسعه داشت، اما نهادها را با افزایش هزینه ها مواجه ساخت حجم افزونتر فعالیت های مربوط به ایجاد آرشیو و نیاز آن به پشتیبانی از مشارکت کنندگانی که از مکان های دور در این طرح شرکت می کردند، ضرورت ایجاد زیر ساخت های پیچیده فنی متشکل از سیستم ارسال و تحویل مدارک، مدیریت آرشیو، و ذخیره سازی را ایجاب می کرد این، امر به مفهوم وجوب هزینه های ضروری در امر توسعه بود، هر چند که این هزینه ها همچنان از طریق بودجه های موجود پرسنلی برای توسعه مجموعه سازی و بخش فناوری اطلاعات فراهم می شد.

کتابخانه مذکور نمی توانست نظام مدیریت آرشیو مناسبی برای خریداری پیدا کند، بنابراین، نظام ایجاد آرشیو دیجیتالی پاندورا (پانداس) (3) به همراه نظام تحویل مدارک پاندورا، پا به عرصه وجود

گذاشتند. همچنین کتابخانه اقدام به خرید سیستم ذخیره سازی شیء دیجیتالی (دی او اس اس) کرد که زمینه اشتراک منابع را میان پاندورا و دیگر مجموعه های دیجیتالی کتابخانه را فراهم می کرد.

یکی از چیزهایی که در 9 سال گذشته تغییر نکرده این است که: کتابخانه هنوز هیچ بودجه اضافی برای انجام این فعالیت ها در اختیار ندارد به طرز عجیب شواهد موجود بر سودآور بودن چنین اقدامی دلالت دارد؛ لذا چنین مسئله ای ما را بر آن می دارد که بودجه ای را برای تحقق این امر از محل بودجه دولت استرالیا که هر ساله رقم آن رو به تزاید است تعریف نماییم با این حال هیچ گونه طرح کوتاه

ص: 207

مدت خاصی که بتواند بودجه ریزی این فعالیت را به گونه ای توجیه کند تا هنگام [طرح] اتمام به ثبات این فعالیت بی انجامد و آن را در حالت بی ثباتی و ناپایداری رها نکند، وجود ندارد.

آرشیو پاندورا در حال حاضر

آرشیو پاندورا مجموعه قابل توجهی از نشریات برخط و وب گاه های استرالیاست که توسط کتابخانه ملی این کشور و شرکای آن ایجاد شده و گسترش یافته است. آرشیوی که در کتابخانه ملی کانبرا (1)، به شیوه ای متمرکز ذخیره، اداره، و حفظ و نگهداری می شود. از 30 آوریل 2005، که آرشیو مشتمل بر 8235 عنوان، بود، تا کنون نرخ متوسط رشد مجموعه سالیانه در حدود 2400 عنوان می باشد. این عناوین شامل فایل مجزایی از اسناد متنی نظیر فایل های پی دی اف (2)، یا اقلام تحت، وب، نظیر وب گاه های حجیم و پیچیده، متشکل از هزاران فایل در اشکال و فرمت های، مختلف شامل، متن، صدا، تصویر و یا ویدئوست که بر پایه نظامی مشخص و به منظور ایجاد «شاخصی» جدید از عناوین آرشیوی ای که در حال حاضر رقم آن ها در آرشیو موجود به معادل 16/736 مورد بالغ می گردد گردآوری شده است. هدف این بود که عناوین منتخب به شیوه ای ایمن و صحیح در آرشیو ذخیره سازی و به طور مستمر برای همگان دسترس پذیر شود.

آرشیو مزبور شامل هر دو نوع انتشارات برخط ایستا و پویا (3) بوده و وب گاه ها را نیز در بر می گیرد، و نشان دهنده طیف گسترده ای از انواع انتشارات و فرمت های استفاده شده توسط ناشران و پدید آورندگان بر روی وب است. افزون بر این نشریات برخط و وب گاه هایی را که هم اکنون در محیط زنده وب ناپدید گردیده و در هیچ کجای دیگر نیز در دسترس نمی باشند نیز در بر می گیرد.

بسیاری از عناوین موجود در این آرشیو می تواند با اتصال به اینترنت به راحتی و آزادانه در دسترس تمامی افراد در هر نقطه از جهان قرار گیرد تقریباً پنج سال است - اندکی کمتر - که به دلایل تجاری دسترسی به بخش بسیار کوچکی از آرشیو محدود شده است. با وجود این، اطلاعات مربوط به این عناوین دارای این محدودیت را می توان بر روی کامپیوتری واحد در اتاق مطالعه اصلی کتابخانه مشاهده کرد.

دسترسی به محتویات آرشیو را می توان یا از طریق ایجاد پیوندی مناسب به پیشینه فهرست نویسی منبعی خاص یا از طریق موضوع و فهرست عناوین موجود در وب گاه پاندورا امکان پذیر کرد. موتور های جست و جوی متعارف مانند گوگل و یاهو، به نمایه سازی منابع آرشیوی بر حسب آن چه که در عنوان آن ها آمده است می پردازند.

محدوده وظایف

وظایفی که کارکنان کتابخانه ملی و دیگر شرکای آن ها به عنوان بخشی از برنامه ایجاد آرشیو تحت وب پاندورا انجام می دهند متأثر از مجموعه عوامل متعددی نظیر تصمیمات مربوط به سیاست گذاری ها و شرایط محیطی است؛ مجموعه سیاست ها و عواملی که نه به شکلی، حداقل بلکه به شیوه ای کاملاً مؤثر بر

ص: 208

هزینه های انجام طرح اثر گذار هستند.

افراد شرکت کننده در طرح پاندورا به بالاترین شکل بر حفظ اصل به طرفه العینی به نتایج مناسب (رسیدن خصوصیات ظاهری و کارآمد پذیری منبع) در باب منبع منتشره و یا وب گاه، و نیز توجه همزمان بر محتویات آن تأکید می ورزند به محض آن که فرد گردآوری کننده نسخه ای از یک منبع را بر روی سرور کتابخانه ملی قرار دهد افراد همکار در سازمان های دیگر نیز نسخه مزبور را جهت کسب اطمینان از جامعیت و کارآمدی آن و قبل از ارسال آن منبع به آرشیو جهت استفاده عموم، مورد بازبینی قرار می دهند. این تضمین کیفیت فرآیندی بسیار وقت گیر و در نتیجه گران قیمت و پرهزینه است که در واقع، گران ترین جنبه فرآیند فراهم آوری منابع نیز محسوب می شود.

هر عنوان در آرشیو بر مبنای رکوردهای موجود در کتابخانه ملی فهرستگان پیوسته کتابخانه های، دیگر و بر مبنای استفاده از اطلاعات پایگاه کتابشناختی ملی (فهرستگان مشترک اطلاعات کتاب شناختی منابع مربوط به بیش از 850 کتابخانه استرالیا به همراه فراهم آوری امکان دست یابی به منابع توسط خدمات مرورگر کین نکا (1)، فهرست نویسی گردیده است) تصمیم در اتخاذ چنین سیاستی از آن جهت صورت گرفته که امکان شناخت و استخراج منابع برخط یکپارچه سازی شده با دیگر منابع کتابخانه ای، امری بسیار مهم تشخیص داده شده است با وجود، این محرز است که انجام چنین اقداماتی بر هزینه های مربوط به ایجاد آرشیو منابع تحت وب می افزاید.

عامل مهم و یاری کننده دیگری که به صرفه جویی زمانی و استفاده بهینه از اوقات کاری کارکنان منجر می شود شرایط مطرح در قوانین و اسپاری کشور استرالیا است. در بسیاری از ایالت های استرالیا، حتی آن هایی که مشترک المنافع هستند نیز قانون و اسپاری منابع و از آن جمله منابع برخط همچنان مهجور مانده و به تصویب نرسیده است تنها در بخش های شمالی این کشور اخیراً قانونی به تصویب رسیده است که به وضوح امکان و اسپاری منابع برخط را میسر می سازد بدان معنا که تمامی کتابخانه های همکار در این طرح از جمله کتابخانه ملی باید قبل از تهیه نسخه ای از یک منبع نسبت به کسب اجازه از ناشر آن، اقدام سپس آن را برای استفاده همگان در آرشیو دسترس پذیر کنند.

نیروی انسانی و افراد شاغل در پاندورا

ایجاد آرشیو پاندورا اجرا و گسترش برنامه های کاربردی وابسته به آرشیو حفظ و نگهداری سیستم های مربوط، تخصص گرایی و برنامه ریزی به منظور محافظت بلند مدت از آن مستلزم به کارگیری نیروهای کاری شش واحد در دو بخش کتابخانه ملی است

کارکنان شاغل در آرشیو دیجیتال بخش مجموعه سازی مسئول انتخاب منابع موجود در آرشیو و محتوای آن می باشند واحد برنامه های کاربردی بخش فناوری اطلاعات نیز مسئول توسعه زیر ساخت های فنی آرشیو هستند بخش مربوط به خدمات وبگاه مسئولیت توسعه رابط کاربر را برای هر دو سیستم پاندورا و پانداس بر عهده دارد و نظام های حمایتی سیستم های تجاری خود عاملی جهت حفظ تولیدات

ص: 209

آزمایش آموزش و ارزیابی محسوب می شوند و موجب می گردند اهداف مذکور به خوبی اجرا و محقق شود. افزودن بر، این کارکنان واحد خدمات حفاظتی بخش مجموعه سازی مسئولیت دسترس پذیر سازی طولانی مدت محتوای آرشیو را بر عهده دارند.

در سال 2004، در زمان تخمین و برآورد هزینه های زیر سهم کمک های کتابخانه ملی به آرشیو پاندورا در قالب استفاده از نیروهای تمام وقت به شرح سطوح زیر بود:

شعبه آرشیو دیجیتال:

• 1 نفر نیروی تمام وقت در سطح مدیر - کتابدار

• 2 نفر نیروی تمام وقت در سطح یک نفر سرپرست یک نفر کارشناس پروژه های خاص -

• کتابداران

• 4 نفر نیروی تمام وقت در سطح کارکنان ستادی - کتابداران؛ و

• 17 نفر نیروی تمام وقت در سطح کارشناس رفع مشکلات فنی - دارای 5 سال سابقه کار در حوزه فناوری اطلاعات

• 2 نفر مدیر تمام وقت که 10 درصد از وقت خود را هر کدام صرف توسعه و نگهداری سیستم کنند؛

و

• 2 نفر مدیر تمام وقت که 25 درصد از وقت خود را هر کدام صرف توسعه و نگهداری سیستم نمایند. همان طور که در بالا ذکر شد نیروی انسانی شاغل در بخش حفاظت از منابع آرشیوی، به صورت تمام وقت در این طرح حضور داشتند، با وجود این هزینه های مربوط در این برآورد منظور نگردیده است.

هزینه دستیابی به وب گاه ها و نشریات برخط

گرچه برای کتابخانه کاملاً مشخص بود که هزینه سرانه ایجاد واحد آرشیو منابع برخط در مقایسه با منابع چاپی نظیر کتاب ها و پیاپی ها بالا می باشد؛ با این حال، تا همین اواخر هیچ گونه اطلاعات جامع و دقیقی در خصوص هزینه های ایجاد آرشیو تحت وب در اختیار نداشتیم. در سال 2004، کتابخانه تصمیم گرفت تا هزینه های مربوط به فعالیت های واسپاری منابع پیاپی و تک نگاشت ها را برآورد کند. افزون بر این، از آن جا که به نظر می رسید گردآوری منابع برخط استرالیا نیز امتداد بسط همان مسئولیت های مذکور باشد، تصمیم گرفته شد هزینه اجرای این کار نیز به همان روال برآورد شود.

محدوده برآورد هزینه

مطالعات مربوط به برآورد هزینه به بررسی و مطالعه هزینه های مربوط به فراهم آوری منابع کتابخانه ملی به عنوان یک «نمونه» (1) پرداخته سپس آن را به هزینه های پاندورا می افزاید افزون بر این مرز میان هزینه هایی که باید لحاظ یا حذف شوند نیز مشخص می گردد.

هزینه های مستقیم عبارت بودند از:

• هزینه های کارکنان بخش آرشیو دیجیتال کتابخانه ملی؛

• هزینه های اداری بخش آرشیو دیجیتال از قبیل، مسافرت، آموزش شرکت در کنفرانس و

• تجهیزات اداری (شامل هزینه های فرد تأمین کننده)؛ و

• هزینه های مربوط به توسعه زیر ساخت ها و نگهداری از آن ها نظیر هزینه کارکنان بخش فناوری اطلاعات و خرید تجهیزات سخت افزاری و نرم افزاری

هزینه هایی که باید حذف شوند عبارت بودند از:

• هزینه های غیر مستقیم مانند تأمین ایستگاه های کاری کارکنان

• نورپردازی و محافظت از ساختمان و

• هزینه های مربوط به محافظت از محتوای منابع موجود در آرشیو

در برآورد هزینه، حاضر فقط هزینه های انجام شده توسط کتابخانه ملی در نظر گرفته شد بدین معنی که هزینه های مربوط به استخدام کارکنان آرشیو که سازمان های شریک در نظر داشتند، حذف گردید.

روش شناسی انجام کار

جهت برآورد هزینه کارکنان اقدام به تهیه نموداری از تمامی فعالیتهای و فرآیندهای (وظایف صورت گرفته توسط بخش آرشیو دیجیتال) مورد نیاز جهت کسب شاخصی در زمینه ایجاد آرشیو کردند. در نمودار مزبور، هزینه های (فعالیت های) اساسی مشخص شده بود، سپس کارکنان اقدام به تخمین مدت زمانی کردند که به طور متوسط هر یک از کارکنان در هر روز صرف می کردند یک روز کاری شامل 441 دقیقه است

• شناسایی و انتخاب مواد و منابع - 30 دقیقه

• تماس با ناشر، مذاکره برای کسب اجازه جهت افزودن منبعی به آرشیو و بایگانی مکاتبات - 30 دقیقه

• گردآوری تضمین کیفیت و آرشیو منابع - 210 دقیقه

• فهرست نویسی - 81 دقیقه

• فعالیت های دیگر (شامل انجام مکاتبات با افراد و نهاد های مسئول در امر نمایه سازی و چکیده نویسی، سوالات مرجع آرشیو مدار و همکاری کارکنان آرشیو دیجیتال برای بسط و توسعه پانداس) - 60 دقیقه

● برقراری ارتباط با شرکا و نیز پشتیبانی - 30 دقیقه (این فعالیت در برآورد هزینه منظور نشده است) البته تمامی کارکنان کلیه وظایف فوق را به یک نسبت انجام نمی دهند برخی بیش تر و برخی کم تر کار می کنند برای مثال سرپرست بخش عهده دار انجام وظایف اداری و نظارت بر کارکنان نیز می باشد این بدان معنی است که او کم تر از دیگران به کارهایی چون گردآوری، آوری کنترل کیفیت، آرشیو سازی،

ص: 211

و فهرست نویسی می پردازد.

از میان مجموعه وظایفی که در بالا اشاره شد، مدیر بخش ایجاد آرشیو دیجیتال، تنها متعهد به حمایت از شرکا و برقراری تماس با ناشران است اوقات وی بیش تر صرف، مدیریت اعمال سیاست گذاری های، توسعه تبلیغات و ارتباط با دیگر سازمان های فعال در امر ایجاد آرشیو تحت وب است. با وجود این حقوق و دستمزد وی جزء ضروری هزینه های کلی ایجاد آرشیو دیجیتال کتابخانه محسوب می شود و باید در برآورد کلی هزینه ها منظور شود.

در این میان کارمندانی نیز عهده دار انجام اموری هستند که کاملاً از حیطه وظایف آن ها خارج است. برای مثال، دو نفر از کتابداران هر هفته باید وقت خود را صرف میز مرجع موجود در اتاق مطالعه کند، و باید این مدت زمان را از مجموع مدت زمان حضور آن ها در آرشیو کسر کرد.

کتابخانه ملی به عنوان شریک اصلی پاندورا و تأمین کننده زیر ساخت های فنی دارای نقشی حمایتی نسبت به دیگر شرکاست و این، امر به نوبه خود به صرف هزینه های بیش تری می انجامد. تمامی کارکنان بخش آرشیو دیجیتال مقادیر قابل توجهی از وقت خود را صرف برقراری ارتباط با دیگر همکاران و ارائه خدمات فنی پشتیبانی کننده به آن ها می کند کارمندان، ستادی هر روز به طور متوسط 30 دقیقه از وقت خود را صرف انجام این کار، یعنی برقراری ارتباط با دیگر همکاران و حمایت از آن ها می کند، در حالی که همزمان مدیر آن ها تنها 10 دقیقه از وقت خود را در هر روز صرف انجام این مهم می کند. لازم به ذکر است که هزینه های مربوط به موارد فوق در صورت هزینه های برآورد شده منظور نشده است.

نحوه محاسبه هزینه ها

پس از برآورد هزینه ها در مرحله بعد با استفاده از برنامه اکسل روشی برای محاسبه هزینه های مربوط به هر یک از فعالیت های انجام گرفته در آرشیو، طراحی و تدوین شد.

سپس در هر روز مقادیر مربوط به حقوق و دستمزد کارکنان آرشیو دیجیتال و نیز مقادیر زمانی صرف شده توسط آن ها در خصوص انجام هر کار وارد برنامه صفحه گستر اکسل شد و در نهایت، کل تعداد دقیق صرف شده برای انجام هر کار و نیز هزینه های مربوط به هر کارمند، محاسبه گردید.

سپس هزینه های تأمین کنندگان مختلف به شرحی که پیش تر گفته شد نیز اضافه گردید، و هزینه های مربوط به توسعه زیر ساخت های نگهداری نیز در این مرحله لحاظ شدند.

در مجموع می توان گفت که از کل 937 موردی که توسط کتابخانه ملی در ماه های ژوئیه تا اکتبر سال 2004 وارد آرشیو شد به طور متوسط در هر روز 13 عنوان به ثبت رسید.

هزینه های فراهم آوری منابع آرشیوی

با ورود تمامی این اطلاعات در برنامه صفحه گستر اکسل مشخص شد که هزینه های تمام شده برای هر منبع آرشیو شده، صرف نظر از هزینه های مربوط به انجام فعالیت هایی چون برقراری ارتباط با همکاران و حمایت از آن ها رقمی معادل 178/68 دلار استرالیا است.

اطلاعات مربوط به هر یک از اجزای تشکیل دهنده

ص: 212

هزینه های مربوطه در زیر آمده است:

• هزینه های مربوط به استفاده از نیروی انسانی جهت آرشیو هر منبع در آرشیو دیجیتال - به ازای هر مورد 168/36 دلار استرالیا؛

• هزینه های فراهم کنندگان اقلام آرشیوی به ازای هر مورد 3/41 دلار استرالیا؛ و

• هزینه های مربوط توسعه زیرساخت ها و تعمیر و نگهداری اقلام آرشیوی به ازای هر مورد 6/791 دلار استرالیا

همان گونه که ملاحظه می کنید، حقیقت تلخ پر هزینه بودن این کار سخت و پرمشقت، در رویکرد گزینشی ایجاد آرشیو وب به خوبی نمایان و قابل مشاهده است. هزینه های مربوط به نیروی انسانی، 94 درصد کل هزینه ها را تشکیل می دهد

مقایسه با نوع چاپی

تفاوت هزینه های بالای دستیابی به انتشارات، وبی در مقایسه با نشریات چاپی را همچنین می توان از قیاس میان هزینه های مربوط به تهیه این گونه مواد (وبی) و موارد مشابه چاپی (شامل پیابند ها و تک نگاشت ها) که به طور همزمان از طریق قانون و اسپاری فراهم گردیده اند، احراز نمود.

• هزینه فراهم آوری تک نگاشت ها با استفاده از قانون و اسپاری منابع چاپی به ازای هر مورد 12/29 دلار استرالیا؛ و

• هزینه فراهم آوری پیابندهای چاپی با استفاده از قانون و اسپاری منابع چاپی به ازای هر مورد 11/29 دلار استرالیا.

تلاش برای مقایسه هزینه های مزبور اندکی شبیه تلاش برای مقایسه هزینه های مربوط به حمل و نقل هندوانه موز و انگور در بازار در حالت عمده در قیاس با قیمت تمام شده برای هر یک از این نوع اقلام است. در حالت عادی قیمت تمام شده هر یک از این اقلام با قیمت عمده آن ها برابر نیست. در قیاس با این مسئله، قیمت و هزینه تمام شده برای تهیه منابع تحت وب، با قیمت تمام شده تهیه هر یک از این مواد در حالت چاپی نیز برابری ندارد گر چه اختلاف فاحشی میان ماهیت مواد مزبور (منابع تحت وب) و فرآیند تهیه آن ها با نوع چاپی وجود دارد کاملاً آشکار است که در کل، هزینه تهیه و فراهم آوری منابع تحت وب در مقایسه با منابع چاپی بالا تر می باشد

باید این جا به این نکته توجه داشت که در میان اعداد و ارقام موجود به هیچ عنوان هزینه های مربوط به خرید منابع منظور نشده است. کتابخانه منابع چاپی را بر اساس طرح و اسپاری و بدون پرداخت هزینه ای دریافت می کند منابع تحت وب نیز بدون پرداخت هزینه از طریق وبگاه ناشران قابل دستیابی است. این هزینه ها صرفاً مربوط به فرآیند خرید منابع است هزینه های مربوط به قفسه آرایی و آماده سازی اقلام بر روی قفسه ها، منظور، اما هزینه های بعد آن نظیر و جین منابع و مراقبت از مجموعه در محاسبات مربوط به برآورد هزینه ها لحاظ نمی شود.

افزون بر این کتابخانه ملی علاقه مند به تحلیل هزینه های مربوط به ارائه خدمات خاص نیز بود از این رو جزئیات هزینه های مربوط به استفاده از افراد به ازای انجام هر یک از موارد کاری زیر عبارت است از:

• شناسایی و انتخاب مواد و منابع - 10/16 دلار استرالیا

• تماس با ناشر، مذاکره برای کسب اجازه جهت افزودن منبعی به آرشیو و بایگانی مکاتبات - 10/34 دلار استرالیا

• گردآوری تضمین کیفیت و آرشیو منابع - 7/09 دلار استرالیا،

• فهرست نویسی - 27/42 دلار استرالیا، و

• دیگر فعالیت ها (نظیر بسیاری از فعالیت های مدیریتی انجام مکاتبات با افراد و نهاد های مسئول در امر نمایه سازی و چکیده نویسی پرسش های مرجع آرشیو مدار و همکاری کارکنان آرشیو دیجیتال برای بسط و توسعه پانداس - 59,67 دلار استرالیا.

توجه داشته باشید که این برآورد هزینه چندان نیز مفید، نیست زیرا این برآورد هزینه تنها به شکل موردی و در سطح فهرست نویسی ساده و ابتدایی (توصیفی) صورت می پذیرد نه در سطحی تحلیلی و محتوایی در حالی که هر عنوان موجود در آرشیو دارای حداقل دو سطح فهرست نویسی است. به این ترتیب قیمت واقعی تر رقمی معادل 54/84 دلار استرالیا می باشد.

امکان کاهش هزینه ها

با اطلاعات حاصل از این روش (تخمین هزینه فعالیت های صورت گرفته) توانستیم دریابیم که چگونه می توان هزینه ها را کاهش داد ما توانستیم هزینه های خود را با ایجاد تغییر در رویکرد و اقدام جهت ایجاد آرشیو تحت وب (تغییر سیاست های کاری خود) با یافتن شیوه هایی جهت انجام مؤثرتر وظایف خود در این زمینه و با اعمال همزمان هر دو روش کاهش دهیم برای دستیابی به آینده ای قابل پیش بینی، تلاش خواهیم کرد تا تدابیر و سیاست های موجود در این عرصه برای مثال تدوین فهرستگان اطلاعات کتاب شناختی منابع یا تضمین کیفیت را حفظ کنیم. البته موقعیت های دیگری نیز برای صرفه جویی در هزینه ها وجود دارد.

شناسایی و انتخاب عناوین، اقدامی مهم و اساسی برای کاربرد رویکرد گزینشی در ایجاد آرشیو دیجیتال است. در نگاه اول و در نخستین گام مشکل به نظر می رسد که در یابیم چگونه می توان مدت زمان صرف شده برای انجام این قبیل فعالیت ها را کاهش داد به هر حال ما با ناشران دولتی از آن جهت همکاری می کنیم که بتوان از طریق آن ها به ابر داده های مربوط به اطلاعات کتاب شناختی منابع برخط که به صورت گروهی توسط آن ها در پانداس بارگذاری شده و امکان دستیابی خودکار به منابع برخط را فراهم می سازد - دست یابیم در واقع ناشران آن چه را که پاندورا آرشیو خواهد شد - به واسطه تأمین ابر داده های مربوط به، آن ها مشخص می نمایند. پانداس همچنان در صدد است تا از این طریق امکان پردازش و دستیابی به منابع اطلاعاتی را به شکلی انبوه برای همگان محقق سازد با وجود این، این امر زمانی محقق

خواهد شد که قیمت متوسط فراهم آوری منابع اطلاعاتی توسط سازمان های همکار کاهش یابد. هم اکنون فراداده ها توسط تعداد اندکی از ناشران به طور خودکار به رکوردهای مارک ارسال می شوند تا ضمن افزوده شدن به پایگاه داده های کتابشناختی، ملی به طور همزمان در فهرستگان برخط کتابخانه نیز بارگذاری شوند این امر به کاهش هزینه ها خواهد انجامید.

محاسبات بیشتر در خصوص جزئیات مربوط به هر فعالیت (جمع آوری، تضمین کیفیت، و آرشیو منابع) بیان گر آن بود که تضمین کیفیت 86 درصد کل هزینه ها را تشکیل می دهد. چنان چه می توانستیم از نرم افزارهای قابل اعتماد در زمینه ارزیابی کیفیت در پانداس استفاده کنیم صرفه جویی در هزینه ها بسیار با ارزش تر و معنی دارتر می شد

افزون بر این کتابخانه مزبور تلاش می کند با انجام واسطه گری و لابی های مناسب، دولت استرالیا را متقاعد به بسط و گسترش امکان استفاده از قانون واسطه گری منابع برخط، کند تا از این طریق بتواند نسبت به رفع موانع موجود در زمینه ضرورت اخذ مجوز از ناشران به منظور تهیه نسخه منابع آرشیوی اقدام نماید. بدیهی است تحقق این امر موجب کاهش زمان کاری کارکنان آرشیو خواهد شد.

بهره گیری از فناوری به پیشرفت این امر کمک خواهد کرد، و ما را قادر خواهد ساخت که هزینه های مربوط به انجام این فعالیت را کاهش دهیم کنسرسیوم بین المللی محافظت از منابع اینترنتی که کتابخانه ملی نیز از اعضای فعال آن می باشد، در صدد تهیه مجموعه ای از ابزار های لازم جهت ایجاد آرشیو وب نظیر میانجی کاربر به طور اخص طراحی شده برای کتابخانه ملی. است انتظار می رود تحقق این امر نیز به کارآمد شدن این فعالیت بیانجامد.

نتیجه گیری

محاسبه هزینه ها، دانسته های کنونی ما را در خصوص ایجاد آرشیو تحت وب و این که این نوع آرشیو در مقایسه با شیوه فراهم آوری منابع چاپی فعالیتی پر هزینه بوده (حتی اخیراً هزینه های آن به دلیل افزایش هزینه های نیروی کار نسبت به گذشته فزونی یافته است) تأیید کرد. علاوه بر این، مطالعه حاضر اطلاعات سودمندتری را درباره چگونگی ارتقای برنامه ها در اختیار ما قرار داد آیا ما باید از سیاست واحدی در زمینه فهرست نویسی عناوین آرشیوی استفاده کنیم؟ بله شاید از این طریق بتوانیم راه های کم زحمت تری را جهت ثبت پیشینه های کتاب شناختی منابع به دست آوریم تا چه میزان باید وقت و انرژی خود را صرف فرآیندگران قیمت و پرهزینه تضمین کیفیت نماییم؟ مسلماً، خیلی اما بهره گیری از یک نرم افزار ارزیابی مناسب می تواند بسیار سودمند باشد.

بدیهی است افزایش پختگی در زمینه بهره گیری از فناوری های مربوط به ایجاد آرشیو تحت وب، منجر به افزایش کارآمدی تمامی پروژه های مربوط به این زمینه (ایجاد آرشیو تحت وب)، نظیر، آرشیوهای که اخیراً اقدام به استفاده از رویکرد پر زحمت گزینشی کرده و موفق شده اند از هزینه های سرانه پایین تری در یک بازه زمانی برخوردار گردند خواهد شد.

مایلم در اینجا اعلام کنم که بیشترین زحمات در برآورد هزینه‌ها بر عهده ناظم یوسف، مدیر مالی کتابخانه استرالیا بوده است. روش شناسی این کار شامل تخمین هزینه‌ها و تفسیر نتایج مربوط به آن‌ها نیز بر عهده وی بوده است.

منابع

1. Library and Archives Canada. Electronic Collection: a Virtual Collection of Monographs and Periodicals. .1 <http://www.collectionscanada.ca/electroniccollection/> (accessed 22 March 2005); National Diet Library of Japan. Web Archiving Project (WARP). <http://warp.ndl.go.jp/> (accessed 22 March 2005); National Library of Australia. www.nla.gov.au 9|10 Creative Commons Attribution-NonCommercial-ShareAlike 2.1 Australia 19 March 2009 Staff paper Selective Archiving of Web Resources 10 | 10 www.nla.gov.au 19 March 2009 Creative Commons Attribution-NonCommercial-ShareAlike 2.1 Australia PANDORA, (Australia's Web Archive. <http://pandora.nla.gov.au/index.html> (accessed 22 March 2005
2. (National Library of Sweden. Kulturawa3. <http://www.kb.se/kw3/> (accessed 22 March 2005
3. (Library of Congress. MINERVA. <http://www.loc.gov/minerva/> (accessed 22 March 2005
4. State Library of Tasmania. STORS: Long Term Storage of Tasmanian Electronic Documents.. 4 <http://www.stors.tas.gov.au/> (accessed 22 March 2005); Oltmans, E. and H. van Wijngaarden. 2004. Digital preservation in practice: the e-Depot at the Koninklijke Bibliotheek. *Vine* 34 (1): 21-26
5. Information about PANDORA, Australia's Web Archive, and access to its contents are available at .5 <http://pandora.nla.gov.au/index.html>
6. Information about PANDAS is available at <http://pandora.nla.gov.au/pandas.html> .6
7. (National Library of Australia. Kinetica. <http://www.nla.gov.au/kinetica/> (accessed 22 March 2005 .7
8. Levels and pay scales are explained in Attachment A – Salary Table of the .8

.National Library of Australia Certified Agreement 2004–2007 available here

An "instance" is a single gathering of a title. It includes the gathering of a monograph that has been archived .9 once only, the first gathering of a serial or integrating title (for example, a website that changes over time), .and all subsequent gatherings

International Internet Preservation Consortium <http://netpreserve.org/about/index.php> (accessed 22 . 10 .(March 2005

ص: 217

بایگانی اینترنت اقدامی پیشگامانه و تلاشی ابتکاری برای بایگانی منابع وب انجام داده است. معمولاً جست و جوی این نوع بایگانی ها با استفاده از فراداده های توصیفی یا خدمات نمایه سازی متن، کامل امکان پذیر نیست. مسائل بسیاری وجود دارد که بایگانی وب را مشکل می کند و اگر برنامه ریزی سازمانی، برای بایگانی وب از محدوده کوچک تری آغاز شده باشد و بیش تر محدودیت های آن محدودیت های بشری یا محدودیت های منابع سخت افزاری باشد، به برخی مسائل باید با رویکردی متفاوت نگاه کرد. برای نشان دادن برخی مسائل اصلی، بایگانی رقومی برای مطالعه زبان چینی (SHCAD) به عنوان مطالعه موردی در مقاله حاضر مورد بررسی قرار گرفته است.

1. چرایی بایگانی علمی در مقیاس کوچک

با توجه به پیچیدگی های بایگانی وب و نیاز های سخت افزاری و نرم افزاری و همچنین تخصص و پرسنل آیا چنین پروژه هایی فقط برای مؤسسه هایی در مقیاس بزرگ مانند کتابخانه های ملی امکان پذیر است؟ یا اینکه مؤسسه های کوچک تر مانند موزه ها دانشگاه ها و مانند آن ها نیز قادر به انجام وظایف مورد نیاز برای بایگانی وب- با چشم انداز بلند مدت قادر به این کار خواهند بود؟

حتی اگر پاسخ مثبت باشد این سؤال باقی می ماند که آیا در واقع این کار ضروری است یا خیر؟ به هر حال می توان فکر کرد که در حال حاضر بایگانی اینترنت همراه با تلاش برای افزایش تعداد کتابخانه های ملی بسیاری از منابع وب را تحت پوشش دارد و بایگانی می کند.

اجازه دهید با سؤال دیگری شروع کنم بایگانی اینترنت اقدامی پیشگامانه را به عنوان نخستین تلاش ابتکاری، جامع جهت بایگانی منابع وب انجام داده است. موفقیت در انجام این کار، انقلابی بوده و بنیاد و اساس پروژه های بسیاری را پی ریزی کرده است با این حال، بررسی آن چه که بایگانی اینترنت و پروژه های جامع دیگر می توانند به آن دست یابند جهت کشف برخی محدودیت ها، آسان است. از آن جا که اساس

ص: 219

Small Scale Academic Web Archiving: DACHS: in Masanes, Julien (ed.), Web Archiving. Berlin. -1
Heidelberg New York: Springer.pp.213-224

Hanno E. Lecher -2

3- استادیار دانشگاه شاهد

کار مجموعه، بسیار گسترده است، باید برای بخش بزرگی از فعالیت‌ها برای جمع‌آوری اطلاعات وب‌ها به صورت خودکار، تا آن‌جا که امکان پذیر است به روبات‌ها تکیه کرد این نوع جمع‌آوری اطلاعات غالباً بسیار سطحی است بخش‌های زیادی از دست‌رفته بسیاری از صفحات به طور ناقص بارگذاری می‌شوند و بعضی از انواع فایل و نیز وب‌های مخفی به طور کلی نادیده گرفته خواهند شد.

علاوه بر این، از آن‌جا که دریافت به طور خودکار انجام می‌شود در فواصل زمانی نامنظم نمی‌توان انتخابی آگاهانه از منابع داشته باشیم امکان در نظر گرفتن یا تشخیص مطالب مهم موجود - که ممکن است عمر کوتاهی داشته و یا به سختی شناسایی شوند - وجود ندارد، البته بایگانی اینترنت و سایر پروژه‌های بایگانی وب در مقیاس بزرگ برخی امور خاص را انجام داده‌اند که در آن تلاش زیادی صرف توسعه مجموعه‌ها در محدوده خاصی از موضوع‌هایی شده است که انتخاب کرده‌اند. در هر صورت، تعداد این موضوع‌ها بسیار محدود است و طرح‌های پژوهشی زیادی باید آرشیوهای خود را توسعه دهند.

بیش‌ترین مشکل در پروژه‌های جامع بایگانی، وب دسترسی محدود به محتواست. مهم‌ترین اصل شناسایی نشانی دقیق اینترنتی یک سند یا وبگاه است تا بازیابی اطلاعات مقدور شود. معمولاً جست‌وجوی این نوع بایگانی‌ها با استفاده از فراداده‌های توصیفی یا خدمات نمایه‌سازی متن کامل، امکان پذیر نیست. حتی اگر گزینه جست‌وجوی متن کامل در دسترس باشد عدم انتخاب آگاهانه منابع، تنها نتایج نامنظمی شبیه جست‌وجوی وب امروزی ارائه می‌دهد.

با نگاه به این محدودیت‌ها بدیهی است که بایگانی اینترنت و برخی منابع دیگر، نه بایگانی اینترنت را به طور کامل انجام می‌دهند و نه راه‌های دسترسی مناسبی برای بسیاری از اهداف علمی و یا پژوهشی دیگر فراهم می‌کنند. در نتیجه بایگانی علمی در مقیاس کوچک به یک نیاز مهم تبدیل شده است. اما آیا بازگشت به پرسش مطرح شده در ابتدا امکان پذیر است؟

به برخی مسائل باید با رویکردی متفاوت نگاه کرد اگر برنامه ریزی سازمانی، برای بایگانی وب از محدوده کوچک تری آغاز شده باشد و بیش‌تر محدودیت‌های آن محدودیت‌های بشری/محدودیت‌های منابع سخت‌افزاری باشد برای نشان دادن برخی مسائل اصلی بایگانی دیجیتال برای مطالعه زبان چینی (DACHS) به عنوان مطالعه موردی بررسی قرار خواهد گرفت. هر چند باید در ذهن داشته باشید که پروژه‌های متفاوت غالباً به روش‌ها یا راه‌حل‌های متفاوتی نیاز دارند.

2. بایگانی دیجیتال برای مطالعات زبان چینی

اشاره

اهداف اصلی بایگانی دیجیتال مطالعات زبان چینی (DACHS)، اطمینان از دسترسی درازمدت برای شناسایی و بایگانی منابع اینترنت مربوط به مطالعات زبان چینی است انتخاب نقش مهمی در این فرآیند دارد و به طور ویژه در گفتمان سیاسی و اجتماعی بر آن تأکید و در اینترنت چینی منعکس شده است.

در حال حاضر پروژه DACHS توسط کتابخانه‌های دو مؤسسه چین شناسی با نام‌های «مؤسسه چین» در شهر هایدلبرگ آلمان و «مؤسسه چین شناسی» در دانشگاه لیدن هلند به کار گرفته شده است. بدین ترتیب، زیرساختی کاملاً متفاوت با کتابخانه‌های بزرگ ملی دارد (جدول 1)

جدول 1- حجم ذخیره مجموعه Dchs

نخستین بار هنگامی که اندیشه بارگذاری منابع برخط در مورد چین در اواخر 1999 در هایدلبرگ مطرح شد، به هیچ وجه وضعیت آن برای کسی روشن نبود این اندیشه هنوز توسعه نیافته بود و به عنوان بخشی از برنامه ای بزرگ تر برای ایجاد مرکز اروپایی منابع دیجیتال در مطالعات چینی، با هدف بهبود شرایط مربوط به تحقیقات چین و دسترسی به اطلاعات در اروپا، معرفی شد. این پروژه، شامل فعالیت هایی همچون خرید طیف وسیعی از پایگاه داده های تجاری تمام متن حمایت از توسعه پروژه های پایگاه داده دانشگاهی و توسعه امکانات بود، همچنین توسعه برای یافتن بیماری ایدز در منابع چاپی و غیر چاپی در چین و فراهم آوردن امکان دسترسی آزاد به تمام متن منابع به طور گسترده را شامل می شد. این پروژه به مدت پنج سال در نظر گرفته شد و امکانات مالی آن برای بهبود وسایل سخت افزاری موجود و برای استخدام برخی کارکنان - به عنوان دستیاران دانشجویی - صرف شد.

دستور عمل اصلی این پروژه «حداکثر انعطاف پذیری با حداکثر پاسخگویی» بود. این امر به این معنی است که فضایی برای توسعه طرح های تفصیلی بسیاری از پروژه های فرعی و حتی ایده هایی برای پروژه های فرعی جدید مورد نیاز وجود دارد.

در سال 2000، بیشتر به جذب پروژه های فرعی برای اجرای برنامه ریزی پیشرفته و سازماندهی مجدد زیر ساخت های فناوری اطلاعات مؤسسه پرداخته شد، و برنامه ریزی واقعی برای بایگانی وب در سال 2001 آغاز گردید.

2-1. گام های اولیه

با نگاهی به زیر ساخت های موجود که بایگانی وب باید از طریق آن توسعه داده شود، محدودیت ها قابل مشاهده است. بایگانی دیجیتال برای مطالعات زبان چینی توسط کتابخانه مؤسسه به اجرا در آمد. در آن زمان، کتابخانه با استفاده از فناوری اطلاعات با چهار سرور و نزدیک به 100 ایستگاه کاری [جریان کار] نظارت داشت. مسئول نگهداری این زیر ساخت های فناوری اطلاعات کتابداران با کمک یک یا دو دستیار دانشجویی پاره وقت بودند و در صورت نیاز از سوی بخش آی. سی. تی دانشگاه پشتیبانی می شدند.

برای پروژه بایگانی وب باید یک دستیار دیگر دانشجویی به صورت پاره وقت برای تهیه برخی امور (بارگذاری کردن بایگانی ایجاد ابر داده ها و مانند آن) استخدام می شد؛ زیرا کتابدار - علاوه بر مسئولیت های خود برای کتابخانه و محیط فناوری اطلاعات - باید کارهای مدیریت پروژه و توسعه سایر امور نیز در نظر داشته باشد در واقع سهم دستیار پروژه برای توسعه DACHS مهم است. مدرک تحصیلی بالای دستیار برای جلوگیری از تمرکز زدایی دانش در مورد چارچوب نظری پروژه ها اهمیت دارد به خصوص اینکه وقتی کتابدار مؤسسه را جهت کار دیگری محل را ترک کند، اهمیت آن روشن خواهد شد.

البته برخی مسائل، با توجه به اندازه مؤسسه و امکانات آن باید در آغاز پروژه در نظر گرفته شود و به برخی سؤال ها نیز پاسخ داده شود مانند هدف دسترسی دراز مدت منابع بایگانی شده چیست؟ الزامات و نیازهای سخت افزاری و نرم افزاری برای ایجاد و نگهداری بایگانی چه چیزهایی است؟ چگونه انتخاب منابع باید به عنوان یک کار در حال انجام سازماندهی شود و چگونه باید اطلاعات ایجاد شده در دسترس قرار گیرد؟ و مهم تر از همه چه چیز دیگری نیاز است که برای یک برنامه ریزی مناسب در نظر گرفته شود و کجا در جست و جوی پاسخ باشیم؟

پاسخ به آخرین سؤال را می توان در سندی که در سال 2003 به چارچوب استاندارد بایگانی وب تبدیل شده پیدا کرد. 14721:2003 ISO که با عنوان OAIS (سیستم اطلاعات آرشیوی باز) شناخته می شود ثابت کرده است که از اهمیت حیاتی برخوردار است زیرا در آن پس زمینه های تئوری بسیار مهمی تعیین شده است و کمک می کند تا بسیاری از مسائل مهم بایگانی وب را در آن ببینیم.

این سند مفید یک اشکال عمده نیز دارد: همان طور که از نام آن پیداست، فقط چارچوبی برای راهنمایی های تئوری در مسائل مختلف است و پیاده سازی واقعی مهارت های کاربر را ارائه نمی دهد، بنابراین لازم است جست و جوی دیگری برای درک چگونگی تبدیل چارچوب به عمل انجام شود.

اطلاعات بسیاری می توان در جاهایی مانند، RLG PADI و مانند آن یافت اما مهم تر از آن مشارکت فعالانه در کارگاه ها و همایش هایی است که در برخورد با مسائل مرتبط با بایگانی وب ارائه می دهند.

2-2. توسعه پایدار سازمانی

یکی از مهم ترین سؤال هایی که کل پروژه با آن روبه روست چگونگی فراهم کردن توسعه پایدار مؤسسه برای بایگانی دراز مدت است. هیچ یک از سه عامل اساسی در این پرسش نمی توانند به صورت دراز مدت در نظر گرفته شوند تأمین مالی پروژه نمی تواند محدود شود و به پایان برسد؛ علاقه مندی های کاربران در سازمان می تواند تغییر کند و منجر به غفلت از پروژه شود؛ و حتی وضعیت دراز مدت سازمان به هیچ وجه تضمین شده نیست. بنابراین واضح است که سازمان آن جایی با قابلیت اعتماد به مراتب کم تری از کتابخانه یا بایگانی ملی برای بایگانی دراز مدت است که در آن برخی مفاد قانونی می تواند مؤسسه را برای تحقق مسئولیت ها نسبت به مجموعه مجبور می کند.

بنابراین، راهبرد توسعه باید برای بقای بایگانی، فنی در صورت توقف کار مؤسسه یا عدم توانایی آن برای حمایت از پروژه باشد بقا را می توان به دوروش تعریف کرد: یا باید بایگانی را فعال نگه داشت بدین معنی که تمام فعالیت های انتخاب منابع برای داده های در دسترس ادامه داشته باشد؛ یا این که بایگانی باید حداقل در یک وضعیت رکودی حفظ شود، بدین معنی که با وجود این که منابع جدید به بایگانی اضافه نمی شوند حداقل دسترسی به آن چه که در حال حاضر موجود است تضمین شود.

دوره برای این کار وجود دارد نخست و مهم تر از همه بایگانی باید ویژگی های اولیه یک منبع قابل اعتمادی را محقق کند، که به معنای پایبندی به استانداردهای اعلام شده است همان گونه که در مدل OAIS توصیف شد. در عمل استفاده از آن را برای مؤسسه های دیگر ممکن سازد - حالت ایده ال کتابخانه های ملی - البته تا زمانی که قادر به همسانی آن با سایر بایگانی ها نباشیم

راه دوم توسعه بایگانی است که برای آن تلاش شده است. اگر تعدادی از مؤسسه ها به طور فعال و تعاملی در این پروژه مشارکت داشته باشند در زمان توقف کار برای یکی از شرکا برای دیگری این امکان به وجود می آید که کار بایگانی را ادامه داده و از امکانات آن استفاده کند. در عمل، پایبندی به استانداردهای برقرار شده در این نوع کار ضروری به نظر می رسد برای این پروژه نیز تصمیم گرفته شد که به صورت مشارکتی کار شود [پروژه] در جست و جوی شرکایی است که امکان همکاری در این زمینه را داشته باشند و مسائل مربوط را حل کنند و [در نهایت] بایگانی به عنوان یک کمک اساسی به حوزه مطالعات چین وارد شده و شناخته شود.

3-2. سخت افزار

برای اطمینان از کارکرد مناسب DACHS در سطح محلی باید محیط سخت افزاری مناسبی را اندازی شود بعد از تدارک سخت افزاری برای سرور مناسب و ایستگاه های کاری اختصاص داده شده به امور روزمره مانند بارگذاری و مقاصد، مدیریتی لازم است به امکانات پشتیبانی و امنیت و محافظت از این سرور در مقابل ویروس های کامپیوتری توجه شود.

مرکز کامپیوتر دانشگاه هایدلبرگ با استفاده از روش ذخیره سازی توزیع شده IBM ADSTAR، سیستم پشتیبان گیری مهمی را فراهم کرده است. با استفاده از این سیستم هر شب یک نسخه پشتیبان از کل بایگانی بر روی نوار های مغناطیسی در کامپیوتر مرکزی ذخیره می شود کپی های پشتیبانی از این نوارها نیز به طور مرتب در دانشگاه کالسروهه ذخیره می شود بنابراین داده های بایگانی در سه مکان مختلف که امنیت مناسبی را ارائه می دهند نگهداری می شوند.

دستگاه منبع تغذیه اضطراری (UPS) نیز به عنوان یک سیستم RAID (آرایه افزونه دیسک های مستقل ارزان) (سطح اول) جهت تأمین امنیت اولیه برای دسترسی بی وقفه نصب شد. برای محافظت در مقابل ویروس از نرم افزار ویروس یاب McAfee استفاده شد. از طریق تعاریفی که در این ویروس یاب وجود دارد، به صورت ساعتی از سرور McAfee استفاده می شود تمام داده های ورودی به طور مداوم چک شده و به طور خودکار فرآیندهای اسکن منظم از کل بایگانی انجام می پذیرد.

بدینوسیله می‌توانیم زیر ساخت‌های فناوری اطلاعات موجود در مؤسسه را که حمایت بخش ICT دانشگاه را نیز دارد در این بخش مورد استفاده قرار دهیم در تمامی موارد این تجهیزات باید از نظر دانشگاه‌ها استاندارد، باشد به همین دلیل این تجهیزات باید در حوزه آی.سی.تی. مرکزی بوده و در آن جا نگهداری شود. این موضوع می‌تواند از لحاظ حرفه‌ای نوعی مزیت به حساب آید. از طرفی هم می‌تواند یک نقص باشد زیرا به نوعی می‌تواند محدودیت‌هایی بر روی سخت‌افزار و نرم‌افزار اعمال کند.

4-2. نرم افزار

برای شروع کار بایگانی دیجیتال برای پروژه مطالعات چینی به نرم‌افزار مناسب نیاز داریم که باید برخی شرایط این کار را داشته باشد.

اینترنت مجموعه عظیمی از داده‌های به هم پیوسته در فرمت‌ها و سیستم‌های کدگذاری مختلف است و بایگانی چنین داده‌هایی برای محافظت از آن‌ها به نحوی که محتوا و عملکرد آن به صورت اصلی باشد، قدری مشکل به نظر می‌رسد. از آن جا که کارکرد محتوای وب وابسته به نرم‌افزار مرورگر است، هیچ داده قابل اعتمادی برای محافظت از آن وجود ندارد با این حال چنان چه اساس آن محافظت شود، بسیاری از موارد فوق برآورده می‌شود به این معنی که اگر کل وبگاه بارگذاری شد، باید ساختار فایل اصلی دست نخورده نگه داشته شود البته برای نگهداری اسناد بارگذاری شده و ارتباط آن با باید ارتباطات مناسبی در نظر گرفته شود برای این منظور به نرم‌افزاری نیاز است که مناسب این کار بوده و قابلیت تنظیم را داشته باشد این نرم‌افزار باید توانایی اداره طیف گسترده‌ای از فرمت‌های مختلف را داشته و مقرون به صرفه نیز باشد بعد از آزمایش‌هایی که بر روی نرم‌افزارها صورت گرفت به این نتیجه رسیدیم که نرم‌افزار آفلاین اکسپلورر، تمام قابلیت‌های مورد نیاز را برای شروع کار دارد.

5-2. فراداده

مسئله‌ای که موجب دغدغه فکری شده ایجاد فراداده است از یک طرف فراداده اطلاعات لازم را برای مشاهده و توصیف محتوای سند ایجاد می‌کند و از طرفی هم موجب دسترسی کاربران به اطلاعات می‌شود. فراداده‌هایی مانند نویسنده عنوان و موضوع تا زمانی که جست‌وجوی متنی وجود نداشته باشد این نوع فراداده به کاربر کمک می‌کند از طرفی هم فراداده مسائل فنی را برای کمک به حفاظت درازمدت فراهم می‌کند؛ زیرا اطلاعات لازم را برای مدیریت مناسب را فراهم و بررسی را آن در آینده میسر می‌کند.

ایجاد فراداده پر هزینه است؛ با این حال با وجود فراداده نیمه خودکار، جمع‌آوری برخی داده‌ها باید به صورت دستی انجام شود. بنابراین بدیهی است که ایجاد فراداده برای صدها هزار سند با همان سرعتی که از اینترنت بارگذاری شده غیر ممکن است (شکل 1 را ببینید). حتی اگر این کار ممکن بود یافتن آن چه که دقیقاً این فراداده باید حاوی آن باشد، باز هم یافتن جزئیات و فرمت‌های آن، بسیار سخت است.

یکی از پرسش‌های مهمی که در این مورد وجود دارد این است که آیا وجود فراداده ضروری است یا بهتر است بگوییم: آیا صرفاً تکیه بر الگوریتم‌های جست‌وجو برای بازیابی تمام داده‌های مورد نیاز

برای دسترسی و نیز برای اهداف بلند مدت حافظت از داده‌ها امکان پذیر است؟ از دیدگاه کاربر می‌توان استدلال کرد که جست و جوی متن، کامل در مقایسه با سرفصل‌های موضوعی خام یا عناوین حاوی فراداده‌های ناقص جهت پیدا کردن اسناد ابزار قابل اعتماد تری است از نظر فنی بسیاری از اطلاعات ضروری از جمله فرمت فایل، سرعت بارگذاری، رمز گذاری و غیره برای راحتی در بازیابی داده‌هاست یا حتی می‌تواند برای ساختار نام گذاری فایل باشد

عکس

بایگانی وب علمی در مقیاس کوچک DACHS ۲۲۵

برای دسترسی و نیز برای اهداف بلند مدت حافظت از داده‌ها امکان پذیر است؟ از دیدگاه کاربر می‌توان استدلال کرد که جست و جوی متن کامل، در مقایسه با سرفصل‌های موضوعی خام یا عناوین حاوی فراداده‌های ناقص جهت پیدا کردن اسناد ابزار قابل اعتماد تری است. از نظر فنی، بسیاری از اطلاعات ضروری از جمله فرمت فایل، سرعت بارگذاری، رمز گذاری، و غیره برای راحتی در بازیابی داده‌هاست یا حتی می‌تواند برای ساختار نام گذاری فایل باشد.



شکل ۱-۱. رابط کاربری جست و جوی فراداده DACHS که امکان جست و جوی پیشرفته در بایگانی را دارد.

در پایان این ایده رد شد. منابع دیجیتال نوعی اطلاعات هستند که زود از بین می‌روند. بنابراین امکان دارد تا در اسناد جداگانه‌ای نیز ذخیره شوند تا این منابع را با اطمینان و به صورت یکپارچه در آینده و به صورت طولانی مدت قابل دسترسی کرد. با توجه به نوع داده‌ها و اطلاعات، امکان نبودن آنها وجود دارد. برای مدیریت و ادغام بهتر مواد بایگانی شده در موجودی کتابخانه‌های بزرگ، تصمیم گرفته شد تا فراداده به‌عنوان بخشی از کاتالوگ به صورت منظم ایجاد و تهیه شود. بنابراین، کاتالوگ با فراداده تطبیق داده شد تا بخش‌هایی برای مدیریت درست، تاریخچه مبدأ، سابقه مدیریت، نوع فایل، شناسه، و سایر موارد تعیین شود. بسته به پیچیدگی منابع، در حال حاضر سرگرم تهیه رکوردهای فراداده‌ای هستیم که یا بتواند تک‌فایل را توصیف کند که فقط متن را شامل می‌شود یا شامل مجموعه‌ای کامل از فایل‌ها مانند وبگاه، انجمن و یا روزنامه باشد.

۲-۶. سیاست گذاری و خط‌مشی مجموعه

کاملاً بدیهی است که هدف از بایگانی دیجیتال مطالعات چینی نمی‌تواند شامل تمامی موارد موجود در اینترنت باشد و از آن محافظت نماید. نه از لحاظ فنی چنین امکانی وجود دارد و نه مفید به‌نظر می‌رسد. به‌عنوان مؤسسه‌ای تحقیقاتی، علاقه‌مند به بخشی از اینترنت چین هستیم که منعکس کننده جنبه‌های

شکل 101 رابط کاربری جست و جوی فراداده DACHS که امکان جست و جوی پیشرفته در بایگانی را دارد.

در پایان این ایده رد شد منابع دیجیتال نوعی اطلاعات هستند که زود از بین می روند. بنابراین امکان دارد تا در اسناد جداگانه ای نیز ذخیره شوند تا این منابع را با اطمینان و به صورت یکپارچه در آینده و به صورت طولانی مدت قابل دسترسی کرد. با توجه به نوع داده ها و اطلاعات امکان نابودی آن ها وجود دارد.

برای مدیریت و ادغام بهتر مواد بایگانی شده در موجودی کتابخانه های بزرگ، تصمیم گرفته شد تا فراداده به عنوان بخشی از کاتالوگ به صورت منظم ایجاد و تهیه شود. بنابراین کاتالوگ با فراداده تطبیق داده شد تا بخش هایی برای مدیریت درست، تاریخچه مبدأ، سابقه مدیریت، نوع فایل، شناسه، و سایر موارد تعیین شود. بسته به پیچیدگی، منابع در حال حاضر سرگرم تهیه رکوردهای فراداده ای هستیم که با بتواند تک فایل را توصیف کند که فقط متن را شامل می شود یا شامل مجموعه ای کامل از فایل ها مانند، وبگاه انجمن و یا روزنامه باشد.

2-6. سیاست گذاری و خط مشی مجموعه

کاملاً بدیهی است که هدف از بایگانی دیجیتال مطالعات چینی نمی تواند شامل تمامی موارد موجود در اینترنت باشد و از آن محافظت نماید نه از لحاظ فنی چنین امکانی وجود دارد و نه مفید به نظر می رسد. به عنوان مؤسسه ای، تحقیقاتی علاقه مند به بخشی از اینترنت چین هستیم که منعکس کننده جنبه های

ص: 225

خاصی از جامعه چینی است که موارد زودگذری هم هستند بنابراین فراهم کردن انتخاب آگاهانه منابع یک سرمایه است که در حال حاضر به ما کمک می کند و نیز به کاربران کمک می کند در آینده تا جهت شناسایی موادی که ما نیاز به آن داریم و مربوط به هدف ما می شود یاری می رساند. البته ارزش ها در گذر زمان تغییر خواهد کرد و نسل بعد ممکن است انتخاب های مختلفی داشته باشد. در این مورد، بایگانی اینترنت هنوز جایگزینی بسیار غنی را فراهم میکند که در آن گزینه های دیگری در دسترس خواهد بود و ترکیبی از دو رویکرد گزینشی و جامع است که کاربر در آینده با گسترده ترین آرایه از امکانات برای تحقیقات خود قادر به تأمین آن است.

با این حال همزمان میخواهیم خطر محدودیت بیش از حد را در انتخابمان کاهش دهیم. همانطور که در قبلا نیز گفته شد و در ادامه بحث خواهد شد DACHS در حال تبدیل شدن به پروژه تعاونی بزرگ تر است شرکای، مختلف معیارهای مختلفی دارند و در نتیجه مقادیر مختلفی را در روند انتخاب خود اعمال می کنند این امر تنها می تواند محتوای بایگانی را غنی، کند در حالی که در سیاست انتخاب آگاهانه منابع تغییری رخ نمی دهد

با این حال به منظور استفاده بهتر از منابع محدود و در دسترس باید راهبردهای هوشمندانه ای را برای توسعه مجموعه به کار گرفت. با توجه به تعداد زیاد سیاست های اینترنتی در مورد چین و سرعت زیاد آن و توسعه آن و همچنین ناپدید شدن این مطالب در شبکه چین وظیفه ساختن بایگانی منابع برخط که بتواند مهم و ارزشمند باشد کاری سخت و دلهره آور است.

برای رفع این مشکل شروع به ساخت یک شبکه اطلاعاتی از افراد (محققان بومی و شهروندان) کردیم که به طور فعال یا غیر فعال بخشی از این کار را بر عهده گیرند. این کار زمانی عملی می شود که با استفاده از دانشی باشد که برای شناسایی سیاست های خاص و متناسب با کار باشد با این حال این شبکه اطلاعاتی بزرگ می تواند خیلی متنوع تر از فرآیند انتخاب باشد.

معمولاً شبکه اطلاع رسانی نیز هدف دیگری را دنبال می کند. از آن جا که اعضای آن بخشی از فرهنگ اینترنت معاصر است آن ها می توانند اطلاعات زمینه ارزشمندی در مورد منابع بایگانی شده ارائه دهند تا آن جا که ممکن است این اطلاعات باید در فراداده و یا صفحات وب تخصیص داده شده به عنوان بخشی از مجموعه ایجاد و حفظ گردد.

از آن جا که چند ماه نیز بر روی تعدادی پروژه خاص مانند شعر معاصر، سارس (1) یا موضوع همجنس بازی در چین کار کرده ایم، محققان و دانشجویان کارشناس ارشد برای ایجاد آرشیوهای جامع بر روی این موضوع ها کار می کنند، از جمله این وب گاه ها، آن هایی هستند که در معرض خطر بوده و یا شامل موارد دیگر مانند عکس و پوستر می شود برای بایگانی این منابع می توان متن ساده ای را در زمان بایگانی به آن اضافه کرد این نوع عملکرد را برای درک بهتر این منابع که زود از بین می روند، به ویژه در آینده دور، ضروری می دانیم.

اما این کار صرفاً با تکیه بر آگاهان و دانشمندان انجام نمی شود. تأثیر رویدادهای خاص بین المللی مانند حمله تروریستی 11 سپتامبر و یا بازی های المپیک چین در سال 2008، غالباً باعث بحث های داغ بر

روی اینترنت می شود برای جلوگیری از شیوع افکار عمومی ما در حال کار بر روی سیاهه واریسی تالار گفتگو و روزنامه ها هستیم که نوع گزارش ها به آن مربوط می شود تا بتوان فاصله زمانی چند هفته قبل و بعد از چنین رویدادهایی را پوشش دهد.

مجموعه های مشابه دیگری هم توسط اشخاص، خصوصی، پژوهشگران گروه های پژوهشی و یا سایر مؤسسه های به DACHS اهدا و یا فروخته شده اند. گاهی اوقات برای کمک به ناشران وب گاه ها در معرض خطر برای حفظ محتوای وبگاه شان نزدیک شدیم (و یا با آن ها تماس گرفتیم). به هر حال، این مجموعه ها برای DACHS طراحی و ایجاد نشده اند (یا حتی ممکن است همه استانداردهای کیفیت مد نظر ما را دنبال نکنند) اما DACHS به آن ها کمک می کند تا این منابع در درازمدت در دسترس باشند.

2-7. مشارکت

ملاحظه کردیم که اجرای بایگانی تلاشی مشارکتی است و مشارکت راهبردی مهمی برای بقای آرشیوهای وب در مقیاس کوچک است مشارکت می تواند به نگهداری طولانی مدت از طریق عرضه خدمات به اعضا کمک کند. بنابراین، شرکا نیز اجازه توزیع کار را می دهند که این کار به معنای کاهش هزینه و امکان انتخاب گسترده تر از منابع بایگانی است استانداردهای سخت افزار تجربه و کیفیت نیز می تواند به طور وسیعی به اشتراک گذاشته شود و عملکرد کلی بایگانی را بهبود بخشد. مشارکت به عنوان بخشی از راهبرد سیاسی مهم است. پروژه بین المللی مشارکتی قطعاً ساده تر از شناخت جامعه علمی و کمپین موفق برای تأمین بودجه است.

شناخت مسائل DACHS چیزی است که توسعه دستور عمل ها برای مشارکت ها را ممکن می نماید این دستور عمل ها توسط دانشگاه لیدن، هلند مورد بررسی قرار گرفت که در انتهای سال 2003 مشارکت آن عملی شد. هدف اصلی این بود که شرکا مستقل مانده و تا حد امکان هویت خود را حفظ کنند ولی برخی استانداردها و خدمات وجود دارد که باید به اشتراک گذاشته شود.

یکی از مسائل اصلی در همکاری ایجاد یافته های ایدز، است از جمله راهنمای موضوعی پیونددار (یا جدول محتوا) و همچنین متن کامل و یکپارچه سازی امکانات جست و جو در فراداده و بایگانی آن مسائلی است که همکاری شرکای فعلی و آینده پروژه را می طلبد این امر نه تنها به دسترسی متمرکز به فرم صفحه اصلی به صورت اشتراکی باید دارد بلکه باید در مورد چگونگی گزینه های جست و جو برای رسیدن به هدف مطلوب نیز مذاکره شود جست و جوی متن کامل باید در تمام حوزه ها امکان پذیر باشد و برای فراداده یا یک فهرستگان لازم است که همه داده های فیزیکی با هم ذخیره شود که بتواند از طریق فراداده ها منابع مختلف حتی منابع محلی را جست و جو کند و آن ها را به شیوه ای منسجم ارائه دهد. در هر مورد به ایجاد استانداردهای مختلف جهت ایجاد فراداده نیاز داریم.

مسائل دیگری که باید مورد بحث قرار گیرند عبارت اند از سیاست محدود کردن دسترسی مشترک، مقررات تقسیم کار و امور روزمره ای که از دوباره کاری منابع بایگانی جلوگیری می کند.

نکته مهم دیگر برای همکاری این است که کدام یک از شرکا در آینده قادر هستند زیر ساخت بایگانی به صورت محلی ایجاد نمایند و یا اینکه از امکانات خود برای این پروژه استفاده کنند. از آن جا که

تمامی کارهای فوق هنوز در شرف انجام می باشد به اشتراک گذاشتن تجربه در این زمان غیر ممکن است.

بیش تر مسائل در فرآینده مذاکره و پیاده سازی واقعی برای امور مشارکتی پدیدار خواهد شد. در هر صورت، مشارکت برای پروژه های بایگانی وب در مقیاس کوچک در بسیاری از سطوح ضروری است و از امور مهم و بدیهی است.

3. درس های آموخته شده: جمع بندی

در حدود 4 سال پس از اجرای DACHS، در حال حاضر بسیاری از پرسش ها و مشکلات هنوز حل نشده باقی مانده و یا در حال برطرف شدن و در حال توسعه است. با نگاهی به گذشته، می توان گفت که مسائل کمی وجود دارد که لاینحل باقی مانده است که سعی خواهیم کرد آن ها به طور متفاوتی نگاه کنیم.

ضرورت امر در این زمینه تخصیص سمت ها برپروژه است. علاوه بر روال کار روزانه باید به کار توسعه و مدیریت پروژه های بایگانی وب نیز توجه کرد و آن را دست کم نگرفت. تصمیم به انحصاری کردن آن و سپردن همه چیز به یک کتابدار که در حال حاضر مجموعه ای از وظایف دیگر را انجام می دهد، قدری جای تأمل دارد.

در هر صورت نقش او در توسعه پروژه مزایای بسیاری دارد و مهم است سمت مدیریتی تخصیص داده شده برای DACHS بسیار مناسب بوده و منجر به اجرای موفق پروژه می شود. جا داشت تا بسیاری از مسائل خیلی زود تر و مؤثر تر به کار گرفته میشد تا تأثیر مثبتی در توسعه پروژه داشته باشد.

مسئله دوم که - حداقل از امروز - باید به طور ویژه ای به پرداخته شود و مسائل آن حل شود انتخاب نرم افزار جمع آوری داده و بایگانی می باشد. در سال 2001 که DACHS شروع به بایگانی وب کرد هنوز در مرحله ابتدایی بود و بسیاری از دست اندکاران بزرگ امروز تنها فقط شروع به توسعه و انتشار تلاش های خود کردند. اگر چه نرم افزار امروز ما که برای اهداف مورد نظر در آن زمان انتخاب کردیم بیش تر نیازهای ما را برآورده می کند ابزار های امروزی نیز برای این کار بسیار مناسب هستند. لازم است ذکر گردد که یک انتخاب جدید باید روند استفاده و ایجاد فراداده را ساده کند.

مطالبی که عنوان شد نقطه نهایی بوده است که من می خواستم آن را ایجاد کنم. مسائل بسیاری هم وجود دارد که بایگانی وب را مشکل می کند و در برخی موارد به مسائلی بر می خورید که به دفعات سراغ آن رفتید ولی خودتان نتوانستید آن را حل کنید بدون اینکه نگران نوع فرمت فایل یا توسعه نرم افزار جمع آوری داده باشیم این امر امکان پذیر است - حتی ضروری است - که بر تلاش دیگران تکیه کنیم. همکاری نه تنها در میان شرکای پروژه بایگانی است بلکه همکاری در سایر مؤسسه های در زمینه بایگانی وب، می تواند راه حل های مهم و ابزارهایی را فراهم کند که شما به تنهایی قادر به ایجاد آن نخواهید بود.

4. منابع مفید

فهرست زیر تنها انتخاب بسیار کوچکی از منابع است که برای کار مفید است اگر شما وب سایت های زیر را مشاهده کنید منابع بسیار بیش تری را خواهید یافت به خصوص این که این موضوع بسیار عالی با

عنوان دروازه دیجیتالی PADI. نام دارد که همه منابع در ماه می 2006 قابل دسترس بود.

Websites 4.1

Council on Library and Information Resources (CLIR)

<http://www.clir.org/>

Electronic Resource Preservation and Access Network (erpaNet)

<http://www.erpanet.org/>

International Internet Preservation Consortium

<http://netpreserve.org/>

Internet Archive

<http://www.archive.org/>

Networked European Deposit Library (NedLib)

<http://www.kb.nl/coop/nedlib/>

PADI Preserving Access to Digital Information (National Library of Australia)

<http://www.nla.gov.au/padi>

PADI: Web archiving

<http://www.nla.gov.au/padi/topics/92.html>

Mailing Lists 4.2

Archivists

<http://groups.yahoo.com/group/archivists/>

DigiCULT

<http://www.digicult.info/pages/subscribe.php>

DIGLIB - Digital Libraries Research mailing list (IFLA)

<http://infoserv.inist.fr/wwsympa.fcgi/info/diglib/>

OAIS Implementers (RLG)

<http://lists2.rlg.org/cgi-bin/lyris.pl?enter=ois-implementers>

PadiForum (National Library of Australia) <http://www.nla.gov.au/padi/forum/> Web-Archive/

<http://listes.cru.fr/wws/info/web-archive>

Newsletters and Magazines 4.3

CLIR Issues

<http://www.clir.org/pubs/issues/> DigiCULT Newsletter

<http://www.digicult.info/pages/newsletter.php>

D-Lib Magazine

<http://www.dlib.org/DPC/PADI> What is new in digital preservation

<http://www.nla.gov.au/padi/qdiges>

RLG DigiNews <http://www.rlg.org/en/page.php?Page-ID=12081>

ص: 229

هدف این پژوهش بررسی قابلیت های قالب های یونی مارک و مارک 21 برای سازماندهی منابع اطلاعاتی وب و مقایسه آن ها با یکدیگر است تا از این طریق بتوان به ارزیابی این دو قالب در تهیه پیشینه های اطلاعاتی برای منابع اطلاعاتی وب دست یافت روش پژوهش تحلیل محتوا، است به صورتی که ابتدا به شناسایی عناصر توصیف مرتبط با منابع اطلاعاتی وب در قالب کتاب شناختی یونی مارک و مارک 21، و سپس به تحلیل و مقایسه آن ها با یکدیگر پرداخته شد. از یافته های این پژوهش می توان چنین دریافت که قالب مارک می تواند به عنوان استاندارد برای ذخیره و بازیابی منابع اطلاعاتی وب به کار رود با این حال مارک راه حلی برای چگونگی روزآمد نمودن و تعیین سطح توصیف منابع اطلاعاتی وب تهیه نکرده است. کلید واژه ها منابع اطلاعاتی، وب، یونی مارک مارک، 21 قالب کتاب شناختی، سازماندهی

رقیه حجازی، (1) دکتر مرتضی کوکبی (2)

1-مقدمه و بیان مسئله

منابع اطلاعاتی وب بر اساس ویراست 2005 قواعد فهرست نویسی انگلوامریکن، شامل منابع (داده ها و یا برنامه) کد گذاری شده به منظور کاربرد با دستگاه های رایانه ای است که نیاز به برقراری ارتباط با یک شبکه رایانه ای دارد (Weitz 2006) این منابع محسوس و شامل محمل فیزیکی نیستند، و همچنین برای استفاده از آن ها نیاز به برقراری ارتباط با یک دستگاه رایانه (مانند شبکه) یا منابع ذخیره شده بر روی دیسک سخت یا دیگر ابزارهای ذخیره سازی است (Miller 2008).

با ظهور محمل های جدید، اطلاعات از جمله منابع اطلاعات موجود در محیط وب نظام های سازماندهی اطلاعات نیازمند روش های نوینی در ذخیره و بازیابی اطلاعات شده اند. نمایه سازی موتورهای کاوش را نمی توان روش مناسبی در ذخیره و بازیابی اطلاعات دانست زیرا با توجه به گفته فتاحی (1380) آشفتگی اینترنت در حال حاضر بیشتر ناشی از پراکندگی و ناکارآمدی روش های سازماندهی اطلاعات در این محیط است بنابراین برای سازماندهی اطلاعات نیاز به رویکردهایی استاندارد با توجه به ویژگی های محمل های اطلاعاتی است

ص: 231

1- کارشناس ارشد علوم کتابداری و اطلاع رسانی 2 r.hejazi86@yahoo.com

2- استاد کتابداری و اطلاع رسانی دانشگاه شهید چمران اهواز kokabi80@yahoo.com

نظام سازماندهی دانش برای سازگاری با تغییر ساختار منابع اطلاعاتی همواره دورویکرد را برگزیده است. نخست هماهنگی سازگاری و تطابق نظام های سنتی با محیط و رسانه های جدید و دیگر طراحی و ایجاد نظام های جدید به منظور حداکثر بهره وری از امکانات و قابلیت های محیط جدید (طاهری 1387). با توجه به اینکه در محیط های کتابخانه ای حجم عظیمی از منابع اطلاعاتی با استفاده از استانداردهای موجود سازماندهی شده اند و تبدیل و تغییر در آن ها نیازمند صرف وقت و هزینه است، به نظر می رسد استفاده از رویکرد نخست حداقل در محیط های کتابخانه ای مناسب تر باشد یکی از استانداردهای مورد استفاده در نظام های سنتی قالب مارک است قالب مارک از استانداردهای مورد استفاده برای ماشین خوان کردن داده های کتابشناختی است. قالب های مختلفی از مارک همچون مارک 21 و یونی مارک وجود دارند. از سوی دیگر، بیانیه کتابخانه کنگره برای چارچوب کتابشناختی برای دوران دیجیتال (2011 در عمرانی 1390) به این نکته اشاره می کند که «استاندارد مارک مسئولیت تولید میلیون ها پیشینه کتاب شناختی در گوشه و کنار دنیا را بر عهده دارد و نیاز است که در تمام دوره نقل و انتقال به پشتیبانی از مارک ادامه دهیم. علاوه بر این بیانیه عنوان می کند که سامانه ها و خدمات مبتنی بر پیشینه های مرتبط با قالب مارک تا سال های زیادی بخش مهمی از زیر ساخت ها خواهند بود». بنابراین توجه به این قالب در سازماندهی اطلاعات حائز اهمیت است زیرا علاوه بر استفاده کنونی از این، قالب الگوها و استانداردهای جایگزین آینده نیز سازگار با آن طراحی خواهند شد.

همان طور که بیان شد قالب های مختلفی از مارک موجود است که می توان قالب مارک 21 (به دلیل کاربرد آن در پیشینه های کتاب شناختی کتابخانه کنگره آمریکا، کتابخانه بریتانیا، کتابخانه ملی کانادا و برخی کتابخانه های دیگر) و قالب یونی مارک (به دلیل بین المللی بودن و مبنای مارک ایران) را از مهم ترین آن ها دانست. هر دو قالب بر اساس قواعد فهرست نویسی انگلومریکن و تأکید بر منابع چاپی تهیه شده اند. با ظهور محمل های اطلاعاتی، جدید تغییراتی در قالب ها ایجاد شد تا بتوانند با محمل های جدید سازگار و در جهت سازماندهی آن ها کارآمد باشند که از آن جمله می توان به افزودن فیلدها و فیلدهای فرعی، نشان گرها و نویسه های موجود در برجسب پیشینه اشاره کرد.

با توجه به ویژگی منابع اطلاعاتی وب از یک سو و تفاوت های موجود میان قالب مارک 21 و یونی مارک این پرسش پیش می آید که آیا قالب هایی همچون یونیمارک و مارک 21 توانایی سازماندهی منابع اطلاعاتی وب را دارند؟ علاوه بر این تفاوت های موجود در میان این دو قالب تأثیری بر کیفیت سازماندهی منابع اطلاعاتی وب دارد؟ پژوهش حاضر سعی دارد با بررسی دو قالب یونی مارک و مارک 21، به بررسی عناصر تعریف شده بر اساس خصوصیات منابع اطلاعاتی، وب و همچنین تبیین کارایی و قابلیت دو قالب در سازماندهی بهینه منابع اطلاعاتی وب بپردازد. بنابراین به ارزیابی توان و کارایی هر دو قالب در سازماندهی اطلاعات وب و مقایسه آن ها با یکدیگر پرداخته شد.

2- پرسش های پژوهش

این پژوهش در پی پاسخ به پرسش های زیر انجام شد:

1-1 هر یک از قالب های یونیمارک و مارک 21 شامل چه عناصری (اعم از فیلد فیلد، فرعی نشان گر و غیره) برای سازماندهی منابع اطلاعاتی وب هستند؟

1-2 آیا عناصر تعریف شده در دو قالب یونی مارک و مارک 21 نشان دهنده خصوصیات منابع اطلاعاتی وب هستند؟

1-3 آیا تفاوتی میان عناصر تعریف شده در قالب یونی مارک و مارک 21 برای سازماندهی منابع اطلاعاتی وب وجود دارد؟

1-4 اگر پاسخ پرسش سوم مثبت است تفاوت های موجود چه تأثیری بر سازماندهی منابع اطلاعاتی وب می گذارد؟

3-پیشینه پژوهش

بررسی پیشینه های مرتبط با پژوهش نشان داد که پژوهش های صورت گرفته بر اساس دو موضوع کلی قابل تقسیم بندی هستند. موضوع نخست در رابطه با ماهیت منابع اطلاعاتی وب و شناسایی تفاوت های آن ها با دیگر قالب های اطلاعات بود که از آن جمله می توان به پژوهش وارد (2001) (1) اشاره نمود. وی در این پژوهش به اهمیت سازماندهی منابع اینترنتی اشاره کرده و معتقد است که مهارت های فهرست نویسی برای سازماندهی منابع موجود بر روی وب مناسب است در این مقاله به برخی از فیلدهای مارک 21 که در سازماندهی منابع اینترنتی کارآمد هستند نیز اشاره شده است. همچنین حاجی زین العابدینی (1381) به بررسی و تحلیل آخرین فعالیت ها و تحقیقات در زمینه سازماندهی اطلاعات در اینترنت و ارائه یک الگوی مناسب برای فهرست نویسی منابع فارسی پرداخت. نتایج به دست آمده نشان داد که دو روش مهم برای سازماندهی اطلاعات در اینترنت وجود دارد. روش نخست ایجاد پیشینه های کتاب شناختی با استفاده از قواعد عام فهرست نویسی و قواعد خاص منابع اینترنتی و روش دوم ایجاد پیشینه های کتاب شناختی با استفاده از روش های ابر داده عنوان شده است. او به دلیل ویژگی های منابع اینترنتی فارسی روش دوم را برای منابع اینترنتی فارسی مناسب ندانسته است.

موضوع دیگر در رابطه با بررسی انواع فراداده ها در سازماندهی منابع اطلاعاتی وب بود که از آن جمله می توان به پژوهش فارد و ریگیو (2004) (2) اشاره نمود که به بررسی چند استاندارد فراداده از جمله قالب مارک برای مدیریت منابع الکترونیکی (که منابع اطلاعاتی وب زیر مجموعه ای از آن است) پرداختند تحلیل های این پژوهش گران مشخص می کند که در زمان انجام پژوهش هیچ استاندارد و الگوی فراداده ای توانایی نمایش پیچیدگی های منابع الکترونیکی را نداشته است. طاهری (1387) نیز به مقایسه کارایی طرح فراداده ای هسته دوپلین و قالب فراداده ای مارک 21 در سازماندهی منابع اطلاعاتی وب پرداخت نتایج پژوهش وی نشان داد که قالب مارک 21 برای ذخیره پردازش و مبادله اطلاعات محیط وب مناسب تر است. علاوه بر این پژوهش های دیگری همچون مرور قالب های ابر داده (Heer)

ص: 233

Ward -1

Fard and Riggio -2

1996)، کنترل مستند در زمینه کنترل کتاب شناختی در محیط الکترونیک (Gorman 2004)، بررسی و مقایسه قواعد فهرست نویسی انگلومریکن و عناصر هسته دوبلین برای سازماندهی منابع اینترنتی (کوکبی و آخشیک 1385 و سازماندهی صفحات وب با استفاده از نظام های رده بندی کتابخانه (ملای مقدم و نعیم آبادی 1385) نیز با هدف بررسی و مقایسه دو رویکرد سنتی و نوین فهرست نویسی منابع اینترنتی صورت گرفته است

مطالعه پژوهش های مرتبط نشان داد که سازماندهی منابع وب با استفاده از استانداردها و قواعد فهرست نویسی مورد توجه بوده است. علاوه بر این پژوهش گران همواره در پی یافتن استانداردهای مناسب با ویژگی های منابع و بی در جهت سازماندهی هر چه بهتر آن ها بوده اند به همین جهت به ارزیابی و بررسی آن ها پرداخته اند.

4- روش شناسی پژوهش

روش این پژوهش تحلیل محتوا بود به صورتی که در ابتدا به شناسایی عناصر توصیف مرتبط با منابع اطلاعاتی وب در قالب کتاب شناختی یونی مارک و مارک 21 و سپس به تحلیل و مقایسه آن ها با یکدیگر پرداخته شد. برای انجام این پژوهش از قالب کتاب شناختی یونی مارک (International Federation of 2008 (Library Associations and Institutions (IFLA) و قالب کتاب شناختی مارک 21 (Library 2012 of Congress) استفاده شد. قابل ذکر است که شماری از عناصر دادهای موجود در قالب های کتاب شناختی یونی مارک و مارک 21 می توانند در پیشینه های منابع اطلاعاتی وب کاربرد داشته باشند، همان طور که در پیشینه های قالب های دیگر هم کاربرد دارند اما به دلیل تمرکز این پژوهش بر روی منابع اطلاعاتی، وب تنها عناصری که به طور خاص به این نوع از منابع اختصاص داشتند، مورد بررسی قرار گرفتند.

5- تجزیه و تحلیل یافته ها

برای پاسخ به پرسش نخست پژوهش به شناسایی فیلدها فیلدهای فرعی نشان گرها و نویسه های موجود در برچسب پیشینه مرتبط با منابع اطلاعاتی وب پرداخته شد برای پاسخ به این پرسش، علاوه دست نامه کامل هر دو قالب دست نامه آموزشی میلر برای فهرست نویسی منابع اینترنتی (Miller 2008)، راهنمای یونی مارک برای منابع الکترونیکی (International Federation of Library Associations and Institutions (IFLA) 2000). و جدول های تبدیل یونی مارک به مارک 21 (International Federation of Library Associations and Institutions (IFLA) 2001 of Library Associations and Institutions) مرجع در نظر گرفته شد. در جدول 1 تنها فیلدها فیلدهای فرعی و نشان گرهایی که با محتوای منابع اطلاعاتی وب (به طور خاص) سازگار بودند، نشان داده شده است علاوه بر این معادل هر یک از عناصر در دو قالب در یک سطر قرار گرفت.

بررسی و مقایسه قابلیت های یونی مارک ... ۲۳۵

جدول ۱: جدول تطبیقی عناصر مربوط به منابع اطلاعاتی وب در قالب های یونی مارک و مارک ۲۱

یونی مارک				مارک ۲۱			
کد A شامل سیستمها و خدمات پیوسته هم می شود		نویسه ششم: نوع و کدورد / کد A: فایبل رایانه ای		برجسب پیشینه		کد A شامل سیستمها و خدمات پیوسته هم می شود	
توضیحات	فیلدهای فرعی	نشانه ها		نام فیلد	شماره فیلد	توضیحات	فیلدهای فرعی
		اول	دوم				
عبارت "پیوسته" در فیلد فرعی Sb نشان می دهد که مدرک پیوند شده به شکل پیوسته می باشد.	Sb نام عام مواد	—	—	بلوک شناسه رابط	4—		
عبارت "پیوسته" در فیلد فرعی Sb نشان دهنده منابع وبی است	Sb نام عام مواد	—	—	عنوان و نام پدیدآور	200	عبارت "پیوسته" در فیلد فرعی Sh نشان دهنده منابع وبی است.	Sh رسانه
	Sb نام عام مواد	—	—	عنوان قراردادی	500		Sh رسانه
	Sb نام عام مواد	—	—	عنوان مشترک قراردادی	501		Sh رسانه
عبارت "پیوسته" در فیلد فرعی Sn نشان دهنده منابع اطلاعاتی وب است	Sn اطلاعات مترقه	—	—	عنوان به منزله موضوع	605	عبارت "پیوسته" در فیلد فرعی Sg نشان دهنده منابع اطلاعاتی وب است	Sg اطلاعات مترقه
عناصر داده ای یا طول ثابت- ویژگی مواد اضافی						کد 0 نشان دهنده شکل مدرک است	نویسه 06: شکل مدرک
							نویسه 09: فایبل نشان دهنده خدمات و سیستم های رایانه ای
توضیحات فیزیکی فیلدهای ثابت- اطلاعات عمومی				فیلد داده- های کد شده: منابع الکترونیکی	135	کد F نشان دهنده دسترسی از راه دور است	نویسه 01: تعیین مواد خاص
عناصر داده ای یا طول ثابت				فیلد داده- های کد شده: منابع الکترونیکی	135	کد J نشان دهنده خدمات و سیستم های پیوسته است	نویسه 26: نوع فایبل رایانه ای
							نویسه 23: شکل مدرک
				Sa داده های کد شده برای منابع الکترونیکی نویسه 0: نوع منبع الکترونیکی / کد J			

جدول 1 جدول تطبیقی عناصر مربوط به منابع اطلاعاتی وب در قالب های یونی مارک و مارک 21

۲۳۶ مدیریت منابع اطلاعاتی وب

256	مشخصات قابل رایانه- ای	—	—	—	متنقه خاص مواد مشخصات منابع الکترونیکی	230	اطلاعاتی همچون نوع قابل و اندازه آن در این فیلد فرعی وارد می‌شود	Sa مشخصات قابل رایانه‌ای	—	—	اطلاعاتی همچون نوع قابل و اندازه آن در این فیلد فرعی وارد می‌شود
338	نوع حامل	—	—	—	—	—	بر اساس آردی‌ای عبارت "منابع" پیوسته "یکی از حامل‌های رایانه‌ای است.	Sa اصطلاح نوع حامل	—	—	—
516	نوع قابل رایانه‌ای یا پادداشت داده	—	—	—	پادداشت نوع منابع الکترونیکی	336	عبارت "سیستم‌های پیوسته" در این فیلد فرعی نشان دهنده نوع منبع الکترونیکی است.	Sa نوع قابل رایانه‌ای یا پادداشت داده	—	—	عبارت "خدمات و سیستم‌های پیوسته" در این فیلد فرعی نشان دهنده نوع منبع الکترونیکی است.
538	پادداشت جزئیات سیستم	—	—	—	پادداشت سیستم مورد نیاز	337	اطلاعات در مورد شکل دسترسی به منابع ویس در فیلد فرعی Sd آدرس دسترسی خودکار به مدرک در فیلد فرعی Su و بخشی از مواد شرح داده شده در S3 وارد می‌شود	Sa پادداشت جزئیات سیستم Si متن نمایش Su شناسگر متحدالشکل منبع S3 مواد معین	—	—	اطلاعات در مورد شکل دسترسی به منابع ویس در فیلد فرعی Sa و آدرس دسترسی خودکار به مدرک در فیلد فرعی Su وارد می‌شود
753	جزئیات دسترسی سیستم به قابل رایانه- ای	—	—	—	جزئیات دسترسی فنی منابع الکترونیکی	626	—	Sa ساختار و الگوی دستگاهها Sb زبان‌های برنامه‌نویسی سیستم Sc عامل	—	—	Sa ساختار و الگوی دستگاهها Sb زبان‌های برنامه- نویسی سیستم عامل
856	دسترسی و تعیین محل الکترونیکی	روشن- های دسترسی	روشن- های دسترسی	تعریف نشده	دسترسی و تعیین محل الکترونیکی	856	Sa نام میزان؛ Sb شماره دسترسی؛ Sc فهرده‌سازی اطلاعات؛ Sd مسیر؛ Sf نام الکترونیکی؛ Sh پردازشگر درخواست؛ Si دستورالعمل؛ Sj بیت در هر تابه؛ Sk کلمه عبور؛ Sl ورود به سیستم؛ Sm برقراری تماس برای دریافت کمک؛ Sn نام محل میزان؛ So سیستم عامل؛ Sp درگاه؛ Sq نوع قالب الکترونیکی؛ Sr تنظیمات؛ Ss اندازه قابل؛ St تکرار پایانه‌ای؛ Su شناسگر متحدالشکل منبع؛ Sv ساعاتی که روش دسته‌ای قابل دستیابی است؛ Sw شماره کنترل پیشینه؛ Sx پادداشت غیر عمومی؛ Sy متن پیونده؛ Sz روش‌های دسترسی؛ S3 مواد مشخص شده؛ S6 پیونده؛ S8 پیونده فیلد و شماره توالی	روابط	—	—	Sa نام میزان؛ Sb شماره دسترسی؛ Sc فهرده‌سازی اطلاعات؛ Sd مسیر؛ Sf نام الکترونیکی؛ Sh پردازشگر درخواست؛ Si دستورالعمل؛ Sj بیت در هر تابه؛ Sk کلمه عبور؛ Sl ورود به سیستم؛ Sm برقراری تماس برای دریافت کمک؛ Sn نام محل میزان؛ So سیستم عامل؛ Sp درگاه؛ Sq نوع قالب الکترونیکی؛ Sr تنظیمات؛ Ss اندازه قابل؛ St تکرار پایانه‌ای؛ Su شناسگر متحدالشکل منبع؛ Sv ساعاتی که روش دسته‌ای قابل دستیابی است؛ Sw شماره کنترل پیشینه؛ Sx پادداشت غیر عمومی؛ Sy متن پیونده؛ Sz روش‌های دسترسی؛ S3 مواد مشخص شده؛ S6 پیونده؛ S8 پیونده فیلد و شماره توالی

همان طور که در جدول 1 مشاهده می شود در هر دو قالب یکی از نویسه های برجسب پیشینه می تواند نشان دهنده منابع الکترونیکی باشد و در توضیحات ذیل کد آمده که شامل منابع پیوسته نیز هست. «برجسب پیشینه در ابتدای هر پیشینه قرار می گیرد و حاوی داده های مربوط به پردازش آن پیشینه است» (مارک ایران 1381 ، 28) و به طور غیر مستقیم برای کاربرد در تشخیص خود مدرک کتاب شناختی به کار می رود (International Federation of Library Associations and Institutions (IFLA). 2008 12).

همچنین نشان گرهای یکی از فیلدها (856) برای نمایش خصوصیات منابع وبی تعریف شده است. داده های جدول 1 نشان می دهد که در قالب مارک، 21 می توان در سیزده فیلد، و در قالب یونیمارک، می توان در دوازده فیلد اطلاعات مربوط به منابع وبی را وارد کرد.

برای پاسخ به پرسش دوم، پژوهش از خصوصیات تعریف شده برای منابع اطلاعاتی وب توسط هیری (1) (1996) استفاده و عناصر تعریف شده در دو قالب یونی مارک و مارک 21 مقایسه شد که نتایج آن در جدول 2 قابل مشاهده است.

عکس

همان‌طور که در جدول ۱ مشاهده می‌شود، در هر دو قالب، یکی از نویسه‌های برچسب پیشنهادی می‌تواند نشان دهنده منابع الکترونیکی باشد و در توضیحات ذیل کد آمده که شامل منابع پیوسته نیز هست. «برچسب پیشنهادی در ابتدای هر پیشنهاد قرار می‌گیرد و حاوی داده‌های مربوط به پردازش آن پیشنهاد است» (مارک ایران، ۱۳۸۱، ۲۸) و به‌طور غیرمستقیم برای کاربرد در تشخیص خود مدرک کتاب‌شناختی به‌کار می‌رود (International Federation of Library Associations and Institutions (IFLA). 2008. 12).

همچنین، نشانگرهای یکی از فیلدها (۸۵۶) برای نمایش خصوصیات منابع وبی تعریف شده است. داده‌های جدول ۱ نشان می‌دهد که در قالب مارک ۲۱، می‌توان در سیزده فیلد، و در قالب یونی‌مارک، می‌توان در دوازده فیلد اطلاعات مربوط به منابع وبی را وارد کرد.

برای پاسخ به پرسش دوم پژوهش، از خصوصیات تعریف شده برای منابع اطلاعاتی وب توسط هییری^۱ (۱۹۹۶) استفاده، و عناصر تعریف شده در دو قالب یونی‌مارک و مارک ۲۱ مقایسه شد که نتایج آن در جدول ۲ قابل مشاهده است.

جدول ۲: مقایسه خصوصیات منابع اطلاعاتی وب با عناصر موجود در قالب‌های یونی‌مارک و مارک ۲۱

ردیف	ویژگیها	مارک ۲۱	یونی‌مارک
۱	اطلاعات مربوط به مکان‌های مختلف یک منبع	فیلد ۸۵۶ و قابلیت تکرارپذیری آن برای ثبت مکان‌های مختلف یک منبع	فیلد ۸۵۶ و قابلیت تکرارپذیری آن برای ثبت مکان‌های مختلف یک منبع
۲	شیوه (های) دسترسی به منبع	نشانگر اول در فیلد ۸۵۶ و قابلیت تکرارپذیری برای بیش از یک روش دسترسی فیلد فرعی S۲ در فیلد ۸۵۶ فیلد ۷۵۳	نشانگر اول در فیلد ۸۵۶ و قابلیت تکرارپذیری برای بیش از یک روش دسترسی فیلد فرعی Sy در فیلد ۸۵۶ فیلد ۶۲۶
۳	قالب‌های مختلف نسخه‌های مربوط به یک منبع (PDF, HTML, XML, ...)	فیلد فرعی Sq در فیلد ۸۵۶ فیلد فرعی Sy در فیلد ۸۵۶ فیلد ۳۳۸	فیلد فرعی Sq در فیلد ۸۵۶ فیلد فرعی Sb در بلوک ۴
۴	ثبت تغییرات ناشی از نبود ثبات	تکرارپذیری بودن فیلد ۸۵۶ در فیلدهای فرعی Sa, Sb, Sc	تکرارپذیری بودن فیلد ۸۵۶ در فیلدهای فرعی Sa, Sb, Sc
۵	روزآمد سازی منابع	_____	_____
۶	سطح توصیف	_____	_____
۷	اطلاعات مربوط به دسترسی و شرایط آن	فیلد ۵۳۸ فیلد ۲۵۶	فیلد ۳۳۷ فیلد ۲۳۰

داده‌های جدول ۲ نشان می‌دهد که برای ۵ مورد از خصوصیات بیان شده توسط هییری، حداقل یک عنصر در یونی‌مارک و مارک ۲۱ تعریف شده است. همچنین، برای نمایش دو مورد از خصوصیات نیز

1. Heery

جدول ۲: مقایسه خصوصیات منابع اطلاعاتی وب با عناصر موجود در قالب‌های یونی‌مارک و مارک 21

داده‌های جدول 2 نشان می‌دهد که برای 5 مورد از خصوصیات بیان شده توسط هییری حداقل یک عنصر در یونی‌مارک 21 تعریف شده است. همچنین، برای نمایش دو مورد از خصوصیات نیز

برای پاسخ به پرسش سوم پژوهش با استفاده از داده‌های به دست آمده در جدول 1، به مقایسه عناصر تعریف شده در مارک 21 و یونیمارک پرداخته شد که نتایج این بررسی در جدول 3 نمایش داده شده است.

عکس

۲۳۸ مدیریت منابع اطلاعاتی وب

عنصری در مارک ۲۱ و یونی‌مارک یافت نشد.

برای پاسخ به پرسش سوم پژوهش، با استفاده از داده‌های به دست آمده در جدول ۱، به مقایسه عناصر تعریف شده در مارک ۲۱ و یونیمارک پرداخته شد که نتایج این بررسی در جدول ۳ نمایش داده شده است.

جدول ۳. مقایسه عناصر موجود در مارک ۲۱ و یونیمارک برای نمایش خصوصیات منابع اطلاعاتی وب

یونیمارک	مارک ۲۱
نمایش خصوصیات منابع اطلاعات وب در ۱۲ فیلد	نمایش خصوصیات منابع اطلاعات وب در ۱۳ فیلد
نمایش شکل مدرک در فیلدهای شناسه رابط	عدم نمایش شکل مدرک در فیلدهای شناسه رابط
عدم تعیین نوع مدرک در فیلدهای داده‌ای	تعیین نوع مدرک در فیلد داده‌ای (۰۰۸)
عدم تعیین نوع حامل بر اساس استاندارد آر.دی.ای.	تعیین نوع حامل بر اساس استاندارد آر.دی.ای.
در فیلد ۸۵۶ نشانگر دوم برای نمایش روابط تعریف شده است.	در فیلد ۸۵۶ نشانگر دوم برای نمایش روابط تعریف شده است.

داده‌های موجود در جدول ۳ نشان می‌دهد که قالب‌های یونی‌مارک و مارک ۲۱ برای نمایش خصوصیات منابع وبی در ۵ مقوله با یکدیگر متفاوت هستند. در قالب مارک ۲۱، در ۴ مقوله به نمایش خصوصیات و اختصاص عناصر جزئیتر در نمایش آنها توجه بیشتری در مقایسه با یونی‌مارک شده است. همین تفاوت در قالب یونی‌مارک در یک مقوله نسبت به مارک ۲۱ دیده می‌شود.

برای پاسخ به پرسش چهارم، به تجزیه و تحلیل نتایج به دست آمده در پرسش‌های اول و سوم پژوهش پرداخته شد که نتایج آن به قرار زیر است:

- فراوانی فیلدهای اختصاص داده شده به نمایش خصوصیات منابع وبی: با توجه به داده‌های موجود در جدول ۳، تفاوت مارک ۲۱ و یونی‌مارک در این مورد یک فیلد است. جدول ۱ نشان می‌دهد که فیلد فرعی \$b\$ در فیلدهای شناسه رابط یونی‌مارک نشان دهنده قالب مدرک پیوندی است. با توجه به ویژگی بلوک شناسه رابط در قالب یونی‌مارک، با نمایش قالب پیوسته در این فیلدها، علاوه بر مشخص شدن قالب پیوسته پیشینه‌های مرتبط با پیشینه اصلی، در هنگام بازیابی اطلاعات، قالب پیوسته مرتبط با مدرک بازیابی شده نیز مشخص خواهد شد.

در مارک ۲۱، فیلد ۰۰۶ مربوط به نمایش ویژگی مواد اضافی است. تفاوت این فیلد با فیلد ۰۰۸ این است که این فیلد برای مواردی که نمی‌توان آنها را در فیلد ۰۰۸ وارد کرد، آورده می‌شود. فیلد ۰۰۶ بیشتر برای مواردی که دارای چند نوع از یک گروه خصوصیت هستند استفاده می‌شود. به نظر می‌رسد فیلد ۰۰۶ جایگزینی برای فیلد فرعی \$b\$ در بلوک --۴ یونی‌مارک باشد، اگرچه فیلد فرعی \$b\$ علاوه بر نمایش قالب‌های دیگر مدرک مربوط، نمایش دهنده قالب تمام پیشینه‌های مرتبط (بر اساس انواع روابط موجود در میان آثار) است. علاوه بر این، در مارک ۲۱ بر اساس استاندارد آر.دی.ای، فیلد ۳۳۸ نشان دهنده نوع حامل است که در یونی‌مارک تعریف نشده است.

- تعیین نوع مدرک در فیلد داده‌ای: در مارک ۲۱، علاوه بر تعیین نوع فایل رایانه‌ای در فیلد داده‌ای،

مارک 21

داده های موجود در جدول 3 نشان می دهد که قالب های یونی مارک و مارک 21 برای نمایش خصوصیات منابع وبی در 5 مقوله با یکدیگر متفاوت هستند در قالب مارک ، 21، در 4 مقوله به نمایش خصوصیات و اختصاص عناصر جزئیتر در نمایش آن ها توجه بیشتری در مقایسه با یونی مارک شده است. همین تفاوت در قالب یونی مارک در یک مقوله نسبت به مارک 21 دیده می شود.

برای پاسخ به پرسش چهارم به تجزیه و تحلیل نتایج به دست آمده در پرسش های اول و سوم پژوهش پرداخته شد که نتایج آن به قرار زیر است:

- فراوانی فیلهای اختصاص داده شده به نمایش خصوصیات منابع وبی: با توجه به داده های موجود در جدول 3 تفاوت مارک 21 و یونی مارک در این مورد یک فیلد است جدول 1 نشان می دهد که فیلد فرعی b در فیلهای شناسه رابط یونی مارک نشان دهنده قالب مدرک پیوندی است. با توجه به ویژگی بلوک شناسه رابط در قالب یونی مارک با نمایش قالب پیوسته در این فیلهای، علاوه بر مشخص شدن قالب پیوسته پیشینه های مرتبط با پیشینه اصلی در هنگام بازیابی اطلاعات، قالب پیوسته مرتبط با مدرک بازیابی شده نیز مشخص خواهد شد.

در مارک ، 21 فیلد 006 مربوط به نمایش ویژگی مواد اضافی است تفاوت این فیلد با فیلد 008 این است که این فیلد برای مواردی که نمی توان آن ها را در فیلد 008 وارد کرد آورده می شود. فیلد 006 بیشتر برای مواردی که دارای چند نوع از یک گروه خصوصیت هستند استفاده می شود. به نظر می رسد فیلد 006 جایگزینی برای فیلد فرعی b در بلوک - یونی مارک ، باشد اگر چه فیلد فرعی b علاوه بر نمایش قالبهای دیگر مدرک مربوط نمایش دهنده قالب تمام پیشینه های مرتبط (بر اساس انواع روابط موجود در میان آثار) است. علاوه بر این در مارک 21 بر اساس استاندارد آر.دی.ای فیلد 338 نشان دهنده نوع حامل است که در یونی مارک تعریف نشده است.

- تعیین نوع مدرک در فیلد داده ای: در مارک 21، علاوه بر تعیین نوع فایل رایانه ای در فیلد داده ای

شکل مدرک نیز با کدی مشخص در نویسه 23 مشخص می شود باید توجه داشت که در مارک 21 مشخصات تمام انواع قالب ها در فیلد داده ای 008 وارد می شود، در صورتی که در یونی مارک بلوک --1 به اطلاعات کد شده اختصاص یافته و برای انواع مدارک دارای فیلد خاصی است. فیلد 135 مختص داده های کد شده برای منابع الکترونیکی است به همین دلیل نیازی به تعریف کدی برای شکل مدرک نبوده است.

- تعریف نشان گر برای نمایش: روابط در مارک ، 21 نشانگر دوم برای نمایش رابطه بین اطلاعات موجود در فیلد 856 و منبع توصیف شده در پیشینه است. این نشان گر ممکن است برای تولید نمایشی پیوسته یا نظم همگانی فیلدهای 856 استفاده شود (Network Development and MARC 2003 .Standards Office, Library of Congress) فیلد فرعی 13(1) نشان می دهد. به نظر می رسد نمایش روابط در فیلد 856 (حداقل برای منابع وبی) مناسب تر، باشد زیرا با توجه به این که تمام مشخصات منبع الکترونیکی (مانند مسیر شناس گر متحد الشكل منبع، و غیره) در فیلد 856 وارد می شود می توان قالب وبی هر منبع را (در صورت وجود) در پیشینه اصلی منبع وارد کرد و دیگر نیازی به تهیه پیشینه جدید برای قالب وبی آن نیست.

6- بحث و نتیجه گیری افزایش روزافزون منابع اطلاعاتی وب و آشفته گی موجود در بازیابی اطلاعات جامع و مانع، لزوم سازماندهی منابع اطلاعاتی وب را پر رنگ تر می سازد از آن جایی که منابع وبی سازماندهی شده جزئی از نظام های ذخیره و بازیابی کتابخانه ای خواهند بود لازم است برای سازماندهی این گونه منابع استانداردهای موجود در سازماندهی اطلاعات کتابخانه ای مورد بررسی قرار گیرند.

نتایج پژوهش حاضر نشان داد که قالب های مارک 21 و یونی مارک عناصری را به نمایش خصوصیات منابع وبی اختصاص داده اند این عناصر خصوصیات هم چون، قالب نوع دسترسی، ملزومات دسترسی، و نشانی دسترسی را نمایش می دهند که با نتایج طاهری (1387) مبنی بر توانایی های مارک همخوانی داشت. علاوه بر این عناصر تعریف شده در این دو قالب توانایی نمایش سطح توصیف و روزآمدی را برای منابع اطلاعاتی وب. نداشتند نتایج این بخش از پژوهش با نتایج پژوهش وارد (2001) همخوانی نداشته و با نتایج پژوهش فارد و ریگو (2004) مبنی بر ناتوانی ابر داده ها در سازماندهی پیچیدگی های منابع الکترونیکی همخوانی دارد دو مورد بیان شده از مسائل اصلی در سازماندهی منابع وبی هستند زیرا همین خصوصیات منابع اطلاعاتی وب را از دیگر منابع اطلاعاتی متمایز و تصمیم گیری در چگونگی

ص: 239

1- نیز برای نمایش اطلاعات بیش تر در مورد وضعیت است که در آن یک رابطه یک به یک وجود نداشته باشد. برای مثال نشان گر 1 نمایش گر وضعیت است که مدرک توصیف شده در پیشینه کتاب شناختی الکترونیکی نیست اما یک نسخه الکترونیک از آن موجود و اطلاعات ثبت شده در فیلد 856 نشان دهنده نسخه الکترونیکی است (Library of Congress. 2003) همان طور که در بالا اشاره شد یونی مارک برای نمایش روابط از بلوک شناسه رابط استفاده می کند و قالب مدرک را در فیلد فرعی b

سازماندهی آن‌ها را چالش برانگیز نموده است. مقایسه دو قالب کتاب شناختی یونی مارک و مارک 21 نیز نشان داد که تفاوت چندانی در میان تعریف عناصر مربوط به قالب وبی وجود ندارد و تنها تفاوت‌هایی در شکل تعریف آن‌ها مشاهده شد به نظر می‌رسد در انتخاب یکی از دو قالب باید به وضعیت کنونی دو قالب در ایران میزان انطباق هر یک با استانداردها و الگوهای مطرح در سازماندهی اطلاعات وب، و همچنین وضعیت هر یک از آن‌ها در سازماندهی دیگر قالب‌های اطلاعاتی توجه شود.

از یافته‌های این پژوهش می‌توان چنین دریافت که قالب مارک می‌تواند به عنوان استاندارد برای ذخیره و بازیابی منابع اطلاعاتی وب به کار رود، با این حال مارک راه حلی برای چگونگی روزآمد نمودن و تعیین سطح توصیف منابع اطلاعاتی وب تهیه نکرده است پژوهشگران برای ویژگی سطح توصیف، استفاده از فیلدهای رابطه‌ای را پیشنهاد می‌کنند به این صورت که برای نمایش اجزای یک منبع وبی (مانند صفحه‌های یک وبگاه یا همان پیوندهای درونی) از فیلدهای رابطه‌ای کل و جزء و برای نمایش پیوندهای خارجی از فیلدهای رابطه‌ای که نشان دهنده روابطی همچون هم‌ارز هستند، استفاده شود. پژوهشگران برای ویژگی روزآمدی استفاده از فیلدی خاص برای ثبت تغییرات را پیشنهاد می‌دهند. این فیلد شامل فیلدهای فرعی باشد که هر یک نشان دهنده یکی از تغییرات احتمالی (تغییر در آدرس، عنوان و غیره) در منابع اطلاعاتی وب و تعیین نشانگرها برای نمایش نوع تغییر (روزآمد شده، حذف شده، جایگزین شده و غیره) باشد با توجه به این که پژوهش حاضر به بررسی قالب مارک برای سازماندهی منابع اطلاعاتی وب به طور کلی پرداخته است انجام پژوهش‌هایی برای بررسی قالب مارک برای انواع منابع اطلاعاتی وب و ارزیابی توانایی‌های آن ضروری به نظر می‌رسد.

منابع

حاجی زین العابدینی محسن 1381 بررسی مسائل فهرست نویسی منابع اینترنتی و ارائه دست‌نامه پیشنهادی برای کتابخانه‌های ایران پایان‌نامه کارشناسی ارشد دانشگاه علوم پزشکی ایران. دانشکده مدیریت و اطلاع‌رسانی پزشکی

طاهری، سید مهدی 1387. مقایسه کارآیی طرح فراداده‌ای هسته‌دوئین و قالب فراداده‌ای مارک 21 در سازماندهی منابع اطلاعاتی شبکه جهانی وب فصلنامه کتابداری و اطلاع‌رسانی، 3(11): 139-158.

عمرانی ابراهیم 1390 چارچوب کتاب شناختی برای دوران دیجیتال خبرنامه انجمن کتابداری و اطلاع‌رسانی (27) 21-31

فتاحی، رحمت‌الله. 1380. چالش‌های سازماندهی دانش در قرن بیستم فصلنامه کتاب، 12(4): 59-83 کمیته ملی مارک ایران 1381 مارک ایران تهران کتابخانه ملی جمهوری اسلامی ایران

کوکبی، مرتضی، سمیه سادات آخشیک 1385 سازماندهی منابع اینترنتی: قواعد فهرست نویسی انگلو-امریکن یا عناصر فراداده‌های هسته دوئین؟ سازماندهی اطلاعات: رویکردها و راهکارهای نوین (مجموعه) مقالات اولین همایش انجمن کتابداری و اطلاع‌رسانی، ایران 16 و 17 اسفند 1385 (صص. 333-344) تهران: کتابدار

ملای، مقدم، گلناز محمد نعیم آبادی 1385 سازماندهی صفحات وب با استفاده از نظام های رده بندی کتابخانه ای سازماندهی اطلاعات: رویکرد ها و راهکار های نوین مجموعه مقالات اولین همایش انجمن کتابداری و اطلاع رسانی ایران، 16 و 17 اسفند 1385 (صص. 345-364). تهران: کتابدار.

Fard, S.E., A. Riggio. 2004. Medium or message? A new look at standards, structures, and schemata for managing electronic resources. *Library Hi Tech*, 22(2): 144-152

Gorman, M. 2004. Authority Control in the Context of Bibliographic Control in the Electronic Environment. *Cataloging Classification Quarterly*, 38(3-4). Retrieved on September 19, 2011, from <http://www.sba.unifi.it/ac/relazioni/gorman-eng.pdf>

Heery, R. 1996. Review of metadata formats. *Program: Electronic Library and Information Systems*, 30(4): 345-373

International Federation of Library Associations and Institutions (IFLA). 2000. UNIMARC guideline no.6: electronic resources. Retrieved on September 20, 2012, from <http://archive.ifla.org/VI/3/p1996-1/guid6.htm>

International Federation of Library Associations and Institutions (IFLA). 2001. UNIMARC to MARC 21 conversion specifications. Retrieved on September 12, 2012, from <http://www.loc.gov/marc/unimarc21.html>

International Federation of Library Associations and Institutions (IFLA). 2008. UNIMARC manual. München: K. G. Saur

Library of Congress 2003. MARC 21 format for bibliographic data: 856: Electronic location and access. Retrieved on September 19, 2012, from <http://www.loc.gov/marc/bibliographic/bd856.html>

Library of Congress 2012. MARC 21. Retrieved on September 19, 2012, from www.loc.gov/marc/translations.html

Miller, S. 2008. Rules and tools for cataloging internet resources (trainee manual). Retrieved on September 29, 2012, from <http://www.loc.gov/catworkshop/courses/cataloginginternet/pdf/ceig1-IM-FINAL.pdf>

Network Development and MARC Standards Office, Library of Congress 2003. MARC 21 formats: Guidelines for the use of field 856. Retrieved on September 29, 2012, from <http://www.loc.gov/marc/856guide.html>

- Ward, D. 2001. Internet resource cataloging: the SUNY Buffalo Libraries' response. OCLC Systems Services, 17(1): 19-26
- Weitz, J. 2006. Cataloging electronic resources: OCLC-MARC coding guidelines. Retrieved on September/ 20, 2012, from <http://www.oclc.org/support/documentation/worldcat/cataloging/electronicresources>

هدف پژوهش حاضر سنجش کیفیت محیط های رابط کاربر پایگاه های اطلاعاتی مجلات پیوسته تمام متن فارسی «سید»، «مگ ایران»، «نما متن»، «نورمگز»، «پژوهشگاه علوم و فناوری اطلاعات» و «مرکز منطقه ای» از طریق بررسی و مقایسه آن ها با معیارها و استانداردهای رعایت شده در سیاهه واری رابط کاربر پایگاه اطلاعاتی می باشد. روش پژوهش حاضر پیمایشی-توصیفی است گردآوری داده ها با استفاده از یک سیاهه واری محقق ساخته و از طریق مشاهده ی مستقیم انجام گرفته است. یافته ها نشان داد که پایگاه اطلاعاتی «پژوهشگاه علوم و فناوری اطلاعات» در مجموع با رعایت 35 مؤلفه (55/55 درصد) از مؤلفه های سیاهه واری برترین محیط رابط کاربری را از میان پایگاه های اطلاعاتی پیوسته مجلات تمام متن فارسی دارا می باشد محیط های رابط کاربری پایگاه های اطلاعاتی نور، مگز، مرکز منطقه ای و نامتن به ترتیب با میزان رعایت 34 مؤلفه (53/96 درصد)، 32 مؤلفه (50/79 درصد) و 28 مؤلفه (44/44 درصد) رتبه های دوم سوم و پنجم را به خود اختصاص دادند همچنین محیط کاربری دو پایگاه مگ ایران و سید به طور مساوی با میزان رعایت 31 مؤلفه (49/20 درصد) در جایگاه چهارم قرار گرفتند نتایج پژوهش نشان داد که محیط رابط کاربری پایگاه های تمام متن فارسی در امکانات نوینی که طی سال های اخیر پدیدار شده اند دچار کمبود می باشند از این رو پیشنهاد می شود طراحان پایگاه های مذکور به طراحی مجدد پایگاه های خود پرداخته و این امکانات و قابلیت ها را به پایگاه های خود بی افزایند.

کلیدواژه: رابط کاربر پایگاه های اطلاعاتی پیوسته مجلات تمام متن مجلات فارسی

اشاره

صدیقه جعفر زاده (1) معصومه، پیروزفر، (2) عبدالحسین فرج پهلوی (3)

مقدمه

امروزه با رشد و گسترش شبکه جهانی اینترنت پایگاه های اطلاعاتی الکترونیک به عنوان مهم ترین و پر کاربرد ترین ابزار فراهم کننده ی دسترسی به اطلاعات تبدیل شده اند به طوری که پژوهشگران و محققان می توانند اطلاعات مورد نیاز خود را در میان انبوه اطلاعات از طریق این محمول های اطلاعاتی، به دقت و سرعت جستجو و بازیابی کنند نخستین نقطه برخورد کاربر با پایگاه های اطلاعاتی، محیط رابط کاربر است. در پایگاه های اطلاعاتی رابط کاربر عامل مهمی در تسهیل دستیابی کاربران به اطلاعات مورد نیاز خود محسوب می شود از این رو میزبان ها و تولید کنندگان این پایگاه ها می کوشند تا ضمن در نظر گرفتن نیاز کاربران و با پیروی از اصول موجود عوامل و خصیصه های ضروری را جهت طراحی رابط کاربر پایگاه های اطلاعاتی خود شناسایی کرده و از این طریق دسترسی مؤثر کاربران نهایی را به اطلاعات موجود تضمین کنند (مهرداد و زاهدی 1386).

یکی از نمودهای عینی کیفیت استفاده کاربر از پایگاه های اطلاعاتی به ویژه پایگاه های اطلاعاتی پیوسته

ص: 243

1- دانشجوی کارشناسی ارشد کتابداری و اطلاع رسانی دانشگاه شهید چمران sedigheh.jafarzadeh88@gmail.com

2- دانشجوی کارشناسی ارشد کتابداری و اطلاع رسانی دانشگاه شهید چمران piruzfar@gmail.com

3- عضو هیئت علمی و مدیر گروه کتابداری و اطلاع رسانی دانشگاه شهید چمران farajpahlou@gmail.com

را می توان در قالبی با عنوان رابط کاربر مشاهده کرد (یمین فیروز 1382). در واقع رابط کاربر نخستین نقطه ی برخورد کاربر با پایگاه های اطلاعاتی است و به عنوان پلی ارتباطی بین انسان و سامانه ی اطلاعاتی عمل می کند. به همین دلیل مهم ترین هدف از طراحی رابط کاربر برآوردن رضایت کاربران و ایجاد تعامل بیش تر و بهتر بین کاربر و محیط های رایانه ای است (Shneiderman 1998 نقل در اعظمی و فتاحی 1388).

با توجه به این که رابط کاربر خوب باعث می شود کاربران مسیر خود را در پایگاه اطلاعاتی بهتر شناسایی کنند و تأثیر به سزایی در عملکرد آنان خواهد داشت (Bates 2002 ، نقل در علیجانی و دهقانی 1386)؛ بنابراین محیط رابط کاربر باید به کاربران کمک کند تا کلیدواژه ها و عبارت های مناسبی برای جستجوهای خود به کار ببرند منبع مورد نظر را از میان منابع اطلاعاتی انتخاب کنند، نتایج جستجو ها را درک و از چگونگی پیشرفت کارشان آگاه شوند تا بتوانند آسان تر و بهتر به اطلاعات مورد نیاز خود دست یابند (Hirst 1999)

در کشور ما نیز در دو دهه ی اخیر تعدادی پایگاه اطلاعاتی پیوسته ی تمام متن طراحی شده است که پژوهش گران و محققان بسیاری برای انجام کارهای علمی و پژوهشی و رفع نیازهای اطلاعاتی خود به آن ها مراجعه می کنند با توجه به نقش مهم رابط کاربر در دستیابی به محتوای پایگاه های اطلاعاتی و بالا بردن کیفیت خدمات ، آنان محیط رابط این پایگاه ها باید به گونه ای طراحی شده باشد که دسترسی مناسب و مؤثر به اطلاعات را برای کاربران تسهیل بخشد؛ به طوری که آنان بتوانند با درک روشنی از نحوه ی طراحی پایگاه اطلاعاتی ساختار اطلاعات و چگونگی ارائه اطلاعات به ، کاربران اطلاعات مورد نیاز خود را در پایگاه مورد نظر به سرعت جستجو و بازیابی کنند. بنابراین بررسی و مقایسه محیط رابط کاربر پایگاه های اطلاعاتی پیوسته مجلات تمام متن فارسی با معیارهای رعایت شده در پایگاه های اطلاعاتی معتبر پر کاربرد و استاندارد جهانی امری ضروری است این امر ضمن نشان دادن فاصله بین وضع موجود تا وضع مطلوب سبب بهبود کیفیت ارائه ی اطلاعات در پایگاه های اطلاعاتی مجلات فارسی می شود.

بیان مسأله و ضرورت پژوهش

پایگاه های اطلاعات علمی پیوسته مهم ترین و پرکاربرد ترین منابعی هستند که امروزه پژوهشگران جهت دسترسی به اطلاعات مورد نیاز خود به آن ها مراجعه می کنند در ایران نیز در همین راستا طراحی پایگاه های اطلاعاتی الکترونیکی آغاز شده و فعالیت های پراکنده ای توسط برخی سازمان ها در خصوص ارائه ی شکل الکترونیکی مجلات به صورت پیوسته یا ناپیوسته انجام گرفته است. از بین پایگاه هایی که به طور پیوسته شکل الکترونیکی نشریه های ایرانی را ارائه می دهند چهار پایگاه اطلاعاتی «سید» (1)، «مگ ایران»، «نما متن» و «نور مگز» جامعیت و پوشش موضوعی گسترده تری دارند. از طرفی، مدت زمان بیش تری از طراحی آن ها می گذرد بنابراین امکانات بالاتری نسبت به سایر پایگاه ها دارند و برای کاربران شناخته شده تر هستند (اسداللهی و نوکاریزی 1389) در پژوهش حاضر علاوه بر این چهار پایگاه که صرفاً پایگاه های اطلاعاتی نشریات می باشند به بررسی محیط کاربری پایگاه های اطلاعاتی مجلات

ص: 244

فارسی «پژوهشگاه علوم و فناوری اطلاعات» و «مرکز منطقه ای» (1) که دارای پایگاه مجلات و پایگاه های گوناگون دیگری نیز می باشند پرداخته می شود.

با توجه به استفاده گسترده کاربران از پایگاه های اطلاعاتی به عنوان منابع بازیابی اطلاعات و علی رغم تأکید متخصصان بر اهمیت رابط کاربر پایگاه ها و وب سایت ها پژوهش های انجام شده در این زمینه نشان می دهد که طراحان در پاره ای از موارد به دلیل در نظر نگرفتن شرایط و ویژگی های لازم از این مسئله غفلت کرده و در نتیجه مشکلاتی را در خصوص استفاده بهینه کاربران از خدمات خود به وجود آورده اند (Hansen 1998) همچنین بررسی ها نشان داد که به نظر می آید محیط رابط کاربری برخی پایگاه های اطلاعاتی پیوسته مجلات فارسی موجود در ایران در پاره ای موارد با استاندارد ها هماهنگی ندارند و دارای برخی کمبود ها می باشند از این رو وجود این مسئله پژوهش گران را به این امر واداشت تا بررسی و ارزیابی کنند که وضعیت محیط رابط کاربر پایگاه های اطلاعاتی پیوسته تمام متن فارسی در مقایسه با معیارهای رعایت شده در محیط رابط کاربری پایگاه های اطلاعاتی معتبر بین المللی چگونه است؟ و در این راستا نقاط قوت و ضعف این پایگاه ها کدامند؟ نتایج حاصل از این پژوهش می تواند سبب استفاده، بهتر افزایش موفقیت در جستجو ها و انجام بهتر و سریع تر بازیابی اطلاعات شود.

هدف پژوهش

هدف پژوهش حاضر سنجش کیفیت محیط های رابط کاربر پایگاه های اطلاعاتی مجلات پیوسته تمام متن فارسی «پژوهشگاه علوم و فناوری اطلاعات»، «سید»، «مرکز منطقه ای»، «مگ ایران»، «نما متن» و «نور مگز» از طریق بررسی و مقایسه آن ها با معیارهای رعایت شده در محیط رابط کاربر پایگاه های اطلاعاتی مجلات پیوسته تمام متن معتبر و برجسته بین المللی می باشد. به منظور دستیابی به این هدف نقاط قوت و ضعف موجود در محیط رابط کاربری این پایگاه ها مشخص می شود

با توجه به اهداف فوق سعی می شود به پرسش های زیر پاسخ داده شود:

1. وضعیت محیط رابط کاربری هر یک از پایگاه های اطلاعاتی مورد بررسی از لحاظ مقوله جستجو چگونه است؟
2. وضعیت محیط رابط کاربری هر یک از پایگاه های اطلاعاتی مورد بررسی از لحاظ مقوله نمایش اطلاعات چگونه است؟
3. وضعیت محیط رابط کاربری هر یک از پایگاه های اطلاعاتی مورد بررسی از لحاظ مقوله خدمات چگونه است؟
4. وضعیت محیط رابط کاربری هر یک از پایگاه های اطلاعاتی مورد بررسی از لحاظ مقوله پیوندها چگونه است؟
5. وضعیت محیط رابط کاربری هر یک از پایگاه های اطلاعاتی مورد بررسی از لحاظ مقوله راهنمایی چگونه است؟

ص: 245

6. کدام یک از مقوله های کلی رابط کاربری پایگاه اطلاعاتی بیش ترین میزان رعایت را در پایگاه های اطلاعاتی مجلات الکترونیک پیوسته تمام متن فارسی دارد؟

7. کدام یک از پایگاه های اطلاعاتی مجلات الکترونیک پیوسته تمام متن فارسی بیش ترین انطباق را با معیارهای موجود در سیاهه واریسی محیط رابط کاربری پایگاه اطلاعاتی دارد؟

پیشینه ی پژوهش

سابین (2001) (1) رابط کاربر میزبان های پایگاه های اطلاعاتی را بر اساس پنج طبقه کلی عملکرد ارزیابی کرد و سپس خصوصیات بیش تری را که مربوط به هر طبقه بود به عنوان خصیصه های ضروری، مطلوب یا مورد نیاز مشخص کرد یافته های این مطالعه نشان داد که دایالوگ (2) با 32 امتیاز از مجموع 41 امتیاز بالاترین رتبه را دارد. او سی ال سی (3)، اوید (4) با 27 امتیاز در رتبه دوم پروکوئست و سیلور پلاتر (5) با 26 امتیاز در مقام سوم قرار دارند. نکسیس (6) با کسب 25 امتیاز در رتبه هفتم و ویلسون (7) با 20 امتیاز مقام آخر را کسب کرد. در پایان این پژوهش خصیصه های سیستم ایده آل مشخص شد بر این اساس دایالوگ به این سیستم بسیار نزدیک است اما هنوز فاقد تعدادی از خصیصه های مهم می باشد.

کراستینی (2004) (8) یک رابط کاربر گرافیکی جدید برای بازیابی مدارک به صورت سلسله مراتبی ارائه کرده و سپس به نحوی طراحی، اجرا و ارزیابی رابط کاربر پیشنهادی پرداخته است. نتایج این پژوهش حاکی از آن است که رابط کاربر پیشنهادی ابزارهای قدرتمند و مؤثری را برای جستجوی مدارک، رهیابی فهرست بازیابی شده و تصحیح جستجو فراهم می آورد.

شاید بتوان نخستین تلاش در ارزیابی رابط کاربر پایگاه های اطلاعاتی کتاب شناختی فارسی را پژوهش مهرداد و زاهدی (1386) دانست. آنان در پژوهش خود به بررسی و مقایسه محیط رابط کاربر دو میزبان داخلی (کتابخانه منطقه ای علوم و تکنولوژی (9) و پژوهشگاه اطلاعات و مدارک علمی ایران) با چهار میزبان خارجی (الزویر، امرالد، ابسکو و پروکوئست) ارائه دهنده ی پایگاه های اطلاعاتی به روش پیمایش تطبیقی پرداختند. با استفاده از سیاهه واریسی جامع با پنج خصیصه (خصیصه های کلی، جستجو، بازیابی، نمایش و کاربر پسندی) تلاش شد تا ضمن شناخت ویژگی ها و نقاط قوت و ضعف هر یک رابط کاربر میزبان های داخلی و خارجی با یکدیگر مقایسه شوند یافته های این پژوهش نشان داد در میزبان های، داخلی به ترتیب کتابخانه منطقه ای علوم و تکنولوژی و پژوهشگاه اطلاعات و مدارک علمی ایران و در میزبان های خارجی به ترتیب ابسکو، پروکوئست امرالد و الزویر در پنج خصیصه مورد بررسی دارای

ص: 246

Sabin -1

Dialogweb -2

OCLC First Search -3

Ovid -4

Silver Platter -5

Nexis -6

WilsonWeb -7

Crastini -8

9- منظور «مرکز منطقه ای اطلاع رسانی علوم و فناوری» می باشد

اعظمی و فتاحی (1388) به تعیین همخوانی محیط رابط پایگاه‌های اطلاعاتی اِبِسکو، امرالد پروکوئست و ساینس دایرکت با عناصر رفتار اطلاع‌یابی مدل «الیس» (1) که عبارتند از: شروع پیوندیابی، مرور تمایز، بازیابی و استخراج پرداختند نتایج نشان داد که عناصر «شروع» «پیوندیابی» و «تمایز» تا حدودی به وسیله‌ی برخی از محیط‌های رابط کاربر پایگاه‌های مورد بررسی حمایت می‌شوند، اما دیگر عناصر رفتار اطلاع‌یابی (تورق‌بازنگری و استخراج) در ساختار رابط کاربر این پایگاه‌ها لحاظ نشده‌اند. به طور کلی میزان مطابقت و همخوانی رابط کاربر پایگاه‌های اطلاعاتی با عناصر رفتار اطلاع‌یابی الیس در حد متوسط است بنابراین استفاده از این عناصر در طراحی و ارزیابی محیط رابط کاربری می‌تواند تأثیر زیادی بر بهینه‌شدن محیط رابط پایگاه‌های اطلاعاتی و در نتیجه بر فرآیند جستجو و بازیابی داشته باشد.

انتظاریان و فتاحی (1389) به تحلیل تبیین و شناسایی نقاط قوت و ضعف عناصر و ویژگی‌های مهم در محیط رابط پایگاه‌های اطلاعاتی مقاله‌های الکترونیکی مرکز منطقه‌ای اطلاع‌رسانی علوم و فناوری و پژوهشگاه اطلاعات و مدارک علمی پرداختند آن‌ها در این بررسی میزان همخوانی محیط رابط پایگاه‌های مورد بررسی با 10 مؤلفه نیلسن (2)، مشکلات اساسی محیط رابط این پایگاه‌ها و نیز تفاوت بین میزان درک کاربران متخصص و مبتدی را مورد سنجش قرار دادند. یافته‌ها نشان داد میزان همخوانی محیط رابط پایگاه‌های پژوهشگاه با 10 مؤلفه نیلسن به طور کلی در حد متوسط و در پایگاه‌های مرکز منطقه‌ای کمی بیش از حد متوسط است هر دو پایگاه در برخی از مؤلفه‌های مدل نیلسن دارای مشکلات اساسی هستند

اسداللهی و نوکاریزی (1389) طی پژوهشی به ارزیابی ساختار و محتوای پایگاه‌های اطلاعاتی الکترونیکی، سید مگ ایران و نما متن پرداختند. میزان مطابقت پایگاه سید، مگ ایران و نما متن با معیارهای ساختار به ترتیب 50/43 درصد، 54/70 درصد و 41/88 درصد برآورد شد. پایگاه نما متن بر خلاف دو پایگاه سید و مگ ایران به دلیل استفاده از زبان نمایه‌سازی کنترل شده، روند مشخص و منسجمی برای نمایه‌سازی مجله‌ها داشت. نتایج حاصل از بررسی ویژگی‌های منحصر به فرد ساختاری محتوایی پایگاه‌ها مبین آن بود که پایگاه سید از این نظر در رتبه‌ی نخست قرار دارد. به طور کلی هر سه پایگاه می‌توانند با توجه بیش‌تر به معیارهای مطرح شده در سیاهه واریسی بر قابلیت‌های ساختاری خود بیفزایند در بخش محتوا نیز افزایش جامعیت و توسعه‌ی پوشش، موضوعی استفاده از رویکرد ترکیبی یعنی زبان آزاد و کنترل شده، عدم گزینش در نمایه‌سازی افزایش پوشش گذشته‌نگر و کاهش دوره‌ی، روزآمد سازی در هر سه پایگاه توصیه می‌شود.

جمع‌بندی پیشنهادی

مرور پیشنهادها نشان می‌دهد که اهمیت بازیابی اطلاعات باعث گردیده پژوهش‌هایی بر روی پایگاه‌های

اطلاعاتی نشریات الکترونیک انجام گیرد. این پژوهش ها ابعاد مختلف این پایگاه ها از جمله رفتار اطلاع یابی در این پایگاه ها، گرافیک ساختار و محتوای پایگاه ها و به ویژه رابط کاربری این پایگاه ها را مورد بحث قرار داده اند اما نکته ای که حائز اهمیت است این است که تاکنون پژوهشی که رابط کاربری تمام پایگاه های شناخته شده مجلات فارسی را به صورت یک جا با معیارهای رعایت شده در چندین پایگاه بین المللی مورد ارزیابی قرار دهد انجام نشده است. از این رو در این پژوهش رابط کاربری پایگاه های مجلات فارسی با معیارهای پایگاه بین المللی مورد بررسی قرار می گیرد.

روش شناسی پژوهش

روش پژوهش حاضر پیمایشی- توصیفی است. جامعه ی مورد مطالعه در این پژوهش پایگاه های اطلاعاتی مجلات پیوسته تمام متن فارسی «پژوهشگاه علوم و فناوری اطلاعات»، «سید»، «مرکز منطقه ای» «مگ ایران»، «نما متن» و «نور مگز» می باشد با توجه به اهداف پژوهش ابزار گردآوری داده ها با استفاده از یک سیاهه واریسی محقق ساخته است بدین ترتیب که محتوای سیاهه واریسی از معیارهای رعایت شده در محیط رابط کاربری پایگاه های اطلاعاتی مجلات پیوسته تمام متن معتبر و برجسته دنیا از قبیل «ساینس دایرکت» (1)، «پروکوئست» (2)، «امرالده» (3)، «ابسکو» (4) و «اشپرینگر» (5) و همچنین بررسی متون مرتبط استخراج و طراحی شد. سپس روایی آن توسط متخصصان این حوزه بررسی گردید آن گاه با مراجعه مستقیم به وب سایت پایگاه های اطلاعاتی مورد نظر بر اساس سیاهه واریسی به ارزیابی آن ها پرداخته شد. بر اساس معیارهای مطرح شده در سیاهه واریسی به ازای دارا بودن آن معیار عدد یک و در صورت نداشتن آن معیار عدد صفر منظور گردید برای تجزیه و تحلیل داده های گردآوری شده از طریق سیاهه واریسی نیز از آمار توصیفی (فراوانی درصد فراوانی و میانگین) استفاده شد.

یافته های پژوهش

در پاسخ به پرسش اول پژوهش باید گفت که پایگاه های اطلاعاتی سید و نور مگز به طور مساوی با میزان رعایت 11 مؤلفه (50 درصد) بیش ترین میزان رعایت معیارهای موجود در سیاهه واریسی را نسبت به پایگاه های اطلاعاتی دیگر داشته اند همچنین پایگاه های اطلاعاتی پژوهشگاه علوم و فناوری اطلاعات و مگ ایران و مرکز منطقه ای به طور مساوی 10 مؤلفه (45/45 درصد) و پایگاه اطلاعاتی نامتن 8 مؤلفه (36/36 درصد) از مؤلفه های مربوط به مقوله ی جستجو را در خود رعایت کرده اند. همچنین از مجموع 22 مؤلفه ی مربوط به مقوله ی جستجو مؤلفه های «جستجوی مجاورتی»، «ریشه سازی»، «توانایی اصلاح استراتژی جستجوی قبلی»، «محدود گره های مناسب جهت اصلاح جستجو قابلیت جستجوی متن آزاد»، «عملکردهای اصطلاحنامه» و «عملکرد مقالات مرتبط» در محیط رابط کاربری هیچ کدام از پایگاه های

ص: 248

ScienceDirect -1

ProQuest -2

Emerald -3

Ebsco -4

Springer -5

مورد بررسی رعایت نشده است. مؤلفه های «ذخیره جستجوها در حساب کاربری خود» و «تاریخچه جستجو» تنها در پایگاه اطلاعاتی پژوهشگاه علوم و فناوری اطلاعات رعایت شده است. به طور کلی می توان گفت محیط رابط کاربری پایگاه های اطلاعاتی مورد بررسی در این پژوهش به طور میانگین 45/45 درصد از مؤلفه های مربوط به جستجو را رعایت کرده اند (جدول 10)

عکس

سنجش رابط کاربر پایگاه های اطلاعاتی... ۲۴۹

مورد بررسی رعایت نشده است. مؤلفه های «ذخیره جستجوها در حساب کاربری خود» و «تاریخچه جستجو» تنها در پایگاه اطلاعاتی پژوهشگاه علوم و فناوری اطلاعات رعایت شده است. به طور کلی می توان گفت محیط رابط کاربری پایگاه های اطلاعاتی مورد بررسی در این پژوهش به طور میانگین ۴۵/۴۵ درصد از مؤلفه های مربوط به جستجو را رعایت کرده اند (جدول ۱۰).

جدول ۱. وضعیت رابط کاربری پایگاه های اطلاعاتی از لحاظ مقوله جستجو

مؤلفه های جستجو پایگاه های اطلاعاتی	مرکز منطقه ای	پژوهشگاه علوم و فناوری اطلاعات	نورمگز	مگ ایران	نماتن	سید	میانگین
جمع امتیازها	۱۰	۱۰	۱۱	۱۰	۸	۱۱	---
درصد	۴۵/۴۵	۴۵/۴۵	۵۰	۴۵/۴۵	۳۶/۳۶	۵۰	۴۵/۴۵

در پاسخ به پرسش دوم پژوهش می توان این طور بیان نمود که از لحاظ مقوله نمایش، محیط های رابط کاربری پایگاه پژوهشگاه علوم و فناوری اطلاعات با رعایت ۱۸ مؤلفه (۸۱/۸۱ درصد) بیشترین میزان رعایت و پایگاه مگ ایران با رعایت ۱۰ مؤلفه (۴۵/۴۵ درصد) کمترین میزان رعایت مؤلفه های مربوط به مقوله نمایش را داشته اند. همچنین از میان مؤلفه های مربوط به مقوله نمایش اطلاعات مؤلفه «نمایش آمار تعداد دانلود هر مقاله» در رابط کاربری هیچ کدام از پایگاه ها رعایت نشده است. به طور کلی یافته ها حاکی از آن است به طور میانگین پایگاه ها از لحاظ مقوله نمایش ۶۴/۳۸ درصد معیارها را رعایت کرده اند. نتایج در جدول ۲ قابل مشاهده است.

جدول ۲. وضعیت رابط کاربری پایگاه های اطلاعاتی از لحاظ مقوله نمایش اطلاعات

مؤلفه های نمایش پایگاه های اطلاعاتی	مرکز منطقه ای	پژوهشگاه علوم و فناوری اطلاعات	نورمگز	مگ ایران	نماتن	سید	میانگین
جمع امتیازها	۱۵	۱۸	۱۶	۱۰	۱۲	۱۴	---
درصد	۶۷/۱۸	۸۱/۸۱	۷۲/۷۲	۴۵/۴۵	۵۴/۵۴	۶۳/۶۳	۶۸/۳۸

پاسخ به پرسش سوم پژوهش: همان طور که از جدول ۳ مشخص است پایگاه های مگ ایران و نورمگز با رعایت ۵ مؤلفه (۶۲/۵۰ درصد)، در میزان هماهنگی با مؤلفه های مقوله خدمات در رتبه نخست قرار دارند. پایگاه نماتن و پژوهشگاه علوم و فناوری اطلاعات نیز با رعایت ۳ مؤلفه (۳۷/۵۰ درصد)، پایگاه مرکز منطقه ای با رعایت ۲ مؤلفه (۲۵ درصد)، پایگاه سید به ترتیب در رتبه های دوم، سوم و

جدول ۱. وضعیت رابط کاربری پایگاه های اطلاعاتی از لحاظ مقوله جستجو

در پاسخ به پرسش دوم پژوهش می توان این طور بیان نمود که از لحاظ مقوله نمایش محیط های رابط کاربری پایگاه پژوهشگاه علوم و فناوری اطلاعات با رعایت 18 مؤلفه (81/81 درصد) بیش ترین میزان رعایت و پایگاه مگ ایران با رعایت 10 مؤلفه (45/45) درصد کم ترین میزان رعایت مؤلفه های مربوط به مقوله نمایش را داشته اند همچنین از میان مؤلفه های مربوط به مقوله نمایش اطلاعات مؤلفه ی «نمایش آمار تعداد دانشجو هر مقاله در رابط کاربری هیچ کدام از پایگاه ها رعایت نشده است. به طور کلی یافته ها حاکی از آن است به طور میانگین پایگاه ها از لحاظ مقوله نمایش 64/38 درصد معیارها را رعایت کرده اند. نتایج در جدول 2 قابل مشاهده است

جدول 2. وضعیت رابط کاربری پایگاه های اطلاعاتی از لحاظ مقوله نمایش اطلاعات

پاسخ به پرسش سوم پژوهش همان طور که از جدول 3 مشخص است پایگاههای مگ ایران و نورمگز با رعایت 5 مؤلفه (62/50 درصد) در میزان هماهنگی با مؤلفه های مقوله خدمات در رتبه نخست قرار دارند پایگاه نما متن و پژوهشگاه علوم و فناوری اطلاعات نیز با رعایت 3 مؤلفه (37/50 درصد)، پایگاه مرکز منطقه ای با رعایت 2 مؤلفه (25 درصد) پایگاه سید به ترتیب در رتبه های دوم، سوم و

ص: 249

چهارم قرار گرفتند همچنین از میان مؤلفه های مقوله خدمات مؤلفه های «امکان ساخت پرونده شخصی (شخصی سازی محتوی)» و سهولت خروج از سیستم به طور مساوی با فراوانی 4 و درصد فراوانی 80 بیشترین میزان رعایت را داشته اند مؤلفه های «استخراج اطلاعات کتابشناختی برای نرم افزارهای مدیریت ارجاعات» و «ارائه خدمات آگهی رسانی جاری (آلرت نشریات، آلرت جستجو)» در محیط های رابط کاربری هیچ کدام از پایگاه های اطلاعاتی مورد بررسی رعایت نشده است. میانگین امتیاز پایگاه ها در این مقوله نیز 37/49 درصد به دست آمد (جدول 3)

عکس

۲۵۰ مدیریت منابع اطلاعاتی وب

چهارم قرار گرفتند. همچنین از میان مؤلفه های مقوله خدمات، مؤلفه های «امکان ساخت پرونده شخصی (شخصی سازی محتوی)» و «سهولت خروج از سیستم» به طور مساوی با فراوانی ۴ و درصد فراوانی ۸۰ بیشترین میزان رعایت را داشته اند. مؤلفه های «استخراج اطلاعات کتابشناختی برای نرم افزارهای مدیریت ارجاعات» و «ارائه خدمات آگهی رسانی جاری (آلرت نشریات، آلرت جستجو)» در محیط های رابط کاربری هیچ کدام از پایگاه های اطلاعاتی مورد بررسی رعایت نشده است. میانگین امتیاز پایگاه ها در این مقوله نیز ۳۷/۴۹ درصد به دست آمد (جدول ۳).

جدول ۳. وضعیت رابط کاربری پایگاه های اطلاعاتی از لحاظ مقوله خدمات

مؤلفه های خدمات	میانگین	سید	نماتن	مگ ایران	نورمگز	پژوهشگاه علوم و فناوری اطلاعات	مرکز منطقه ای	پایگاه های اطلاعاتی
جمع امتیازها	---	۰	۳	۵	۵	۳	۲	۲
درصد	۳۷/۴۹	۰	۳۷/۵۰	۶۲/۵۰	۶۲/۵۰	۳۷/۵۰	۲۵	۲۵

در پاسخ به سؤال چهارم می توان گفت پایگاه های مرکز منطقه ای، مگ ایران، نماتن و سید هر یک با رعایت ۳ مقوله (۶۰ درصد) مشترکاً در رتبه نخست و پایگاه های پژوهشگاه علوم و فناوری اطلاعات و نورمگز به طور مساوی با رعایت ۲ مقوله (۴۰ درصد) در رتبه دوم قرار دارند. پایگاه ها از لحاظ مقوله پیوندها به طور میانگین ۵۲/۳۳ درصد از مؤلفه ها رادر محیط رابط کاربری خود رعایت کرده اند. نتایج در جدول ۴ آمده است.

جدول ۴. وضعیت رابط کاربری پایگاه های اطلاعاتی از لحاظ مقوله پیوندها

مؤلفه های پیوندها	میانگین	سید	نماتن	مگ ایران	نورمگز	پژوهشگاه علوم و فناوری اطلاعات	مرکز منطقه ای	پایگاه های اطلاعاتی
جمع امتیازها	---	۳	۳	۳	۲	۲	۳	۳
درصد	۳۷/۴۹	۶۰	۶۰	۶۰	۴۰	۴۰	۶۰	۶۰

پاسخ سؤال پنجم: همان گونه که از جدول ۵ مشاهده می شود پایگاه مرکز منطقه ای با رعایت ۴ مقوله (۶۶/۶۶ درصد) در رتبه نخست، پایگاه های نورمگز، مگ ایران و سید با رعایت ۳ مقوله مشترکاً در رتبه دوم، پایگاه پژوهشگاه علوم و فناوری اطلاعات با رعایت ۲ مقوله (۳۳/۳۳ درصد) در رتبه سوم و پایگاه نماتن تنها با رعایت ۱ مقوله (۱۶/۶۶ درصد) در رتبه چهارم قرار دارد. همچنین از لحاظ مقوله راهنمایی مؤلفه های «پیام های اختطاری قابل درک» و «نقشه سایت» به ترتیب با درصدهای فراوانی ۱۰۰ و

جدول 3. وضعیت رابط کاربری پایگاه های اطلاعاتی از لحاظ مقوله خدمات

در پاسخ به سؤال چهارم می توان گفت پایگاه های مرکز منطقه ای مگ، ایران، نامتن و سید هر یک با رعایت مقوله (60 درصد) مشترکاً در رتبه نخست و پایگاه های پژوهشگاه علوم و فناوری اطلاعات و نورمگز به طور مساوی با رعایت 2 مقوله (40 درصد) در رتبه دوم قرار دارند پایگاه ها از لحاظ مقوله پیوندها به طور میانگین 53/33 درصد از مؤلفه ها را در محیط رابط کاربری خود رعایت کرده اند. نتایج در جدول 4 آمده است.

جدول 4. وضعیت رابط کاربری پایگاه های اطلاعاتی از لحاظ مقوله پیوندها

پاسخ سؤال پنجم همان گونه که از جدول 5 مشاهده می شود پایگاه مرکز منطقه ای با رعایت 4 مقوله (6666) درصد در رتبه نخست پایگاه های نورمگز مگیران و سید با رعایت 3 مقوله مشترکاً در رتبه دوم پایگاه پژوهشگاه علوم و فناوری اطلاعات با رعایت 2 مقوله (33/33 درصد) در رتبه سوم و پایگاه نامتن تنها با رعایت 1 مقوله (1666) (درصد در رتبه چهارم قرار دارد هم چنین از لحاظ مقوله راهنمایی مؤلفه های پیام های اختطاری قابل درک و نقشه سایت به ترتیب با درصدهای فراوانی 100 و

ص: 250

بیشترین و کمترین میزان رعایت را در رابط کاربری پایگاه‌های مورد بررسی داشته‌اند. به طور میانگین 44/44 درصد مؤلفه‌های مربوط به مقوله راهنمایی توسط پایگاه‌های مورد بررسی رعایت شده‌اند.

عکس

سنجش رابط کاربری پایگاه‌های اطلاعاتی... ۲۵۱

• بیشترین و کمترین میزان رعایت را در رابط کاربری پایگاه‌های مورد بررسی داشته‌اند. به طور میانگین 44/44 درصد مؤلفه‌های مربوط به مقوله راهنمایی توسط پایگاه‌های مورد بررسی رعایت شده‌اند.

جدول ۵. وضعیت رابط کاربری پایگاه‌های اطلاعاتی از لحاظ مقوله راهنمایی

مؤلفه‌های راهنمایی پایگاه‌های اطلاعاتی	مرکز منطقه‌ای	پژوهشگاه علوم و فناوری اطلاعات	نورمگز	مگ‌ایران	نماتن	سید	میانگین
جمع امتیازها	۴	۲	۳	۳	۱	۳	---
درصد	۶۶/۶۶	۳۳/۳۳	۵۰	۵۰	۱۶/۶۶	۵۰	۴۴/۴۴

پاسخ به پرسش ششم پژوهش: همان‌گونه که از جداول ۱ تا ۵ ملاحظه می‌گردد بیشترین امتیاز و درصد فراوانی میزان رعایت مقوله‌های مورد بررسی در محیط رابط کاربری پایگاه‌های اطلاعاتی مرکز منطقه‌ای، پژوهشگاه علوم و فناوری اطلاعات، نورمگز و سید به ترتیب با امتیازهای ۱۵، ۱۸، ۱۶، ۱۴ و درصدهای ۶۷/۱۸، ۸۱/۸۱، ۷۲/۷۲ و ۶۳/۶۳ مربوط به مقوله‌ی نمایش و می‌باشد. همچنین بیشترین امتیاز و درصد فراوانی میزان رعایت مقوله‌های مورد بررسی در محیط رابط کاربری پایگاه‌های اطلاعاتی مگ‌ایران و نماتن به‌طور مساوی با ۳ امتیاز و ۶۰ درصد مربوط به مقوله پیوندها است.

پاسخ به پرسش هفتم پژوهش: همان‌طور که در جدول ۶ مشاهده می‌شود محیط رابط کاربری پایگاه اطلاعاتی پژوهشگاه علوم و فناوری اطلاعات در مجموع با میزان رعایت ۳۵ مؤلفه (۵۵/۵۵ درصد) بیشترین میزان رعایت مقوله‌های موجود در سیاهه‌وارسی را داشته است و رتبه نخست را به خود اختصاص داده است. سپس محیط‌های رابط کاربری پایگاه‌های نورمگز، مرکز منطقه‌ای، سید و مگ‌ایران به‌طور مساوی و نماتن در رتبه‌های دوم تا چهارم قرار داشتند.

جدول ۶. وضعیت کلی رابط کاربری پایگاه‌های اطلاعاتی

مؤلفه کلی پایگاه‌های اطلاعاتی	۱ بسیار زیاد	۲ زیاد	۳ متوسط	۴ کم	۵ بسیار کم	مجموع	درصد
مرکز منطقه‌ای	۸	۱۵	۲	۳	۴	۳۲	۵۰/۷۹
پژوهشگاه علوم و فناوری اطلاعات	۱۰	۱۶	۲	۲	۳	۳۵	۵۵/۵۵
نورمگز	۱۱	۱۶	۵	۲	۳	۳۷	۵۳/۹۶
مگ‌ایران	۱۰	۱۰	۵	۳	۳	۳۱	۴۹/۲۰
نماتن	۸	۱۲	۳	۳	۱	۲۷	۴۲/۸۵
سید	۱۱	۱۴	۰	۳	۳	۳۱	۴۹/۲۰

جدول ۵. وضعیت رابط کاربری پایگاه‌های اطلاعاتی از لحاظ مقوله راهنمایی

پاسخ به پرسش ششم پژوهش: همان گونه که از جداول 1 تا 5 ملاحظه می گردد بیش ترین امتیاز و درصد فراوانی میزان رعایت مقوله های مورد بررسی در محیط رابط کاربری پایگاه های اطلاعاتی مرکز منطقه ای، پژوهشگاه علوم و فناوری اطلاعات نورمگز و سید به ترتیب با امتیازهای 15، 18، 16، 14 درصدهای 18/18، 81/81، 72/72 و 63/63 مربوط به مقوله ی نمایش و می باشد. همچنین بیشترین امتیاز و درصد فراوانی میزان رعایت مقوله های مورد بررسی در محیط رابط کاربری پایگاه های اطلاعاتی مگ ایران و نما متن به طور مساوی با 3 امتیاز و 60 درصد مربوط به مقوله پیوندها است.

پاسخ به پرسش هفتم پژوهش: همان طور که در جدول 6 مشاهده می شود محیط رابط کاربری پایگاه اطلاعاتی پژوهشگاه علوم و فناوری اطلاعات در مجموع با میزان رعایت 35 مؤلفه (55/55 درصد) بیش ترین میزان رعایت مقوله های موجود در سیاهه واری را داشته است و رتبه نخست را به خود اختصاص داده است. سپس محیط های رابط کاربری پایگاههای نورمگز مرکز منطقه ای، سید و مگایران به طور مساوی و نامتن در رتبه های دوم تا چهارم قرار داشتند.

جدول 6. وضعیت کلی رابط کاربری پایگاه های اطلاعاتی

ص: 251

این پژوهش با هدف سنجش رابط کاربری پایگاه های اطلاعاتی پیوسته مجلات تمام متن فارسی انجام شد. نتایج نشان می دهد که از لحاظ مقوله ی جستجو پایگاه ها در سطح مشابه ای می باشند و در حد متوسط (میانگین 45/45 درصد) معیارها را رعایت کرده اند. از لحاظ مقوله ی نمایش نیز پایگاه ها در وضعیت خوبی میانگین 64/38) قرار دارند از این لحاظ رابط کاربری پایگاه پژوهشگاه علوم و فناوری اطلاعات با میزان رعایت 81/81 درصد و پایگاه مگ ایران با میزان رعایت 45/45 درصد به ترتیب بیش ترین و کم ترین میزان هماهنگی را با مقوله های موجود در سیاهه واری داشته اند. رابط کاربری پایگاه ها در میزان رعایت مقوله خدمات کم ترین هماهنگی را نسبت به سایر مقوله ها داشته اند؛ به طوری که میانگین رعایت معیارهای این مقوله 37/49 درصد و در حد ضعیف محاسبه شده است. در ارزیابی میزان رعایت مؤلفه های مقوله پیوندها نیز رابط کاربری پایگاه ها در سطح متوسطی قرار داشتند و به طور میانگین 53/33 درصد معیارها را رعایت نموده اند از لحاظ میزان رعایت مؤلفه های مقوله راهنمایی به جز پایگاه نما متن که تنها 16/66 درصد معیارهای این مقوله را رعایت نموده است؛ رابط کاربری سایر پایگاه های مورد بررسی در وضعیت مشابه و متوسطی (میانگین 44/44 درصد) قرار داشتند.

همچنین نتایج دیگری که از این پژوهش حاصل شد این است که پایگاه اطلاعاتی «پژوهشگاه علوم و فناوری اطلاعات» با بیشترین میزان رعایت، مؤلفه ها برترین محیط رابط کاربری را از بین محیط رابط کاربری پایگاه های اطلاعاتی پیوسته مجلات تمام متن فارسی دارا می باشد در نهایت می توان اظهار نمود که محیط رابط کاربری پایگاههای مورد بررسی از لحاظ میزان رعایت معیارهای موجود در سیاهه واری در حد «متوسطی» می باشند و این نتیجه یافته های پژوهش های انتظاریان و فتاحی (1389)، اعظمی و فتاحی (1388) و اسداللهی و نوکاریزی (1389) را تأیید می نماید به طور کلی بهتر است هر شش پایگاه نسبت به گنجاندن آن دسته از معیارهای مطرح شده در سیاهه واری که فاقد آن هستند، اقدام کنند.

در نهایت پژوهشگران بر اساس نتایجی که از ارزیابی رابط کاربری پایگاه ها به دست آمد به طراحان محیط رابط کاربری پایگاه های اطلاعاتی مجلات تمام متن فارسی پیشنهاد می دهند که جهت استفاده بهتر، افزایش موفقیت در جستجوها و انجام بهتر و سریعتر بازیابی اطلاعات از طریق محیط رابط کاربری و در نتیجه جذب کاربران بیشتر برخی نکات را که در محیط رابط کاربری آن ها به ندرت استفاده شده یا اصلاً استفاده نشده است رعایت کنند. این نکات عبارتند از: در مقوله جستجو به طراحی آیتم های «محدودگرهای مناسب»، «عملکرد اصلاح نامه»، «عملکرد مقالات مرتبط»، «عملکرد پیشنهاد کلیدواژه های مرتبط»، «ذخیره جستجوها در حساب کاربری» و «تاریخچه جستجو»؛ در زمینه نمایش اطلاعات به طراحی آیتم های «نشانه دار کردن نتایج» و «ارسال نتایج به ایمیل»؛ از لحاظ مقوله خدمات پایگاه ها «امکان استخراج اطلاعات کتابشناختی به نرم افزارهای مدیریت ارجاعات» «نشان دادن مقالات استناد کننده» و «خدمات آگاهی رسانی جاری»؛ در زمینه پیوندها نیز «امکان فرایبوند به عناوین مشابه» و «فرایبوند به دیگر مقالات یک نویسنده» پردازند. بنابراین پیشنهاد می شود که طراحان پایگاه ها به طراحی مجدد پایگاه های خود پرداخته و به ویژه خدمات جدید معرفی شده در این پژوهش را در محیط رابط کاربری پایگاه های اطلاعاتی خود اعمال نمایند.

پیشنهاد هایی برای پژوهش های آینده:

پیشنهاد می شود پایگاه های مذکور را از نظر جنبه های خاص شامل «دامنه پوشش موضوعی و زمانی مجلات»، «نحوه سازمان دهی اطلاعات» و غیره مورد پژوهش قرار دهند همچنین مطالعاتی نیز بر پایگاه های تهیه شده در کشور جز پایگاه های مجلات از قبیل پایگاه های اطلاعات کتابشناختی پایگاه های مقالات کنفرانس ها انجام شود.

منابع

اسد اللهی، زهرا، محسن نو کاریزی. 1389. ارزیابی ساختار و محتوای پایگاه های اطلاعاتی الکترونیکی نشریات ایرانی کتابداری و اطلاع رسانی. 13(2)

در [http://www.aqlibrary.ir/index.php?module=TWArticlesfile=indexfunc=view&pubarticlesdid=882pid=10.\(91/5/3](http://www.aqlibrary.ir/index.php?module=TWArticlesfile=indexfunc=view&pubarticlesdid=882pid=10.(91/5/3) - (دسترسی در

اعظمی، محمد رحمت الله فتاحی 1388. تطابق رابط گرافیکی کاربر پایگاه های اطلاعاتی با مدل رفتار ابی الیس علوم و فناوری اطلاعات. 25 (2): 247-264.

انتظاریان، ناهید، رحمت الله فتاحی. 1389. مبانی طراحی رابط کاربر مبتنی بر شناخت ویژگی ها، ادراک و رفتار کاربران. کتابداری و اطلاع رسانی. 13 (2).

در [http://www.aqlibrary.ir/index.php?module=TWArticlesfile=indexfunc=view&pubarticlesdid=878pid=10.\(91/4/24](http://www.aqlibrary.ir/index.php?module=TWArticlesfile=indexfunc=view&pubarticlesdid=878pid=10.(91/4/24) (دسترسی در

علیجانی، رحیم، لیلا دهقانی 1386 مقایسه ی رابط کاربر پایگاه های اطلاعاتی کتاب مدار بین المللی فصلنامه کتاب 72: 233-252
مهرداد، جعفر، لیلا- دهقانی 1385 معیارهای ارزیابی رابط های کاربر نسخه های مختلف پایگاه های اطلاعاتی مجله کتابداری بهار و تابستان 77-95.

مهرداد جعفر زهره زاهدی. 1386 بررسی و مقایسه رابط کاربر دو میزبان داخلی کتابخانه منطقه ای علوم و تکنولوژی و پژوهشگاه اطلاعات و مدارک علمی ایران با چهار میزبان خارجی Emerald, Ebsco, Proquest Elsevier کتابداری و اطلاع رسانی 10 (3): 107-124.

یمین فیروز، موسی 1382 ویژگی ها و عناصر تشکیل دهنده رابط کاربر در وب سایت ها فصلنامه کتاب. 14: 159-168

Cracstini, F. 2004. A Graphical user interface for the retrieval of hierarchically structured documents. Information Processing Management. 40 (2): 269-289

Hansen, P. 1998. Evaluation of IR user interface implications for user interface design. www.)

Hb.Se/bhs/ith/2-98/ph.htm (Available at 12/6/2012

Hirst, S.J. 1999. HyperLib Deliverable 2.1.1: The Use of Icons in a Multilingual OPAC Interface. Hyperlib
.Electronic Document Store, University of Antwerp–University of Loughborough

<http://lib.ua.ac.be/MAN/WP211/root.htm>(Available at 8/5/2012)

Sabin–Kidiss, L. 2001. Assessing the functionality of web–based versions of traditional search engines
..Online. 25 (2): 18–24

ص: 253

طبق قانون حق مؤلف کشور فرانسه مصوب اول اوت 2006، کتابخانه ملی فرانسه (BNF از این پس کتابخانه) وظیفه گردآوری و حفاظت اینترنت فرانسه را بر عهده دارد. این کتابخانه «مدل تلفیقی» را برای آرشیو وب ایجاد کرده است. این مدل، متشکل از خزش های فراگیر دامنه .fr و خزش های کانونی و واسپاری های الکترونیکی است.

کتابخانه ملی فرانسه، به برکت همکاری پژوهشی با آرشیو اینترنت از سال 2004 هر سال به اجرای چهار خزش فراگیر مبادرت کرده است. آخرین آن ها، با ویژگی های متفاوت چشمگیری ایجاد شده است و از مهم ترین این ویژگی ها کاربرد فهرست همه جانبه ای از اسامی دامنه های .fr، بود که توسط آفنیك انجمن همکاری نام گذاری اینترنت، فرانسه برای ثبت (.fr) بعد از امضای توافق نامه ای میان دو سازمان در سپتامبر 2007 به کتابخانه ملی فرانسه ارائه شد.

گزینش های ماهرانه قبل و حین، خزش نقش تعیین کننده های در شکل آینده مجموعه خواهند داشت، بنابراین تصمیم هایی که باید بر طبق ساختار قانونی و معنوی در مدت انجام خزش اتخاذ شوند عبارت اند از: برای بیان اف این مجموعه شامل 5 قرن سنت گذشته و اسپاری قانونی است. برای تعیین پیامدها و نتایج راه حل های فنی موجود، قصد داریم نتایج آخرین خزش این کتابخانه را تحلیل و با برداشت های سال های گذشته مقایسه کنیم به علاوه این مطالعات سودمندی تلاش های ما را برای توصیف وب 2007 فرانسه تأیید می کند.

کلید واژه ها: آرشیو سازی، وب قانون و اسپاری اینترنت کتابخانه ملی فرانسه بیان اف، کنسرسیوم بین المللی حفاظت، اینترنت IIPC آرشیو اینترنت میراث دیجیتال

اشاره

راهبردهایی برای گردآوری دامنه ملی (1)

نوشته: فرانس لاس فارگوس کلمنت کیوری، برت وندلاند (2)

ترجمه: سودابه نوذری (3)

1. قلمرو فرانسه

1,1 تعریف حوزه قانون و اسپاری

در اول اوت 2006، قانون حق مؤلف جدید در مجلس فرانسه به رأی گذاشته شد. مفاد این قانون، که کتابخانه مدت ها انتظارش را می کشید، قانون واسپاری را به اینترنت کشاند. قانون واسپاری، به هر ناشر تکلیف می کند نسخه هایی از تولیداتش را به کتابخانه ملی ارسال کند این، قانون که نخستین بار در 1537 برای منابع چاپی تصویب شده بود با گذشت قرن ها به اشکال متفاوت تولیدات فرهنگی انتشاراتی جدید از منقوشات گرفته تا نرم افزارها و بازی های ویدئویی تعلق گرفت به موازات گسترش شبکه جهانی وب، به عنوان مکانی مطلوب برای ایجاد دانش و اطلاعات لازم بود برای نهادهای حافظ میراث ملی فرانسه، چارچوبی قانونی برای سازماندهی منابع مورد حفاظت شان ارائه شود.

قانون فوق درباره قلمرو اینترنت فرانسه شفاف نیست و برای روشن شدن این مسئله در آینده ای

ص: 255

1- Legal deposit of the French Web: harvesting strategies for a national domain

2- France Lasfargues, Clément Oury, and Bert Wendland

3- عضو هیئت علمی سازمان اسناد و کتابخانه ملی ایران

نزدیک صدور فرمانی انتظار می رود در این اثنا، کتابخانه خط مشی خود را درباره دامنه ملی برخط به پیش می برد که احتمالاً با، فرمان زمانی که تأیید شود سازگار باشد، کتابخانه به منظور پالایش حوزه سیاست آرشیوسازی وب خود پنج قرن و اسپاری قانونی و تجربه جدیدتر چالش های فنی گردآوری وب را با یکدیگر تلفیق کرده است بنابراین روش ما هم عملی است (یعنی موافق با ابزار اکتشافی برداشت حجمی) و هم سازگار (با قانون و اسپاری گذشته فرانسه).

این سنت بر اساس سه معیار قرار داشته است:

- انتشارات: اسناد گردآوری شده برای یک مخاطب و در خدمت یک هدف عمومی هستند، آن ها نباید در حد حوزه های خصوصی یا ارتباطات درون شرکتی تنزل کنند؛

- رسانه: همه ترتیبات قانون و اسپاری قبلی بر اساس موجودیت فیزیکی یک رسانه قرار داشتند: منابع، چاپی متون موسیقی (1) تصاویر نوارهای صدا دیسکت و جز آن؛ و

- قلمرو: اسناد باید در محدوده مرزهای قلمرو ملی منتشر یا توزیع می شدند.

به طور خلاصه، قانون و اسپاری گذشته برای تمام انتشارات موجود در رسانه های تولید شده یا توزیع شده در قلمرو کشور فرانسه قابل اجرا بوده است.

متأسفانه هیچ یک از این معیارها برای گردآوری [منابع] وبی به راحتی قابل اجرا نیست، زیرا وب گاه ها:

- به ترکیب و ادغام ارتباطات شخصی و عمومی و ترکیب مبتکرانه آن ها تمایل دارند؛

- به معنای دقیقه کلمه، رسانه نیستند، بلکه بیش تر سکویی هستند که سایر رسانه ها تمایل به انتقال به آن ها دارند (می توان کتاب، تصاویر، فیلم، متون موسیقی، و جز آن را در وب یافت)؛ و

- به سهولت، به یک سرزمین خاص - دست کم در مقیاسی وسیع - محدود نمی شوند؛

به علاوه، دیگر نمی توان زبان فرانسه را معیار انتخاب قرار داد زیرا قانون و اسپاری بدون توجه به زبان انتشاراتی اعمال می شود: مجموعه های و اسپاری شده به کتابخانه ملی فرانسه شامل اقلامی به زبان های خارجی می شود که در کشور فرانسه منتشر چاپ یا توزیع شده باشند به همین ترتیب، نمی توانیم از راهبرد قانون و اسپاری خود انتظار تمرکز بر موضوع نویسنده یا سطوح انتشاراتی خاص داشته باشیم جنبه مهم این قانون آن است که مجموعه ها باید منعکس کننده جامعه و فرهنگ فرانسه به هر شکل و صرف نظر از ارزش علمی یا محبوبیت آن ها باشند.

اگر مجبور به انتخاب هستیم منظور بیش تر تهیه نمونه هایی از منابع است تا انتخاب نسل های بعدی تصمیم می گیرند که چه چیزی ارزشمند است نه کتابخانه ملی در قفسه های کتابخانه، نویسندگان ناشناخته و مجله های مستهجن در کنار متفکران بزرگ قرار گرفته اند، انتظار داریم این فلسفه در مورد منابع و بی هم اعمال شود برداشت، فل های فرصت بزرگی را برای عرضه این رویکرد در مقیاس وب فراهم می کند قانون و اسپاری درباره محتوا و قالب - یا محمل - است به این مفهوم که کتابخانه ملی فرانسه به ساخت مجموعه هایی که

منعکس کننده گرایش هایی از نمونه های انتشاراتی است توجه می کند تلاش می کنیم طیف گسترده اشیا و قالب هایی را برای ارائه فراهم کنیم که به طور عملی توسط افراد

ص: 256

scores -1

بهره مند از اطلاعات تولید می شود.

در نتیجه، برای دامنه ملی خود نیازمند تعریفی بودیم که با وجود انعطاف پذیری و سادگی اجرا منعکس کننده «روح» این گذشته باشد. تعریف یک «کانون» فرانسه در عین داشتن قابلیت انعطاف پذیری، تنها راه تسهیل شیوه های گردآوری اکتشافی در مقیاس کلان بود ثابت شده است، زبان، مکان جغرافیایی اسامی یا موضوع برای تعیین تفاوت در مقیاس کلان بسیار بی ربط و چالش برانگیز هستند، بنابراین، آن چه به عنوان نقطه آغاز احتمالی برای اکتشاف دامنه ملی باقی ماند استفاده از دامنه سطح بالای ملی (تی. ال دی) Top Level domain (1) (TLD) بود این دامنه به سبب تهیه فهرست آغازین اصلی برای اکتشاف با سایر راهبردها تلفیق شده است. از این رو خط مشی واسپاری برخط ما حین انتظار آشکار شدن در حکم صادر شده، به شکل زیر تعریف شد. ما «فرانسه» را چنین در نظر گرفتیم:

- به عنوان یک اصل هر وبگاهی که در دامنه سطح بالای fr. یا هر دامنه سطح بالای مشابه دیگر، با ارجاع به قلمرو رسمی فرانسه ثبت می شود (برای، نمونه re برای جزیره فرانسوی لاری نیو (2))؛

- هر وب گاهی (شاید خارج از فرانسه) که تولید کننده اش در سرزمین فرانسه قرار دارد (معمولاً این مسئله را می توان در وب گاهی یا با استفاده از (3) Domain Name System تعیین کرد)؛

- هر وبگاهی (شاید خارج از فرانسه) که بتواند برای نمایش محتواهای تولید شده در سرزمین فرانسه تأیید شده باشد (این معیار آخر چالش برانگیزتر از آن است که بررسی شود اما مجالی را برای تفسیر و گفت و گو میان کتابخانه ملی فرانسه و تولیدکنندگان وب به وجود می آورد).

البته انتظار نمی رود هیچ یک از این معیارها قبل از این که ما فهرست هسته (4) را بسازیم، به طور کامل تأمین شوند (این امر مانع اکتشاف می شود و در بررسی وب گاه ها، قبل از خزش (5)، ایجاد اخلاص می کند، که به سادگی قابل اندازه گیری نیست). با وجود، این آن ها قصد شان خدمت در چارچوبی قانونی و فرهنگی است به منظور:

- تعریف و توصیف شیوه کلی سیاست آرشیو سازی وب ملی برای مردم صاحبان محتوا، و کتابداران درگیر در پروژه؛

- تبیین وظایف و دستورالعمل های مورد نیاز برای خود تا وقت نظارت بر خزش ها به یاد داشته باشیم؛ تعیین عناصر عینی تصمیم گیری وقتی ما نیازمندیم بدانیم آیا وب گاهی به طور کامل در حوزه کار ماست یا خیر این امر به ویژه زمانی که وب مستر یک وب گاه خاص از کتابخانه ملی می خواهد کار خزش را متوقف کند یا حتی هنگام رویارویی با یک دعوی، قانونی سودمند است: اگر وب گاه در هیچ یک از معیارهای فوق، ننگنجد باید از خزش های آینده و از مجموعه کنار گذاشته شود.

بنابراین، رویکرد دامنه ملی مصالح های را میان قانون واسپاری گذشته و ویژگی های چالش برانگیز

ص: 257

fr-1

La Ré union-2

DNS-3

seed list -4 (فهرست یو. آر.ال).

crawl -5

وب نشان می دهد، همچنین مصالح های میان رویکرد کاملاً باز و نامعقولانه ای که به احتمال می توانست ما را طوری هدایت کند تا کل شبکه جهانی وب را به عنوان فرانسه بالقوه در نظر بگیریم، یا برعکس، رویکردی محدود کننده، که فقط به دامنه سطح بالای fr تنزل می یافت در حالی که معروف است این دامنه تنها بخش محدودی از وب گاه های فرانسوی را شامل می شود در یک کلام، هدف این رویکرد، تأمین تمرکز و انعطاف پذیری است.

2.1. در حال حاضر کجا هستیم

کتابخانه اجازه دارد راه های متفاوتی را برای گردآوری اینترنت فرانسه به کار ببرد: «مؤسسه های قیم ممکن است منابع اینترنت را با به کارگیری فنون خودکار یا با تنظیم موافقت نامه های خاص و فرآیندهای واسپاری با همکاری تولیدکنندگان گردآوری کنند» (ماده 41 II) به همین منظور کتابخانه ملی فرانسه مدلی تلفیقی متشکل از سه راهبردی زیر را تعریف کرده است.

- برداشت فله ای اینترنت فرانسه هدف گردآوری دامنه fr، دست کم به صورت سالانه است. این خزش های فراگیر (1) امکان آرشیو تصاویر وب فرانسه را به کتابخانه می دهند. این رویکرد، در مقایسه با هزینه های ماشینی و نیروی انسانی برداشت [و حجم] اطلاعات بازیابی، شده باصرفه تر است. با وجود، این به خاطر محدودیت منابع و ابزاری با چنین خزش گر (2) هایی امکان ندارد بتوان وب عمیق (وبگاه های بسیار بزرگ پایگاه های اطلاعاتی و جز آن) را گردآوری کرد.

- خزش های کانونی (3) برای وب گاه ها به تعداد محدود این سایت ها در داخل یا خارج از کتابخانه ملی، فرانسه با همکاری شبکه ای از کتابداران و پژوهش گران کشف شده اند خزش گر های کانونی به وبگاه های بزرگ و به وب گاه های غالباً در حال اصلاح، اختصاص دارند.

- واسپاری های الکترونیکی ویژه شمار محدودی انتشارات الکترونیکی

کتابخانه، منتظر تصویب قانونی که فنون خزش را بررسی کند نشد پروژه آرشیو سازی وب در 1999 آغاز شد. نخستین خزش گر کانونی رویداد محور حدود سال 2002 آزمایش شد در آن زمان کتابخانه نزدیک به دو میلیون وب گاه وابسته به انتخابات کشور (انتخابات ریاست جمهوری و مجلس) را گردآوری کرد این کار برای اروپا و برای انتخابات محلی دو سال تمدید شد کتابخانه ملی فرانسه 1162 وب گاه را گردآوری کرد (4)

با وجود این، تجهیزات فنی (سخت افزاری و نرم افزاری) مهارت ها و تجربه لازم برای تحقق خزش گر های بزرگ مقیاس اینترنت، فرانسه در کتابخانه ملی هنوز کافی نبود این دلیل چرایی موافقت نامه همکاری کتابخانه ملی فرانسه با آرشیو اینترنت (IA)، بنیاد غیرانتفاعی در آرشیوسازی شبکه جهانی وب، از 1996 است. در نوامبر 2004 این دو مؤسسه موافقت نامه تحقیقاتی را با نام «پروژه پژوهشی: انتخاب

ص: 258

broad crawls -1

crawler -2

Focused crawls -3

4- خزش گرہا برای انتخابات 2002 و 2004 استفاده شدند برای اطلاعات بیش تر درباره این دو خزش گر به منبع 19 مراجعه کنید.

دامنه ملی برای آرشیوسازی وب)) امضا کردند هدف این پروژه تعیین شیوه ها و ابزارهایی برای استفاده در یک خزش دامنه ملی وب بود موافقت نامه تصریح کرد برای خزش های فراگیر، لازم است بررسی این ابزارها و شیوه ها توسط آرشیو اینترنت (IA) انجام و داده هایی که طی این خزش ها گردآوری می شد به قفسه های ذخیره سازی بی.ان.اف تحویل شود.

نخستین خزش فراگیر در چارچوب این موافقت، نامه در پایان سال 2004 اتفاق افتاد. از آن، در سال های 2005 2006 و 2007 سه خزش فراگیر fir انجام شد (این آخرین خزش که تا سال 2009 ادامه خواهد داشت مرهون بسط موافقت نامه تحقیقاتی است).

خزش های با مقیاس کوچک تر نیز - مستقیم یا غیر مستقیم - توسط کتابخانه ملی فرانسه اجرا شده است. دو خزش کانونی در شمار محدودی وب گاه (در حدود 4000 عدد) توسط IA برای پروژه تحقیقاتی انجام شد [15] از 2007، این وب گاه ها توسط بی.ان.اف با استفاده از امکانات خود، برداشت می شوند. در همان سال سایر خزش های کانونی موضوعی یا رویداد محور به اجرا درآمدند، مانند وب گاه های مربوط به انتخابات ملی 2007 فرانسه.

به این ترتیب بی.ان.اف تاکنون - مستقیم یا به لطف همکاری با IA- چهار خزش فراگیر به علاوه حجم فراوانی خزش های کانونی به اجرا در آورده است، دیگر آرشیوسازی، وب یک پروژه در کتابخانه ملی فرانسه نیست بلکه فعالیتی روزمره و واحدی دائمی در اداره قانون واسپاری (1) این کتابخانه است. از آن جا که این امر بهترین رویکرد در مواجهه با چالش گردآوری حجم روزافزون اشیای دیجیتال در وب است برداشت فله ای هنوز اولویت نخست به شمار می آید.

با وجود این برداشت فله ای به معنای برداشت کورکورانه نیست حتی وقتی روبات ها سعی دارند حداکثر فایل های وبی را کشف کنند خود را به رعایت قوانین و تنظیمات ملزم می کنند [20]. تصمیمات فنی، قبل حین و بعد از خزش نقش قطعی در نتیجه برداشت دارد.

این مقاله راهبردهایی که بی.ان.اف با همکاری آرشیو اینترنت (IA) برای اجرای خزش های بزرگ مقیاس ارائه کرده است توصیف میکند این خزش ها می تواند با دیدگاه دامنه ملی فرانسه سازگار باشد.

2. طراحی خزش

2.1. هدف چیست؟

به نظر می رسد، اجرای خزش فراگیر در حدود صدها میلیون فایل با طیفی از مشکلات فنی اجتناب ناپذیر همراه باشد. با این حال، نخستین سؤالی که قبل از آغاز یک خزش باید پاسخ داده شود، سؤالی فنی نیست: هدف این خزش گر چیست؟ پاسخ ها بسته به اینکه از یک شرکت بزرگ نهادی پژوهشی، یا مؤسسه ای میراثی باشد متفاوت خواهد بود. اگر لازم است داده ها در مدت طولانی نمایه سازی (توسط موتور جست و جو)، تحلیل (به طور مثال برای شناسایی دامنه های وبی) و یا آرشیو و نمایش داده شوند، خزش فراگیر شیوه یکسانی را اجرا نمی کند.

ص: 259

به عبارت دیگر توجه به محدودیت های فنی، هنگام تعریف اهداف خزش فراگیر ضروری است. این نکته دلیل لزوم گفت و گوی دائمی میان کتابداران (که وظیفه تعریف خط مشی مجموعه سازی را بر عهده دارند) و مهندسان (که وظیفه اجرای خزش ها را بر عهده دارند) را به خوبی بیان می کند متدولوژی ای که در این بخش توصیف می شود نقطه تلاقی دغدغه های مربوط به این دو گروه است

خزش های کتابخانه در چارچوب قانون و اسپاری به اجرا درآمده اند سودمندی این چارچوب تنها برای پرداختن به مسائل حفاظت مالکیت معنوی نیست بلکه آرشیوسازی وب را به عنوان وظیفه ای مستمر مطرح می کند. مطابقت خط مشی مجموعه سازی در آرشیوهای وب، با قالب های قدیمی، انتشارات امری ضروری است.

کشف خودکار وب گاه ها با رویت راه کاری است برای سازگاری با ویژگی «غیر تبعیض آمیزانه» قانون واسپاری، فرانسه تا هم «بهترین» (انتشارات ادبی علمی) و هم «بدترین» (از آگهی ها گرفته تا پورنوگرافی) انتشارات فرانسوی را گردآوری کند با وجود، این حتی رویت ها در معرض پیش داوری هستند. ساختار فرایبندی، وب نخست به کشف و گردآوری پر استناد ترین وب گاه منتهی شد اما با رویت های آرشیوسازی ما وب گاه های کم طرفدار فراموش نمی شوند برای اجتناب از این پیش داوری BNF در سپتامبر 2007، موافقت نامه ای را با انجمن همکاری نامگذاری اینترنت فرانسه (آفنیک) (1)، سازمان مسئول دامنه های fr و re. امضا کرد. طبق شرایط، موافقت نامه این سازمان باید هر دو سال یک بار فهرست کاملی از اسامی دامنه ثبت شده موجود در دامنه های fr و re را ارائه دهد (هم اکنون بیش از یک میلیون اسامی دامنه دارد) به عبارت دیگر BNF باید اعتبار این اطلاعات ارزشمند را تضمین کند.

بنابراین، هدف یک خزش دامنه، بزرگ گردآوری نمونه ای جامع از دامنه ملی و ارائه تصویری از تولیدات [فکری] فرانسه در زمان این گردآوری است در بیش تر موارد نمونه به عنوان یک تصویر در نظر گرفته می شود- راه کاری برای حفظ و تثبیت یک فضای در حال تحول از آن جا که گردآوری همه چیز ممکن نیست به بهای [از دست دادن] یکی برداشت چند سند از هر وب گاه را به گردآوری کل چند وب گاه ترجیح می دهیم

به این سبب خزش های عمیق تری را در مهم ترین وب گاه ها تجربه نکردیم مانند کتابخانه ملی استرالیا که با گذاشتن اولویت بالا در مورد وب گاه های دولتی و دانشگاهی این امر را انجام داد فهرستی از این وب گاه ها توسط کتابداران منتشر شده است [17]) این کار لازم هم نبود زیرا خزش های فراگیر ما با همراهی خزش های کانونی به صورت کامل به اجرا در می آیند.

از سوی دیگر، احتمال تأثیر سریع وب از تحولات فناورانه وجود دارد. قالب های انتشاراتی نوین پدیدار می شوند و در عرض چند ماه گسترش می یابند خزش فراگیر باید بازتاب دهنده این تحولات باشد به عنوان نهاد قانون واسپاری یکی از اهداف، ما روشن کردن قالب های نوین انتشاراتی است و در نتیجه کسب اطمینان از اینکه رویت توانایی برداشت این اسناد را داشته باشد به این سبب، در سال 2006 ، به وب نوشت ها و وب گاه های شخصی توجه بیش تری داشتیم و در سال 2007 بر ویدیو ها تأکید

ص: 260

داریم (در ادامه می آید).

طبق اهداف از پیش تعریف شده در آینده احتمال دارد مجموعه قبل و در زمان خزش شکل گیرد. قبل از آغاز کار، برداشت درباره دو عامل مهم که سهم بزرگی در ساخت مجموعه های آینده دارند تصمیم گیری شد طراحی فهرست هسته، و تنظیمات خزش.

2.2 فهرست هسته

خزش گر وظایفش را با فهرستی از یو. آر. ال آغاز می کند که فهرست هسته نامیده می شود. هسته ها، در هایی برای دستیابی به وب گاه ها به شمار می آیند؛ به این سبب کیفیت برداشت تا حد زیادی به کیفیت این فهرست بستگی دارد.

از آن جا که همکار ما در اجرای خزش فراگیر برای چهار سال همان [آرشیو اینترنت (IA)] باقی ماند، امکان غنی کردن تدریجی فهرست هسته وجود داشت منابع مختلف هسته ها سال به سال افزایش می یافتند:

- 2004: فهرست هسته [یو. آر. ال]، برای نخستین خزش فراگیر از استخراج دامنه های fr. آخرین خزش الکسا به وجود آمد. (1)

- 2005: هسته های حاصل از استخراج دامنه های fr آخرین خزش گر و میزبان آکسا، هنگام خزش فراگیر قبلی کتابخانه ملی فرانسه کشف شدند هدف فرصت دادن به خزش گر برای پیش رفتن و کشف میزبان های جدید بود

- 2006: فهرست هسته به شیوه سال گذشته آغاز شد.

- 2007: به لطف امضای توافق نامه با آفنیک فهرست جامع اسامی دامنه fr و ie. به عنوان فهرست هسته به کار رفت برای اطمینان از تداوم و سازگاری با خزش های فراگیر، قبلی این فهرست با استخراج هایی از خزش آکسا و از خزش های میزبان قبلی ادغام شد.

آفنیک شامل:

- 890064 دامنه fr

- 1516 دامنه re، و

- 21 / 753 اسم دامنه سطح دوم دامنه سطح دوم (2, 3, 3) را ببیند.

با وجود این استفاده از فهرست، آفنیک به عنوان فهرست هسته، ساده امکان پذیر نبود این فهرست شامل اسامی دامنه می شد نه یو آر ال آن ها به عبارتی پشت نام هر دامنه یک وب گاه وجود نداشت

بنابراین در برخورد با فهرست، آفنیک چند تحلیل مورد پردازش قرار گرفتند.

ص: 261

هایی درباره سایت های مورد بررسی فراهم کند و بر اساس داده های گردآمده از دیگر کاربران ، صفحه های مرتبطی را که ممکن است مورد علاقه آن ها باشند، توصیه می کند این شرکت از ، 1996 آرشیوهای خزش خود را در اختیار آرشیوهای اینترنتی قرار داده است [16]

نخستین آن‌ها کمی کردن تعداد دامنه‌هایی بود که هنوز فعال بودند و باید به‌طور عملی، به‌عنوان هسته، مورد استفاده قرار می‌گرفتند. این بررسی سنگین توسط آرشیو اینترنت (IA) انجام شده است. برای کسب اطمینان از برخط بودن پاسخ آن‌ها هر دامنه با دو نشانی مختلف بررسی می‌شد (یعنی بررسی 1/780/128 یو.آر.ال):

<http://domainname.fr> و <http://www.domainname.fr>.

عکس

۲۶۲ مدیریت منابع اطلاعاتی وب

نخستین آن‌ها، کمی کردن تعداد دامنه‌هایی بود که هنوز فعال بودند و باید به‌طور عملی، به‌عنوان هسته، مورد استفاده قرار می‌گرفتند. این بررسی سنگین، توسط آرشیو اینترنت (IA) انجام شده است. برای کسب اطمینان از برخط بودن پاسخ آن‌ها، هر دامنه با دو نشانی مختلف بررسی می‌شد (یعنی بررسی ۱/۷۸۰/۱۲۸ یو.آر.ال):

<http://domainname.fr> و <http://www.domainname.fr>.

هر دو صورت دامنه دارای دی.ان.اس است	۷۹ درصد
یک صورت دامنه دارای دی.ان.اس است	۱۴ درصد
هیچ یک از دو صورت دی.ان.اس ندارد	۷ درصد

شکل ۱. پاسخ به وجود دی.ان.اس^۱ در اسامی دامنه آف‌نیک

بر اساس نتایج، فهرستی از یو.آر.ال‌های معتبر به دست آمد، تا آنجا که، اگر نسخه WWW پاسخ نمی‌داد، از «domainname.fr»، یا <http://www.domainname.fr> یا <http://domainname.fr> استفاده می‌شد. هسته‌های بدون دی.ان.اس در فهرست هسته به‌طور تصادفی وارد می‌شدند. ما اسامی دامنه‌های (و به‌ویژه در مجموعه خزش فراگیر ۲۰۰۶) موجود در فهرست آف‌نیک را که در آرشیوهای وب نیز موجود بودند، بررسی کردیم. فهمیدن این نکته که فقط ۳۰ درصد از دامنه‌های آف‌نیک در مجموعه ما وجود داشت، بسیار حائز توجه بود.

این امر شاید به افزایش چشمگیر اندازه دامنه .fr مربوط باشد. به یمن ساده شدن پی در پی قوانین مختص .fr، از سال ۲۰۰۴، به‌طور مداوم افزایش یافته است. عمده‌ترین آنها، اجازه ایجاد دامنه سطح بالا برای اشخاص، در سال ۲۰۰۶ بود (تا این تاریخ، تنها اداره‌ها و انجمن‌های خصوصی مجوز ایجاد دامنه .fr را داشتند): با گذشت یکسال، .fr، بیش از ۶۳ درصد افزایش یافت. این امر باید مربوط به ایجاد تصور مثبت از ccTLD در کاربران اینترنت باشد. اشخاص ۳۰ درصد دامنه‌های .fr ثبت شده، و ۵۰ درصد ثبت‌نام‌های جدید را تشکیل می‌دهند [۲].

همچنین، وجود پدیده فوق، در فهرست آف‌نیک، ممکن است تاحدی به خاطر حضور اسامی دامنه‌هایی باشد که یا پیوندی به وبگاه‌های دیگر ندارند و یا پیوندشان ضعیف است، و توسط رویات‌ها در سال‌های گذشته کشف نشده بودند.

این رشد فوق‌العاده، تا حد زیادی گسترش چشمگیر فهرست هسته را از ۲۰۰۶ تا ۲۰۰۷ توصیف می‌کند.

۱. DNS [سروری که دامنه‌های آدرس سایت‌ها را پشتیبانی می‌کند، مترجم].

شکل 1. پاسخ به وجود دی. ان. اس (1) در اسامی دامنه آفنیک

بر اساس نتایج، فهرستی از یو. آر. ال های معتبر به دست آمد تا آن جا که اگر نسخه www پاسخ نمی داد، از «domainname.fr» یا http://www.domainname.fr یا http://domainname.fr استفاده می شد.

هسته های بدون دی. ان. اس در فهرست هسته به طور تصادفی وارد می شدند.

ما اسامی دامنه های (و به ویژه در مجموعه خزش فراگیر 2006) موجود در فهرست آفنیک را که در آرشیوهای وب نیز موجود بودند بررسی کردیم فهمیدن این نکته که فقط 30 درصد از دامنه های آفنیک در مجموعه ما وجود داشت بسیار حائز توجه بود

این امر شاید به افزایش چشمگیر اندازه دامنه .fr مربوط باشد به یمن ساده شدن پی در پی قوانین مختص .fr. از سال 2004 .fr به طور مداوم افزایش یافته است. عمده ترین، آن ها اجازه ایجاد دامنه سطح، بالا برای اشخاص در سال 2006 بود (تا این، تاریخ تنها اداره ها و انجمن های خصوصی مجوز ایجاد دامنه .fr را داشتند): با گذشت یکسال .fr. بیش از 63 درصد افزایش یافت این امر باید مربوط به ایجاد تصور مثبت از ccTLD در کاربران اینترنت باشد. اشخاص 30 درصد دامنه های .fr ثبت شده و 50 درصد ثبت نام های جدید را تشکیل می دهند [2].

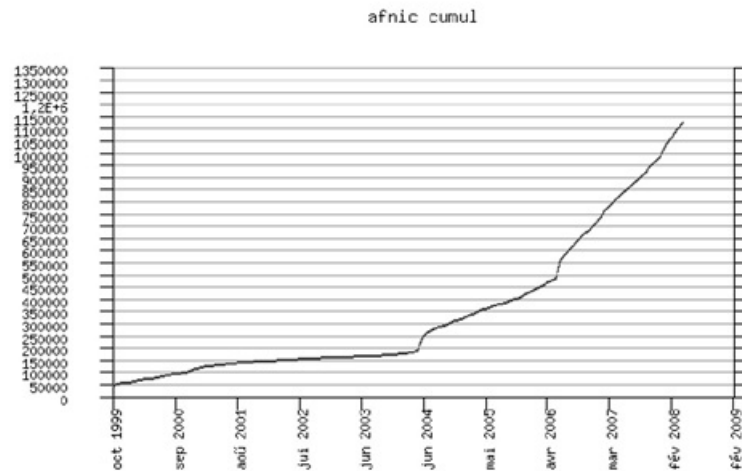
همچنین، وجود پدیده فوق در فهرست، آفنیک ممکن است تا حدی به خاطر حضور اسامی دامنه هایی باشد که یا پیوندی به وب گاه های دیگر ندارند و یا پیوندشان ضعیف است و توسط روایات ها در سال های گذشته کشف نشده بودند.

این رشد فوق العاده تا حد زیادی گسترش چشم گیر فهرست هسته را از 2006 تا 2007 توصیف می کند.

ص: 262

1-DNS [سروری که دامنه های آدرس سایت ها را پشتیبانی می کند مترجم].

قانون واسپاری وب فرانسه... ۲۶۳



شکل ۲. تغییر حجم دامنه .fr از ۲۰۰۶ تا ۲۰۰۷

۲۰۰۷	۲۰۰۶	
۵۸۱/۲۲۴	۲۰۷/۰۴۶	هسته‌های استخراج شده از خزش‌های آلکسا
۲/۲۹۵/۸۹۰	۴۲۷/۴۷۶	هسته‌های استخراج شده از خزش‌های قبلی
۸۹۰/۰۶۴	-	فهرست آف‌نیک
۲/۸۸۸/۷۲۳	۵۶۲/۶۳۷	تعداد کل پس از کاهش تکراری‌ها

شکل ۳. تغییر اندازه فهرست هسته از ۲۰۰۶ تا ۲۰۰۷

۳.۲. تنظیمات خزشگر

برداشت‌ها توسط هریتریکس^۱ انجام شده بود؛ خزشگر معتبر آرشیوی وب منبع باز توسط آرشیو اینترنت (IA) با کمک‌هایی از اعضای کنسرسیوم آی آی پی سی^۲ (به‌ویژه کتابخانه‌های ملی کشورهای حوزه شمال اروپا)^۳ [۱۱ و ۲۱]. این خزشگر، در کتابخانه ملی فرانسه برای خزش‌های فراگیر، از ۲۰۰۴ و برای خزش‌های کانونی از ۲۰۰۶، مورد استفاده قرار گرفته است.

۱. قابل دسترس در: <http://www.afnic.fr/actu/stats/evolution> (دیده شده: ۱۶ مارس ۲۰۰۸).

۲. Heritrix

۳. IIPC

۴. Nordic national libraries

۵. کنسرسیوم بین‌المللی حفاظت اینترنت در ۲۰۰۳ توسط ۱۲ سازمان (آرشیو اینترنت و چند کتابخانه ملی) برای یافتن راه حل‌های مشترک به‌منظور آرشیوسازی و اطمینان از دستیابی بلندمدت به انتشارات الکترونیکی در وب تأسیس شد. از ۲۰۰۷، عضویت برای سازمان‌های جدید باز است. برای اطلاعات بیشتر در باره اهداف و فعالیت‌های کنسرسیوم منبع [۱۳] را ببینید.

afnic cumu1

شکل ۲. تغییر حجم دامنه .fr از ۲۰۰۶ تا ۲۰۰۷ (۱)

شکل ۳. تغییر اندازه فهرست هسته از ۲۰۰۶ تا ۲۰۰۷

3.2. تنظیمات خزش گر

برداشت ها توسط هریتریکس (2) انجام شده بود؛ خزش گر معتبر آرشیوی وب منبع باز توسط آرشیو اینترنت (IA) با کمک هایی از اعضای کنسرسیوم آی. آی. پی. سی (3) (به ویژه کتابخانه های ملی کشورهای حوزه شمال اروپا) (4) [11 و 21] (5) این خزش گر در کتابخانه ملی فرانسه برای خزش های فراگیر از 2004 و برای خزش های کانونی از 2006 مورد استفاده قرار گرفته است.

ص: 263

1- قابل دسترس در <http://www.afnic.fr/actu/stats/evolution> (دیده شده: 16 مارس 2008)

Heritrix -2

IIPC -3

Nordic national libraries -4

5- کنسرسیوم بین المللی حفاظت اینترنت در 2003 توسط 12 سازمان (آرشیو اینترنت و چند کتابخانه ملی) برای یافتن راه حل های مشترک به منظور آرشیوسای و اطمینان از دستیابی بلند مدت به انتشارات الکترونیکی در وب تأسیس شد از 2007 عضویت برای سازمان های جدید باز است. برای اطلاعات بیش تر درباره اهداف و فعالیت های کنسرسیوم منبع [13] را ببینید.

با قابلیت تنظیم، بالا هریتریکس تغییر مقادیر زیادی از تنظیمات شامل، دامنه، اولویت های خزش، فیلترها، آداب رویات ها و جز آن را مقدور می سازد.

تنظیمات فوق به خزشگر فرمان می دهد چه چیزی را و چگونه باید برداشت کند. زیرا تنظیمات تأثیر زیادی در مجموعه سازی دارند و رویات لازم است مطابق با اهداف خزش گر قالب بندی شود.

2.3.1. حوزه

حوزه تعیین شده برای خزشگر تعریف می کند که کدام یک از یو. آر. ال های کشف شده باید در برداشت وارد و کدام یک کنار گذاشته شود.

از این رو، حوزه خزش کتابخانه ملی، فرانسه شامل هر وب گاهی می شود که در هر نام دامنه متعلق باشد به:

- دامنه سطح بالای fr؛

- دامنه سطح بالای re؛ و

- هر دامنه دیگری که خزش گر با آن مواجه شود زیرا از یک دامنه fr. یا از یک re تغییر مسیر داده است (طی بررسی خزش، 398,548 تغییر مسیر از اسامی دامنه آفنیك مورد توجه قرار گرفتند). به سبب تغییر مسیر از <http://yahoo.fr> به <http://fr.yahoo.com> لازم است خزش گر هر چیزی را که بخشی از <http://fr.yahoo.com> است گزینش کند اگر چه در چنین موردی، خزش گر باید در همان میزبان باقی بماند به این ترتیب خزش گر هر چیزی را که بخشی از <http://fr.yahoo.com> است، بر می گزیند، نه از <http://de.yahoo.com> یا حتی <http://fr.news.yahoo.com>.

از یک سو این حوزه بسیار بزرگ به نظر می رسد - شامل بیش از یک میلیون دامنه - و از دیگر سو محدود است زیرا وب فرانسه تنها در دامنه های fir میزبانی نمی شود به نقل از گزارش منتشر شده، آفنیك کمتر از 30 درصد وب گاه های فرانسه در fir میزبانی می شود [2] این رقم با تحلیل مجموعه سایت های انتخابات، 2007 مورد تأیید قرار گرفته است. تنها 36 درصد یو. آر. ال های این مجموعه برای میزبانی در دامنه fr پذیرفته شده اند حتی اگر از قبل چنین تعیین کنیم که وب گاه های سیاسی نماینده تام الاختیار کل وب فرانسه نیستند به طور مثال کاربرد وسیع .org. کل توسط وب گاه های احزاب و اتحادیه های تجاری به بازنمون بیش از حد این دامنه سطح بالا منجر می شود (باید بدانیم fir در بخش عمده وب گاه های فرانسه استفاده نمی شود).

درصد	URLs	TLD
36.11	22938947	.fr
28.93	18373574	.com
25.12	15955225	.org
5.81	3690655	.net
1.39	882656	.de
1.16	733634	.info
257769	257769	.eu
0.17	110944	.tv
0.17	106753	.us
0.16	99012	.re
0.58	368148	سایر TLDs

شکل ۴. تعداد یو.آر.ال‌های هر دامنه سطح بالا (TLD)، خزش‌های کانونی انتخابات ۲۰۰۷

با وجود این، هدف ما برداشت دامنه‌های دیگری غیر از .fr و .re از طریق دنبال کردن پیوندهای تغییر مسیر داده شده بود. این موضوع، راهکاری برای اختیار کردن حوزه‌ای تغییرپذیر است. شاید رویکرد اکتشافی‌تر به گردآوری وبگاه‌های خارجی‌تر منتهی شود: این امر مشکلات قانونی (کتابخانه ملی فرانسه را به خارج از چارچوب قانون واسپاری می‌برد) و اقتصادی را پدید می‌آورد: گردآوری وبگاه‌های غیرمرتبط (مطابق با مأموریت ما) باعث اشغال فضایی می‌شود که از آن می‌توان برای برداشت سایت‌های معتبر فرانسه استفاده بهتری کرد. بنابراین، تأکید بر .fr یک انتخاب واقع بینانه بود. به علاوه، از آنجا که .fr به‌طور چشمگیری در حال رشد است، می‌توان امیدوار بود که این دامنه، بخش بزرگ‌تری از وب فرانسه را سال به سال نمایندگی کند. از دیگر سو، ممکن است در آینده، به‌طور مثال با استفاده از مراجعه به دی.ان.اس خودکار (محل جغرافیایی)، راه‌های نوینی را برای کشف سایت‌های فرانسوی خارج از دامنه .fr بررسی کنیم.

۲.۳.۲. اولویت‌های خزش

هر URL که معلوم شود در حوزه خزش است، در «صف» خزشگر قرار می‌گیرد، یعنی در فهرست فایل‌های در انتظار خزش از آنجا که ممکن است روبات نتواند همه یو.آر.ال‌هایی که زمان اجرای یک خزش فراگیر در سر راهش می‌یابد، برداشت کند، برای خزشگر، مدیریت این صف و تعیین اولویت‌های خزش شدن مسائل اساسی هستند.

نخستین تصمیم مهم، انتخاب میان رویکرد «به ازای هر دامنه»^۱ و «به ازای هر میزبان»^۲ بود. با رویکرد به ازای هر میزبان، که در خزش‌های فراگیر قبلی استفاده شده بود، URLها در صف در حال انتظار خزش شدن در هر میزبان دسته‌بندی می‌شدند، و هر میزبان به‌صورت جداگانه تلقی می‌شد. این دسته‌بندی، به

1. per-domain
2. per-host

شکل 4. تعداد یو.آر.ال‌های هر دامنه سطح بالا (TLD) خزش‌های کانونی انتخابات 2007

با وجود این هدف ما برداشت دامنه‌های دیگری غیر از .fr و .re از طریق دنبال کردن پیوندهای تغییر مسیر داده شده بود این موضوع راه کاری برای اختیار کردن حوزه ای تغییر پذیر است. شاید رویکرد اکتشافی تر به گردآوری وب گاه های خارجی تر منتهی شود: این امر مشکلات قانونی (کتابخانه ملی فرانسه را به خارج از چارچوب قانون واسپاری می برد) و اقتصادی را پدید می آورد گردآوری وب گاه های

غیر مرتبط (مطابق با مأموریت ما) باعث اشغال فضایی می شود که از آن می توان برای برداشت سایت های معتبر فرانسه استفاده بهتری کرد بنابراین تأکید بر `fr` یک انتخاب واقع بینانه بود. به علاوه، از آن جا که `fr` به طور چشم گیری در حال رشد است می توان امیدوار بود که این دامنه، بخش بزرگ تری از وب فرانسه را سال به سال نمایندگی کند. از دیگر سو ممکن است در آینده به طور مثال با استفاده از مراجعه به `fr`، اس خودکار (محل جغرافیایی) راه های نوینی را برای کشف سایت های فرانسوی خارج از دامنه `fr` بررسی کنیم

2.3.2. اولویت های خزش

هر URL که معلوم شود در حوزه خزش است، در «صف» خزشگر قرار می گیرد، یعنی در فهرست فایل های در انتظار خزش از آن جا که ممکن است رویات نتواند همه یو. آر. ال هایی که زمان اجرای یک خزش فراگیر در سر راهش می یابد برداشت کند برای خزش گر مدیریت این صف و تعیین اولویت های خزش شدن مسائل اساسی هستند.

نخستین تصمیم مهم انتخاب میان رویکرد «به ازای هر دامنه» (1) و «به ازای هر میزبان» (2) بود. با رویکرد به ازای هر میزبان که در خزش های فراگیر قبلی استفاده شده بود URL ها در صف در حال انتظار خزش شدن در هر میزبان دسته بندی می شدند و هر میزبان به صورت جداگانه تلقی می شد این دسته بندی به

ص: 265

per-domain -1

per-host -2

خزش وب گاه های بیش تری با چندین میزبان منجر می شد. وب نوشت های میزبانی شده در سکوها های تجاری یا صفحه های، شخصی به عنوان میزبان های متمایز مقادیر زیادی از فضای مجموعه را اشغال می کنند به این ترتیب با استفاده از رویکرد به ازای هر دامنه با این وب گاه ها به عنوان هویتی مستقل برخورد می شود.

برای خزش سال 2007 تصمیم گرفتیم رویکرد به ازای هر دامنه را اتخاذ کنیم. علت اصلی این امر تبعیت از فهرست آف نیک، بود که فقط بر دامنه ها نظارت می. کند همچنین به عنوان یک نهاد و اسپاری می خواستیم به هر دامنه «شانس / فرصت» یکسانی برای برداشت شدن بدهیم. به علاوه، از آن جا که فهرست هسته 2007 بزرگ تر از فهرست های قبلی بود و مقادیر داده هایی که باید توسط آرشیو اینترنت (IA) بازیابی می شد افزایش نمی یافت (قسمت 2، 5 را ببینید) بیم داشتیم، که عمق خزش کم تری را کسب کنیم. بنابراین به وب گاه های تجاری یا سازمانی که اغلب چندین میزبان دارند و با خزش های کانونی بهتر برداشت می شوند، اهمیت زیادی بدهیم

در آغاز خزش لازم بود اسامی بزرگی مانند free.fr skyblog.fr orange.fr همچون سایر دامنه ها مدیریت شوند. تصمیم گیری برای اقدام به رفتار خاص با این وب گاه ها به تعویق افتاد: در طول خزش، هنگام تحلیل گزارش های ارسالی توسط مهندسان آرشیو اینترنت (IA) امکان انتخاب های مرتبط تر به نظر آسان تر می نمود.

همچنین برای جلوگیری از نمایندگی بیش از حد وب گاه های بزرگ برای هر سایت بیش ترین «بودجه» را تعیین کردیم برای روبات از یک دامنه برداشت بیش از 10,000 URL ممنوع بود. معنای این، محدودیت این نبود که خزش گر وقتی به این مقدار می رسید به طور کامل متوقف شود؛ اگر کل بودجه خزش صرف نشده بود روبات اجازه داشت مقدار بیش تری برداشت کند

سرانجام برای اطمینان از اینکه روبات منابع کافی برای برداشت کل فهرست هسته در اختیار داشته، باشد، سطح «میزان دوباره پر کردن» (1) را پایین تر انتخاب کردیم زمان اتصال به هسته هر «ریسمان» روبات (2) دستور برداشت صد URL که با دامنه مشابه را دریافت می کرد و سپس به هسته دیگر می رفت. پس از اتمام برداشت صد URL نخست هر هسته روبات به همان هسته باز می گشت.

این، رویکرد متوسطی میان رویکرد نخست پهنا (3) با هدف گردآوری وب گاه های مختلف تا جای ممکن و رویکرد آرشیو سازی ناب (4) است که در آن خزش عمیق یک سایت قبل از [آغاز] برداشت سایت بعدی انجام می گیرد

تصمیم گیری های فوق خطری را برای کتابخانه به همراه خواهد داشت در حقیقت، مجموعه فراهم شده نسبت به مجموعه های سال های قبل بسیار متفاوت است. با وجود این، برای ما بسیار اهمیت داشت که برای کشف راه های جدید در برآوردن ضرورت های وظیفه واسپاری قانونی خود، این شیوه گردآوری را تجربه کنیم.

ص: 266

replenishamount -1

thread" of the robot" -2

breadth-first -3

pure archiving approach -4

نقض رویکرد در هر دامنه این بود که: برای رویات ممکن نبود به طور خودکار دامنه های سطح دوم (اس. ال. دی) (1) را کشف کند در مورد کار ما دامنه های سطح، دوم بخش های فرعی تخصصی fr هستند. این دامنه ها یا «اسامی بخش سطح دوم» (2) مختص شناسایی یک صنعت یا یک بخش متعارف (3) (مانند aeroport.fr برای فرودگاه ها یا gov.fr برای وب گاه های دولتی هستند یا «دامنه های توصیف سطح دوم» (4) که برای شناسایی یک فعالیت یا عنوان چند نوع، اختصاص داده شده اند (به طور مثال، asso.fr. برای فدراسیون ها یا tm.fr برای دارندگان نشان های تجاری)

بدون هیچ تنظیم خاصی وب گاه های متفاوت که در دامنه های سطح دوم مشابه میزبانی شده اند، از نظر رویات ها به عنوان یک هویت محسوب شده و بودجه یکسانی را دریافت نمی کنند برای اجتناب از این مسئله، فهرست کلی دامنه های میزبانی شده در دامنه های سطح دوم را با دستور برخورد با این وب گاه ها به عنوان وب گاه های منفرد به رویات دادیم (این فهرست توسط آفیک نیز برای کتابخانه ملی فرانسه تهیه شده است).

به محتوایی که در صفحه ای جای داده شده بود اما در دامنه ای متفاوت از خود آن صفحه میزبانی می شد توجه خاصی مبذول شد به رویات اجازه داده شد از سه جهش انتقالی حداکثر (5) پیروی کند که پیوندهایی جاسازی شده در صفحه های وبی هستند. انجام این امر برای برداشت تعدادی زیادی فایل، ویدئویی که هدف مهم این خزش بود ضرورت داشت.

در واقع ما سعی داشتیم راه حل های بیشتری را برای برداشت فایل های ویدئویی بیابیم. تعداد ویدئو در وب به طور مداوم رو به افزایش است و این مسئله برای رویات هایی که آن ها را گردآوری می کنند چالش برانگیز است. مشکلات برداشت فایل های ویدئویی به اندازه و شیوه پخش آن ها مربوط می شود (برای تحلیل مشکلات گردآوری رسانه های کنونی، منابع [4] و [7] را ببینید). کمبود فایل های ویدئویی، به سبب وجود محدودیت های فنی در خزش گر، با هدف ما از آرشیو سازی یعنی تهیه تصویری «بازنمون» از وب تناقض دارد ما تصمیم گرفتیم تلاش های خود را بر دو سکوی اصلی پخش ویدئو مورد استفاده (6) netsurfers بر اینترنت فرانسه متمرکز سازیم یوتوب (7) YouTube و معادل فرانسوی آن دیلیموشن (8) Dailymotion ما از مهندسان آرشیو اینترنت خواستیم مهندسی فنی دیلیموشن را مطالعه کنند (آن ها قبلا یوتوب را به شکل بسیار خوبی بررسی کرده بودند) و نسخه های هریتریکس بهبود بخشند تا خزش گرها امکان گردآوری ویدئو فایل های این وب گاه ها را داشته باشند.

ص: 267

1- Second Level Domains (SLD)

2- second level sector names

3- regulated sector

4- second level descriptive domains

5- max-transshops

6- موجسواران

7- [25]

۲۰۰۷	۲۰۰۶	۲۰۰۵	۲۰۰۴	
دامنه	میزبان	میزبان	میزبان	قلمرو
۱۰۰	۱۰۰	۱۰۰	---	جهش - حداکثر
۳	۳	۳	---	جهش انتقالی حداکثر
۱۰۰	۵۰۰	۱۰۰۰	---	مقدار بازذخیره (در URLها)
۱۰۰۰۰	۲۰۰۰۰۰	۲۰۰۰۰۰	---	بودجه (در URL)
۴	۴	۵	---	عامل تأخیر
۵۰۰۰	۵۰۰۰	۵۰۰۰	---	حداقل تأخیر (میان دو درخواست برای یک میزبان، به میلی ثانیه)
۱۰۰۰۰	۱۰۰۰۰	۱۰۰۰۰	---	حداکثر تأخیر (میان دو درخواست برای یک میزبان، به میلی ثانیه)
مناسب	مناسب	مناسب	مناسب	تغییر مسیر
100 Mo	100 Mo	100 Mo	---	حداکثر اندازه URLهای بارگذاری شده

شکل ۵. تنظیمات خزشگر، از ۲۰۰۴ تا ۲۰۰۷

۲. ۴. پروتکل حذف روبات‌ها

قانون حق مؤلف ۲۰۰۶، به کتابخانه ملی فرانسه، اجازه سرپیچی از پروتکل حذف روبات‌ها (آر.ای.پی)^۱ را می‌دهد: «تولیدکنندگان یا ناشران] نباید در وبگاه‌های خود، رمزگذاری کنند یا مانع دستیابی مؤسسه‌های مسئول برای برداشت شوند^۲». به همین سبب، این کتابخانه، هنگام اجرای خزش‌های کانونی داخلی^۳، معمولاً به robots.txt توجه نمی‌کرد. در واقع، اغلب، حذف روبات‌ها برای جلوگیری از خزشگرهای گردآورنده صفحه‌هایی که قرار نیست نمایه‌سازی شود، توسط مدیران وب به‌کار می‌رود: راهنماهای تصویری یا صفحه‌های CSS. اما، این اسناد می‌توانند برای خزشگرهای آرشیوسازی مهم باشند، زیرا برای نمایش شکل اصلی صفحه‌های آرشیوی وب در آینده ضرورت دارند.

با وجود این، کتابخانه تصمیم گرفت از قوانین robots.txt برای خزش fr. سال ۲۰۰۷ خود اطاعت کند، همانطور که برای خزش‌های فراگیر قبلی نیز انجام می‌داد. تجارب خزش کانونی نشان داده بود که مدیران وب وقتی می‌فهمیدند که یک روبات، با شکستن قوانین robots.txt مربوط، مشغول خزش در سایت آنهاست، خوششان نمی‌آمد. به‌طور مثال، یک بلاگر فرانسوی که مورد خزش کتابخانه ملی فرانسه قرار گرفته بود، نامه‌های الکترونیکی خشونت‌آمیزی به کتابخانه فرستاد و یک تله - خزشگر برای پایین آوردن سرعت روبات کتابخانه ایجاد کرد، و پیامی در پست وبلاگش گذاشت تا دیگر مدیران وب را

1. robots exclusion protocol (REP)

۲. همچنین این جمله به این معناست که بی‌ان اف، به‌طور مثال زمان دستیابی به داده‌هایی که در وبگاه رایگان نیست، به صورت قانونی حق درخواست رمز عبور و کدهای لازم را از مالک وبگاه برای خزش سایت او دارد.

۳ in-house . یا درون سازمانی

شکل ۵. تنظیمات خزش گر، از ۲۰۰۴ تا ۲۰۰۷

۲. ۴. پروتکل حذف روبات‌ها

قانون حق مؤلف ۲۰۰۶، به کتابخانه ملی فرانسه اجازه سرپیچی از پروتکل حذف روبات‌ها (آر.ای.پی) (1) را می‌دهد «تولیدکنندگان یا

ناشران] نباید در وب گاه های خود رمزگذاری کنند یا مانع دستیابی مؤسسه های مسئول برای برداشت شوند». (2) به همین سبب این کتابخانه هنگام اجرای خزش های کانونی داخلی (3)، معمولاً به robots.txt توجه نمی کرد. در واقع، اغلب، حذف روبات ها برای جلوگیری از خزش گر های گردآورنده صفحه هایی که قرار نیست نمایه سازی شود، توسط مدیران وب به کار می رود: راهنماهای تصویری یا صفحه های CSS، اما این اسناد می توانند برای خزش گر های آرشیو سازی مهم باشند زیرا برای نمایش شکل اصلی صفحه های آرشیوی وب در آینده ضرورت دارند.

با وجود این کتابخانه تصمیم گرفت از قوانین robots.txt برای خزش fr. سال 2007 خود اطاعت کند، همان طور که برای خزش های فراگیر قبلی نیز انجام می داد. تجارب خزش کانونی نشان داده بود که مدیران وب وقتی می فهمیدند که یک روبات با شکستن قوانین robots.txt مربوط، مشغول خزش در سایت آن هاست خوش شان نمی آمد به طور مثال یک بلاگر فرانسوی که مورد خزش کتابخانه ملی فرانسه قرار گرفته بود نامه های الکترونیکی خشونت آمیزی به کتابخانه فرستاد و یک تله - خزش گر برای پایین آوردن سرعت روبات کتابخانه ایجاد کرد و پیامی در پست وبلاگش گذاشت تا دیگر مدیران وب را

ص: 268

1- robots exclusion protocol (REP)

2- همچنین این جمله به این معناست که بی ان اف، به طور مثال زمان دستیابی به داده هایی که در وبگاه رایگان نیست به صورت کانونی حق درخواست رمز عبور و کدهای لازم را از مالک وبگاه برای خزش سایت او دارد.

3- in-house . یا درون سازمانی

به گذاشتن تله - خزش گر تشویق کند (1) در واقع صفحه های و بی محدود شده توسط فایل های robots.txt اغلب تله های خزش گر را در خود جای می دهند به علاوه، گاهی هدف قوانین حذف جلوگیری خزش گرها برای جست و جوی URL های است که می توانستند موجب اضافه بار در وب گاه ها شوند (به طور، مثال پیشنهادهایی به یک گروه بحث (2)). وقتی امکان نظارت بر برداشت هر وب گاه شخصی وجود داشته، باشد مسائل فنی و مشکلات تولیدکنندگان وبگاه به راحتی در طول خزش کانونی مدیریت خواهند شد. اما کتابخانه تمایلی برای مدیریت آن ها در مقیاس کلان نداشت و نمی خواست مهندسان آرشیو اینترنت زیر فشار مخالفت های مدیران وب عصبانی غرق شوند (به هر حال، سیاست IA احترام به حذف های robots.txt است).

اغلب مؤسسه های مجری خزش های فراگیر در وب به ویژه برداشت دامنه ملی سیاست احترام به حذف روایات ها را انتخاب کردند. [17]، Nearchive.dk [10] مرکز مجازی (کتابخانه سلطنتی و کتابخانه ایالتی و دانشگاهی) با وظیفه آرشیوسازی دامنه دانمارک ترجیح می دهد از robots.txt چشم پوشی کند، زیرا این روایات بیش تر برای سایت های واقعا مهم شبکه مورد استفاده قرار می گیرد [3] (3) در واقع وب گاه ها با میزبانی محتوای بسیار ارزشمند مانند روزنامه ها یا - سایت های احزاب سیاسی رایج ترین وب گاه هایی هستند که از محدودیت های robots.txt استفاده می کنند [24].

2.5. برنامه ریزی با همکاری آرشیو اینترنت (IA)

برای خزش گر فراگیر سال 2007 آرشیو اینترنت و کتابخانه ملی، فرانسه بر اندازه تعیین شده مجموعه توافق دارند به نظر می رسد حدود 300,000,000 URL برای هماهنگ کردن نیازهای خزش منطقی باشد. در صورت نیاز IA می توانست برای 10-15 درصد افزایش تصمیم بگیرد.

هر دو نهاد با سازماندهی روز به روز خزش موافق هستند حوزه و تنظیمات اصلی برای برداشت باید میان IA و کتابخانه ملی فرانسه به بحث گذاشته شود و کتابخانه درباره آن ها تصمیم گیری کند. مهندسان IA لازم است بر کار ماشین های خزش نظارت داشته باشند و هفته ای دو بار گزارش هایی به کتابخانه ملی ارسال کنند (در ادامه می آید) برداشت باید قبل از پایان سال خاتمه یابد و تا ماه های نخست سال 2008، داده ها (برای دسترسی ماشین های Wayback و NutchWAX) نمایه سازی شوند؛ و پس از آن، آرشیوها به قفسه های ذخیره خود در کتابخانه ارسال شوند (4) دو مهندس IA باید برای کمک به کتابخانه برای نصب، قفسه ها مشاوره به تیم BNF و اطمینان از کیفیت مجموعه به پاریس بیابند.

ص: 269

1- کتابخانه به سرعت به اعتراض های بلاگرها پاسخ داد بعد از تبادل چندای میل با کتابخانه بلاگر با کاربرد قانون واسپاری وب آشنا شد و تصمیم گرفت تله اش را جمع کند

forum-2

3- توجه کنید که این تصمیم راهنمای متخصص آرشیو وب دانمارک برای انتخاب قوانین مؤدبانه محدود کننده روایات است تا مانع اضافه بار مورد درخواست سرور ها و پرونده های دعوی بالقوه شود.

4- این قفسه های Petaboxes، دارای ظرفیت بالا- هزینه پایین سخت افزار ذخیره قدرت پایین، هستند به فناوری های Capricorn اختصاص دارند. <http://www.capricorn-tech.com> [تاریخ دسترسی: 17 آوریل 2008]

3. خزش گر در حال کار

1.3. آزمایش خزش ها

قبل از اقدام به، خزش باید به منظور بررسی روایات و عکس العمل های ماشین برای پیش بینی مشکلات آزمایش خزش ها انجام شود. این، وظایف همراه با نظارت بر روایات ها در طول خزش (2/3 را ببینید) و توصیف مجموعه بعد از برداشت، عملاً مهم ترین بخش فرآیند تضمین کیفیت برای یک خزش هستند (بخش 4 را ببینید). متدولوژی، فوق برای چهار خزش فراگیر ما مورد استفاده قرار گرفت.

همان طور که قبلاً توضیح دادیم وظیفه اصلی، قبل از آغاز خزش ادغام فهرست های هسته متفاوت و بررسی آن ها بود دیگر گام مهم اجرای «آزمایش خزش» بود: برای تجزیه و کشف حجم بالای URL در فهرست هسته روایاتی راه اندازی شد اجرای خزش از آغاز تا پایان مورد نظر نیست، بلکه خزش در زمان کافی برای شناسایی اسامی دامنه غیر مرتبط یا URL های خطرناک به اجرا در می آید.

به طور مثال URL های ارسال کننده کد «خطای 404» حذف شدند به علاوه اسامی دامنه هایی که به شمار کمی میزبان عمومی (در اکثر موارد ثبت کنندگان یا مالکان تصروف کننده دامنه (1)) تغییر مسیر داده بودند نیز از خزش حذف شدند. همچنین، اسامی دامنه برای کشت دامنه (یعنی با هدف بالا بردن رتبه یک وب گاه، با استفاده از اسامی چند دامنه ای برای یک IP واحد) شناسایی شدند.

2.3. ارتباط کتابخانه ملی فرانسه و آرشیو اینترنت در طول خزش

خزش فراگیر در 11 اکتبر 2007 راه اندازی شد و در 29 نوامبر به طور کامل پایان یافت (یعنی بعد از «خزش تکه ای») در زمان خزش رابطه IA/BNF بر اساس تحلیل «گزارش یافته ها» (2) بود. این گزارش، همه دامنه های موجود در صف را با نشان دادن URL هایی که قبلاً برای هر یک از آن ها برداشت شده، مقدار بودجه صرف شده و URL های خزش شده فهرست می کند. هدف از این کار تکمیل تخصص سنتی مهندسان IA در این زمینه با دانش تیمی هسته های فرانسوی BNF است. این کار در BNF توسط یک کتابدار و یک مهندس هدایت گردید.

با افزایش بودجه به دامنه هایی که بالای 10,000 URL داشتند توجه خاصی مبذول شد. صرف نظر از اینکه داده ها مرتبط بودند یا نه این تعداد را به عنوان آستانه ای برای آن چه که باید بررسی میشد تعیین کردیم. در تعریف ما داده های غیر مرتبط اسناد با ارزش علمی اندک نیستند بلکه فایل های اضافی ای هستند که در اثر خصیصه های مربوط به وب گاه های پاتولوژی ایجاد شده اند به طور مثال، تله های روایات (به خاطر ایجاد تقویم یا نسخه های جاوا) شمار نامحدودی URL برای صفحه های ویب یکسان ایجاد کردند. وب گاه های میروور (3) نیز مشکل ساز هستند زیرا چندین دامنه محتوای یکسان را با اسامی مختلف میزبانی می کنند در اغلب موارد این دامنه ها حذف شدند.

ص: 270

1- بیش از 100000 اسامی دامنه به تنها سه وب گاه - دو ثبت نام کننده و یک متصرف دامنه تغییر مسیر داده بودند.

2- frontier report

3- Mirror websites

با استفاده از کشت دامنه وب گاه های اصلی را شناسایی و حذف کردیم. بررسی های برخط نیز برای تعیین مطابقت رشته های حروف با تقویم ها پردازش شد تا آن ها را فیلتر و از ایجاد تله های روباتی جلوگیری کند.

هدف این کنترل، هفتگی اطمینان از کسب 100 درصد کیفیت خزش نبود - به هر حال یک خزش فراگیر بود. هدف مدیریت بزرگ ترین صف ها به بهترین شکل ممکن و جلوگیری از هدر رفتن زمان و منابع زیادی توسط خزش گر بود.

3.3. «خزش تکه ای»

*«خزش تکه ای» (1)

بعد از سه هفته خزش مداوم شبانه روزی مهندسان IA تصمیم گرفتند خزش را متوقف کنند. آن ها به تحلیل QA در داده های بازیابی شده پرداختند تا دامنه هایی را مشخص کند که روبات به چهارمین سطح عمق آن ها (یعنی سه جهش نسبت به صفحه هسته) نرسیده بود این دامنه ها بار دیگر (با عنوان «خزش تکه ای») هر زمان که ممکن بود خزش شدند.

خزش تکه ای دیگری برای بازیابی فایل های ویدئویی شناخته شده در طول مدت برداشت راه اندازی شد اما بنا به دلایل فنی بارگذاری نشد (مشکلات فایل های میزبانی شده در دامنه های مختلف و جز آن).

4. پیامدهای خزش

تحلیل های زیادی درباره مجموعه برداشت شده صورت گرفت مهندسان، IA که برای کمک به ما در زمان نصب قفسه ها به پاریس آمده بودند بسیار کارآمد بودند. نخستین هدف این تحلیل ها کنترل کیفیت داده های دریافتی بود همچنین مایل بودیم مجموعه ها را در مقیاس وسیع شناسایی کنیم: طیف اسناد موجود در قفسه ها در عمل چه بود شکل و عمق وب گاه های برداشت شده چگونه بود و جز آن. از این رو، هدف، تعیین کمیت و کیفیت مجموعه سال 2007 ما بود. سرانجام تحلیل نتایج خزش فراگیر 2007، و مقایسه این نتایج با خزش های قبلی (به ویژه خزش سال 2006) برای تعیین تأثیر تنظیمات جدید خزش و تصمیم گیری درباره این که آیا این خزش ها در جهت وظیفه واسپاری قانونی ما قرار دارند، ضرورت داشت.

1.4. اشکال اصلی

عکس

با استفاده از کشت دامنه، وبگاههای اصلی را شناسایی و حذف کردیم. بررسی‌های برخط نیز برای تعیین مطابقت رشته‌های حروف با تقویم‌ها، پردازش شد تا آنها را فیلتر و از ایجاد تله‌های روباتی جلوگیری کند. هدف این کنترل هفتگی، اطمینان از کسب ۱۰۰ درصد کیفیت خزش نبود - به هر حال یک خزش فراگیر بود. هدف، مدیریت بزرگترین صف‌ها به بهترین شکل ممکن و جلوگیری از هدر رفتن زمان و منابع زیادی توسط خزش گر بود.

۳.۳. «خزش تکه‌ای»^۱

بعد از سه هفته خزش مداوم شبانه‌روزی، مهندسان IA تصمیم گرفتند خزش را متوقف کنند. آنها، به تحلیل QA در داده‌های بازبایی شده پرداختند، تا دامنه‌هایی را مشخص کند که روبات، به چهارمین سطح عمق آنها (یعنی، سه جهش نسبت به صفحه هسته) نرسیده بود. این دامنه‌ها، بار دیگر (با عنوان «خزش تکه‌ای»)، هر زمان که ممکن بود خزش شدند.

خزش تکه‌ای دیگری برای بازبایی فایل‌های ویدئویی شناخته شده در طول مدت برداشت راه‌اندازی شد، اما بنا به دلایل فنی بارگذاری نشد (مشکلات فایل‌های میزبانی شده در دامنه‌های مختلف و جز آن).

۴. پیامدهای خزش

تحلیل‌های زیادی در باره مجموعه برداشت شده صورت گرفت. مهندسان IA، که برای کمک به ما در زمان نصب قفسه‌ها به پاریس آمده بودند، بسیار کارآمد بودند. نخستین هدف این تحلیل‌ها، کنترل کیفیت داده‌های دریافتی بود. همچنین، مایل بودیم مجموعه‌ها را در مقیاس وسیع شناسایی کنیم: طیف اسناد موجود در قفسه‌ها، در عمل، چه بود، شکل و عمق وبگاه‌های برداشت شده چگونه بود، و جز آن. از این رو، هدف، تعیین کمیت و کیفیت مجموعه سال ۲۰۰۷ ما بود. سرانجام، تحلیل نتایج خزش فراگیر ۲۰۰۷، و مقایسه این نتایج با خزش‌های قبلی (به‌ویژه خزش سال ۲۰۰۶)، برای تعیین تأثیر تنظیمات جدید خزش و تصمیم‌گیری در باره اینکه آیا این خزش‌ها در جهت وظیفه واسپاری قانونی ما قرار دارند، ضرورت داشت.

۱.۴. اشکال اصلی

تعداد	۲۰۰۶	۲۰۰۷
URLها	۲۷۱۶۹۷۴۵۶	۳۳۷۳۲۲۲۰۰
میزبان‌ها	۲۹۲۸۳۶۴	۱۵۸۹۴۵۸
دامنه‌ها	۳۸۲۵۴۰	۱۰۶۲۳۱۷
(از دامنه های .fr)	۱۳۱۱۳۶	۷۹۱۹۴۰

1. patch-crawl

۳۱۸	۷۱۰	URL های هر دامنه
۲۱۲	۹۳	یو.آر.ال های هر میزبان
91745	73073	فایل های ARC یکتا
8,8	7,2	اندازه فشرده داده های یکتا (در Tb)

شکل ۶. اشکال اصلی خزش فراگیر fr ۲۰۰۷

به خاطر افزایش چشمگیر فهرست هسته، رشد تعداد دامنه های برداشت شده پیش بینی می شد. همکاری با آفنیک (انجمن نامگذاری اینترنت فرانسه برای ثبت fr) در سال ۲۰۰۶، به کتابخانه اجازه داد دامنه های fr گردآوری شده را در شش مرحله کشف و برداشت کند. به عبارت دیگر، در سال ۲۰۰۷ نسبت به ۲۰۰۶، با وجود افزایش یو.آر.ال های برداشت شده توسط IA، میزبان های کمتری خزش شدند. این امر، احتمال دارد به خاطر رویکرد در ازای هر دامنه، مربوط به آخرین خزش فراگیر ما باشد. شمار یو.آر.ال ها در ازای هر دامنه یا در ازای هر میزبان راهکاری ساده برای ارزشیابی «عمق متوسط» یک خزش است. این شکل، تفاوت های معنی دار فراوانی میان وبگاهها را پنهان می کند: شکل در بخش ۶.۴ با جزئیات بیشتر بررسی شده است.

۶.۴. توزیع هر سرآیند^۱

گزارش سرآیند در باره خزش فراگیر سال ۲۰۰۷، حدود ۱۶۰۴ نوع مختلف را نشان می دهد. تعجب آور نبود که یک سرآیند text/html، تنها دو سوم از فایل های برداشت شده را نمایش می دهد. به علاوه، ۹۷ درصد URL های بارگذاری شده یکی از پنج نوع سرآیند پر استفاده: HTML، JPEG، GIF، PNG و PDF، را دارند. اگر فردی سرآیندهای اسناد برداشت شده طی خزش فراگیر ۲۰۰۷ را نگاه کند، تصور می کند وب فرانسه ۲۰۰۷ را بیشتر شامل متن و تصویر است. با وجود این، باید در باره شکل ها بسیار محتاط باشیم. لازم است محدودیت های فنی رویت ها را نیز در نظر بگیریم؛ زیرا حتی اگر عملکرد خزشگر به طور مداوم پیشرفت کند، قادر به تجزیه و گردآوری قالب های مختلف فایل که در وب می باید، نیست. قالب های فایلی پیچیده، بدون نمایش بوده یا به راحتی در مجموعه غایب هستند.

دلیل مناسب دیگر برای رعایت احتیاط این است که اطلاعات سرآیند مورد نیاز محاسبه، همانی است که توسط سرور فرستاده می شود. در واقع، این اطلاعات قابل اعتماد نیست. گاهی، حتی سرور سرآیندی می فرستد که وجود ندارد (با کمال تعجب، «نرم افزار / X- چیزی»^۲ را در مجموعه خود یافتیم). به جز ۱۶۰۴ سرآیند گوناگون، ۱۴۰۰ سرآیند با کمتر از ۵۰۰ فایل مربوط هستند- می توان چنین استنباط کرد

1. MIME type
2. application/x-something

شکل 6. اشکال اصلی خزش فراگیر 2007 fr

به خاطر افزایش چشم گیر فهرست، هسته رشد تعداد دامنه های برداشت شده پیش بینی می شد. همکاری با آفنیک (انجمن نام گذاری اینترنت فرانسه برای ثبت fr) در سال 2006، به کتابخانه اجازه داد دامنه های fr گردآوری شده را در شش مرحله کشف و برداشت کند.

به عبارت دیگر در سال 2007 نسبت به 2006 با وجود افزایش یو. آر.ال های برداشت شده توسط IA میزبان های کمتری خزش شدند این، امر احتمال دارد به خاطر رویکرد در ازای هر دامنه مربوط به آخرین خزش فراگیر ما باشد.

شمار. یو. آر.ال ها در ازای هر دامنه یا در ازای هر میزبان راه کاری ساده برای ارزشیابی «عمق متوسط» یک خزش است این، شکل تفاوت های معنی دار فراوانی میان وب گاه ها را پنهان می کند: شکل در بخش 4. 6 با جزئیات بیشتر بررسی شده است.

2.4. توزیع هر سرآیند

*توزیع هر سرآیند (1)

گزارش سرآیند درباره خزش فراگیر سال 2007، حدود 1604 نوع مختلف را نشان می دهد. تعجب آور نبود که یک سرآیند text/html تنها دو سوم از فایل های برداشت شده را نمایش می دهد. به علاوه، 97 درصد URL های بارگذاری شده یکی از پنج نوع سرآیند پر استفاده HTML، JPEG، GIF، PNG و PDF را دارند اگر فردی سرآیندهای اسناد برداشت شده طی خزش فراگیر 2007 را نگاه کند، تصور می کند وب فرانسه 2007 را بیشتر شامل متن و تصویر است.

با وجود، این باید درباره شکل ها بسیار محتاط باشیم. لازم است محدودیت های فنی روبات ها را نیز در نظر بگیریم؛ زیرا حتی اگر عملکرد خزش گر به طور مداوم پیشرفت کند قادر به تجزیه و گردآوری قالب های مختلف فایل که در وب می یابد نیست قالب های فایلی، پیچیده بدون نمایش بوده یا به راحتی در مجموعه غایب هستند.

دلیل مناسب دیگر برای رعایت احتیاط این است که اطلاعات سرآیند مورد نیاز محاسبه، همانی است که توسط سرور فرستاده می شود در واقع این اطلاعات قابل اعتماد نیست گاهی، حتی سرور سرآیندی می فرستد که وجود ندارد (با کمال تعجب «نرم افزار / X-چیزی» (2) را در مجموعه خود یافتیم). به جز 1604 سرآیند گوناگون 1400، سرآیند با کمتر از 500 فایل مربوط هستند- می توان چنین استنباط کرد

ص: 272

MIME type -1

application/x-something -2

که سرآیندها به طرز نامناسبی تعیین شده اند.

به نظر می رسد اهمیت این مسئله سال به سال بیشتر شود. یو.آر.ال های خزش فراگیر 2004، دارای 554 سرآیند گوناگون هستند؛ این رقم به 1024 در سال 2006 و به 1604 در 2007 تغییر یافت.

عکس

قانون واسپاری وب فرانسه... ۲۷۳

که سرآیندها به طرز نامناسبی تعیین شده اند.

به نظر می رسد اهمیت این مسئله سال به سال بیشتر شود. یو.آر.ال های خزش فراگیر ۲۰۰۴، دارای ۵۵۴ سرآیند گوناگون هستند؛ این رقم به ۱۰۲۴ در سال ۲۰۰۶ و به ۱۶۰۴ در ۲۰۰۷ تغییر یافت.

درصد	یو.آر.الها	MIME-type
67.96	229257942	متن / html
19.04	64222287	تصویر / jpeg
7.52	25376262	تصویر / gif
1.17	3955885	تصویر / png
1.17	3955463	نرم افزار / pdf
۰,۶۷	۲۲۵۶۷۵۹	متن / ساده
۰,۴۷	۱۵۹۴۳۴۲	برنامه / فلش - shockwave-x
۰,۴۲	1432809	متن / css
۰,۴۲	۱۴۱۵۲۳۰	برنامه / javascript-x
0.32	1083991	برنامه / XML
0.82	2771213	سایر

شکل ۷. ده رتبه نخست سرآیند خزش فراگیر ۲۰۰۷

اما اگر نمی توانیم به سرآیند یک فایل منفرد اطمینان کاملی داشته باشیم، توزیع فراگیر تعیین شده برای صدها میلیون سند به احتمال زیاد قابل اعتماد است. تحول سرآیند می تواند به عنوان راهی برای تحلیل تغییرات و تمایلات در مقیاسی وسیع، دیده شود. به طور مثال می توان، از سال ۲۰۰۴ تا ۲۰۰۷، کاهش در استفاده از قالب GIF را، به نفع JPEG و فورمت باز PNG^۱ را مشاهده کرد (میزان تصاویر GIF در عرض این چهار سال تقریباً نصف شده است). بر همین اساس، میزان اسناد XML پنج برابر شده است. حتی وقتی به حجم اسناد پرداخت شده نگاه می کنیم رشد فایل های XML در وب و وضوح بیشتری می یابد: از ۸۷۰۰۰ در سال ۲۰۰۴ به یک میلیون در ۲۰۰۷. شاید این رشد، تا اندازه ای، به خاطر کاربرد فزاینده آر.اس.اس خوانها^۲ باشد (سرآیند صحیح برای یک فایل آر.اس.اس، «application/rss» است، اما غالباً «application/xml» یا حتی «text/xml» به جای آن به کار می روند).

۱. توجه کنید که این اشکال گاهی به چند سرآیند دسته بندی می شوند: به طور مثال، تعدادی اسناد JPEG معین با افزایش شماری اسناد دارای سرآیند «image/jpeg»، «Image/jpeg» یا «image/JPEG».

۲. توجه کنید که نرخ های مشابهی برای دامنه au مشاهده شده اند (استرالیا، ۲۰۰۴ تا ۲۰۰۵): درصد تصاویر GIF نصف شده است (از ۱۰ تا ۵ درصد)، با وجود این، تصاویر png در استرالیا (۵۶ درصد) نسبت به فرانسه (۱,۱۷ درصد) کمتر مورد استفاده قرار گرفته اند [۱۸].

3. RSS feeds

اما اگر نمی توانیم به سرآیند یک فایل منفرد اطمینان کاملی داشته باشیم توزیع فراگیر تعیین شده برای صد ها میلیون سند به احتمال زیاد قابل اعتماد است. تحول سرآیند می تواند به عنوان راهی برای تحلیل تغییرات و تمایلات در مقیاسی وسیع دیده شود. به طور مثال می توان از سال 2004 تا 2007 کاهش در استفاده از قالب GIF را به نفع JPEG و فورمت باز (2) PNG را مشاهده کرد (میزان تصاویر GIF در عرض این چهار سال تقریباً نصف شده است). بر همین اساس میزان اسناد XML پنج برابر شده است. حتی وقتی به حجم اسناد برداشت شده نگاه می کنیم رشد فایل های XML در وب وضوح بیش تری می یابد: از 88/000 در سال 2004 به یک میلیون در 2007 شاید این، رشد تا، اندازه ای به خاطر کاربرد فزاینده آر.اس. اس خوان ها (3) باشد (سرآیند صحیح برای یک فایل آر.اس. اس، «application/rss» است، اما غالباً «application/xml» یا حتی «text/xml» به جای آن به کار می روند).

ص: 273

1- توجه کنید که این اشکال گاهی به چند سرآیند دسته بندی می شوند به طور مثال تعدادی اسناد JPEG معین با افزایش شماری اسناد دارای سرآیند «image/jpeg»، «Image/jpeg» یا «image/JPEG».

2- توجه کنید که نرخ های مشابهی برای دامنه au. مشاهده شده اند (استرالیا 2004 تا 2005): درصد تصاویر GIF نصف شده است (از 10 تا 5 درصد). با وجود، این تصاویر png در استرالیا (0,56 درصد) نسبت به فرانسه (1,17 درصد) کمتر مورد استفاده قرار گرفته اند

[18]

RSS feeds -3

قالب فایلی دیگر رو به افزایش در مجموعه های خزش های فراگیر ، نرم افزار فلش شاک ویو (1) است. این رشد می تواند دو علت داشته باشد: محبوبیت در حال افزایش قالب فلش در وب توانایی بیش تر خزنده هریتریکس برای برداشت این نوع فایل

عکس

۲۷۴ مدیریت منابع اطلاعاتی وب

قالب فایلی دیگر رو به افزایش در مجموعه های خزش های فراگیر ، نرم افزار فلش شاک ویو^۱ است. این رشد می تواند دو علت داشته باشد: محبوبیت در حال افزایش قالب فلش در وب، توانایی بیشتر خزنده هریتریکس برای برداشت این نوع فایل.

MIME Type evolution	2004	2005	2006	2007
text/html	68.11	67.22	70.15	67.96
image/jpeg	14.04	15.79	15.13	19.04
image/gif	12.70	11.09	8.05	7.52
application/pdf	1.36	1.39	1.19	1.17
image/png	0.79	0.73	0.87	1.17
text/plain	1.0833	1.19	1.01	0.67
application/x-shockwave-flash	0.2488	0.34	0.35	0.47
application/xml	0.07	0.16	0.50	0.32

جدول ۸. تحول در چند سرآیند از ۲۰۰۴ تا ۲۰۰۷

از منظر حفاظت بلندمدت، این اطلاعات بسیار ارزشمند است و قالبی را که - در مقیاس ملی و نیز بین المللی، باید تلاش های خود را بر آن متمرکزسازیم، ارائه می کند. کاربرد رو به افزایش قالب های باز، مانند PNG یا XML، اخبار خوبی از این نقطه نظر است.

۳. ۴. فایل های ویدئویی

افزایش چهار قالب ویدئویی پر استفاده (Quicktime, Windows media video, Flash video و MPEG video) منجر به ایجاد حدود ۴۰۰۰۰ فایل ویدئویی برداشت شده در ۲۰۰۴ گردید، در مقابل ۱۲۰۰۰۰ فایل در چهار سال بعد (یعنی ۰/۰۴ درصد مجموعه). ما متوجه کاهش قالب ویدئویی MPEG در مقابل فلش شده ایم. در ۲۰۰۶ خزشگر هریتریکس فقط صد فایل ویدئویی فلش را برداشت کرده بود، یک سال بعد، نسخه مجوزدار خزشگر ما، برای برداشت محتوای میزبانی شده در سکویهای پخش ویدئویی به اجرا درآمد: هریتریکس سی هزار سند را گرد آورد. اراده ما برای تمرکز بر برداشت فایل های ویدئویی به هدفش رسیده بود: اگرچه اغلب آرشیوهای ما «ضعف هایی»^۲ در ویدئو داشته است، در ۲۰۰۷، نمونه ویدئویی حجیم تری را به دست آوردیم.

1. application/x-shockwave-flash
2. holes

از منظر حفاظت بلندمدت این اطلاعات بسیار ارزشمند است و قالبی را که - در مقیاس ملی و نیز بین المللی باید تلاش های خود را بر آن متمرکز سازیم ارائه می کند کاربرد رو به افزایش قالب های باز مانند PNG یا XML، اخبار خوبی از این نقطه نظر است

3.4. فایل های ویدئویی

افزایش چهار قالب ویدئویی پر استفاده (Flash video Quicktime Windows media video و MPEG video) منجر به ایجاد حدود 40000 فایل ویدئویی برداشت شده در 2004، گردید در مقابل 120000 فایل در چهار سال بعد (یعنی 0/04 درصد مجموعه) ما متوجه کاهش قالب ویدئویی MPEG در مقابل فلش شده ایم. در 2006 خزش گر هریتریکس فقط صد فایل ویدئویی فلش را برداشت کرده بود، یک سال بعد، نسخه مجوزدار خزش گر ما، برای برداشت محتوای میزبانی شده در سکو های پخش ویدئویی به اجرا درآمد هریتریکس سی هزار سند را گرد آورد اراده ما برای تمرکز بر برداشت فایل های ویدئویی به هدفش رسیده بود اگر چه اغلب آرشیوهای ما «ضعف هایی» (2) در ویدئو داشته است، در 2007، نمونه ویدئویی حجیم تری را به دست آوردیم

ص: 274

application/x-shockwave-flash-1

holes-2

قانون واسپاری وب فرانسه... ۲۷۵

MIME-type	2004	2005	2006	2007
ویدئو /x-ms-wmv	4 408	7 705	33 936	39 218
ویدئو /quicktime	22 020	26 687	39 073	36 294
درخواست /x-flv	0	0	104	31 556
ویدئو /mpeg	11 408	17 304	28 413	14 992
جمع	39 840	53 701	103 532	124 067

شکل ۹. تحول سرآیند فایل‌های ویدئویی از ۲۰۰۴ تا ۲۰۰۷

۴.۴. توزیع هر TLD

دیگر کشف سه چهارم اسناد خزش شده متعلق به دامنه سطح بالای fr تعجب‌آور نبود (می‌توان به این رقم، دامنه re را که توسط آف‌نیک مدیریت می‌شود، افزود). با وجود این، شکل ۱۰ تأیید می‌کند که چنین چیزی با تنظیماتی که ما به کار گرفتیم، از فهرست هسته a.fr آغاز و تا مرزهای آن ادامه دادیم، ممکن است. دو نوع دامنه سطح بالا (تی.ال.دی.) ارائه شده در مجموعه نیز دامنه‌های سطح بالای کدهای عمومی و ملی^۱ هستند. اغلب سایت‌های میزبانی شده تحت دامنه‌های سطح بالای عمومی (.org، .net، .com، .info) به احتمال، توسط مدیران وب فرانسوی تولید شده‌اند. کد کشوری مربوط به رتبه‌بندی دامنه‌های سطح بالا در شکل ۱۰ یا متعلق به کشورهای فرانسوی زبان (بلژیک و سوئیس) است یا به همسایگانی که (آلمان و انگلیس) فرانسه با آنها ارتباط تجاری عمده دارد مربوط می‌شود: این پدیده برای اسپانیا نیز مورد توجه قرار گرفت [۵]. دامنه eu. ویژه اروپاست.

TLD	تعداد یوآرآلها	درصد
fr	259 869 452	77.12
com	59 843 624	17.76
net	4 951 932	1.47
org	3 171 196	0.94
de	2 808 359	0.83
eu	993 456	0.29
info	900 544	0.27
be	660834	0.20
ch	461 021	0.14
uk	434 315	0.13
re	381 746	0.11
Other TLDs	2 471 064	0.73

جدول ۱۰. تعداد یوآرآل‌های هر TLD، خزش‌گر فراگیر ۲۰۰۷

این ارقام، کاملاً با داده‌های مجموعه‌های قبلی، از ۲۰۰۴ تا ۲۰۰۶ مشابه است. به سقوط بالای biz. (که

1. general and country codes

شکل ۹. تحول سرآیند فایل‌های ویدئویی از ۲۰۰۴ تا ۲۰۰۷

4.4. توزیع هر TLD

دیگر کشف سه چهارم اسناد خزش شده متعلق به دامنه سطح بالای fr. تعجب‌آور نبود (می‌توان به این رقم دامنه re. را که توسط آف‌نیک

مدیریت می شود افزود). با وجود، این شکل 10 تأیید می کند که چنین چیزی با تنظیماتی که ما به کار گرفتیم، از فهرست هسته a.fr آغاز و تا مرزهای آن ادامه دادیم، ممکن است دو نوع دامنه سطح بالا (تی.ال.دی.) ارائه شده در مجموعه نیز دامنه های سطح بالای کدهای عمومی و ملی (1) هستند اغلب سایت های میزبانی شده تحت دامنه های سطح بالای عمومی (org .net .com. info) به احتمال، توسط مدیران وب فرانسوی تولید شده اند کد کشوری مربوط به رتبه بندی دامنه های سطح بالا در شکل 10 یا متعلق به کشورهای فرانسوی زبان (بلژیک و سوییس) است یا به همسایگانی که (آلمان و انگلیس) فرانسه با آن ها ارتباط تجاری عمده دارد مربوط می شود: این پدیده برای اسپانیا نیز مورد توجه قرار گرفت. [5] دامنه eu. ویژه اروپاست

جدول 10. تعداد یو. آر. ال های هر، TLD، خزش گر فراگیر 2007

این ارقام، کاملاً با داده های مجموعه های قبلی از 2004 تا 2006 مشابه است. به سقوط بالای biz. (که

ص: 275

سال ها قبل در 10 رتبه نخست دامنه سطح بالا جای داشت)، و ظهور ناگهانی eu. توجه کنید.

با وجود این، اگر به تعداد دامنه های میزبانی شده در یک دامنه سطح بالای خاص نگاه کنیم توزیع های متفاوت را متوجه می شویم.

عکس

۲۷۶ مدیریت منابع اطلاعاتی وب

سالها قبل در ۱۰ رتبه نخست دامنه سطح بالا جای داشت)، و ظهور ناگهانی eu. توجه کنید.
با وجود این، اگر به تعداد دامنه های میزبانی شده در یک دامنه سطح بالای خاص نگاه کنیم،
توزیع های متفاوت را متوجه می شویم.

TLD	دامنه های ۲۰۰۵ (درصد)	دامنه های ۲۰۰۶ (درصد)	دامنه های ۲۰۰۷ (درصد)
com	44.59	42.78	17.10
fr	26.86	34.28	74.55
net	5.37	5.95	2.06
org	4.39	4.69	1.46

جدول ۱۱. درصد دامنه هر تی.ال.دی (منحصر به com, fr, net, و org)، از ۲۰۰۵ تا ۲۰۰۷.

مجموعه های ۲۰۰۵ و ۲۰۰۶، که با تنظیمات مشابه (فهرست هسته و تنظیمات خزش) شکل گرفته اند برتری بیشتر دامنه های سطح بالای عمومی را بر fr. نشان می دهند. این ارقام می تواند توسط حجم زیاد دامنه های غیرموجود در فهرست هسته، که توسط رویات لمس شده بود توصیف شود (یعنی جایی که یک یو.آر.ال. پیوند یافته به یک یو.آر.ال. درون حوزه، مورد خزش قرار گرفته است). برخی پیوندها به دامنه های com یا net. در صفحه های fr. قابل دسترس بودند، و به همین دلیل تا حدودی گردآوری شدند. به خاطر افزایش چشمگیر اندازه فهرست هسته^۱، میزان زیادی از دامنه های قابل دسترس در مجموعه ۲۰۰۵ و ۲۰۰۶ - اما نه در خزش فراگیر ۲۰۰۷ - عرضه می شوند.

۴. Robots.txt. ۵.

در سال ۲۰۰۷، پروتکل حذف رویات^۲ مانع ما در آرشوسازی ۱۵ میلیون فایل بود، یعنی، ۴.۵ درصد فایل های کشف شده - و به وضوح از آرشوسازی تمام اسنادی که رویات می توانست از آغاز تشکیل این فایل ها کشف کرده باشد جلوگیری کرد. این ارقام درکل، زمانی که فایل های robots.txt بارگذاری ۶ درصد اسناد کشف شده را متوقف کردند نسبت به خزش فراگیر ۲۰۰۶ بسیار پایین تر هستند. تفسیر این میزان مشکل است، زیرا آنها با آخرین مطالعات این پروتکل، که کاربرد در حال رشد robots.txt را شناسایی می کند، مغایرت دارند [به طور مثال منبع ۲۳].

زمان تحلیل آنچه خزش نشده بود، استفاده از سرآیند فایل ها امکان نداشت، زیرا سرور آنها را ارسال نکرده بود. با وجود این، می توانیم از پسوند فایل URL های خواسته شده استفاده کنیم. تقریباً ۴۰ درصد این فایل ها تصویر هستند^۳. ممکن است برخی مدیران وب می خواستند مانع خزش رویات در این فایل ها

۱. توجه کنید برخی دامنه های برداشت شده com یا net بین ۲۰۰۶ و ۲۰۰۷ تفاوت خیلی زیادی با یکدیگر ندارند. ۱۶۳۶۳۲ دامنه com در ۲۰۰۶ نسبت به ۱۸۱۶۲۶ در ۲۰۰۷؛ ۲۲۷۴۶ دامنه net در ۲۰۰۶ نسبت به ۲۱۸۵۳ در ۲۰۰۷.
2. robots Exclusion Protocol
۳. jpg. (۲۶ درصد)، gif. (۸ درصد)، JGP. (۲ درصد)، png. (۲ درصد). به علاوه توجه کنید که ۱۰۴،۱۰۹ فایل (۱ درصد) از فایل های CSS بودند.

جدول 11. درصد دامنه هر تی.ال.دی (منحصر به com, fr, net, و org)، از 2005 تا 2007.

مجموعه های 2005 و 2006، که با تنظیمات مشابه (فهرست هسته و تنظیمات خزش) شکل گرفته اند برتری بیش تر دامنه های سطح بالای عمومی را بر fr نشان می دهند این ارقام می تواند توسط حجم زیاد دامنه های غیر موجود در فهرست، هسته که توسط روبات لمس شده بود توصیف شود (یعنی جایی که یک یو.آر.ال. پیوند یافته به یک یو.آر.ال درون حوزه مورد خزش قرار گرفته است). برخی پیوندها به دامنه های .com یا .net در صفحه های fr قابل دسترس بودند و به همین دلیل تا حدودی گردآوری شدند به خاطر افزایش چشمگیر اندازه فهرست هسته (1) میزان زیادی از دامنه های قابل دسترس در مجموعه 2005 و 2006 - اما نه در خزش فراگیر 2007 عرضه می شوند.

Robots.txt.5.4

در سال 2007، پروتکل حذف روبات ها (2) مانع ما در آرشیو سازی 15 میلیون فایل بود، یعنی، 4,5 درصد فایل های کشف شده - و به وضوح از آرشیو سازی تمام اسنادی که روبات می تواند از آغاز تشکیل این فایل ها کشف کرده باشد جلوگیری کرد این ارقام در کل زمانی که فایل های robots.txt بارگذاری 6 درصد اسناد کشف شده را متوقف کردند نسبت به خزش فراگیر 2006 بسیار پایین تر هستند. تفسیر این میزان مشکل است زیرا آن ها با آخرین مطالعات این پروتکل که کاربرد در حال رشد robots.txt را شناسایی می کند مغایرت دارند [به طور مثال منبع 23].

زمان تحلیل آن چه خزش نشده بود استفاده از سرآیند فایل ها امکان نداشت زیرا سرور آن ها را ارسال نکرده بود. با وجود این می توانیم از پسوند فایل URL های خواسته شده استفاده کنیم. تقریباً 40 درصد این فایل ها تصویر هستند (3) ممکن است برخی مدیران وب می خواستند مانع خزش روبات در این فایل ها

ص: 276

1- توجه کنید برخی دامنه های برداشت شده .com یا net بین 2006 و 2007 تفاوت خیلی زیادی با یکدیگر ندارند. 163632 دامنه .com در 2006 نسبت به 181626 در 2007؛ 22746 دامنه .net در 2006 نسبت به 21853 در 2007

2- robots Exclusion Protocol

3- jpg (26 درصد)، gif (8 درصد)، png (2 درصد)، JGP (2 درصد)، png (2 درصد) به علاوه توجه کنید که 104,109 فایل (1 درصد) از فایل های CSS بودند.

شوند، زیرا آن‌ها توسط موتورهای جست‌وجو نمایه‌سازی نشده بودند این فرض با کمک شیوه‌ای که روبات‌های ما این فایل‌ها را کشف کردند مورد تأیید قرار گرفته است: نیمی از آن‌ها (7/378/578)، هنگام پیگیری یک پیوند جاسازی شده یافت شدند.

بنابراین، تمکین از robots.txt مانع برداشت بسیار مرتبط داده‌هایی می‌شود که برای روبات‌های موتورهای جست‌وجو، غیر ضروری اما برای روبات‌آرشیوسازی ما بسیار مفید است. در برخی مواقع، پروتکل حذف روبات‌ها (آر.ای.پی) از دسترسی ما به کل یک سایت جلوگیری می‌کرد بیش از URL 150/000 که به عنوان هسته‌هایی که به وسیله robots.txt محافظت می‌شدند به کار رفتند.

6.4 عمق خزش

برای تعیین عمق خزش، ممکن است تعداد یو.آر.ال‌ها برای هر دامنه .fr را محاسبه کنیم.

عکس

شوند، زیرا آنها توسط، موتورهای جست‌وجو نمایه‌سازی نشده بودند. این فرض با کمک شیوه‌ای که روبات‌های ما این فایل‌ها را کشف کردند مورد تأیید قرار گرفته است: نیمی از آنها (۷/۳۷/۵۷۸)، هنگام پیگیری یک پیوند جاسازی شده، یافت شدند.

بنابراین، تمکین از robots.txt مانع برداشت بسیار مرتبط داده‌هایی می‌شود که برای روبات‌های موتورهای جست‌وجو، غیرضروری، اما برای روبات آرشیوسازی ما بسیار مفید است. در برخی مواقع، پروتکل حذف روبات‌ها (آر.ای.پی) از دسترسی ما به کل یک سایت جلوگیری می‌کرد: بیش از ۱۵۰/۰۰۰ URL که به‌عنوان هسته‌هایی که به وسیله robots.txt محافظت می‌شدند به کار رفتند.

۶.۴ عمق خزش

برای تعیین عمق خزش، ممکن است تعداد یو.آر.ال‌ها برای هر دامنه .fr را محاسبه کنیم.

تعداد دامنه‌ها	تعداد یو.آر.ال‌ها
498777	10<
146356	10-100
103370	100-1000
43101	1000-10000
334	10000>

جدول ۱۲. تعداد یو.آر.ال‌ها برای هر دامنه .fr، خزش فراگیر ۲۰۰۷

تقریباً ۵۰ درصد دامنه‌های برداشت شده شامل ۱۰ یو.آر.ال یا کمتر هستند. این امر می‌تواند به سبب عدم دسترسی به سرور دائمی در طول خزش باشد. با وجود این، ما چندین آزمایش «دستی» انجام دادیم، یعنی برای کشف آنچه به‌صورت برخط قابل دسترس بود، در یک دو جین نام دامنه با آستانه ۱۰ یو.آر.ال. کلیک کردیم. این آزمایش‌ها نشان دادند که این وبگاه‌ها خالی بودند (مالکان آنها را خریداری کردند تا مطمئن شوند که آنها استفاده نخواهند شد، اما خودشان هم از آنها استفاده نمی‌کنند) یا اینکه آنها برای کشت پیوند استفاده شدند.

این ارقام با ارقام متعلق به خزش فراگیر قبلی بسیار تفاوت داشتند.

تعداد دامنه‌ها	تعداد یو.آر.ال‌ها
38 439	10<
32 258	10-100
41 352	100-1000
15 159	1000-10000
3 928	10000>

شکل ۱۳: تعداد یو.آر.ال‌ها برای هر دامنه .fr، خزش فراگیر ۲۰۰۶

جدول 12. تعداد یو.آر.ال‌ها برای هر دامنه .fr، خزش فراگیر 2007

تقریباً 50 درصد دامنه‌های برداشت شده شامل 10 یو.آر.ال یا کم تر هستند. این امر می‌تواند به سبب عدم دسترسی به سرور دائمی در طول خزش باشد با وجود این ما چندین آزمایش «دستی» انجام دادیم یعنی برای کشف آن چه به صورت برخط قابل دسترس بود در یک دو جین نام دامنه با آستانه 10 یو.آر.ال کلیک کردیم این آزمایش‌ها نشان دادند که این وب‌گاه‌ها خالی بودند (مالکان آن‌ها را خریداری کردند تا مطمئن شوند که آن‌ها استفاده نخواهند شد، اما خودشان هم از آن‌ها استفاده نمی‌کنند) یا این که آن‌ها برای کشت پیوند استفاده شدند.

این ارقام با ارقام متعلق به خزش فراگیر قبلی بسیار تفاوت داشتند.

شکل 13: تعداد یو.آر.ال ها برای هر دامنه .fr. خزش فراگیر 2006

ص: 277

سه تفاوت عمده میان سال های 2006 و 2007 درباره عمق دامنه های .fr. عبارت اند از:

- دامنه های بسیار زیاد تقریباً خالی در 2007 که می توانست ناشی از افزایش اندازه فهرست هسته باشد.

- برخی دامنه ها در جایی که بیش از 10,000 یو آر ال آرشیو شده بود میان دو خزش به عدد 10 تقسیم شدند این امر پیامد محدود کردن «بودجه» برای 10,000 یو.آر.ال بود؛ همچنین شاید به خاطر رویکرد «به ازای هر دامنه»: دامنه های چندین میزبانه باز نمون نمی شوند.

- از سوی دیگر این رویکرد خزش عمیق تر تعداد بیش تر دامنه های «کوچک» یا «متوسط»، بین 100 و 10,000 یو.آر.ال را اجازه می داد.

ممکن است این ارقام تحت تأثیر تله های رویات، باشند اما در حال حاضر نمی دانیم با کدام پسوند. سایر کنترل کیفیت ها، بیش از همه کنترل بصری وب گاه های شخصی برای تخمین این ارقام ضرورت دارد. با وجود این ارقام از دید قانون واسپاری فرانسه و اهداف اولیه این خزش، رضایت بخش به نظر رسد: با تضمین اینکه همه وب گاه های .fr موجود در مجموعه وظیفه سنگین برداشت وب گاه های کوچک و متوسط شاید به بهای وب گاه های بزرگ تر را داشته باشند

7.4. وب گاه های بزرگ

برای پالا-پیش تحلیل عمق وب گاه ممکن است بر بزرگترین وب گاه ها تمرکز کنیم. هدف، بررسی علت چرایی بیش تر خزش شدن این وب گاه ها نسبت به دیگران است - آیا به این دلیل است که آن ها نسبت به دیگران بیش تر برخط هستند؟

اطلاعات مختلفی از نمایه خزش (سی.دی.ایکس) (1) اخذ شده بود: فهرست 50 دامنه بزرگ (مربوط به سال 2005 تا 2007)؛ و فهرست 1000 دامنه بزرگ در سال های 2006 و 2007.

1.7.4. دامنه ها

از این گذشته میان مجموعه های مختلف تناقض های زیادی وجود دارد نخستین تناقض، اندازه دامنه های بزرگ است تنها سه دامنه دارای بیش از 1,000,000 یو آر ال در مجموعه در 2007 نسبت به 25 یو.آر.ال در مجموعه 2006 - در بخش قبلی ارقام متفاوتی را دیدیم و می توانیم با همان دلایل این ارقام را توصیف کنیم، هم چنین محتوای 50 سایت بزرگ نخست فهرست بسیار متفاوت است: تنها 20 درصد 50 دامنه بزرگ نخست مربوط به خزش گرهای فراگیر سال 2007 در وضعیتی معادل آن در 2006 ارائه شده اند.

بزرگ ترین وبگاه 2006 (free.fr) که دارای بیش از هفت میلیون یو.آر.ال بود، در سال 2007 فقط 40000 یو.آر.ال «وزن» دارد!

ص: 278

قانون واسپاری وب فرانسه... ۲۷۹

نام دامنه ۲۰۰۷	تعداد یو.آرال ها	نام دامنه ۲۰۰۶	تعداد یو.آرال ها
asso.fr	3 984 821	free.fr	7 405 987
com.fr	1 760 957	amiz.fr	5 194 030
tm.fr	1 270 244	asso.fr	4 036 224
gouv.fr	534 599	lrencontre.fr	3 482 657
cci.fr	495 753	sportblog.fr	2 547 231
co.uk	408 881	gouv.fr	2 360 960
nom.fr	179 256	promovacances.fr	2 113 314
presse.fr	144 207	football.fr	1 895 302
dailymotion.com	108 222	mbpro.fr	1 885 549
notaries.fr	102 081	com.fr	1 856 720

شکل ۱۴. ده دامنه بزرگ اول، از ۲۰۰۵ تا ۲۰۰۷

به علت انجام تنظیمات خاص در مجموعه ۲۰۰۷، بالاترین رتبه دامنه‌ها تقریباً مربوط به دامنه‌های سطح دوم است. از سوی دیگر، در خزش‌های ۲۰۰۵ و ۲۰۰۶، دو نوع وبگاه کشف شد: سکوهایی با میزبانی و بنوشت‌ها و وبگاه‌های شخصی (free.fr, sportblog.fr) که توسط رویکرد به ازای هر میزبان پشتیبانی شده‌اند، و وبگاه‌های تجاری با آگهی در صفحه‌های متعدد (به‌طور مثال promovacances.fr، آژانس مسافرتی برخط). خزش فراگیر ۲۰۰۵ نیز چند وبگاه دانشگاهی را نشان می‌دهد (۱۶ وبگاه دانشگاهی در ۵۰ دامنه بزرگ اول)، مانند jussien.fr یا cnrs.fr. در سال‌های بعد، این وبگاه‌ها تقریباً در فهرست ۵۰ وبگاه بزرگ اول وارد شدند.

حضور گسترده وبگاه‌های تجاری توصیف کننده تعداد دامنه‌های سطح بالای عمومی در ۵۰ دامنه بزرگ اول هستند: حتی برای مجموعه ۲۰۰۷، تنها ۳۲ درصد در fr هستند. به این ترتیب، توصیف مجموعه منعکس کننده وب فرانسه در سال ۲۰۰۷ چنین است: عمدتاً، فضایی برای تجارت، خدمات و روابط اجتماعی.

۴. ۷. ۲. دامنه‌های سطح دوم

	2006	2007	Evolution
asso.fr	4 036 224	3 984 821	↙
com.fr	1 856 720	1 760 957	↙
tm.fr	1 150 555	1 270 244	↗
gouv.fr	2 360 960	534 599	↙

شکل ۱۵. تحول دامنه‌های سطح دوم از ۲۰۰۶ تا ۲۰۰۷

شکل ۱۴. ده دامنه بزرگ اول، از ۲۰۰۵ تا ۲۰۰۷

به علت انجام تنظیمات خاص در مجموعه ۲۰۰۷، بالاترین رتبه دامنه‌ها تقریباً مربوط به دامنه‌های سطح دوم است از سوی دیگر در خزش‌های ۲۰۰۵ و ۲۰۰۶، دو نوع وبگاه کشف شد سکوهایی با میزبانی و بنوشت‌ها و وبگاه‌های شخصی (free.fr, sportblog.ir) که توسط رویکرد به ازای هر میزبان پشتیبانی شده‌اند و وبگاه‌های تجاری با آگهی در صفحه‌های متعدد (به‌طور مثال

promovacances.fr، آژانس مسافرتی برخط) خزش فراگیر 2005 نیز چند وب گاه دانشگاهی را نشان می دهد (16 وب گاه دانشگاهی در 50 دامنه بزرگ اول)، مانند jussien.fr یا cnrs.fr. در سال های بعد این وب گاه ها تقریباً در فهرست 50 وبگاه بزرگ اول وارد شدند.

حضور گسترده وب گاه های تجاری توصیف کننده تعداد دامنه های سطح بالای عمومی در 50 دامنه بزرگ اول هستند: حتی برای مجموعه 2007 تنها 32 درصد در fr. هستند به این ترتیب، توصیف مجموعه منعکس کننده وب فرانسه در سال 2007 چنین است: عمدتاً فضایی برای تجارت، خدمات و روابط اجتماعی

2.7.4. دامنه های سطح دوم

شکل 15. تحول دامنه های سطح دوم از 2006 تا 2007

ص: 279

توجه خاص به دامنه های سطح دوم در طول برداشت سال 2007 به کتابخانه ملی فرانسه خزش مقادیر داده اندکی پایین تر را نسبت به مشابه آن ها در برداشت سال 2006 اجازه داد. مهم ترین استثنا دامنه سطح دوم .gou.fr است از 2/3 میلیون به 500,000 یو.آر.آل سقوط کرده است. یک راه برای توصیف این رقم کاربرد مشترک دامنه های سطح سوم و حتی چهارم در وب گاه های دولتی (به طور مثال، www.rhone.pref.gouv.fr یا www.auvergne.culture.gouv.fr) است.

تأکید بر فایل های ویدئویی مجموعه بهتری از وب گاه های پخش ویدئویی (1) به بار آورد دیلی موشن (2) 10 تا از پر رتبه ترین ها دامنه ها را وارد کرد (در سال 2006، تنها 30,000 یو.آر.آل در این دامنه برداشت شده بود). تعداد فایل های گردآوری شده در یوتیوپ دوبرابر شدند.

5. نتیجه گیری

برای برداشت دامنه ملی فقط یک راه وجود ندارد. انتخاب های فنی متفاوتی (قبل و حین خزش، و حتی بعد از حذف URL ها) مجموعه را شکل می دهند سیاست مجموعه سازی حتی در مقیاس وسیع اعمال می شود. برای اجرای یک، خزش دامنه بر طبق قانون و اهداف، آن شناسایی این انتخاب ها و تخمین پیامد های آن ها ضروری است.

به طور سنتی در رسالت قانون و اسپاری برای تصمیم گیری درباره مرتبط بودن سند به حوزه، مجموعه سه معیار به کار می رفته است باید با قالبی، خاص قابل دسترس مردم و داخل مرزهای سرزمین فرانسه موجود باشد، همه این ها باید ویژگی های جدید وب را پذیرفته باشند و باید هنگام انجام خزش فراگیر به حساب آمده باشند.

هدف برداشت وب «فرانسه» تأکید بر .fr را توصیف می کند در واقع این امر برآورده نمی شود اگر تصور کنیم 50 تا 60 درصد وب گاه های فرانسوی خارج از .fir هستند اما در حال حاضر تأکید بر .fr. یک انتخاب واقع گرایانه و اقتصادی است همه محتوای فرانسه در .fr موجود نیست، اما هر چیزی که در .fr است متعلق به فرانسه است. به علاوه این تأکید قابل انعطاف است زیرا روایات اجازه دارد از تغییر مسیرهای .fir به دیگر دامنه ها تبعیت کند. سرانجام دامنه .fr به سرعت با ساده کردن قوانین ویژه این دامنه در حال توسعه است و امید می رود به زودی بخش بزرگ تری از دامنه فرانسه را بنمایاند. انتظار می رود این خواسته با قوانین CCTLD جدید مربوط به ICANN متوقف نگردد.

همچنین تأکید بر تأکید بر .fr بسیار آسان است زیرا مطمئن هستیم به لطف موافقت نامه با آفیک، قادر به برداشت همه جانبه این دامنه سطح بالا هستیم این راه کاری است که با دومین اصل مهم قانون واسپاری منطبق است: گردآوری کل تولید فرهنگی کشور هر آن چه به محض قرار گرفتن در دسترس عموم «کیفیت» می یابد یا «ارزش» دارد. از تعداد بسیار زیاد هسته آغاز کردن ضمانتی است برای فراموش نکردن وب گاه های با پیوند کمتر یا نامشهور

ص: 280

این اصول «غیر تبعیض آمیزانه» با ویژگی دیگر قانون واسپاری: میل به گردآوردن تمام قالب های انتشاراتی در حال ظهور ناسازگار نیستند. کتابخانه در گذشته [از قالب های مختلفی] پشتیبانی می کرد: متون، تصویر، صدا یا ویدئو. با آرشیوسازی وب دیگر رفتار متفاوت با آن ها قابل تصور نیست زیرا این قالب ها به عناصری در شبکه ای یکسان پیوند یافته اند با وجود، این در صورت ضرورت می توان به راه حل های مناسبی برای انواع مختلف رسانه دست یافت گردآوری نمایه سازی، حفاظت و دستیابی به فایل های متنی و دیداری - شنیداری یا تعاملی همواره به مسائل یکسانی ختم نمی شوند.

هنوز پرسش هایی باقی است: به طور مثال دشوار است که بگوییم، در حال حاضر، برای کسب تصویری جامع از وب فرانسه آیا رویکرد تأکید تنها بر دامنه ها بهتر از نگهداری به ازای هر میزبان به طور جداگانه است لازم است تحلیل بیشتری برای پاسخ به این سؤال انجام شود و کتابخانه ملی فرانسه، به شنیدن گزارش های بین المللی درباره این موضوع بسیار علاقه دارد سرانجام، مسئله، چگونگی تعریف یک وبگاه است زیرا این هویت هوشمند اغلب با میزبانی فنی سازگار نیست اگر وب گاه را به عنوان دامنه تعریف کنیم، رویکرد به ازای هر دامنه باید پذیرفته شود، زیرا به خزش بهتر وب گاه های کوچک یا متوسط منتهی می شود. اما اگر وبگاه را به عنوان هویتی فکری ایجاد شده توسط یک نویسنده یا یک ویراستار (یک یا چند شخص یک مؤسسه عمومی یا خصوصی) تعریف کنیم سایت و دامنه دیگر با هم سازگار نیستند. در واقع، تعداد یا حتی حجم عظیمی از وب گاه ها می توانند تحت یک نام دامنه میزبانی شوند، همچون و نوشتن های میزبانی شده در سکوها تجاری. این سایت ها - با وجود مرتبط بودن - با راه کاری که در سال 2007 استفاده کردیم باز نمون نشد «به طور مثال، صفحه های شخصی فراوانی، به کل، از free.fr در 2006 برداشت شدند اما یک سال بعد ناپدید شدند برای اجتناب از این مسئله، در خزش های فراگیر آینده به کارگیری رویکردی خاص برای بسیاری از سکوها عمومی برای میزبانی کردن برخی وب گاه ها یا وب نوشت های شخصی امکان پذیر است زیرا متوجه سکوها پخش ویدئو در 2007 بودیم.

این تصمیم ها نشان دهنده اصلاحاتی است که می توانست هدف خزش هریتریکس باشد. قابلیت های بهتر برای تجزیه و برداشت قالب های پیچیده، فایل یکی از مسائل ضروری است - از این نقطه نظر هریتریکس پیش از این خود را بسیار پیکربندی شده نشان داد. در چارچوب پروژه «خزشگر هوشمند»، تاکنون سه ویژگی، دیگر توسط IA، IIPC، کتابخانه کنگره کتابخانه، بریتانیا و کتابخانه ملی فرانسه - که هدف شان توسعه روایات هریتریکس است - ارائه شده است نخستین، آن ها اجتناب از برداشت محتوایی است که از زمان آخرین خزش تغییری نکرده است: این ویژگی کاهش تکرار به صرفه جویی در تخمین منابع و ذخیره سازی آن ها و بنابراین به خزش عمیق تر وب گاه ها منجر می شود. دومین ویژگی، دادن مجوز به روایات برای اولویت بندی URL های داخل صف است. تلفیق رویکرد تمام گزینشی (در آغاز خزش) با به کارگیری قابلیت های کنترل خزش خودکار بهینه بسیار مفید خواهد بود. سومین توسعه، شناسایی خودکار بسامد تغییر وب گاه ها نیز به روایات اجازه شناسایی سایت هایی که باید به آن ها توجه خاصی مبذول شود، می دهد. از این رو، به طور مثال این امر باعث می شود تاریخ ها و بسامدهای خزش فراگیر را انتخاب کنیم، یا خزش های کانونی را در سایت های پر تغییر اجرا کنیم

در حقیقت اگر قرار است خزش های فراگیر با عمقی متوسط هر وب گاهی را یک یا دوبار در سال برداشت کنند باید هدف از خزش های کانونی را چنین تعریف کنیم: خزش های کانونی باید از ابتدا قصد شان آرشیو وب گاه های بزرگ و عمیق باشد؛ نه وب گاه های فرانسوی fit یا وب گاه های پرتغییر - آرشیو آن ها به بهترین شکل ممکن زیرا حتی خزش های کانونی اغلب برای برداشت کامل وب گاه های عظیم و گردآوری اسناد موجود در وب پنهان کافی نیستند با وجود این باید به خاطر داشته باشیم که تهیه تصاویر تنها راه - و اقتصادی ترین راه - گردآوری حافظه دیجیتال فرانسه است و اتخاذ هر تصمیم در این موضوع در سایر راه های آرشیوسازی در راهبرد تلفیقی باید به حساب آیند.

تشکر و قدردانی

مایلیم از کریس کارپنتر (1) و تیم آرشیو، اینترنت به عنوان همکار ما در این چهار سال پایانی (از ابتدا!) قدردانی کنیم به ویژه ایگور رانیتوویک (2) که بر تمام خزش های BNF مربوط به سال های 2004 تا 2007 نظارت داشت همراه با جان لی (3)، برد تافل (4) و میخائیل ماگین (5) که کمک کردند قفسه ها را نصب کنیم و قالب های نوین مجموعه را در پاریس و سان فراسیسکو تحلیل نماییم. همچنین سپاس بسیار از مارالینو چوونیک (6) برای دقت در ویرایش داریم سرانجام، مراتب سپاس مان را به گیلداس ایلین (7)، رئیس واسپاری دیجیتال به خاطر توصیه ها و حمایت هایش تقدیم می کنیم

ص: 282

Kris Carpenter -1

Igor Ranitovic -2

John Lee -3

Brad Tofel -4

Michael Magin -5

Mireille Chauveinc -6

Gildas Illien -7

- Abiteboul, S., Cobena, Masanès, J. and Sedrati, G. 2002. A First Experience in Archiving the French [1] Web. In Proceedings of the Research and advanced technology for digital libraries: 6th European conference .(Italy, 2002
- AFNIC. 2007. French Domain Name Industry report. 2007 Edition. AFNIC, Saint Quentin en Yvelines. [2] <http://www.afnic.fr/data/actu/public/2007/afnic-frenchdomain-name-report-2007.pdf>
- Andersen, B. 2005. The DK-domain: in words and figures. Netarkivet.dk, Aarhus, Copenhagen. [3] <http://netarchive.dk/publikationer/DFreyv-english.pdf>
- Ashenfelder, M. 2006. Web Harvesting and Streaming Media. In Proceedings of the 6th International [4] Web Archiving Workshop (Alicante, Spain). <http://www.iwaw.net/06/PDF/iwaw06-proceedings.pdf>
- Baeza-Yates, R., Castillo, C. and Lopez, V. 2005. Characteristics of the Web of Spain. In Cybermetrics, [5] 9. <http://www.catedratelefonica.upf.es/webes/2005/Characteristics Web-Spain.pdf>
- Baeza-Yates, R., Castillo, C., Marin, M. and Rodriguez, A. 2005. Crawling a country: Better Strategies [6] than BreadthFirst for Web Page Ordering. In Proceedings of the 14th international conference on World .(Wide Web (Chiba, Japan
- Baly, N. and Sauvin, F. 2006. Archiving Streaming Media on the Web, Proof of concept and Firsts [7] Results. In Proceedings of the 6th International Web Archiving Workshop (Alicante, Spain). <http://www.iwaw.net/06/PDF/iwaw06-proceedings.pdf>
- Brin, S. and Page, L. 1998. The Anatomy of a Large-scale Hypertextual Web Search Engine. In [8] Computer Networks and ISDN Systems, 30 (1-7), 107-117. <http://www7.scu.edu.au/programme/fullpapers/1921/com1921.htm>
- Dailymotion. Dailymotion [Accessed: May 10, 2008]. Partagez vos vidéos. [9] <http://www.dailymotion.com>
- Gomes, D and Silva, M. Characterizing a National Community Web. ACM Transactions on Internet [10] Technology (volume 5, issue 3), New York, 508-531. <http://xldb.fc.ul.pt/daniel/gomesCharacterizing.pdf>
- [Heritrix. Heritrix Home Page. <http://crawler.archive.org> [Accessed: May 22, 2008 [11]

- IIPC. International internet preservation consortium-welcome. <http://www.netpreserve.org>. [13]
.[[Accessed: May 24,2008
- Illien, G., Aubry, S., Hafri Y. and Lasfargues, F 2006. Sketching and checking quality for web archives: [14]
a first stage report from BnF. Bibliothèque nationale de France, Paris.
<http://bibnum.bnf.fr/conservation/index.html>
- Illien, G. 2006. Web archiving at BnF. In International Preservation News, Paris, BnF, 27-34. [15]
<http://www.ifla.org/VI/4/news/ipnn40.pdf>
- Kimpton, M., Braggs, M. and Ubois, J. 2006. Year by Year: From an Archive of the Internet to an [16]
.Archive on the Internet. In Web Archiving, J. Masanè s, Ed, Springer, Berlin, Heidelberg, New York
- Koerbin, P. 2005. Report on the crawl and Harvest of the Whole Australian Web Domain Undertaken [17]
during June and July 2005. National Library of Australia, Canberra. [http://
pandora.nla.gov.au/documents/domain-harvest-report-public.pdf](http://pandora.nla.gov.au/documents/domain-harvest-report-public.pdf)
- Koerbin, P. 2008. The Australian Web domain harvests: a preliminary quantitative analysis of the [18]
archive data. National Library of Australia, Canberra. <http://pandora.nla.gov.au/documents/auscrawls.pdf>
- Masanè s, J. 2002. Towards continuous Web Archiving: First results and an agenda for the future. In D- [19]
Lib Magazine, 8 (12).<http://www.dlib.org/dlib/december02/masanes/12masanes.html>
- Masanè s, J. 2006. Selection for Web Archives. In Web Archiving, J. Masanè s, Ed, Springer, Berlin, [20]
.Heidelberg, New York
- Mohr, G., Kimpton, M., Stack, M, and Ranitovic, I. 2004.Introduction to Heritrix, an archival quality [21]
Web crawler.Paper presented at the 4th International Web Archiving Workshop (Bath, United Kingdom,
2004). <http://www.iwaw.net/04/Mohr.pdf>
- Najork, M. and Wiener, J. L. 2001. Breadth-First Search Crawling Yields High-Quality Pages. In: [22]
Proceedings of the 10th international conference on World Wide Web. Elsevier Science, Hong Kong, 114-
.118
- Sun, Y. Zhuang Z., Council I. and Giles C L. 2007 Determining Bias to Search Engines from [23]
Robots.txt. In Proceedings of the IEEE/WIC/ACM International Conference on Web Intelligence. IEEE
Computer Society Washington, 149-155. <http://www.personal.psu.edu/yus115/docs/sun-robotstxtbias.pdf>

:Sun, Y. Zhuang Z. and Giles C. L..2007. A large-scale study of robots.txt. In WWW '07 [24]

ص: 284

Proceedings of the 16th international conference on World Wide Web, ACM Press, New York 1123-1124.

http://www2007.org/posters/poster_1034.pdf

YouTube. YouTube - Broadcast Yourself. <http://www.youtube.com> [Accessed: May 15, 2008 [25]]

ص: 285

مجموعه های کتابخانه ملی فرانسه (1) بخشی از میراث ملی هستند و به طور تقریبی 31 میلیون سند از همه نوع (کتاب نشریه نسخه های خطی عکس ها، نقشه ها و غیره) را شامل می شوند. چالش های جدید مجموعه با گسترش اینترنت ایجاد شده اند. کتابخانه ملی فرانسه در قالبی بین المللی رهنمون های خط مشی، گردش کارها، و ابزارها را توسعه می دهد تا قسمت های مرتبط و معرف بخش فرانسوی اینترنت را جمع آوری کند و حفاظت و دسترسی به آن ها را سازماندهی کند.

آرشیوهای وب حوزه ملی فرانسه به عنوان خدمتی جدید توسعه یافت به عنوان کاربردی جدید عرضه شد و در آوریل 2008 در دسترس عموم قرار گرفت. از آن، پس راهبردهایی بوده است و توسعه می یابد تا کتابداران را درگیر آن کند و آن را در دسترس کاربران نهایی قرار دهد

این، مقاله تجربه کتابخانه ملی فرانسه را به ویژه با تمرکز بر این چهار موضوع بررسی می کند:

- ساختمان مجموعه آرشیو وب به عنوان مجموعه ای جدید و چالش برانگیز،
 - کشف منبع: خدمات و ابزارهای دسترسی برای کاربران نهایی
 - کاربرد: اطلاعات و ارقام
 - مشارکت: راهبردهایی برای ساخت یک انجمن کتابداران و رسیدن به کاربران نهایی
- کلید واژه ها: آرشیوهای ملی آرشیو کردن وب گاه ها، ساختمان مجموعه، کشف منبع، کاربران نهایی، کاربرد، فرانسه

ص: 286

نوشته: سارا اوبری (1) | ترجمه: زهرا تهوری (2)

ساختمان مجموعه آرشیوهای وب به عنوان مجموعه ای جدید و چالش برانگیز اینترنت نقش مهمی در زندگی های روزانه ما به عهده گرفته است به عنوان مثال، مدیریت الکترونیکی یادگیری الکترونیکی تجارت الکترونیکی انتشارات آنلاین هنرهای دیجیتال بلاگ ها و فضاهای جدید عمومی بحث و گفتگو بسیاری از فعالیت ها تاحدی یا به طور کامل به سمت وب حرکت کرده و تغییر مکان داده و فعالیت های جدیدی ایجاد کرده. اند با افزایش مستمر و رو به رشد تعداد کاربران اینترنت (امروزه، حدود 35 میلیون در فرانسه) و تعداد فزاینده وب گاه های فرانسوی (فقط 1/7 میلیون دامنه .fr. به نام Top Level Domain or TLD ثبت می شود) (3)، حیاتی است که این نوع انتشارات و رسانه ارتباطی را که هنوز تاحدی برای یک کتابخانه ملی جدید است، بررسی کرد.

یک وبگاه در موارد زیر متفاوت از دیگر انواع انتشارات است:

- محدود به شکل خاصی همچون یک تکه کاغذ یا یک قطعه موسیقی روی دیسک نیست بلکه وابسته به زیرساخت شبکه ای پیچیده تری است یک وبگاه یک فایل پی.دی.اف تنها (شبهه پایان نامه های

ص: 287

Sara Aubry, Web Archiving Project Manager, IT Department, National Library of France, Quai François – 1

Mauriac, sara.aubry@bnf.fr

2- کارشناس سازمان اسناد و کتابخانه ملی ایران

Observatoire du marché des noms de domaine en France, FR Network Information Center, – 3

[http://www.afnic.fr/actu/observatoire [last accessed on 2010-06-15]

الکترونیکی) یک JPEG تنها یا یک تصویر TIFF (شبهه تصاویر و عکس ها) نیست بلکه چند وجهی است: یک صفحه وب گردآوری تعدادی عوامل (متون، تصاویر، نوشته ها، شیوه نامه ها (1)، فایل های صوتی و ویدئویی و غیره) است که ممکن است از جاهایی دیگر (سامانه فایل محلی پایگاه داده وبگاه راه دور دیگر و غیره) بیاید و هنگامی که هر کاربر وبگاه را با یک مرورگر مشاهده می کند، جمع آوری شوند. تحلیل گردآوری یک نمونه 2/9 میلیون وبگاه فرانسوی در سال 2007 نشان داد که حدود 1600 نوع رسانه اینترنتی متفاوت وجود دارد. (2)

• یک وبگاه آغازی ندارد و نمی تواند تا آخر خوانده شود موجودیتی عقلانی است که می تواند توسط کاربری متفاوت از کاربر دیگر دیده شود. در ارتباط با شبکه ای از دیگر وب گاه های پیوند یافته توسط پیوندهای، فرامتنی، شبکه ای از پیوندهایی که حتی قوی تر از شبکه یک انتشارات علمی مبتنی بر استناد ها و اطلاعات کتاب شناختی است وجود دارد.

• وب گاه ها بسیار متعدد هستند و وب هیچ مرزی ندارد. بر اساس آخرین نظر سنجی نتکرافت (3)، بیش از 206 میلیون وبگاه وجود دارد. (4) بیش تر آن ها از هر کشوری در دنیا قابل دسترس هستند و ممکن است بخشی از یک مجموعه ملی بر اساس رهنمودهای خط مشی توسعه مجموعه ملی حقوقی یا سازمانی باشند.

• محتوای وب همیشه درست شبیه یک بخار در حرکت است صفحات وب ممکن است در عرض یک روز چندین مرتبه روزآمد شوند (در صفحات وب روزنامه هایی مانند لوموند (5) و لیبراسیون (6) اعلام کوچکی وجود دارد که نشان می دهند محتوی چند دقیقه قبل روزآمد شده است).

• وب گاه ها و محتوای وب نیز گذرا هستند: یک صفحه وب ممکن است در هر زمانی و به چند دلیل ناپدید شود: توقف اختیاری یا غیر عمدی توسط وب مستر عدم تمدید نام دامنه شکستن دیسک یا مشکلات دسترسی شبکه به سرور میزبان و غیره محتوای وب به یک حادثه خاص پیش بینی شده یا غیر مترقبه پیوند دارد و به ویژه در معرض خطر است به مناسبت یک گزینش مشارکتی و طرح گردآوری وب گاه های سیاسی در طول انتخابات ریاست جمهوری 2007 فرانسه کتابخانه لیون (7) دریافت که 52 درصد از 421 وب گاهی که انتخاب شده بودند یا به طور کامل یا به طور تقریبی پنج ماه بعد از رأی گیری بسته شدند. (8)

چارچوب حقوقی

هر زمان که یک نوع جدید ماده نمایش و ایجاد از جمله فناوری های متنوع جدیدی که در فرانسه ظاهر

ص: 288

1- style sheets

2- Legal deposit of the French Web: harvesting strategies for a national domain. France Lasfargues, Clément

Oury, Bert Wendland, IAWW, 2008: <http://iaww.net/08/IWAW2008-Lasfargues.pdf> [last accessed on 15-06-2010].

3- Netcraft

4- Netcraft May 2010 Web Server Survey, <http://news.netcraft.com/archives/category/web-server-survey/>

.[[last accessed on 2010-06-15

Le Monde -5

Libé ration -6

Library of Lyon -7

La netcampagne des lé gislatives 2007 en Rhône-Alpes: la course au Net et après,- 8

.[[http://www.pointsdactu.org/article .php3?id_article=863](http://www.pointsdactu.org/article.php3?id_article=863) [last accessed on 2010-06-15

شدند، اختراع میشد کتابخانه ملی فرانسه نخست آزمایش می کرد، سپس سازمانش را با جمع آوری، حفظ و دسترس پذیری به این انتشارات دیجیتال متولد شده وفق می داد زمان آن رسیده است که پس از، کتاب ها حکاکی ها پارتیسیون های موسیقی عکس ها، پوستر ها مدارک صوتی تصویری و چندرسانه ای وب گاه ها نیز بایگانی شوند.

قانون میراث فرانسه (1)، اکنون حق مالکیت IV (ماده 1-311 L تا پایان (1-133 L) از قانون حق مؤلف و نگرش های حقوقی در جامعه اطلاعاتی 2006-961 را که مطابق دستور عمل EC/2001 / 29 پارلمان اروپا و شورای 22 می 2001 در مورد سازگاری جنبه های خاص کپی رایب و حقوق مرتبط با جامعه اطلاعاتی (2) است به ثبت می رساند.

این قانون که به طور رسمی سوم آگوست 2006 منتشر شد:

• دامنه واسپاری حقوقی (3) اینترنت را به این موارد گسترش می دهد: مضمون واسپاری حقوقی هر علامت، نشانه نوشته، تصویر صدا یا پیغام های هر نوع ارتباط با عموم به وسیله کانال های الکترونیکی نیز می شود (ماده 39) قانون برای همه نوع انتشارات الکترونیکی آنلاین از جمله مجموعه ای از علائم نشانه ها، تصاویر، صداها یا هر نوع پیغامی که در اینترنت در دسترس عموم باشد، قابل اجر است نه تنها وب گاه ها، بلکه خبرنامه ها و رسانه های جاری نیز شامل این تعریف می شوند؛

• چگونگی تقسیم مسئولیت های واسپاری وب بین سازمان های تحت قیمیت را تعریف می کند: مؤسسه خبرگزاری ملی که مسئول حفظ میراث صوتی تصویری فرانسه است وب گاه های مرتبط با ارتباطات صوتی تصویری (به طور عمده رادیو و تلویزیون) و کتابخانه ملی فرانسه تمامی وبگاه ها را جمع آوری خواهند کرد؛ یک حکم در دست اقدام است تا فرایندهای گزینش و دسترسی را به اجرا در آورد؛

• راهبردهای جمع آوری را خاص می کند در کتابخانه ملی فرانسه و اسپاری حقوقی اینترنت نیاز به اجازه از ناشران ندارد و به جمع آوری خودکار بخش عمده اولویت می دهد: سازمان های تحت قیمیت ممکن است مواد را مطابق با فرایندهای خاص واسپاری تولید کنندگان از اینترنت جمع آوری کنند. قانون نیز تصریح می کند که هیچ مانعی همچون برقراری ارتباط (4) رمز عبور یا شکل های دیگر محدودیت دسترسی ممکن نیست توسط ناشران برای محدود کردن این فرایند استفاده شود.

دامنه

اگر چه بیشتر وب عمومی توسط هر کس در فرانسه می تواند ملاحظه شود، از نظر فنی و حقوقی غیر ممکن است کل وب را بایگانی کرد. کتابخانه ملی فرانسه دستور دارد تا وب گاه های دامنه ملی فرانسه را جمع آوری کند، یعنی:

ص: 289

The French Heritage Law, or «Code du patrimoine» – 1

Loi n°2006-961 du 1 août 2006 relative au droit d'auteur et aux droits voisins dans la société de – 2

l'information , <http://www.legifrance.gouv.fr/affchTexte.do?cidTexte> = JORF

[TEXT000000266350dateTexte= [last accessed on 2010-06-15

legal depository –3

• به عنوان یک هسته هر وب گاهی با دامنه fr TLD. یا هر TLD مشابه دیگری که به قلمرو مدیریتی فرانسه مربوط می شود (به عنوان مثال re برای جزیره فرانسوی La Reunion)

• هر وب گاهی (به احتمال خارج از دامنه fr) که تولید کننده اش از نظر جغرافیایی تحت قلمرو فرانسه هست (به طور معمول این می تواند در صفحات وب یا با استفاده از سرورهای خاص بررسی شود)؛

• هر وب گاهی (به احتمال خارج از دامنه fr) که بتوان ثابت کرد محتوای تولید شده اش در قلمرو فرانسه به نمایش گذاشته می شود (بررسی این معیار اخیر چالش برانگیز تر است اما فرصتی برای تفسیر و مذاکره با کتابخانه و تولیدکنندگان اینترنت ایجاد می کند).

ابزارها و روشهای جمع آوری

گر چه ما از یک «واسپاری» حقوقی صحبت می کنیم وب گاه ها در حقیقت توسط ناشران به کتابخانه واسپاری نمی شوند. در عوض توسط تکه هایی از نرم افزارهایی به نام روبات های خزنده آرشیو جمع آوری می شوند یک خزنده آرشیو شبیه خزنده های نمایه سازی موتورهای جستجو عمل می کند. برنامه ای است که وب را به روشی خودکار مطابق مجموعه خط مشی هایی مرور می کند. با فهرستی از نشانی های یو.آر.ال (1) شروع و هر صفحه شناسایی شده توسط یو.آر.ال را ذخیره می کند تمامی فرامتن ها را در صفحه می یابد (به عنوان مثال پیوندهایی به دیگر صفحات، تصاویر نوشته ها یا شیوه نامه ها، ویدئوها و غیره) و آن ها را به فهرست یو.آر.ال ها می افزاید تا به طور مسلسل دیده شوند.

پارامترهای فنی بر هویت و رفتار خزنده (دامنه، عمق سرعت فیلترهای تحریم، و غیره) تأثیر می گذارد اما از آن جا که فنون وب بسیار پیچیده هستند و خیلی به سرعت توسعه می یابند، خزنده با بسیاری موانع فنی مواجه می شود که مانع آن از جمع آوری تمامی عوامل یک وب گاه یا حتی یک صفحه وب می شود. بنابراین آرشیوهای وب اغلب ناقص هستند کتابخانه ملی فرانسه از خزنده منبع باز Heritrix که توسط مؤسسات عضو کنسرسیوم حفاظت بین المللی اینترنت توسعه یافته است، (2) استفاده می کند

از آن جا که ممکن نیست جامعیت را هدف قرار داد یا گزینش دستی وب گاه ها را به عهده گرفت، کتابخانه ملی فرانسه در نظر گرفته است دوروش جمع آوری مکمل را برای رویارویی با چالش های واسپاری حقوقی وب ترکیب کند:

• جمع آوری خودکار بخش عمده ای از وب گاه های فرانسوی خزش های وسیع در جهت ارائه که تصویری کلی از هزاران فایل از تعداد بسیار زیادی از وب گاه ها انجام می شود. به عنوان مثال، برای خزش وسیع 2010 که هنوز در زمان نگارش این مقاله هم فعال است، کتابخانه ملی فرانسه بیشترین حد 10000 یو.آر.ال را برای هر 1/6 میلیون وب گاه جمع آوری می کند خزنده ها محتوی را بدون هیچ تمایزی بین محتوای دانشگاهی سازمانی تجاری یا مبتذل جمع آوری می کنند. این روش به راستی به

ص: 290

URL-1

International Internet Preservation Consortium (IIPC): <http://netpreserve.org> [last accessed on 2010-06-22]

[15] Heritrix crawler: <http://crawler.archive.org> [last accessed on 2010-06-15].

اسلوب واسپاری حقوقی است (به عنوان مثال فرض را بر این قرار نمی دهد آن چه را که مورد علاقه پژوهش گران در 100 سال آینده است، شناسایی کند). با این حال آرشیوهایی که با این روش ایجاد می شوند بسیار سطحی هستند؛ محتوای عمیق یا تحولات و بگاہ را حفظ نمی کنند.

● خزش های متمرکز (1): جمع آوری گزینشی خزش های وسیع را کامل می کند. کتابداران موضوعی وب گاه هایی را برای چنین خزش های متمرکز منطبق با طرح های توسعه مجموعه (همکاری با دیگر کتابخانه ها و پژوهشگران نیز ممکن است) انتخاب می کنند خزش های متمرکز می توانند واقعه محور (انتخابات فرانسه در سال های 2002، 2004، 2007، و 2009) یا موضوعی (خاطرات و بلاگ های شخصی توسعه، پایدار عمل گرایی وب و غیره) باشند خزش های متمرکز ساخت آرشیوهای کاملتر و بیشتر را از تعداد محدودی از وب گاه ها ممکن می سازند.

امروزه آرشیوهای وب کتابخانه ملی فرانسه یا آرشیوهای وب حوزه ملی فرانسه شامل 12/5 میلیارد .یو.آر.ال. می شود و 145 ترابایت فضای دیسک اشغال می کند. تاریخ قدیمی ترین صفحات وب به 1996 باز می گردد و این صفحات به لطف آرشیو ملی به دست آمدند که سازمانی غیر انتفاعی است و در نظر دارد یک کتابخانه اینترنتی بسازد و پیشگام آرشیو وب است آخرین تاریخ صفحات وب به چند ساعت قبل باز می گردد.

کشف منبع: خدمات و ابزارهای دسترسی برای کاربران نهایی

تعیین جای خدمت و محدودیت های دسترسی

دسترسی به این آرشیوها مانند دسترسی به اسناد فیزیکی واقع در ساختمان های کتابخانه نیست. وب گاه های گردآوری شده در فهرست کتابخانه ثبت نمی شوند زیرا مجموعه بسیار بزرگ و بسیار ناهمگن است؛ غیر ممکن خواهد بود فهرستی جامع از وب گاه های آرشیو شده ایجاد کرد؛ عنوان ها و محتوای مفصل شان را دانست در، عوض کتابخانه ملی فرانسه فرایندهای نمایه سازی خودکار را ساخته است تا دسترسی سریع به محتوای جمع آوری شده را ممکن سازد. هر فایلی تاریخ می خورد و توصیف می شود تا فقط اطلاعات ضروری (مکان اصلی در وب شکل، اندازه، تعیین جا در آرشیوها، و غیره) را جمع آوری کند. این فرایند نمایه سازی وب گاه های آرشیو شده را قادر می سازد مجدد در محیط انتشاراتی شان نقش داشته باشند و آن ها را با کلیک روی پیوندها درست شبیه وب در حال فعالیت اما در یک محتوای تاریخی و تاریخ خورده مرور کنند.

از آوریل 2008 آرشیوهای وب برای کاربران مجاز در اتاق های مطالعه کتابخانه پژوهشی، در مکان های متفاوت کتابخانه ملی فرانسه (طبقه همکف در کتابخانه فرانسوا میتران و بخش های مجموعه های تخصصی در ریشولیو، لوووا، اوپرا آرسنال و ژان - ویلار در آوینیون) (2) قابل دسترس هستند. گر چه آرشیوها اساساً شامل وب گاه های با دسترسی عمومی و رایگان هستند، این محدودیت وضع

ص: 291

شده است تا از فراهم آوری های حقوقی که به تمامی مجموعه های واسپاری و میراث حقوقی مربوط می شوند و در نظر دارند کاملاً به کپی رایت و قوانین شخصی احترام بگذارند، تبعیت کند.

برای دسترسی، توافقی کاربران نهایی باید بالاتر از 18 سال باشند و مدرکی جهت نیاز دانشگاهی فعالیت های حرفه ای یا شخصی پژوهشی شان برای دسترسی به این آرشیوها ارائه کنند (اکنون کتابخانه ملی فرانسه آخرین و تنها منبع این نوع سند است زیرا تنها کتابخانه ای است که این خدمت را در فرانسه ارائه می کند). کارت های خوانندگان توسط سرویس راهنمای خوانندگان در کتابخانه های فرانسوا - میتران یا ریشولیو پس از مصاحبه شخصی با یک کتابدار صادر می شود بر اساس نیازهای کاربران این مصاحبه تعیین می کند کاربران اجازه ورود به یک یا چند بخش را دارند و دوره اعتبار کارت (3 روز 15 روز سالانه) چقدر است

آرشیوهای وب همراه با تمامی خدمات الکترونیکی کتابخانه (وب گاه های اطلاعاتی، تسهیلات رزرو، فهرست ها، دسترسی به اینترنت) و منابع الکترونیکی (کتاب ها و تصاویر دیجیتالی، نشریات الکترونیکی پایگاه داده های آنلاین، لوح های فشرده، بوک مارک ها) روی 350 کامپیوتر واقع در اتاق های مطالعه کتابخانه پژوهشی در دسترس هستند این کامپیوترها در دسترس عموم هستند اما کاربران برای یک جا در کتابخانه و استفاده از کامپیوترها نیاز دارند از قبل جا ذخیره کنند

ابزارهای جستجو و دیدن

آرشیوهای وب از طریق برنامه کاربردی اختصاصی ای که مجموعه «آرشیوهای اینترنت» (1) نامیده می شود، قابل دسترس می شوند واژه «آرشیوها» استفاده می شود تا تأکید کند که این مجموعه ها جامع نیستند. برنامه کاربردی توسط یک نوار نارنجی با حروف WWW و یک نشانگر ماوس روی آن ارائه می شود تا نشان دهد گر چه صفحات وب در یک جعبه قرار می گیرند هنوز قابل کلیک هستند کتابخانه ملی فرانسه توسعه یک همانندی دیداری و یک علامت خاص (2) را انتخاب کرده است تا آرشیوهای وب را بیش تر قابل دیدن کند.

عکس

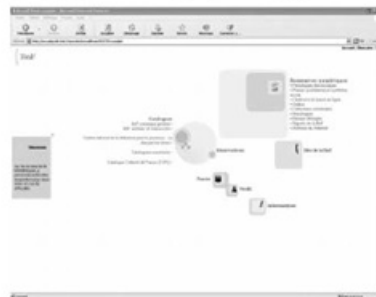
شده است تا از فراهم‌آوری‌های حقوقی که به تمامی مجموعه‌های واسپاری و میراث حقوقی مربوط می‌شوند و در نظر دارند کاملاً به کپی‌رایت و قوانین شخصی احترام بگذارند، تبعیت کند.

برای دسترسی توافقی، کاربران نهایی باید بالاتر از ۱۸ سال باشند و مدرکی جهت نیاز دانشگاهی، فعالیت‌های حرفه‌ای، یا شخصی پژوهشی‌شان برای دسترسی به این آرشیوها ارائه کنند (اکنون کتابخانه ملی فرانسه آخرین و تنها منبع این نوع سند است، زیرا تنها کتابخانه‌ای است که این خدمت را در فرانسه ارائه می‌کند). کارت‌های خوانندگان توسط سرویس راهنمای خوانندگان در کتابخانه‌های فرانسوا-میتران یا ریشولیو پس از مصاحبه شخصی با یک کتابدار صادر می‌شود. براساس نیازهای کاربران، این مصاحبه تعیین می‌کند کاربران اجازة ورود به یک یا چند بخش را دارند و دوره اعتبار کارت (۳ روز، ۱۵ روز، سالانه) چقدر است.

آرشیوهای وب همراه با تمامی خدمات الکترونیکی کتابخانه (وبگاه‌های اطلاعاتی، تسهیلات رزرو، فهرست‌ها، دسترسی به اینترنت) و منابع الکترونیکی (کتاب‌ها و تصاویر دیجیتالی، نشریات الکترونیکی، پایگاه‌داده‌های آنلاین، لوح‌های فشرده، بوک‌مارک‌ها) روی ۳۵۰ کامپیوتر واقع در اتاق‌های مطالعه کتابخانه پژوهشی در دسترس هستند. این کامپیوترها در دسترس عموم هستند، اما کاربران برای یک جا در کتابخانه و استفاده از کامپیوترها نیاز دارند از قبل جا ذخیره کنند.

ابزارهای جستجو و دیدن

آرشیوهای وب از طریق برنامه کاربردی اختصاصی‌ای که مجموعه «آرشیوهای اینترنت»^۱ نامیده می‌شود، قابل دسترس می‌شوند. واژه «آرشیوها» استفاده می‌شود تا تأکید کند که این مجموعه‌ها جامع نیستند. برنامه کاربردی توسط یک نواز نارنجی با حروف www و یک نشانگر ماوس روی آن ارائه می‌شود تا نشان دهد گرچه صفحات وب در یک جعبه قرار می‌گیرند، هنوز قابل کلیک هستند. کتابخانه ملی فرانسه توسعه یک همانندی دیداری و یک علامت خاص^۲ را انتخاب کرده است تا آرشیوهای وب را بیشتر قابل دیدن کند.



تصویر ۱. صفحه اصلی خدمات و منابع الکترونیکی کتابخانه پژوهشی کتابخانه ملی فرانسه

1. Archives de l'Internet
2. logo

تصویر ۱. صفحه اصلی خدمات و منابع الکترونیکی کتابخانه پژوهشی کتابخانه ملی فرانسه

برای مرور آرشیوها کتابخانه ملی فرانسه سه ابزار متفاوت ارائه می کند:

- جستجو با یو آر. ال.
 - جستجو با کلیدواژه
 - مجموعه های مشخصی که در نظر دارند کشف مجموعه های اختصاصی را تسهیل کنند.
- هر سه ابزار در برنامه کاربردی واحدی به نام آرشیو اینترنت منسجم شده است (تصویر 2)

عکس

برای مرور آرشیوها، کتابخانه ملی فرانسه سه ابزار متفاوت ارائه می‌کند:

- جستجو با یو.آر.آل.
 - جستجو با کلیدواژه.
 - مجموعه‌های مشخصی که در نظر دارند کشف مجموعه‌های اختصاصی را تسهیل کنند.
- هر سه ابزار در برنامه کاربردی واحدی به نام «آرشیو اینترنت» منسجم شده است (تصویر ۲).



تصویر ۲. صفحه خانگی عملکرد آرشیو اینترنت

- جستجوی یو.آر.آل.^۱ کاربران را قادر می‌سازد آرشیو یک وبگاه، یک صفحه وب یا حتی فایل را با وارد کردن جایگاه اصلی اینترنتی‌اش جستجو کنند. به عنوان مثال، جستجوی <http://www.lemonde.fr> ۶۱۷ نتیجه می‌دهد (از ۱۵ ژوئن ۲۰۱۰) که در یک نگاه تقویمی از ۱۹۹۶ تا ۲۰۱۰ نمایش داده می‌شوند (نتایج می‌توانند به تاریخ‌های خاصی محدود شوند). هر تاریخ قابل کلیک است و دسترسی به صفحه خانگی روزنامه لوموند می‌دهد که در آن تاریخ قابل دسترس بوده است. این نوع جستجو زمان جستجوی توسعه یک وبگاه و مقایسه ویرایش‌های آن بسیار مفید است، اما درست شبیه جستجوی یک شماره بازیابی در یک فهرست، کاربران باید یو.آر.آل. وبگاه یا صفحه وبی که می‌خواهند جستجو کنند، بدانند (به عنوان مثال، از یک کتابشناسی) یا باید قادر باشند آن را پیدا کنند (با استفاده از راهنماها یا پیوندهای بیرونی^۲ که به وبگاه‌های همکار ارجاع می‌دهند). کاربران نیز ممکن است وب در حال فعالیت و آرشیوهای وب را هم‌زمان در حال استفاده از دو مرورگر با ظاهری متفاوت، به نوبت مرور کنند. در حال حاضر، جستجوی یو.آر.آل. تنها روش جستجوی تمام آرشیوهاست.
- جستجوی کلیدواژه‌ای شبیه یک موتور جستجوی سنتی کار می‌کند: کاربران را قادر می‌سازد اسناد متنی را که شامل یک یا چند واژه است، جستجو نمایند. به لطف گزینه‌های جستجوی پیشرفته نیز کاربران ممکن است عین یک عبارت یا جمله را جستجو کنند یا جستجو را به یک وبگاه خاص

۱. URL کونه‌نوشت Uniform Resource Locator است و اشاره به جایگاه اینترنتی صفحه وب دارد.

2. out-links

تصویر ۲. صفحه خانگی عملکرد آرشیو اینترنت

جستجوی یو.آر.آل. (1) کاربران را قادر می‌سازد آرشیو یک وبگاه یک صفحه وب یا حتی فایل را با وارد کردن جایگاه اصلی اینترنتی‌اش جستجو کنند. به عنوان مثال جستجوی <http://www.lemonde.fr> 617 نتیجه می‌دهد (از 15 ژوئن 2010) که در یک نگاه تقویمی از 1996 تا 2010 نمایش داده می‌شوند (نتایج می‌توانند به تاریخ‌های خاصی محدود شوند). هر تاریخ قابل کلیک است و دسترسی به صفحه خانگی روزنامه لوموند می‌دهد که در آن تاریخ قابل دسترس بوده است. این نوع جستجو زمان جستجوی توسعه یک وبگاه و مقایسه ویرایش‌های آن بسیار مفید است اما درست شبیه جستجوی یک شماره بازیابی در یک فهرست، کاربران باید یو.آر.آل. وب‌گاه یا

صفحه وبی که می خواهند جستجو کنند، بدانند (به عنوان مثال از یک کتابشناسی) یا باید قادر باشند آن را پیدا کنند (با استفاده از راهنماها یا پیوندهای بیرونی (2) که به وب گاه های همکار ارجاع می دهند). کاربران نیز ممکن است وب در حال فعالیت و آرشیوهای وب را همزمان در حال استفاده از دو مرورگر با ظاهری متفاوت به نوبت مرور کنند. در حال حاضر، جستجوی یو. آر. ال. تنها روش جستجوی تمام آرشیوهاست.

● جستجوی کلید واژه ای شبیه یک موتور جستجوی سنتی کار می کند: کاربران را قادر می سازد اسناد متنی را که شامل یک یا چند واژه است جستجو نمایند. به لطف گزینه های جستجوی پیشرفته نیز کاربران ممکن است عین یک عبارت یا جمله را جستجو کنند یا جستجو را به یک وب گاه خاص

ص: 293

1- URL کوتاه نوشت Uniform Resource Locator است و اشاره به جایگاه اینترنتی صفحه وب دارد.

2- out-links

محدود کنند. این جستجو هنوز آزمایشی است و فقط پنج درصد آرشیو های وب کتابخانه ملی فرانسه را شامل می شود نمایه سازی تمام متن میلیاردها فایل بینظم با به حساب آوردن نسخه های تکراری و انسجام موقت که آن ها را به یکدیگر پیوند می دهد، یک چالش فنی است که چندین طرح پژوهشی بین المللی در تلاش اند تا آن را برطرف کنند اما هنوز موفق نشده اند.

• علاوه بر گزینه های جستجو مهمترین کارآمدی های این عملکرد عبارت اند از: توانایی نمایش صفحات وب با جمع آوری تعدادی عوامل که ممکن است در تاریخ های متفاوتی یا دست کم در زمان های مختلف آرشیو شده باشند و در نتیجه دوباره یک وبگاه مصنوعی ایجاد کنند.

• توانایی کلیک روی پیوندها با در نظر گرفتن انسجام موقت (به عنوان مثال اگر ما به وبگاه الیزه ریاست جمهوری فرانسه (1) در 7 می 2007 نگاه کنیم و روی یک پیوند کلیک کنیم تا وبگاه دولت را ببینیم انتظار داریم ویرایش ماه می 2007 را مشاهده کنیم).

این کارآمدی ها و جستجوی یو.آر.آل توسط ماشین وی یک منبع باز (2) پشتیبانی می شوند و توسط آرشیو اینترنت با مشارکت ICP توسعه می یابند جستجوی کلیدواژه ای بر پایه نرم افزار ناچ واکس منبع باز (3) (ناچ با الحاقات آرشیو وب) بنا شده است. همچنین این نرم افزار که توسط آرشیو اینترنت و گردهمایی کتابخانه های ملی نوردیک (4) توسعه یافته ابزاری است برای نمایه سازی و جستجوی آرشیوهای وب با استفاده از موتور جستجوی ناچ و الحاقات که آرشیوهای وب را جستجو می کند. کتابخانه ملی فرانسه جهت ساخت یک برنامه کاربردی بر اساس این ابزارها آن ها را در ابزارها و فرایندهای توسعه اش یکپارچه و سفارشی کرده و به کارآمدی های کوچک تر گسترش داده است (به عنوان مثال، اجرای طرحی برای سؤالات مداوم یو.آر.آل.).

مجموعه های ویژه

کتابداران موضوعی کتابخانه ملی فرانسه جهت جبران ابهام مجموعه ها و فقدان ابزارهای دسترسی (که هنوز توسعه می یابند) و نیز آگاهی در مورد مجموعه های موضوعی غنی، با همکاری پژوهشگران، مسیرهای هدایت شده (به معنی واقعی کلمه: تورهای هدایت شده یا مجموعه های ویژه در فرانسه) را ساخته اند که مجموعه ای از صفحات مصور و منظم ویرایشی است که شامل پیوندهای مستقیم به آرشیوهای وب می شود مقدمه ای ارائه می کند جستجو را شبیه سازی می نماید و ایده هایی در مورد دیگر جستجوهای ممکن می دهد از سال 2008 سه مجموعه ویژه منتشر شده است:

• کلیک کن و رأی بده: وب انتخاباتی: گزینشی از وب گاه ها از انواع تهیه کننده (سازمان ها، کاندیداها، حامیان، تماشاچیان، افراد و غیره) از مبارزات انتخاباتی، 2002، 2004 و 2007 تصویر 3؛

ص: 294

the French Presidency Elysée website – 1

The open source Wayback Machine, <http://archive-access.sourceforge.net/projects/wayback/> [last – 2
[accessed on 2010-06-15].

NutchWAX (Nutch Web Archive eXtensions), <http://archive-access.sourceforge.net/projects/nutch/> – 3
.[[last accessed on 2010-06-15]

● در مورد خود نوشتن در وب خاطرات شخصی و ادبی: در نظر دارد نتایج انتقال از کاغذ را به وب و روشی که بلاگ ها نوشته های شخصی ادبی و انتقادی را تغییر داده اند، بررسی کند؛

● عمل گرایی وب (1): نشان می دهد چطور پس از سال ها وب توسط عمل گراها به عنوان ابزار انتشار و، ارتباط ابزاری که قبول مسئولیت را تشویق می کند و مکانی برای، بحث مشارکت سازماندهی و عمل استفاده شده است.

عکس

معرفی آرشیوهای وب ... ۲۹۵

- در مورد خود نوشتن در وب: خاطرات شخصی و ادبی: در نظر دارد نتایج انتقال از کاغذ را به وب و روشی که بلاگ ها نوشته های شخصی، ادبی، و انتقادی را تغییر داده اند، بررسی کند؛
- عمل گرایی وب: نشان می دهد چطور پس از سال ها، وب توسط عمل گراها به عنوان ابزار انتشار و ارتباط ابزاری که قبول مسئولیت را تشویق می کند و مکانی برای بحث، مشارکت، سازماندهی، و عمل استفاده شده است.



تصویر ۳. مجموعه ویژه «کلیک کن و رأی بده: وب انتخاباتی»

زمانی که این مقاله نوشته می شود، سه مجموعه ویژه تر در دست تهیه است: توسعه پایدار، فیلم های غیر حرفه ای، و سفرنامه نویسی.

به عنوان یک خدمت مکمل، کاربران اجازه دارند صفحات وب را چاپ و کپی کنند و نمونه های متن را در یک فایل نوت پد^۱ درج کنند، محتوای نوت پد را چاپ کنند یا آن را از طریق ایمیل بفرستند. به علت محدودیت های حقوقی، در حال حاضر تسهیلاتی وجود ندارد که کاربران را قادر سازد از صفحه عکس بگیرند یا عناصر وب را از آرشیوها بردارند.

طرح های توسعه آتی، نمایه سازی تمام متن و جستجو به علاوه دسترسی ناپیوسته و بگانه از کتابخانه های منطقه ای را همراه با تسهیم مأموریت و اسپاری حقوقی با کتابخانه ملی فرانسه در اطراف کشور دربرمی گیرد.

1. Web activism
2. note pad

تصویر 3. مجموعه ویژه «کلیک کن و رأی بده: وب انتخاباتی»

زمانی که این مقاله نوشته می شود سه مجموعه ویژه تر در دست تهیه است توسعه پایدار فیلم های غیر حرفه ای، و سفرنامه نویسی.

به عنوان یک خدمت مگمل، کاربران اجازه دارند صفحات وب را چاپ و کپی کنند و نمونه های متن را در یک فایل نوت پد (2) درج کنند محتوای نوت پد را چاپ کنند یا آن را از طریق ایمیل بفرستند. به علت محدودیت های حقوقی در حال حاضر تسهیلاتی وجود ندارد که کاربران را قادر سازد از صفحه عکس بگیرند یا عناصر وب را از آرشیوها بردارند.

طرح های توسعه آتی نمایه سازی تمام متن و جستجو به علاوه دسترسی ناپیوسته وب گاه از کتابخانه های منطقه ای را همراه با تسهیم مأموریت و اسپاری حقوقی با کتابخانه ملی فرانسه در اطراف کشور در بر می گیرد.

ص: 295

Web activism -1

note pad -2

ارزش آرشیوهای وب و علاقه به آن ها فقط با گذشت زمان ثابت می شود پس از اینکه منابع وب از وب ناپدید شدند؛ این بخشی از فلسفه میراث کتابخانه ملی فرانسه و سنتی نیست که انتظار یک بازگشت سرمایه کوتاه مدت را داشته باشد؛ رسالتش فهم اهمیت زمان است.

با این حال نخستین ارزیابی کاربران و کاربرد آرشیوهای وب مرحله ای مهم در اجرای واسپاری حقوقی اینترنت است. اثبات سودمندی عمومی و علمی مجموعه های جمع آوری شده و نیز توسعه و تحلیل نخستین آمارهای کاربرد ما را قادر خواهد ساخت توان بالقوه این مجموعه ها و نیز محدودیت های شان را اندازه گیری کنیم این تحلیل ها در رویارویی با انتظارات پژوهشگران برای توسعه مجموعه و ابزار هر دو مفید خواهند بود

تحلیل کمی

ابزاری تحلیلی به نام AWStats ایجاد شده است تا ترافیک همزمان عملکرد را تحلیل کند. شبیه ماشین وی بک و ناچ واکس درون سازمانی سفارشی ساخته شده است تا بین استفاده عمومی توسط خوانندگان، استفاده مرجع توسط کتابداران در هنگام ورود یا در میز مرجع و استفاده حرفه ای توسط کتابداران موضوعی کتابخانه ملی فرانسه تمایز ایجاد نماید

عکس

کاربرد: اطلاعات و ارقام

ارزش آرشیوهای وب و علاقه به آنها فقط با گذشت زمان ثابت می‌شود، پس از اینکه منابع وب از وب ناپدید شدند؛ این بخشی از فلسفه میراث کتابخانه ملی فرانسه و سستی نیست که انتظار یک بازگشت سرمایه کوتاه‌مدت را داشته باشد؛ رسالتش فهم اهمیت زمان است.

با این حال، نخستین ارزیابی کاربران و کاربرد آرشیوهای وب مرحله‌ای مهم در اجرای واسپاری حقوقی اینترنت است. اثبات سودمندی عمومی و علمی مجموعه‌های جمع‌آوری شده و نیز توسعه و تحلیل نخستین آمارهای کاربرد، ما را قادر خواهد ساخت توان بالقوه این مجموعه‌ها و نیز محدودیت‌هایشان را اندازه‌گیری کنیم. این تحلیل‌ها در رویارویی با انتظارات پژوهشگران برای توسعه مجموعه و ابزار هر دو مفید خواهند بود.

تحلیل کمی

ابزاری تحلیلی به نام AWStats ایجاد شده است تا ترافیک هم‌زمان عملکرد را تحلیل کند. شبیه ماشین وی‌بک و ناچ‌واکس، درون‌سازمانی سفارشی ساخته شده است تا بین استفاده عمومی توسط خوانندگان، استفاده مرجع توسط کتابداران در هنگام ورود یا در میز مرجع، و استفاده حرفه‌ای توسط کتابداران موضوعی کتابخانه ملی فرانسه تمایز ایجاد نماید.

جدول ۱. شاخص‌های کاربرد اصلی در سال ۲۰۰۸ و ۲۰۰۹

۲۰۰۹	۲۰۰۸	
۱۰۶	۳۵	میانگین تعداد جلسه‌ها در هر ماه
۱۲۷۵	۳۱۶	تعداد کل جلسه‌ها
۹۰۰۶۳	۳۵۸۹۱	تعداد کل صفحات مشاهده شده

جدول ۲. کمیت‌های جامع در سال ۲۰۰۹

صفحات مشاهده شده	جلسه‌های طولانی (مشاهده‌های بیشتر از ۱ ساعت)	جلسه‌ها (مشاهده‌های بیشتر از ۵ دقیقه)	مشاهده‌کنندگان	مشاهده‌ها	
۱۴۰۴۵	۷۱	۳۳۸	۴۸۱	۹۰۴	عمومی
۶۲۴۲	۲۳	۹۴	۱۱۸	۲۱۶	مرجع
۶۹۷۷۶	۲۳۹	۸۴۳	۵۵۲	۱۹۴۵	حرفه‌ای
۹۰۰۶۳	۳۳۳	۱۲۷۵	۱۲۳۵	۳۰۶۵	کل

جدول ۱. شاخص‌های کاربرد اصلی در سال ۲۰۰۸ و ۲۰۰۹

جدول ۲. کمیت‌های جامع در سال ۲۰۰۹

دو نتیجه ای که توجه به کتابخانه ملی فرانسه را جلب می کند عبارت اند از:

• تعداد کل جلسه ها بین سال 2008 و 2009 سه بار افزایش یافته است (بعداً خواهیم دید که آموزش اصلی و ابتکار عمل های اطلاعاتی برای خوانندگان و کتابداران مرجع توسعه یافته است)؛

• تعداد فزاینده ای از کاربران نهایی وجود دارند که آرشیوها را برای پژوهش عمیق استفاده می کنند. میانگین جلسه در فوریه، 2009، 13 دقیقه طول کشید و در دسامبر 2009 تا 30 دقیقه افزایش یافت بین فوریه و سپتامبر 2009 بیش از دو جلسه وجود نداشت که بیش از یک ساعت یا بیشتر طول بکشد. سپس، تعداد هر جلسه در اکتبر 2009 به 7 جلسه در نوامبر به 32 جلسه، در دسامبر به 22، جلسه و در ژانویه 2010 به 23 جلسه افزایش یافت این جلسه ها نشان می دهد که اجرای پژوهش وسیع روی آرشیوهای وب آغاز می شود با نگاهی به فهرست بیش ترین صفحات مشاهده شده به نظر می رسد این طرح ها بیش از همه توسط پژوهشگران علوم اجتماعی و علوم سیاسی اجرا می شود.

تحلیل کیفی

در تکمیل تحلیل آماری و فهرستی جامع از یو.آر.ال.ها و کلید واژه های جست و جو شده، کتابخانه ملی فرانسه یک لاگ داخلی و مشترک مورد استفاده کتابدارانی ایجاد کرد که کاربران را به مناسبت های مختلف از جمله مصاحبه پذیرش میز، مرجع درخواست ملاقات، شخصی، نمایش به تهیه کنندگان وب گاه ها و غیره ملاقات می کنند از آوریل 2008 کتابداران 34 کاربر نهایی را با وارد کردن اطلاعاتی همچون تاریخ/ساعت گروه و نام کتابدار شکل در خواست (در محل، تلفنی، ایمیل، ...)، نام نوع، کارت موضوع پژوهش، یادداشت ها /سؤال ها نظرات/ اظهارات، ثبت نام کردند (تصویر 4).

عکس

دو نتیجه‌ای که توجه به کتابخانه ملی فرانسه را جلب می‌کند عبارت‌اند از:

- تعداد کل جلسه‌ها بین سال ۲۰۰۸ و ۲۰۰۹ سه بار افزایش یافته است (بعدها خواهیم دید که آموزش اصلی و ابتکار عمل‌های اطلاعاتی برای خوانندگان و کتابداران مرجع توسعه یافته است)؛
- تعداد فزاینده‌ای از کاربران نهایی وجود دارند که آرشیوها را برای پژوهش عمیق استفاده می‌کنند. میانگین جلسه در فوریه ۲۰۰۹، ۱۳ دقیقه طول کشید و در دسامبر ۲۰۰۹ تا ۳۰ دقیقه افزایش یافت. بین فوریه و سپتامبر ۲۰۰۹، بیش از دو جلسه وجود نداشت که بیش از یک ساعت یا بیشتر طول بکشد. سپس، تعداد هر جلسه در اکتبر ۲۰۰۹ به ۷ جلسه، در نوامبر به ۳۲ جلسه، در دسامبر به ۲۲ جلسه، و در ژانویه ۲۰۱۰ به ۲۳ جلسه افزایش یافت. این جلسه‌ها نشان می‌دهد که اجرای پژوهش وسیع روی آرشیوهای وب آغاز می‌شود. با نگاهی به فهرست بیشترین صفحات مشاهده شده، به‌منظر می‌رسد این طرح‌ها بیش از همه توسط پژوهشگران علوم اجتماعی و علوم سیاسی اجرا می‌شود.

تحلیل کیفی

در تکمیل تحلیل آماری و فهرستی جامع از یو.آر.ال‌ها و کلیدواژه‌های جست‌وجوشده، کتابخانه ملی فرانسه یک لاگ داخلی و مشترک مورد استفاده کتابدارانی ایجاد کرد که کاربران را به مناسبتهای مختلف از جمله مصاحبه پذیرش، میز مرجع، درخواست ملاقات شخصی، نمایش به تهیه‌کنندگان وبگاه‌ها، و غیره ملاقات می‌کنند. از آوریل ۲۰۰۸، کتابداران ۳۴ کاربر نهایی را با وارد کردن اطلاعاتی همچون تاریخ/ساعت، گروه و نام کتابدار، شکل درخواست (در محل، تلفنی، ایمیل، ...)، نام، نوع کارت، موضوع پژوهش، یادداشت‌ها/سؤالها/نظرات/اظهارات، ثبت‌نام کردند (تصویر ۴).

سوالات	پاسخ‌ها
تاریخ و ساعت	می / ژانویه / اوت ۲۰۰۹
بخش یا خدمت و نام کارگزار	کریستین ژن
پشتیبانی از درخواست	در محل
نام کوچک و نام خانوادگی خواننده	
عنوان دسترسی	نمایش‌ها
موضوع پژوهش / رشته	<p>۲۹ می (۱۷ ساعت و نیم): الیزابت لگرو، بلاگر و عضو گروه یار APA http://2009sediments.wordpress.com/ ۲۰ ژانویه (۱۶ ساعت): مارتین سونه، پژوهشگر CNRS، نویسنده و بلاگر www.martinesonnet.fr/blogwp ۱۱ اوت (۱۷ ساعت): سیلور مرسیه که یادداشتی را در بلاگش نوشته است: http://www.bibliosession.net/2009/09/17/archives-de-liternet-demadez-votre-ticket-pour-la-posterite/ ۲۱ اوت (۱۷ ساعت): اوریان دزینی که در Paris XIII تدریس می‌کند و نویسنده رساله‌ای در مورد روزنامه‌های خصوصی آنلاین وی به‌ویژه علاقه‌مند به مسیر هدایت‌شده است و از یک همکار بانک سالن در مورد این موضوع سؤال کرده است که آن را به‌سوی من ارجاع داده است.</p>
مشاهدات	

سؤالات	پاسخها
تاریخ و ساعت	۲۴ نوامبر ۲۰۰۹
بخش یا خدمت و نام کارگزار	SOL کریستوف تربویی
پشتیبانی از درخواست	در یک محل رسمی
نام کوچک و نام خانوادگی خواننده	فابین گرفه
عنوان دسترسی	کارت سالانه (۲۰۰۹/۱۱/۲۴ تحویل شد)
موضوع پژوهش / رشته	مبارزه انتخاباتی اروپایی سال ۲۰۰۹ در اینترنت
مشاهدات	

سؤالات	پاسخها
تاریخ و ساعت	۲۰۰۹/۸/۱۹
بخش یا خدمت و نام کارگزار	SOL
پشتیبانی از درخواست	در محل
نام کوچک و نام خانوادگی خواننده	کارول دافینی
عنوان دسترسی	کارت سالانه
موضوع پژوهش / رشته	آرشیوهای وب (کارآموز محافظه کار)
مشاهدات	

تصویر ۴. لاگ کاربران نهایی آرشیوهای وب

این ابزار ساده، کتابخانه ملی فرانسه را قادر ساخت مجموعه‌هایی را که مورد درخواست هستند، شناسایی کند؛ این مجموعه‌ها در حال حاضر در حوزه علوم اجتماعی و علوم سیاسی هستند. روزی که کتابخانه ملی فرانسه خدمت را عرضه کرد، یک دانشجوی کارشناسی ارشد که در مورد «اینترنت و انتخابات ریاست جمهوری ۲۰۰۷» کار می‌کرد، درخواست دسترسی کرد. یک دانشجوی دکترای زبان‌شناسی که در مورد تحلیل سخنرانی‌های کاندیداهای زن کار می‌کرد از روم، ایتالیا آمده بود تا در مورد بلاگ رهبر کمونیست، ماری - ژرژ بوفه^۱ که بعد از انتخابات ریاست جمهوری بسته شده بود، کار کند. فعالیت وب نیکلا سارکوزی، حزب سوسیالیست، انتخابات اروپایی ۲۰۰۹، استفاده از ویدئوها در فعالیت انتخاباتی، کارتونها و کاریکاتورهای سیاسی، نمونه‌هایی از پژوهش در این حوزه‌ها هستند.

پژوهش‌های دیگر شامل موضوعاتی بوده است همچون جستجوی اطلاعات در اتحادیه اروپا، آرشیوهای وزارت بوم‌شناسی و موجودیت‌های نامتمرکز، وب‌سایت‌های نویسندگان غیر حرفه‌ای، وب‌سایت‌های مدیریت استرس، وب‌سایت‌های شخصی و سؤالاتی درباره اینکه چرا و چگونه

1. Marie-George Buffet

تصویر ۴. لاگ کاربران نهایی آرشیوهای وب

این ابزار ساده کتابخانه ملی فرانسه را قادر ساخت مجموعه‌هایی را که مورد درخواست هستند شناسایی کند؛ این مجموعه‌ها در حال حاضر در حوزه علوم اجتماعی و علوم سیاسی هستند. روزی که کتابخانه ملی فرانسه خدمت را عرضه کرد یک دانشجوی کارشناسی ارشد که در مورد «اینترنت و انتخابات ریاست جمهوری ۲۰۰۷» کار می‌کرد درخواست دسترسی کرد. یک دانشجوی دکترای زبان‌شناسی که

در مورد تحلیل سخنرانی های کاندیداهای زن کار می کرد از، روم ایتالیا آمده بود تا در مورد بلاگ رهبر کمونیست ماری - ژرژ بوفه (1) که بعد از انتخابات ریاست جمهوری بسته شده بود، کار کند. فعالیت وب نیکلا سارکوزی، حزب سوسیالیست انتخابات اروپایی 2009 استفاده از ویدئوها در فعالیت انتخاباتی، کارتون ها و کاریکاتورهای، سیاسی نمونه هایی از پژوهش در این حوزه ها هستند.

پژوهش های دیگر شامل موضوعاتی بوده است همچون جستجوی اطلاعات در اتحادیه اروپا آرشیوهای وزارت بوم شناسی و موجودیت های نامتمرکزش وب سایت های نویسندگان غیر حرفه ای وب سایت های مدیریت استرس وب سایت های شخصی و سؤالاتی درباره این که چرا و چگونه

ص: 298

Marie-George Buffet -1

وب سایت ها جمع آوری می شوند و قوانین رقابتی تا در یک مورد حقوقی کمک کند.

نظر سنجی ها

نخستین نظر سنجی از اکتبر 2006 تا ژوئن 2007 در محتوای یک دوره درسی در مقطع کارشناسی ارشد به نام «اینترنت در طول فعالیت» انجام شد و توسط کتابخانه ملی فرانسه و یک پژوهشگر دانشگاه در رسانه اجتماعی به طور هماهنگ سازماندهی شد. 17 جلسه مشاهده و 5 مصاحبه انجام شد تا نیازهای کاربران به ابزارها و کارآمدی ها و نگرش شان به یک نوع رسانه جدید شناسایی شوند (برای آزمون نهایی آن ها، دانشجویان مجبور شدند مقاله ای بنویسند در موضوعی که مربوط به منبع وب در حال فعالیت و محتوای آرشیو شده اش بود).

دومین نظر سنجی برای نوامبر 2010 برنامه ریزی می شود این نظر سنجی کاربران فعلی و بالقوه را در نظر دارد و قصد دارد نیازهای شان را به محتوای مجموعه تعریف و ابزار توسعه کاربرد را شناسایی کند (جدول 3).

عکس

وبسایت‌ها جمع‌آوری می‌شوند و قوانین رقابتی تا در یک مورد حقوقی کمک کند.

نظرسنجی‌ها

نخستین نظرسنجی از اکتبر ۲۰۰۶ تا ژوئن ۲۰۰۷ در محتوای یک دوره درسی در مقطع کارشناسی ارشد به نام «اینترنت در طول فعالیت» انجام شد و توسط کتابخانه ملی فرانسه و یک پژوهشگر دانشگاه در رسانه اجتماعی به‌طور هماهنگ سازماندهی شد. ۱۷ جلسه مشاهده و ۵ مصاحبه انجام شد تا نیازهای کاربران به ابزارها و کارآمدی‌ها و نگرش‌شان به یک نوع رسانه جدید شناسایی شوند (برای آزمون نهایی آنها، دانشجویان مجبور شدند مقاله‌ای بنویسند در موضوعی که مربوط به منبع وب در حال فعالیت و محتوای آرشیو شده‌اش بود).

دومین نظرسنجی برای نوامبر ۲۰۱۰ برنامه‌ریزی می‌شود. این نظرسنجی کاربران فعلی و بالقوه را در نظر دارد و قصد دارد نیازهایشان را به محتوای مجموعه تعریف و ابزار توسعه کاربرد را شناسایی کند (جدول ۳).

جدول ۳. چارچوبی برای نظرسنجی بعدی کاربرد آرشیو وب در نوامبر ۲۰۱۰

کاربران بالقوه آرشیو وب	کاربران فعلی آرشیو وب
<p>۵-۸ مصاحبه برای هر گروه هدف:</p> <ul style="list-style-type: none"> • دانشمندان در دیگر حوزه‌های پژوهشی (هنرهای دیجیتال، علوم انسانی، پژوهشگران رسانه‌ای علاقه‌مند به خود وب، و غیره). • یک گروه وسیع‌تر علاقه‌مند به حافظه وب، وکلا و متخصصان اطلاعات. 	<p>۵-۸ مصاحبه</p>
<p>سؤالاتی درباره:</p> <ul style="list-style-type: none"> • آگاهی از آرشیوهای وب • علاقه شناسایی شده به خود حوزه پژوهشی • انواع اطلاعات جستجو شده اگر مورد نیاز هستند • ضروریات کاربرد واقعی 	<p>سؤالاتی درباره:</p> <ul style="list-style-type: none"> • خط‌مشی توسعه مجموعه: گزینش، تناوب و کیفیت جمع‌آوری • کاربران هدف و انواع استفاده برای توسعه • علاقه و موضوعاتی برای مجموعه‌های ویژه • اطلاعات و متادیتاهای ارائه شده در نمایش نتایج و نمایش یک صفحه

مشارکت: راهبردهایی برای تأسیس یک انجمن کتابداران و رسیدن به کاربران نهایی

معرفی آرشیوهای وب در کتابخانه به معنی ساخت یک مفهوم کلی از مالکیت این نوع جدید مجموعه توسط هر دو کارکنان و کاربران نهایی است.

مجموعه‌ها به شیوه‌های مختلفی باعث سرخوردگی می‌شوند (ذخیره‌های وب ناقص هستند و برخی وبگاه‌ها در آرشیوها هم‌زمان ظاهر نمی‌شوند). مانند سرخوردگی عملکرد جستجو: برای کارکنان، به‌ویژه، فقدان یک فهرست و این حقیقت که کتابداران هر کار نظام‌مند توصیفی را در مورد وبگاه‌ها انجام

جدول 3. چارچوبی برای نظرسنجی بعدی کاربرد آرشیو وب در نوامبر 2010

مشارکت: راهبردهایی برای تأسیس یک انجمن کتابداران و رسیدن به کاربران نهایی

معرفی آرشیوهای وب در کتابخانه به معنی ساخت یک مفهوم کلی از مالکیت این نوع جدید مجموعه توسط هر دو کارکنان و کاربران نهایی است.

مجموعه ها به شیوه های مختلفی باعث سرخوردگی می شوند (ذخیره های وب ناقص هستند و برخی وب گاه ها در آرشیوها همزمان ظاهر نمی شوند). مانند سرخوردگی عملکرد جستجو: برای کارکنان، به ویژه، فقدان یک فهرست و این حقیقت که کتابداران هر کار نظام مند توصیفی را در مورد وب گاه ها انجام

ص: 299

نمی دهند برای آنان سخت کرده است که آرشیوهای وب را به عنوان مجموعه های میراث با ارزش مورد توجه قرار دهند (چقدر احتمال دارد یک کتابدار بتواند در مورد یک مجموعه ساخته شده توسط نرم افزار خزش گر و نه توسط خودش احساس مثبتی داشته باشد؟)

بنابراین چالش این است که کیفیت های متمایز آرشیوهای وب - به خصوص این حقیقت که آن ها تنها گواه واحد باقیمانده از تغییرات اساسی جامعه ما هستند که در 15 سال گذشته تجربه شده اند و گذارشان از آنالوگ به دیجیتال - را مشخص کرد، اما همچنین مشخص کردن الگوهای این مجموعه که در واقع آن را کمی شبیه به مجموعه های منظمی می کند که کتابداران یاد گرفته اند در دست بگیرند: مربوط می شود به شعار "دیجیتال متفاوت نیست" - و واژه را منتشر می کند.

گفته می شود که تجارت معمول آرشیو وب ساختن ابزار اولیه استفاده از ارتباط استاندارد، سازماندهی و راهبردهای بازاریابی است.

راهبردهای به حساب آوردن آرشیوهای وب به عنوان بخشی از کار روزانه کتابخانه

سازمان. مدیریت آرشیوهای وب نباید فقط به عنوان یک فعالیت فنی در نظام رسمی کتابخانه دیده شود بیش تر سازمان هایی که به آرشیو کردن وب به عنوان یک مورد فنی نگاه کرده اند، موضوعی تحت آی.تی. که نیاز به مهندس و رهبری توسعه نرم افزار دارد به سختی می توانند کتابداران را درگیر ارتقای آرشیوها کنند برای این که کار آرشیو کردن وب کاری معمول به حساب آید کتابخانه ملی فرانسه تصمیم گرفت فعالیت هایش را در قدیمی ترین واحد تولید، کتابخانه بخش واسپاری حقوقی به انجام برساند. وب گاه ها اکنون با منابع چاپی واسپاری حقوقی (کتاب ها و نشریات ادواری) در واحدی که سال 2008 به نام واحد «واسپاری حقوقی دیجیتال» (1) ایجاد شد، مدیریت می شوند. این واحد شامل گروهی پنج نفره می شود که فعالیت ها را با مشارکت متخصصان آی.تی. از یک طرف و متخصصان مجموعه از سویی دیگر اجرا می کنند این اجرا در بخشی بزرگ و قدیمی کمک بسیاری کرده است تا متخصصان و فعالیت های آرشیو کردن وب را در رسالت اصلی و تاریخ کتابخانه ملی فرانسه بگنجانند. به عبارت دیگر، این نوع سازماندهی به اتصال قدیم به جدید کمک کرده است و نمایش می دهد که راهبردهای نمونه برداری با مقیاس بزرگ برای حجم های عظیم داده های دیجیتال متولد شده خیلی متفاوت نیست از راهی که تاریخ سنت واسپاری حقوقی منابع چاپی فرانسوی را طی پنج سال گذشته شکل داده است. این نگرش بخشی از یک تلاش وسیع تر توسط کتابخانه ملی فرانسه است تا سازماندهی مهارت های کارکنان را با تغییر دیجیتال انطباق دهد. (2)

شبکه سازی. ساخت شبکه ای از متخصصان مجموعه موضوعی از سراسر کتابخانه مرحله سرنوشت ساز دیگری جهت تشویق مشارکت های فعال و پذیرش آرشیوهای وب به عنوان مجموعه ای

ص: 300

1- digital legal deposit

2- The Human Face of Digital Preservation: Organizational and Staff Challenges, and Initiatives at the 'Bibliothèque nationale de France', Emmanuelle Bermès, Louise Faudet, iPres 2009, [http://www.escholarship.org/uc/cdl_ipres09 [last accessed on 2010-06-15

توسط کتابداران و مدیران بود. در سال 2005 این شبکه به عنوان گروهی متشکل از 20 پیشاهنگ شروع به کار کرد. تا سال 2010 حدود 80 کتابدار به گونه ای درگیر جمع آوری گزینشی بودند، به عنوان مثال، گزینش وبگاه کنترل کیفی و ارتقاء منابع و خدمات برای عموم نخست برای مدیران بخش مجموعه آگاهی ایجاد شد. در نتیجه، بیش تر بخش های موضوعی کتابداران موضوعی را برای طرح آرشیو کردن وب اختصاص دادند اکنون در هر حوزه اصلی (هنر ادبیات فلسفه علوم و غیره) از جمله موضوعاتی که مجموعه های نادر را در بر می گیرد (مانند نقشه ها یا موسیقی) یک کارگزار واسپاری حقوقی وب (1) وجود دارد، فردی که از آموزش خاص بهره مند بوده و دانش اساسی آرشیو کردن وب را کسب کرده است طوری که بتواند در حوزه تخصصی از آن استفاده کند برخی کتابداران موضوعی کتابخانه ملی فرانسه اکنون کتاب ها و وب گاه ها را به دست می آورند از سال 2008 به بعد تعداد کارکنان این شبکه کارگزاران مرحله واسپاری حقوقی وب که به خصوص به عنوان کتابداران مرجع در نظر گرفته شده بودند، افزایش یافته است. این جا هدف توسعه دانش و علاقه به آرشیوهای وب بود در هر کس که با عامه مردم خواه در اتاق های مطالعه یا به صورت آنلاین در ارتباط است.

ابزارهایی برای اشاعه داخلی. کتابخانه ملی فرانسه خوش شانس است که قادر است از مجموعه منابع مهمی برای ارتقای اطلاعات و مهارت ها به طور بین المللی سود برد: آموزش کارکنان، گفتگوهای داخلی ماهانه، کارکنان و کنفرانس هایی در وقت، ناهار یک مجله چاپی داخلی ماهانه و البته اینترنت همگی، قبل از شروع برنامه آرشیو کردن وب در دسترس بودند. چالش این بود که از تمامی این منابع ارتباطی استفاده و آن ها را به مسیری هوشمند هدایت کند تا کارکنان در مورد آرشیوهای وب در زمان و با روش مناسب بشنوند به عنوان مثال کتابخانه ملی فرانسه درست پایان مبارزات ریاست جمهوری سال 2007 یک نمایش اساسی داخلی از مجموعه وب گاه های منتخب برگزار کرد. همچنین مقالات مجله داخلی یا مقالاتی که در اینترنت کتابخانه ملی فرانسه بود یک طرح ارتباطی ای را دنبال کرد که منافع و علائق کلی کتابخانه را به حساب می آورد به عنوان مثال زمانی که کتابخانه مجموعه فعالیت هایی را که در راستای بالا بردن آگاهی نسبت به توسعه پایدار بود عملی کرد (و کارکنان را تشویق کرد که رفتاری مسئولیت پذیر در این حوزه به عهده بگیرند) گروه آرشیو کردن وب ارتباطات داخلی اش را روی مجموعه های مرتبط با توسعه پایدار نیز افزایش داد در یک کلام طرح ارتباط داخلی جهت ارتقاء آرشیوها هرگز به عنوان یک فرایند مستقل دیده نشد. ما تمامی فرصت های ممکن را به کار گرفتیم تا آرشیوهای وب رل در برنامه کلی کار کتابخانه بگنجانیم بیش از آن که تلاش کنیم این برنامه کار را با تحمیلی دیدن آن بر هم زنیم.

راهبردهایی برای رسیدن به کاربران نهایی

بحث ها: «آیا اصلاً آرشیوهای وب استفاده می شوند؟» «چرا منابع را صرف این داده های بیهوده کنیم در حالی که ما هیچ تضمینی نداریم که آیا این نوع ماده مورد علاقه عموم قرار خواهد گرفت؟» «چرا تلاش های مان را روی فراهم آوری یا دیجیتالی کردن موادی که می دانیم بدون شک میراثی برای نسل های

ص: 301

آینده، هستند متمرکز نمی کنیم؟» - این ها سؤالات معمولی است که خواه مربوط به مدیریت باشند یا مربوط به رسانه ها بیشتر اوقات نیاز به پاسخ دارند. از جهتی با نگاهی به گذشته تاریخ کتابخانه، چنین سؤال هایی جدید نیست در طول، زمان هر رسانه جدیدی سؤالات مشابهی برانگیخته است. دلیل آن این است که این امر زمان می برد - زمانی برای چیزهایی که از عموم فضای کسب و کار رایج ناپدید می شوند - قبل از این که پژوهشگران یا افراد متوجه شوند آن ها را از دست می دهند آنان به این مستندات نیاز دارند تا تاریخ جامعه را تشریح کنند.

روی هم رفته، موقعیت ویژه ای برای یک سازمان میراثی است تا برنامه آرشیو کردن وب خود را اجرا کند: گرایش به این وجود دارد که بروندهای ملموس را با مفید و با ارزش نشان دادن آن ها نمایش دهد در حالی که همزمان کاربرد نمی تواند بلافاصله توسعه یابد هنوز برای این کار خیلی زود است زیرا بیش تر کاربران آرشیوهای وب هنوز دارند متولد می شوند راهبردهای کتابخانه ملی فرانسه جهت کنترل این موقعیت دو جانبه است از یک سو با ایجاد ابزارها و آمارهایی که کتابخانه را قادر می سازد توسعه کاربرد را به روشی مشابه همان گونه که برای مجموعه های عادی انتظار می رفت نمایش دهد (دیجیتالی یا غیر دیجیتالی؛ به بخش بالا مراجعه کنید)، از سوی دیگر توجه عموم را به بحث در مورد قابلیت رؤیت از خارج از کتابخانه جلب کند.

ایجاد یک بحث عمومی. گر چه تعداد زیادی از مردم اکنون از آرشیوهای وب استفاده نمی کنند، در واقع بیشتر آنان علاقه به این مسأله را زمانی که از آن ها در این مورد سؤال شد، نشان داده اند. دلیل این امر این است که اینترنت زندگی های شخصی و عمومی افراد را تحت تأثیر قرار داده است بنابراین هر کس چیزی دارد که درباره آن بگوید. نخستین واکنش به طور معمول این است: «من هرگز در مورد آن فکر نکردم اما اکنون که شما به من می گوئید تصور می کنم مهم است این حافظه ماست و هر روز از یاد می رود» ارتباطات کتابخانه ملی فرانسه این مردم را برای توسعه حمایت عمومی از برنامه آرشیو کردن وب کتابخانه هدف قرار می دهد - آنان که ممکن است امروز برای دسترسی به مجموعه ها سروکله شان در کتابخانه پیدا نشود، اما این آگاهی را به دست آوردند که کاری است که باید انجام شود.

زنجیره کنفرانس حافظه های وب (1) به عنوان یک گردهمایی عمومی جهت فراهم آوردن این نوع حمایت و قابلیت رؤیت طراحی و در مارچ 2010 اجرا شد هر کنفرانسی یک نصف روز طول می کشد و سه نوع شرکت کننده را که در ترکیب مخاطبان منعکس می شوند نیز گردهم می آورد: پژوهشگران، وب ناشران و مسئولان همه کسانی که از آن ها درخواست می شود به موضوع مشابه بپردازند به عنوان مثال، عمل گرایی و سیاست وب خاطرات، وب حمایت از داده های شخصی در مقابل وسعت فضای عمومی، و غیره. کتابخانه ملی فرانسه نزدیک به 100 شرکت کننده را برای هر یک از این رویدادها گردهم آورده است و پوشش رسانه ای مناسبی را (بیشتر بر روی وب) برای دو نشست نخست دریافت کرده است. هدف این است که آرشیو کردن وب یک موضوع بحث عمومی در خارج از کتابخانه اما به همراه کتابخانه شود. این شکل ارتقا (که باز خورد عالی ای را برای ساخت مجموعه و خدمات به پژوهشگران نیز به دنبال دارد)

ص: 302

فعالیت های ارتقاء سازمان یافته مستقیم تر را برای کاربران نهایی امروز کتابخانه ملی فرانسه تکمیل می کند.

بیشتر کاربران آرشیوهای وب هنوز وجود ندارند آنان به احتمال با نسل های آینده ای می آیند که به دنیا آمده اند و با وب به دنیا خواهند آمد و آنان که از آن استفاده می کنند یا آن را به عنوان ابزار اصلی اطلاعاتی و ارتباطی شان استفاده خواهند کرد برای چنین مجموعه ای باید بپذیریم که فقط زمان است که به ما خواهد گفت آیا ما انتخاب درست را انجام می دهیم مدیریت دسترسی و گزینش نیز به معنی مدیریت خطرات است. ما مجبوریم آن چه که امروز غیر منتظره و ناخواسته است اما فردا ممکن است مورد علاقه، باشد بگیریم و ارتقا دهیم. همچنین مجبوریم کاربران و کاربرد فردا را در نظر بگیریم. این مسأله جدید نیست: چالشی است که مؤسسه میراثی قدیمی با آن آشناست. در این میان مؤسسات هنوز به راهبردهای میان مدت برای ایجاد اطمینان و حس مالکیت جمعی نسبت به این مجموعه های جدید در میان کارکنان و سرمایه داران شان نیاز دارند سازماندهی اشاعه داخلی ارتباطات و تلاش های آموزشی همگی در کل، کلید توسعه جوامع جدید آماده انطباق با مجموعه دیجیتالی هستند.

منابع

1. Aubry, Sara (2008): 'Les archives de l'Internet: un nouveau service de la BnF', Documentaliste Sciences .1 de l'Information, 45(4), pp. 12-13
2. Bermès, Emmanuelle and Gildas Illien (2009): 'Metrics and Strategies for Web Heritage Management .2 and Preservation', IFLA, <http://www.ifa.org/files/hq/papers/ifa75/92-bermes-en.pdf>
3. Bermès, Emmanuelle and Louise Faudet (2009): 'The Human Face of Digital Preservation: .3 Organizational and Staff Challenges, and Initiatives at the Bibliothèque nationale de France', iPres, <http://www.escholarship.org/uc/cdl-ipres09>
4. Illien, Gildas (2008a): 'L'archivage d'Internet, un défi pour les décideurs et les bibliothécaires: scénarios .4 d'organisation et d'évaluation, l'expérience du consortium IIPC et de la BnF, IFLA, <http://archive.ifa.org/IV/ifa74/papers/107-Illien-fr.pdf>
5. Illien, Gildas (2008b): 'Re-Inventing Collection Development Policy in the Age of Web Archiving: the .5 Experience of the BnF, LIBER Annual Conference, <http://www.ku.edu.tr/ku/images/LIBER/LIBER-ILLIEN-2008.ppt>
6. Lasfargues, France, Clément Oury and Bert Wendland (2008): 'Legal deposit of the French Web: .6 harvesting strategies for a national domain', IWAW, <http://iwaw.net/08/IWAW2008-Lasfargues.pdf>

WCT نوعی منبع باز برای مدیریت آرشیو وب گزینشی است که به عنوان پروژه ای مشترک بین کتابخانه ملی زلاندنو و کتابخانه بریتانیا توسعه یافته است. از ژانویه 2007 WCT در کتابخانه ملی زلاندنو به طور روزانه استفاده می شود. این مقاله نخستین سال آرشیو وب گزینشی ما را با نرم افزار جدید توصیف می کند. کتابخانه ملی زلاندنو، فواید WCT را توسعه داده و قصد دارد برنامه گزینشی گردآوری با WCT را برای آینده ای قابل پیش بینی ادامه دهد.

اشاره

یک سال آرشیو وب گزینشی با (1) WCT در کتابخانه ملی زلاندنو

نوشته: گوردون پنیتز، (2) سوزانا جو (3)، وانیتا لا لا (4)، گیلیان لی (5)

ترجمه: احترام السادات کیان مهر (6)

مقدمه

WCT نرم افزاری است که منابع برخط گزینش شده، گردآوری شده، و ارزیابی کیفی شده را که توسط کاربران گروهی در محیط کتابخانه ای به کار گرفته می شود پشتیبانی میکند از این نرم افزار، در مواقعی برای گردآوری وب گزینشی استفاده می شود که یک متخصص موضوعی قسمت هایی از یک وبگاه یا کل آن را معمولاً در ارتباط با یک ناحیه موضوعی تمرکز یافته یا یک رویداد مهم شناسایی کرده باشد.

نرم افزار، به عنوان پروژه ای گروهی بین کتابخانه ملی زلاندنو و کتابخانه بریتانیا توسعه داده شده است و زیر نظر کنسرسیوم فراهم آوری اینترنت بین المللی مدیریت می شود.

WCT، نرم افزار منبع باز است و از طریق وبگاه <http://webcurator.sf.net> برای استفاده جامعه آرشیو وب بین المللی آزادانه قابل دسترس است.

کتابخانه ملی زلاندنو (از این پس کتابخانه) از ژانویه 2007 از WCT به عنوان پایه برنامه آرشیو وب

ص: 305

Web Curator Tool -1

Gardon Panynter -2

Susanna Joe -3

Vanita Lala -4

Gillian Lee -5

6- کارشناس ارشد سازمان اسناد و کتابخانه ملی

گزینشی استفاده می کند.

طی سال اول ویرایش جدید نرم افزار توسعه داده شد و به طور شگرفی، افزایش و بهبود کیفیت فعالیت های گردآوری و نیز گردآوری منابع و بی به طور خودکار را میسر ساخته است.

این مقاله تجربه ما را در استفاده از WCT در محیط کار تعریف می کند

بخش بعدی، مقاله پیش زمینه فعالیت های بهره برداری وب کتابخانه و WCT را در اختیار می گذارد بخش های زیر تجربه ما را با نرم افزار جمع بندی می کند و سرانجام با توصیف یک رویداد اجرا شده گردآوری با نرم افزار خاتمه می یابد.

2. آرشیو وب گزینشی در کتابخانه ملی زلاندنو

2-1. انگیزه

کتابخانه ملی زلاندنو حکمی قانونی و مسئولیتی اجتماعی برای محافظت تاریخ فرهنگی و اجتماعی زلاندنو به شکل، کتاب، روزنامه عکس وبگاه بلاگ، یوتیوب، و ویدئو بر عهده دارد.

علاوه بر آن میراث مستند کتابخانه ملی زلاندنو فقط به صورت برخط قابل دسترس است از نظر این محتوای دسترس پذیر با ارزش است اما ناپایداری فقدان مالکیت شفاف و صحیح و طبیعت پویای آن چالش های مهمی برای هر مؤسسه ای است که برای اندوختن و نگهداری آن ها تلاش می کند.

وب گاه ها، WCT برای حل همین مشکلات ارتقا یافت به این ترتیب که به برخی موسسه ها اجازه می داد هر نوع سند برخط شامل صفحه های وب، وب گاه ها و وب نوشت ها و سایر اشکال جاری، شامل صفحه های HTML، تصاویر PDF و مدارک word مانند محتوای چند رسانه ای نظیر فایل های دیداری و شنیداری را دریافت کنند این انواع با مراقبت های، ممکن طوری هماهنگ می شوند که یکپارچگی و اصالت شان حفظ شود

برخط منفعت عمومی از نگهداری و حفاظت دراز مدت میراث برخط زلاندنو غیر قابل محاسبه است. تاریخ اجتماعی برخط ما و اکثر تاریخ دولت و سازمان برای حفظ کردن قابل امکان خواهد بود و برای محققان تاریخ دانان و شهروندان آینده زلاندنو نیز امکان پذیر خواهد بود آن ها خواهند توانست به گذشته مستند مدارک دیجیتال ما در راهی مشابه که زلاندنویی ها تا امروز به واژه های چاپ شده که برای ما از نسل های گذشته باقی مانده است، نگاه کنند.

2-2. تاریخچه درو / هاروستینگ

کتابخانه ملی، زلاندنو از سال 1999 تا پایان سال 2006، برنامه عملی آرشیو وب گزینشی را اجرا کرده است، کتابخانه برای نگهداری منابع از نرم افزار کپی (1) استفاده کرده و منابع را بر اساس مارک (2) هدایت و پایگاه داده ها را فراهم می کند. نرم افزار HTTrack کتابخانه را با یک پس افت منابع گردآوری شده که

HTTrack Website Copier Software -1
(MARC (Machine readable cataloging -2

قابلیت آرشیو شدن برای نگهداری دراز مدت را ندارد رها می کند. در حال حاضر، برای تبدیل این منابع به شکل مناسب قابل آرشیو کردن برنامه انتقال داده در حال اجراست.

2-3. نرم افزار گردآوری وب

WCT، نوعی نمودار کاری گردآوری را پشتیبانی می کند که شامل وظایف تخصصی است: انتخاب منبع برخط، جست و جو، اجازه نگهداری، آن قابل دسترس بودن، قابل توصیف بودن، تعیین دامنه آن، و دسته بندی ها، نمودار کردن گردآوری وب یا یک سری از گردآوری های وب، اجرا کردن گردآوری ها، انجام دادن مرور کیفیت و تأیید یا تکذیب کردن منابع گردآوری شده و ذخیره منابع تأیید شده در یک مخزن دیجیتال یا آرشیو.

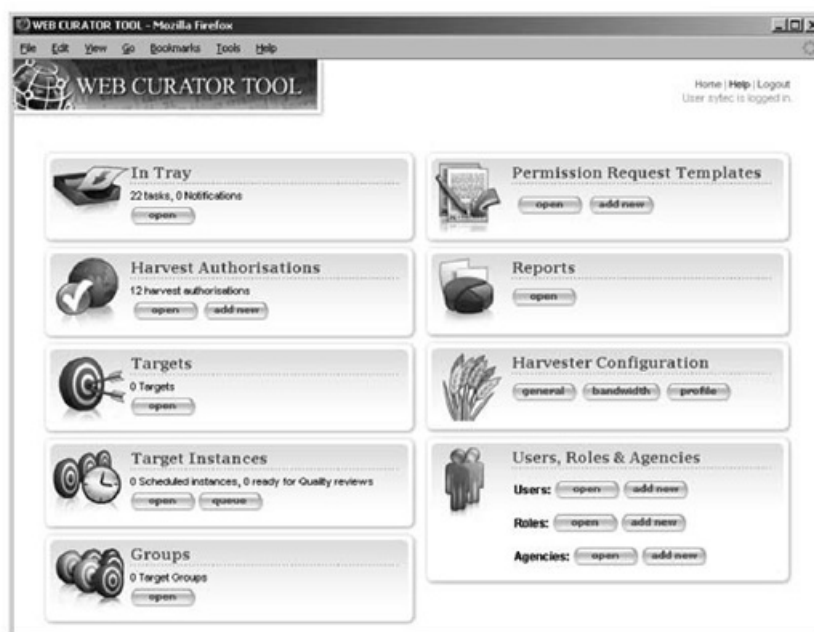
عکس

یک سال آرشیو وب گزینشی ... ۳۰۷

قابلیت آرشیو شدن، برای نگهداری دراز مدت را ندارد، رها می‌کند. در حال حاضر، برای تبدیل این منابع به شکل مناسب قابل آرشیو کردن برنامه انتقال داده در حال اجراست.

۲-۳. نرم افزار گردآوری وب

WCT، نوعی نمودار کاری گردآوری را پشتیبانی می‌کند که شامل وظایف تخصصی است: انتخاب منبع برخط، جست‌وجو، اجازه نگهداری آن، قابل دسترس بودن، قابل توصیف بودن، تعیین دامنه آن، و دسته‌بندی‌ها، نمودار کردن گردآوری وب یا یک سری از گردآوری‌های وب، اجرا کردن گردآوری‌ها، انجام دادن مرور کیفیت و تأیید یا تکذیب کردن منابع گردآوری شده و ذخیره منابع تأیید شده در یک مخزن دیجیتال یا آرشیو.



شکل ۱. فهرست اصلی نرم‌افزار را نشان می‌دهد.

شکل ۱. فهرست WCT

بیشتر فعالیت‌های آرشیو وب به‌طور عمده‌ای بر تخصص‌های فنی اپراتورهای گردآوری تکیه دارد. از سوی دیگر، WCT گردآوری را مسئولیت کاربران و متخصصان موضوعی (ترجیحاً از مهندسان و مدیران سیستم) از طریق آسان و روان نمودن، به‌طور خودکار، جزئیات فنی گردآوری وب می‌داند.

شکل ۱. فهرست اصلی نرم‌افزار را نشان می‌دهد.

شکل ۱. فهرست WCT

بیش تر فعالیت های آرشیو وب به طور عمده ای بر تخصص های فنی اپراتورهای گردآوری تکیه دارد. از سوی دیگر، WCT گردآوری را مسئولیت کاربران و متخصصان موضوعی (ترجیحاً از مهندسان و مدیران سیستم) از طریق آسان و روان نمودن، به طور خودکار، جزئیات فنی گردآوری وب می داند.

نرم افزار برای کار مطمئن و تأثیر گذار در محیط جایی که کارکنان پشتیبانی فنی می توانند آن را نگهداری کنند طراحی شده است.

WCT نرم افزار منبع باز است و به طور آزادانه از طریق وبگاه <http://webcurato.sf.net> تحت گواهینامه Apache public License قابل دسترس است.

وبگاه دسترسی کاربران به، راهنماها دست نامه ها فهرست های، پستی screenshots، سؤال ها، مستندات فنی و اداری کد منبع ردیابی اشتباهات و صفحه پروژه sourceforg را فراهم می کند.

(پینتر و میسون (1)، نرم افزار و توسعه اش را با جزئیات بیش تری در مقاله خود در کنفرانس لیانسا (2) در سال 2006 توصیف می کنند).

2.4. اعضا و منابع

در کتابخانه ملی زلاندنو، WCT نرم افزار اولیه و مسئولیت کتابدارن نشر الکترونیکی در کتابخانه «الکساندر تورن بل» (3) است. در سال 2007 معادل دو و نیم برابر انتخابگر، الکترونیکی به طور مستقیم از این نرم افزار استفاده کردند و نیز همه گزینش ها را مدیریت نموده و گردآوری و مرور کیفی نمودند.

به هر حال نرم افزار با خط مشی های کتابخانه گردش کار ارتباطات راه دور و پشتیبانی سرویس ها یکپارچه و جامع است و بر گروه گسترده تر و بیش تری از اعضا تأثیر می گذارد.

برای مثال، سخت افزار و نرم افزار توسط سرویس های فنی حفظ می شوند و توسط میز پشتیبانی مدیریت می شوند، رده بندی توسط سرویس های محتوا هدایت می شود و آرشیو دیجیتال توسط کتابخانه ملی دیجیتال حفظ و نگهداری می شود.

WCT برای یکپارچه شدن با سیستم های کاری، موجود به طور محکم و استوار - و تا حد ممکن - طراحی شده است.

برای استقرار وضعیت از سخت افزار استاندارد کتابخانه (سرورهای سان اسپارک (4) سیستم های عملیاتی (سولاریس (5))، پایگاه اطلاعاتی (اوراکل (6))، وب سرویس ها (Apache HTTP Server and Tomcat)، و خدمات ثبت کاربر (راهنمای الکترونیکی نوول (7) استفاده می شود.

سیستم تولید بین دو سرور مستقر شده، است یکی برای مدل کور (8) و دیگری برای نرم افزار گردآوری (هم آرایه کردن برای گردآوری هشت منبع و بی همزمان) و اشتراک پایگاه اطلاعاتی موجود و سرورهای فایل با سیستم های دیگر کتابخانه سیستم آزمون جداگانه با ترتیب و وضع مشابه حفظ می شود.

ص: 308

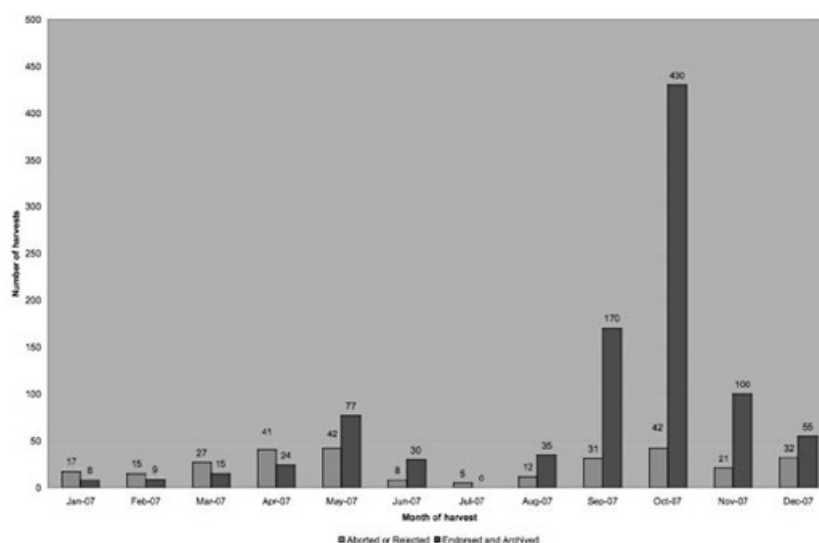
Sun SAARC Servers -4
(Operating System (Solaris -5
Oracle -6
Novell e Directory -7
Core module -8

نمودار 2 تعداد گردآوری های خودکار انجام شده توسط کتابخانه در سال 2007 را نشان می دهد. برای هر ماه، ستون ارغوانی رنگ شماره نتایج را نشان می دهد که در آرشیو تأیید نشده است. این موارد یا متوقف شده‌اند (قبل از آن که بتوانند کامل شوند) یا طی مراحل کنترل کیفیت به دلیل آن که برخی جنبه های اساسی وبگاه را کسب نکرده اند رها شده اند ستونهای، قرمز تعداد نتایج موفقیت آمیزی که در آرشیو کتابخانه دیجیتال آرشیو شده اند نشان می دهد. در مجموع 1249 گردآوری منابع خودکار توسط سیستم محصول ما در سال 2007 کمک شده‌اند. 953 (76 درصد) به طور موفقیت آمیزی گردآوری شده، کنترل کیفی شده و نیز آرشیو شده است.

عکس

۲-۵. سطح گردآوری / درو

نمودار ۲، تعداد گردآوری‌های خودکار انجام شده توسط کتابخانه در سال ۲۰۰۷ را نشان می‌دهد. برای هر ماه، ستون ارغوانی رنگ شماره نتایج را نشان می‌دهد که در آرشیو تأیید نشده است. این موارد یا متوقف شده‌اند (قبل از آنکه بتوانند کامل شوند) یا طی مراحل کنترل کیفیت، به دلیل آنکه برخی جنبه‌های اساسی وبگاه را کسب نکرده‌اند رها شده‌اند. ستون‌های قرمز، تعداد نتایج موفقیت‌آمیزی که در آرشیو کتابخانه دیجیتال آرشیو شده‌اند نشان می‌دهد. در مجموع ۱۲۴۹ گردآوری منابع خودکار توسط سیستم محصول ما در سال ۲۰۰۷ کمک شده‌اند. ۹۵۳ (۷۶ درصد) به‌طور موفقیت‌آمیزی گردآوری شده، کنترل کیفی شده و نیز آرشیو شده است.



نمودار ۲. سطح گردآوری گزینش شده در ۲۰۰۷

باقیمانندگان ۶۹ (۶ درصد) طی گردآوری متوقف شده و ۲۲۴ (۱۸ درصد) رد شده‌اند، زیرا از مرحله کنترل کیفی عبور نکرده‌اند.

سه فاز گردآوری و برداشت منابع وبی به‌طور خودکار در نمودار ۲ قابل رؤیت هستند. نخستین آنها از ژانویه تا می است، یعنی زمانی که انتخابگر الکترونیکی در فعالیتهای گردآوری از ویرایش ۱.۱ WCT استفاده کرده است. گردآوری در این مرحله کم است و در ماه می ۷۵ گردآوری موفق داشته است، یعنی زمانی که گروه در حال تجربه و گردش کار و خروجی‌های مختلف - حتی آنهایی که هنگام کار با نرم‌افزار اولیه کشف نشده بود - بودند.

نمودار 2. سطح گردآوری گزینش شده در 2007

باقیمانندگان 69 (6 درصد) طی گردآوری متوقف شده و 224 (18 درصد) رد شده‌اند، زیرا از مرحله کنترل کیفی عبور نکرده‌اند.

سه فاز گردآوری و برداشت منابع وبی به‌طور خودکار در نمودار 2 قابل رؤیت هستند. نخستین آن‌ها از ژانویه تا می است یعنی زمانی که انتخابگر الکترونیکی در فعالیتهای گردآوری از ویرایش 1.1 WCT استفاده کرده است. گردآوری در این مرحله کم است و در ماه می 75 گردآوری موفق داشته است، یعنی زمانی که گروه در حال تجربه و گردش کار و خروجی‌های مختلف - حتی آن‌هایی که هنگام

کار با نرم افزار اولیه کشف نشده بود - بودند.

ص: 309

از ژوئن تا اواسط سپتامبر، انتخابگران الکترونیکی آزمون ویرایش WCT 1/2 را تقاضا کردند. بنابراین، مقررات فعالیت های منظم گردآوری حتی بیشتر از قبل قطع می شود.

فاز نهایی گردآوری در اواسط سپتامبر شروع می شود و افزایش شگرفی در فعالیت گردآوری را نشان می دهد. عوامل متعددی در این افزایش دخیل هستند که شامل انتخاب ویرایش WCT 2/1 در سیستم، تولید افزایش 50 برابری ظرفیت اینترنت کتابخانه (از 2 مگابایت در ثانیه به 100 مگابایت در ثانیه)، شفافیت بیشتری خروجی های گردش کار پیرامون گردآوری کردن و فهرست کردن و آغاز دورویداد بر اساس گردآوری (در زیر توضیح داده شده).

در نتیجه، تقریباً گردآوری در نیمی از سال فقط در ماه (تعداد وب گاه ها) اکتبر انجام می شود.

نکته جالب در این مرحله این است که اگر چه تعداد گردآوری های موفق به طور چشمگیری افزایش یافته است تعداد «ناتمام» ها و «مردود» ها افزایش نیافته است و عددی که تأیید یا رد می شود، محسوب نمی گردد.

گردآوری با WCT ویرایش 1.1.

زمانی که انتخاب گر الکترونیکی از ویرایش WCT 1.1 از ژانویه تا می استفاده کرد، از تجربه با این نرم افزار سود بردند و به فعالیت های گردآوری و نیز مشکلات مختلفی که از طریق ویرایش 2.1 و کامل شدن آن بود پی بردند

این بخش خلاصه ای از این تجربه است و توصیف بیشتر جزئیات از نویسندگان بر اساس درخواست در دسترس می باشد.

1.3. تجربه اولیه

مهم ترین تغییر ایجاد شده توسط WCT این است که نرم افزار منابع گردآوری شده را در فرمت فایل ARC ذخیره می کند که از لحاظ فراهم آوری سودمند است اما برای فعالیت های کنترل کیفیت.

و این موضوع راه اساسی جدیدی را برای انتخابگران الکترونیکی برای رسیدن و نزدیک شدن به تجدید نظر کیفی را به همراه دارد.

با این شیوه جدید فعل و انفعال، مشکل می توان گفت که آیا خروجی های کیفی، توسط نرم افزارهای گردآوری کننده ایجاد شده یا توسط کمبودهای نرم افزارهای کنترل کیفی

برای مثال اگر یک بازبین کننده یک صفحه وب را در نرم افزار تورق باز کند و نشان داده شود که یک صفحه وب (تصویر) گم شده است انتخابگر الکترونیکی ابتدا باید تعیین کند که آیا به طور موفق، تصویر از منابع و بی گردآوری شده یا آن را به طور موفقیت آمیزی گردآوری کرده، اما مرورگر نمی تواند آن را نشان دهد خیلی زود مشخص شد که برای پشتیبانی اعضا در وجود این تفاوت ها تمرین های اضافی ضروری بوده است

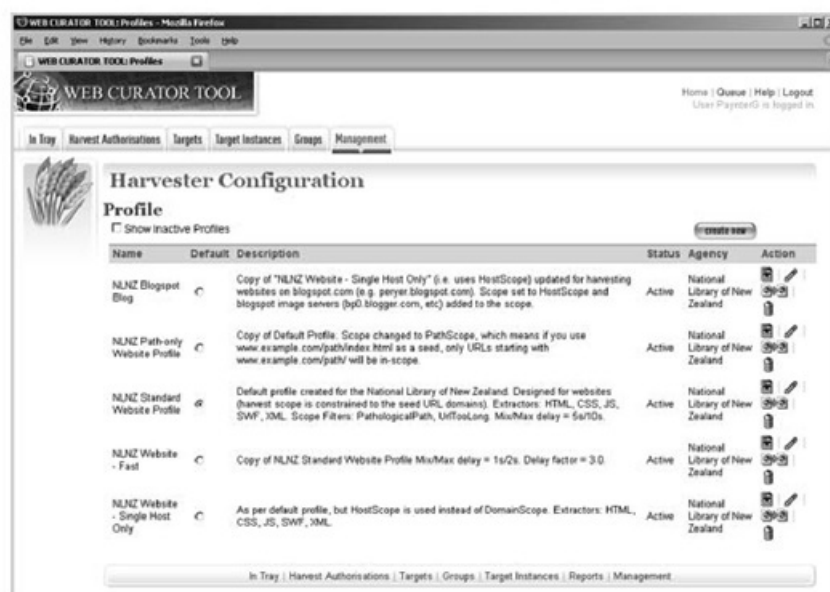
زمانی که کتابداران تجربه مهمی از مرور نتایج گردآوری نرم افزار HTTrack را کسب کردند، نمی دانستند که چگونه از عدم موفقیت گردآوری منابع وبی در WCT گره بکشایند.

تمام مشکلات جدید گردآوری در نسخه جدید ضبط و مشاهده شد که همان نوع اشتباه ها دوباره رخ داده است. بنابراین بخشی از تمرین حفظ و یک روند بازشناختی نیز توسعه و به WCT به طور دستی اضافه شد.

مشکل دیگر که در ماه های نخست نمایان شد این بود که توسعه سودمند پروفایل های گردآوری شده، به خصوص در مورد بلاگ ها زمان بر است توسعه شکل 3. مجموعه ای از پروفایل های گردآوری شده را نشان می دهد که سرانجام به وجود آمده اند.

عکس

زمانی که کتابداران تجربه مهمی از مرور نتایج گردآوری نرم‌افزار HTTrack را کسب کردند، نمی‌دانستند که چگونه از عدم موفقیت گردآوری منابع وبی در WCT گره بگشایند. تمام مشکلات جدید گردآوری در نسخه جدید ضبط و مشاهده شد که همان نوع اشتباه‌ها دوباره رخ داده است. بنابراین، بخشی از تمرین حفظ و یک روند بازشناختی نیز توسعه و به WCT به‌طور دستی اضافه شد. مشکل دیگر که در ماه‌های نخست نمایان شد این بود که توسعه سودمند پروفایل‌های گردآوری شده، به‌خصوص در مورد بلاگ‌ها، زمان بر است. توسعه شکل ۳. مجموعه‌ای از پروفایل‌های گردآوری شده را نشان می‌دهد که سرانجام به‌وجود آمده‌اند.



شکل ۳. پروفایل‌های گردآوری کننده در حال استفاده در کتابخانه در دسامبر ۲۰۰۷

۳.۲. مشکلات گردآوری

مشکلات زیر در ویرایش ۱.۱ وجود داشته و در ویرایش ۲.۱ برطرف شده است.

کمبودهای تورق نرم‌افزار

برخی اوقات محتوای ظاهر شده از دست می‌رود به دلیل اینکه نرم‌افزار تورق (مرور) نمی‌تواند آن را به‌طور صحیحی منتقل کند.

شکل 3. پروفایل‌های گردآوری کننده در حال استفاده در کتابخانه در دسامبر 2007

3.2. مشکلات گردآوری

مشکلات زیر در ویرایش 1.1 وجود داشته و در ویرایش 2.1 برطرف شده است

کمبودهای تورق نرم افزار

برخی اوقات محتوای ظاهر شده از دست می‌رود به دلیل اینکه نرم افزار تورق (مرور) نمی‌تواند آن را به طور صحیحی منتقل کند.

ص: 311

برخی نسخه های کاربر برای انتخاب گران الکترونیکی، سیر کردن سریع در پیرامون نرم افزار را مشکل می ساخت.

اشتباه «در مرحله ایست درگیر شدن»

مشکل دوباره به وجود آمده در ویرایش جدید WCT این بود که گردآوری باید همان گونه که انتظار می رفت با موفقیت به مرحله پایانی می رسید ولی سیستم به وضعیت «گردآوری شده» انتقال نمی یافت و در عوض در وضعیت توقف باقی می ماند مشکلاتی که به این شرایط هدایت می شود نتایجی از نسخه های مورد نظر که در بالا آمده و کامل شده است.

وب گاه های بزرگ

کتابخانه بریتانیا به این موضوع پی برد که نرم افزار از عهده گردآوری های بزرگ و گسترده بر نمی آید و مشکلاتی با نرم افزار گردآوری همزمان وجود دارد.

3.3. خلاصه تجربه با ویرایش 1.1

با وجود برخی ای محدودیت ها ویرایش اول WCT بهبودی وسیع و بزرگی بر نرم افزار HTTrack قبلی بود و آن مرحله مهم و با اهمیتی بود وقتی که ما اولین وب گاه های گردآوری شده را به آرشیو دیجیتال ارائه دادیم

طی سال اول WCT 1250 منابع وبی را به طور خودکار گردآوری کرد و تقریباً 950 مورد دوباره، بررسی تأیید و به آرشیو اضافه شد.

گردآوری با WCT با ویرایش 1.2

ویرایش WCT 1.2.6، در اواسط سپتامبر 2007 به محصول گسترش داده شد و از آن در ماه های بعدی در دو رویداد گردآوری استفاده شد و سپس در پایان اکتبر با ویرایش 1.2.7 جابه جا شد.

1.4. نرم افزار گردآوری تاریخ

یکی از نواحی ای که در آن پیشرفت زیادی به وجود آمد کنترل کیفیت است. افزایش در این ناحیه به طور شگرفی تأثیر کنترل کیفیت را ثابت کرد و باعث پیشرفت بیشتر، سریع تر، ماهرانه تر شد.

نرم افزار گردآوری تاریخ در کنترل کیفی ثابت کرده است که بسیار مفیدتر از تقویت کنترل کیفی است. این نرم افزار همه منابع گردآوری شده را که دارای یک هدف ویژه است با چکیده اطلاعات، نظیر تاریخ شروع تاریخ انتقال URL گردآوری شده و افتاده، زمان صرف شده و وضعیت جاری را فهرست می کند.

این امر بسیار مفید است به طوری که از اطلاعاتی که در روند کنترل کیفیت نیاز بود محکم تر است.

2.4. نرم افزار تورق (مرورگر)

هنوز یک تغییر ساده و تأثیر گذار برای مرور میانجی، بود افزایش پیشنهادها برای دریافت دیداری یک سایت گردآوری شده از سه راه مختلف صورت می گیرد: دریافت دیداری مرحله جاری گردآوری شده در نرم افزار مرورگر دریافت دیداری سایت، جاری یا دریافت دیداری ویرایش آرشیو شده قبلی (در آرشیو اینترنت یا آرشیو مکانی).

این مراحل در Tab دیگر با تورق ویندوز باز می شود و اجازه می دهد که مرورگر یک کپی گردآوری شده را با ویرایش دیگر سایت مقایسه کند.

این پیشرفت بزرگ باعث شد که نرم افزار تورق به ما اجازه دهد که به طور موفقیت آمیزی مرور و تورق سایت های دیگر را - که قبلاً قابل دیدن نبودند - مشاهده کنیم به هر حال، چند خروجی حل نشده باقی می ماند.

وب گاه های تحویل داده شده که از جاوا اسکریپت استفاده می کند باعث مشکلاتی می شود، به ویژه دوباره عناصر وظیفه ای مانند منوی پایین - بالا که می تواند سایت های گردآوری شده را برای سیر کردن مشکل سازد.

گردآوری سبک برگه ها هنوز می تواند مشکل باشد، گرچه نسخه های جمع آوری شده بیش از نسخه های تکراری است نسخه های برجسته، دیگر شامل مشکلات تکرار شده URL ها با فضاها و جمع آوری تصاویر پیش زمینه جاسازی شده مطمئن.

3.4. نرم افزار هرس درخت تصمیم

نرم افزار هرس درخت تصمیم سیستم روز آمد شد اما کتابخانه نرم افزار هرس را در وب گاه های قدیمی تر استفاده نمی کند. بنابراین از این تغییرات استفاده زیادی به عمل نمی آید. به هر حال، ما با نرم افزار تجربه کرده ایم و مدتی می تواند خیلی آهسته برای سایت های بزرگ استفاده شود وظیفه جدید نرم افزار هرس ثابت کرد که می تواند برای انتخاب کردن و نگاه کردن فایل های دل خواه از منابع گردآوری شده بسیار مفید باشد

4.4. وب گاه های بزرگتر جمع آوری می شوند

در ماه های آخر سال ما متوجه شدیم که فعالیت های گردآوری ما شامل سایت های بزرگتر و بزرگ تر هستند. بزرگ ترین گردآوری، کتابخانه به عنوان قسمتی از برنامه گردآوری در اندازه 21 گیگا بایت کامل و کنترل شد گرچه پس از کنترل کیفیت رد شده بود.

وب گاه های متعدد 10 گیگا بایتی به طور موفقیت آمیزی کنترل آرشیو و گردآوری شدند زمانی که اندازه وب گاه ها رشد کرد، انتخاب گران الکترونیکی تکیه بیش تری بر پروفایل های عادی و پروفایل های

برتر کردند به ویژه فیلتر ها که به انتخاب گران الکترونیکی اجازه داد که گردآوری کننده را از جمع آوری قسمت های ویژه وبگاه متوقف کنند.

5.4. گسترش ذخیره دیجیتال با ارزش

نتیجه افزایش تعداد و اندازه وب گاه های گردآوری شده این بود که WCT به ذخیره دیجیتال با ارزش بیش تری نیاز داشت (فضای موقت برای ذخیره گردآوری ها زمانی که کنترل کیفی شده اند و قبل از این که آرشیو شوند).

در اکتبر و نوامبر این مشکلی ویژه و خاص بود زیرا تعداد وب گاه های انتخاب شده گردآوری شده و تأیید شده ناگهان افزایش یافت و موقتاً به طور طولانی تری فهرست شد و نوعی انبار منبع به وجود آمد. در نتیجه فضای اختصاص داده شده برای ذخیره منابع دیجیتال با ارزش که زودتر ظاهر شده از دست رفت و ناگهان به نظر رسید که بسیار کم ذخیره شده است در نتیجه یک دیسک بزرگ تر و جدید برای ذخیره با ارزش منابع دیجیتال بود و ابزارهای گزارش گیری به کتابداران اجازه می دهد که از دیسک به طور موقت استفاده کنند و نمایش دهند.

این وضعیت، همچنین نسخه ای از آرشیو موقتی دیجیتال را فراهم کرد که در آن فضای دیسک با ذخیره ارزشمند دیجیتال مشترک شده مطمئناً تفصیل کنترل های نهایی فایده ای نداشته است.

6.4. ارتباط

ارتباط خوب بین اعضای کتابخانه های مختلف از طریق نرم افزار تحت تأثیر واقع می شود، اما شکاف های مختلفی در ارتباط از زمانی که WCT از ابتدا استفاده می شد در محصول شناسایی شد.

گرچه در ویرایش 1.2 پیشرفت هایی حاصل شد خیلی از قلم افتادگی ها، خارج از نرم افزار از طریق گزارش های دوره ای از طریق پایگاه اطلاعاتی WCT نشانی شده اند (مدرک نرم افزار، شامل یک فرهنگ لغت داده کامل است).

برای مثال، گزارش هفتگی حاصل شده از گردآوری، منابع به فهرست نویس ها برای اعلام آن ها به وبگاه هایی که نیاز است فرستاده شد و تعداد و اندازه وب های گردآوری شده جاداده شده توسط انتشارات الکترونیکی کتابداران برای سرویس های فنی اعضا پیش بینی شد. زمانی که این ها در یک شکل کامل گردآوری شدند ما امیدواریم که به آن ها یک سری گزارش های ساختار WCT را اضافه کنیم.

7.4. فهرست گردش کار و دسترسی

زمانی ممکن است از WCT برای توصیف وب گاه ها استفاده شود و آن خط مشی کتابخانه است برای توصیف مجموعه کامل کتابخانه در فهرستش و فراهم کردن یک پیوند رکورد فهرست کتابخانه به آیتم های دیجیتال نگهداری شده در محزن دیجیتال

(به هر حال انبارش های دیجیتال، پیچیده نظیر سریال ها و وب گاه ها به طور عام قابل دسترس نخواهند)

بود تا آرشیو میراث دیجیتال ملی با مخزن دیجیتال موقت ما در سال 2008 جایگزین شود).

یک مانع و اشکال در جست و جو کردن فهرست کتابخانه برای وب گاه هاست و ناتوانی جمع آوری کامل رویدادها از هر وبگاه که به تنهایی فهرست شده است

نشانی این نسخه کتابخانه بریتانیا توسعه یک نرم افزار میانجی وب را طراحی می کند که دسترسی به منابع گردآوری شده را بر اساس موضوع و رویداد به عنوان یک هدف اضافی برای کاربرانی که می خواهند به طور ویژه برای وب گاه ها جست و جو کنند فراهم می کند ما پیشرفت کتابخانه بریتانیا را با علاقه بسیار زیادی دنبال خواهیم کرد.

5. گردآوری رویداد انتخابات هیئت محلی

کتابخانه، با به کار بردن نرم افزار گردآوری HTTrack در سال 2007 مسئولیت گردآوری برخی رویدادهای عظیم ورزشی را بر عهده داشته است (جام آمریکا) (1) و انتخابات (انتخابات پارلمان ملی).

کتابخانه دو رویداد گردآوری را طراحی کرد که به طور اتفاقی با استقرار ویرایش WCT 1/2 همزمان بود. گزینش های هیئت دولتی محلی، رویدادی سه ساله است و هر قدرت محلی در زلاندنو به نگهداری و هدایت گزینش اعضایش احتیاج دارد کتابخانه برنامه گردآوری رویداد دوازده هفته ای را با تمرکز بر گزینش های هیئت محلی در سال 2007 و شروع شدن در سپتامبر را به عهده گرفت. گردآوری رویداد هیئت محلی اولین گردآوری رویداد عظیمی با استفاده از WCT بود و همچنین هنوز بزرگ ترین تلاش در کتابخانه بود با 238 وب گاه انتخاب شده

این سایت ها شامل وب گاه ها، بلاگ ها، وب گاه های شوراهای منطقه ای، شهری و ناحیه ای، اخبارهای سایت ها و سایت های عمومی یا بلاگ ها با محتوای مربوط و مناسب هم گزینشی یا دولت محلی با این حال طیف وسیعی از وب گاه ها انتخاب شده اند که تفسیرهای رسمی و غیر رسمی را با پراکندگی جغرافیایی وسیعی نمایش می دهند کلیه سایت های انتخاب شده در محدوده وضعیت قانونی زلاندنو بودند، بنابراین نیازی به جست و جوی دستور صریح و واضح و قانون گذاری برای به دست آوردن آن ها نداریم

سایت های انتخاب شده برای تعیین گردآوری، اولویت بندی شدند و در آن موقع یک جدول گردآوری برای پوشاندن کلیه تاریخ ها در مدت 16 هفته طراحی شده بود. جدول با تعداد زیادی سایت انتخاب شده و با قابلیت اداره خاتمه داده شد درگیری مقتضی اولیه برای پایان دادن منابع اعضا نیز وجود داشت. تعداد کمی از گردآوری ها یک فاصله زمان کافی بین هر گردآوری برای 2 عضو برای کامل کردن خزش گرها و روند تکرار کیفیت را مطمئن ساخت.

عامل تأثیر گذار دیگر بر برنامه گردآوری از گردش کاری هایی ریشه می گرفت که در حال حاضر از WCT استفاده می کنند در گردآوری های قبلی ایجاد نمودار گردآوری تکراری عظیمی از کنترل کیفیت امکان پذیر بود که می توان بعد از تکمیل همه گردآوری ها انجام داد.

ص: 315

به هر حال با تعداد زیادی از سایت های انتخاب شده در این دوره این رویکرد بر WCT هم غیر عملی و نشدنی است و هم سنگین و طاقت فرساست.

با گردآوری و سری زمانی، اهداف گروه ها با استفاده از سیستم های گروهی و وظیفه ای به وجود آمد. سایت های برگزیده مرتب شدند و 36 گروه با دو مقوله وبگاه، یا - در مورد سایت های شور- منطقه جغرافیایی تعیین شدند.

این گروه ها غیر رسمی بودند، زیرا استفاده اولیه گروه ها در این رویداد (دوره) - برای آسان تر شدن برنامه گردآوری صورت می گرفت و نه برای گروه بندی رسمی استفاده از گروه در مدیریت گردش کاری در تیم مفید بود و متناوب شدن گردآوری توسط گروه ها نیز مفید بود. همچنین نوعی تعادل فشار در بین سخت افزارهای گردآوری به وجود آورد.

اگر چه هدف، شروع گردآوری در بعد از ظهر بود، حتی با وجود هشت گردآورنده همزمان این کار ممکن نشد برخی گردآوری های برنامه ریزی شده چندین ساعت در نوبت می ماندند تا یک گردآورنده آزاد شود، اما خوشبختانه صف های طولانی برداشت موردی به وجود نیاورد و با کوچک شدن بسیاری از سایت ها صف های برداشت به سرعت ناپدید شد این سیستم با استفاده از هشت گردآورنده همزمان به خوبی با دوره های گردآوری طولانی و فشرده، مقابله کرد

مشکلات وقتی به وجود آمد که خزش های وبگاه شورها در اندازه های بسیار بزرگتر با آرشیو دیجیتال وب گاه هایی که منتظر فهرست نویسی پیش آرشیوی بودند همزمان شدند. نتیجه یک تفسیر مشکل بحرانی در ظرفیت فضای دیسک بود. از آن جا که دیگر مناطق کتابخانه به طور مشترک از همین امکان ذخیره سازی استفاده می کردند عواقب بسیار جدی داشته است. در یک مقطع مجبور شدیم گردآوری را تازمانی که فضای دیسک کافی در دسترس باشد به تعویق بیندازیم.

با وجود پیچیدگی های اخیر، بیش تر وب گاه هایی جمع آوری شده بررسی کیفی و تأیید شده و با موفقیت آرشیو شده اند. در پایان برنامه برداشت، بیش از 600 وبگاه برداشته شده تأیید و یا برای این رویداد خاص آرشیو شده اند که موفقیت (دستاورد) قابل توجهی برای برنامه آرشیو وب کتابخانه است.

به طور کلی انتخابات هیئت، محلی آزمایش مهمی درباره گردآوری رویداد با استفاده از WCT بود. ما به سهولت قادر به برنامه ریزی و مدیریت این رویداد بزرگ آن ها بر روی سیستم بودیم، و تجربه های مفیدی در استفاده از ویژگی های سیستم به دست آوردیم که تاکنون موقعیت استفاده از آن ها را نداشتیم موضوع با اهمیت تر این که گردآوری این رویداد نشان داد که WCT با موفقیت می تواند با تقاضاهای فعالیت گردآوری وب و اندازه افزایش یافته تطبیق دهد مناطقی را که به بررسی و نظارت دقیق تر نیاز داشتند نیز مشخص کند و درس های ارزشمندی آموخته و تجربه ای به دست دهد.

متعاقب انتخابات هیئت محلی سال 2007، رویداد دیگری که بر پایه گردآوری قرار داشت جام جهانی راگی 2007 بود این رویداد در مقیاس بسیار کوچک تر از گردآوری رویداد اول بود، اما شامل وب گاه های خارج از محدوده مقررات واسپاری قانونی زلاندنو می شد که امکان قانونی جمع آوری با استفاده از WCT را برای اولین بار نیاز داشت.

6. نتیجه گیری

انتخاب گران الکترونیکی کتابخانه به طور روز مره برای انتخاب نمودار گردآوری، و مرور وب گاه ها از WCT استفاده می کنند سپس آن ها را برای آرشیو دیجیتال پیشنهاد می دهند.

انتخاب گران الکترونیکی کار را با ویرایش 1.1 شروع کردند و پس از آشنایی بیشتر با این نرم افزار، از نوعی محدودیت روزافزون آگاه شدند که Specification ویرایش جدیدی را اطلاع می داد. نرم افزارهای کنترل کیفی به روز شده تفاوت عظیمی را به وجود آوردند و بسیاری از سایت ها که برای مرور آن ها با ویرایش های اولیه مشکل بودند، جمع آوری و با ویرایش 1.2 مرور شدند.

هنوز برخی مسائل پر دردسر و آزار دهنده وجود دارد و یا در حال ضبط آن ها در bugtraker روی وب گاه های منبع باز هم ضبط کرده ایم.

از آخرین هفته ها ویرایش جدید نرم افزار (ویرایش 3.) قابل بهره برداری شد و انتظار داریم که بتواند در گردآوری گردش کار، نتایج بهتری ایجاد کند.

پس از دو دوره گردآوری کتابخانه، ظرفیت مدل ها برای ذخیره آتی را توسعه داده و نیاز به پهنای باند دارد. ما همچنین دسترسی نرم افزارها و سطح دامنه گردآوری وب گاه های زلاندنورا طراحی می کنیم.

منابع

1. HTTrack Website Copier, (Accessed 2008-05-15).

2. Paynter, G. W. and Mason I. B. (2006) Building a Web Curator Tool for The National Library of New Zealand. New Zealand Library Association (LIANZA) Conference, . October 2006. (Accessed 2008-05-15).

ص: 317

Web Information Management

:Vol. 1

Basics and Worldwide Initiatives

Edited by

Gholam Ali Montazer (Ph.D)

and

Tarbiat Modares University

Farzaneh Shadanpour

National Library and Archives of Islamic Republic of Iran

ص: 318

مشخصات کتاب

مدیریت منابع اطلاعاتی وب جلد دوم

دیدگاه‌های فناورانه اخلاقی و مدیریتی

به کوشش:

دکتر غلامعلی منتظر و

فرزانه شادان پور

زمستان 1391

فهرستتویسی پیش از انتشار کتابخانه ملی جمهوری اسلامی ایران

سرشناسه: منتظر، غلامعلی 1348 - ، گردآورنده

عنوان و نام پدیدآور: مدیریت منابع اطلاعاتی وب / به کوشش غلامعلی منتظر و فرزانه شادان پور.

مشخصات نشر: تهران سازمان اسناد و کتابخانه ملی جمهوری اسلامی ایران 1391

مشخصات ظاهری: 2 ج.

شابک: دوره: 0 - 344 - 446 - 1964 - 978؛ ج. 2: 72 - 345 - 446 - 964 - 978؛

وضعیت فهرستتویسی: فینا

مندرجات: ج. 1 مبانی و تجربه های جهانی - ج. 2 دیدگاه‌های فناورانه، اخلاقی و مدیریتی.

موضوع: وب -- سایت ها -- مدیریت

موضوع: منابع اطلاعاتی -- مدیریت

موضوع: وب -- آرشوسازی

شناسه افزوده: شادان پور، فرزانه 1344 - ، گردآورنده

شناسه افزوده: سازمان اسناد و کتابخانه ملی جمهوری اسلامی ایران

رده بندی کنگره: 4 1391 م TK5105/8888

رده بندی دیویی: 005/72

شماره کتابشناسی ملی: 3077380

خیراندیش دیجیتالی : انجمن مددکاری امام زمان (عج) اصفهان

ویراستار کتاب : خانم مرضیه محمدی سرپیری

ص: 1

اشاره

سازمان اسناد و کتابخانه ملی جمهوری اسلامی ایران

عنوان مدیریت منابع اطلاعاتی، وب جلد اول مبانی و تجربه های جهانی

به کوشش دکتر غلامعلی منتظر (دانشیار دانشگاه تربیت مدرس) و فرزانه شادان پور (مریی سازمان اسناد و کتابخانه ملی جمهوری اسلامی ایران)

ویراستار ادبی آرزو تجلی کارشناس ارشد جامعه شناسی، سازمان اسناد کتابخانه ملی جمهوری اسلامی ایران)

تنظیم و تصحیح مهشید برجیان کارشناس ارشد کتابداری و اطلاع رسانی سازمان اسناد و کتابخانه ملی جمهوری اسلامی ایران)

ویراستار استنادی مقالات تألیفی فروزان رضایی نیا کارشناس کتابداری و اطلاع رسانی سازمان اسناد و کتابخانه ملی جمهوری اسلامی ایران)

نمونه خوانی و اصلاحات مهشید، برجیان فاطمه، رمضانپور آمنه، هزارخانی زهرا، زاهدی محمد رضا، میقانی ملیحه حاجی زاده مقدم

طراحی جلد و صفحه آرایی شهره خوری

ناظر فنی چاپ نصرت الله امیرآبادی

ناشر سازمان اسناد و کتابخانه ملی جمهوری اسلامی ایران

شمارگان: 500 نسخه

بها: 20000 تومان

نشانی تهران بزرگراه شهید حقانی (غرب به شرق)

بعد از ایستگاه مترو، بلوار کتابخانه ملی

تلفن فروشگاه 81623318 - 81623315 - 88941946

دورنگار: 88947496

وب سایت: www.nlai.ir

پست الکترونیک انتشارات Publication@nlai.ir

سازمان اسناد و کتابخانه ملی جمهوری اسلامی ایران

مدیریت منابع اطلاعاتی وب

جلد دوم

دیدگاه های فناورانه اخلاقی و مدیریتی

ص: 3

سخن نخست... یازده

به جای مقدمه ... 1

فصل اول مسائل فناورانه در آرشیو وب...3

آرشیو اشیای داده‌ای با استفاده از فیدهای وب...4

نوشته ماریلنا، اویتا پیرسنلارت / ترجمه لیلی سیفی

آرشیو صفحات وب بر مبنای تحلیل دیداری و 28...DIFF

نوشته میریام بن سعد، استفان گانکارسکی، زینب پهلوان / ترجمه مجیدرضا وحیدی

آرشیو منابع ویدئویی وب...38

نوشته رادو، پاپ گابریل واسیلی ژولین ماسانه / ترجمه فروزان رضائی نیا

استفاده از عاملهای هوشمند نرم افزاری جهت ایجاد قابلیت تعامل پذیری در خدمات محتوایی

و اطلاعاتی سازمانها...52

نوشته محمود خراط مانده مشرف فتانه تقی یاره

تحلیل انسجام و مصورسازی در آرشیو وب...70

نوشته عبدالله حسینیان

بایگانی وب پنهان...90

نوشته ژولین ماسانه / ترجمه افسانه تیموری خانی

بررسی تاثیر بستر نحوی بر میانکنش پذیری استانداردهای فراداده ای گامی در جهت یکپارچه سازی...106

نظامهای اطلاعاتی نوشته سید مهدی طاهری

بهبود سازی کیفیت آرشیوهای وب...126

نوشته میریام بن سعد / ترجمه مهشید برجیان ، ساناز باغستانی

خزش هوشمند در برنامه های کاربردی وب...144

نوشته محمد ، فهیم زیر نظر پیر سنلار ترجمه فرزانه شادان پور

دسته بندی مفهومی صفحات وب با استفاده از اتوماتای یادگیر توزیع شده...158

نوشته روح الله ، گودرزی ، مصطفی پیرهادی

DigiBoard: ابزار افزایش کارایی فعالیتهای پیچیده آرشیو وب در کتابخانه کنگره...168

نوشته آبه ، گروتک جینا جونز ترجمه سعیده اسلامی

رونوشت برداری از وبگاهها...184

نوشته خاویر روش / ترجمه فرزانه شادان پور

رویکرد جدید آرشیو وب مبتنی بر آنالیز بصری صفحه های وب...210

نوشته میریام بن سعد استفان گانچارسکی زینب پهلوان / ترجمه سعیده اسلامی

طراحی درخت تصمیم گیر برای دسته بندی سریع تصاویر در آرشیو وب...222

نوشته سید مجتبی حسینی عالیبه بهرام زاد

فرداده وبگاه پروژه نهایی MIMS

آنورادا ، روی زیر نظر اریک وایلد / ترجمه الهام میرزاییگی کسبی

کاوش مجموعه های وب...268

نوشته آندریاس آشنبرنر آندریاس روبر / ترجمه مهندس سارا کلینی

کیفیت داده در آرشیو وب...294

نوشته پیرسنلار مارک ، اسپانیول دیمیتار ، دنو آرتوراس ، مازیکا گرهارد ویکوم / ترجمه الهام میرزاییگی کسبی

محافظت پویا از وبگاهها...314

نوشته رابرت شارپ / ترجمه آرزو تجلی

مرور و ارزیابی روشهای تشابه سنجی در متن...330

نوشته حمید آهنگر بهان، غلامعلی منتظر

نقش پیوندها و مدیریت آنها در وب سایت کتابخانه های دانشگاهی...358

فائزه دلقتدی فرحناز فتح الله زاده

ص: 5

EverLast یک معماری توزیع شده برای حفاظت از وب...374

نوشته آویشک، آناند سریکانتا، بداتور کلاوس بریچ رالف، اسکنل کریستوس تریفونوپلوس / ترجمه مریم کراری

فصل دوم مسائل اخلاقی و مدیریتی در آرشیو وب...399

اصول اخلاقی حاکم بر آرشیو وب...400

نوشته مگان داگهرتی کرستن ایفوت استون. ام اشنايدر / ترجمه سوده صیرفی

بررسی نحوه سازماندهی منابع اطلاعاتی در کتابخانههای دیجیتال ایران...406

نوشته حامد علیپور، حافظی زهرا، عبداللهی سمیه مجیدی، میترا حیدر تامینی

حفاظت بلندمدت محتوای وب...424

نوشته مایکل دی / ترجمه میترا صمیعی

دسترسی و فهرستهای راهنما...450

نوشته تورستاین هالگریسون / ترجمه امیررضا اصنافی، مریم پاکدا من نائینی

عوامل موثر بر درک حریم خصوصی در شبکههای اجتماعی و راهکارهای پیشنهادی برای

شخصی سازی آن...472

نوشته سعید رضایی شریف آبادی، نسرین علیپور

گزینش آرشیوهای وب...486

نوشته زولین ماسانه / ترجمه دکتر زهرا اباذری

مسائل اخلاقی در ایجاد و به کارگیری آرشیو وب - به سوی یک برنامه پژوهشی...508

نوشته آندریاس، رویر مکس، کیزر برنارد واجر / ترجمه نجلا حریری

حق مؤلف در محیط الکترونیک...524

نوشته داریوش، مطلبی شهمیه السادات حسینی

از ویژگیهای قرون گذشته بی خبری بود و تمایز جدی عصر جدید نسبت به گذشته دسترسی آسان به اطلاعات. است بشر با از سر گذراندن سه موج و پارادایم، کشاورزی صنعت و اطلاعات امروز در قرن بیست و یکم پا در عصر انفجار اطلاعات نهاده است این امر فی نفسه نه مطلوب است نه مذموم، بلکه به نحوه مدیریت ما نسبت به اطلاعات باز میگردد.

بشر امروزی به دلیل رشد روزافزون علم و فناوری در شرایط هشدار آمیز عدم قطعیت بسر میرود و همین مدیریت و تصمیم گیری را با چالش جدی روبرو ساخته است. اگر اطلاعات درست مدیریت شود و در تصمیم گیریها به موقع به کار آید و از دو ویژگی صحت و سرعت برخوردار باشد میتواند منشأ تصمیمهای تحول آفرین شود. ویژگی دیگر این عصر ظهور و حضور همه جانبه اطلاعات دیجیتالی است دورانی فرا رسیده است که در آن بناست دانش مدون و تفکر مضبوط بشر علاوه بر کاغذ و حتی بیش از آن بر محمل «بیت» ها مسیر، تولید نشر و اشاعه و مصرف را پیماید. هم اطلاعات تولید شده تحت وب و هم میزان استفاده از این اطلاعات با سرعت فزاینده ای رو به رشد است. کشور ما بنابر اطلاعات وثیق از حیث تعداد کاربران و میزان حضور و فعالیت آنها در وب جایگاه نخست را در منطقه خاور میانه داراست. این روند رو به رشد با نصب العین قرار دادن آرمانهای بلند انقلاب اسلامی در ترویج تفکر رهایی بخش اسلام ولایت مدار و وظیفه خطیری بر دوش نهادها و دستگاههای مسئول تولید، سیاستگذاری و نشر محتوا در محیط وب قرار می دهد و آن انجام بررسیهای علمی و مستند به منظور ابتنای سیاستگذاریها و عملکردها بر مبانی صحیح و کارآمد و متناسب با نیازهای گوناگون کاربران در این محیط است. اما وجه دیگر، صیانت از این محتوا و انتقال آن به نسلهای آینده است که با توجه به ناپایداری محتوای قرار گرفته بر اینترنت و فناوری پیشرفته ای که برای چنین امر خطیری لازم است

از اهمیت مضاعفی برخوردار می شود.

سازمان اسناد و کتابخانه ملی جمهوری اسلامی ایران بنابر مأموریت خویش دایر بر صیانت از میراث فکری کشور و اشاعه آن عزم راسخ داشته است که برای مدیریت منابع اطلاعاتی مهم و رو به رشد وبی نیز چاره اندیشی نماید؛ بنابر این در سال 1389 نخستین بار در کشور به تهیه ساز و کار لازم برای ایجاد

آرشیو ملی وب همت گماشته است.

از دیگر سو، سازمان با علم به این که مدیریت در این حوزه مشارکت همه صاحبان اندیشه در حوزه تولید سازماندهی و اشاعه اطلاعات تحت وب را می طلبد مصمم شد نخستین کنفرانس ملی مدیریت منابع اطلاعاتی وب را برگزار نماید تا اهل علم و فناوری در این مجمع با هم اندیشی و تضارب آراء همچون گذشته این سازمان را یار و یاور باشند.

این اثر مجموعه ای است فراهم آمده از تلاش پژوهشگرانی که با وجود نبودن مباحث مطرح شده در محورهای موضوعی، کنفرانس به ارائه ثمره پژوهشهای خود همت نمودند؛ که با برگزیده ای از مقالات ترجمه ای در این عرصه پژوهشی ادغام و به طبع رسیده است. رجاء واثق دارم که با الطاف الهی از این پس مدیریت منابع اطلاعاتی وب و آرشیو وب به طور خاص موضوع پژوهش و ابتکار عمل اهل دانش و فناوری در کشورمان قرار خواهد گرفت و در این عرصه نیز فرزندان این مرز و بوم تجسم گفتار نغز رسول اعظم صلی الله علیه و آله خواهند بود که علم اگر تا ثریا، برود مردانی از فارس بدان دست خواهند یافت».

اسحق صلاحی

رئیس کنفرانس و

رئیس سازمان اسناد و کتابخانه ملی

جمهوری اسلامی ایران

ص: 8

گسترش روزافزون اطلاعات در شبکه اینترنت و سادگی بارگذاری انواع داده‌ها بر وب جهان را با شکل جدیدی از تولید انتشار و مصرف اطلاعات مواجه کرده است. تغییر جایگاه شهروندان جامعه از مصرف‌کننده صرف اطلاعات به مولّد و ناشر اطلاعات و فارغ از سازوکارهای موسوم سبب ساز روابطی جدید در عرصه ارتباطات اجتماعی و فرهنگ شده است. از سویی حجم رو به تزاید داده‌ها و چرخه عمر کوتاه اطلاعات موجود در وب موجب شده که گردآوری، «پالایش»، «سازماندهی»، ذخیره‌سازی و «اشاعه» آنها در زمره مسائل پژوهشی در نهادهای علمی و نیز بخشهای پژوهش و نوآوری شرکتها قرار گیرد؛ ضمن اینکه حفظ و دسترسی پایدار به اطلاعات موجود در وب که خود جزئی از میراث فکری ملتها محسوب میشود به دغدغه‌ای جدی برای سازمانها متولی حفظ و اشاعه میراث فکری به ویژه کتابخانه‌های ملی بدل شده است.

این حوزه در جهان موضوعی نسبتاً جدید است و پیشینه آن به کمتر از پانزده سال میرسد لیکن با سرعتی شتابان در حال رشد است و محققان مختلفی را از زوایای مختلف فنی حقوقی اقتصادی و حتی اخلاقی به سوی خود جذب کرده که گواه آن نیز طیف وسیعی از مقاله‌ها کتاب‌ها و گزارشهای سازمانی است که در طی چند سال اخیر در سطح جهانی منتشر شده است به رغم این نکات در ایران همچنان این زمینه حوزه‌ای بکر و کمتر مورد توجه محسوب میشود و در طی سالهای اخیر کمتر تحقیق بدان پرداخته لیکن تزاید اطلاعات فارسی بر روی وب و برنامه‌های ملی کشور مبنی بر توسعه کاربریهای مختلف بر شبکه‌های اطلاعاتی (از جمله توسعه دولت الکترونیکی، یادگیری الکترونیکی و کتابخانه‌های دیجیتال) لزوم توجه به این موضوع را بیش از پیش نمایان میسازد به همین دلیل سازمان اسناد و کتابخانه ملی جمهوری اسلامی همزمان با برگزاری نخستین کنفرانس ملی مدیریت منابع اطلاعاتی «وب» درصدد برآمد تا این حوزه را هرچه بیشتر به متخصصان و پژوهشگران باشناساند کتاب پیش رو حاصل همین نیت

متولیان این موضوع مهم است.

این کتاب مجموعه‌های قریب به 30 مقاله برگزیده از مهمترین منابع علمی منتشر شده در جهان و نیز قریب به 15 مقاله برگزیده از صاحب نظران ایران است که در قالب دو جلد تقدیم

حضور خوانندگان ارجمند می‌شود این مقالات در چهار موضوع اصلی به شرح زیر تقسیم شده اند:

● مبانی مدیریت و آرشیو وب

● تجارب جهانی و مسائل بومی در مدیریت و آرشیو وب

● مسائل فناورانه

● مسائل اخلاقی و مدیریتی

بی‌گمان این مجموعه می‌توانست به افزودنیهای دیگر (هم از منابع خارجی و هم از دیدگاه سایر متخصصان ایرانی) به اثری پربارتر بدل گردد لیک نخستین گامی است که در این حوزه برداشته شده و مطمئناً در مراحل بعدی با همت سایر، اندیشمندان ویراست‌هایی غنی‌تر از آن حاصل خواهد آمد نگارنده امیدوار است این مجموعه به مثابه بذری باشد که در کشتزار ذهن پژوهشگران کاشته شده و ان‌شاء الله در آینده‌ای نه‌چندان دور به نهالی پرتراوت در عرصه علم و عمل در جامعه اسلامیمان مبدل گردد.

در پیدایی این اثر کسان بسیاری همراهی و همکاری داشته‌اند که مقدم بر همه اندیشمندانی است که متن هر مقاله به‌خامه دانش افزای آنان امکان وجود یافته است از این رو نگارنده سپاس فروتنانه خود را نثار نگارندگان و مترجمان ارجمند این اثر مینماید گردآوری، تنظیم و آماده‌سازی مطالب کتاب به همت خانمها فرزانه شادان پور و مهشید برجیان بوده و ویراستاری آن را خانم آرزو تجلی برعهده داشته‌اند. ویراستار استنادی مقالات تألیفی را سرکار خانم فروزان رضایی نیا به انجام رسانده‌اند و سرکار خانم دکتر میترا صمیعی زحمت چکیده‌نویسی شماری از مقالات را که فاقد چکیده بودند متقبل شدند نمونه خوانی و اصلاحات اثر حاصل تلاش خانمها مهشید برجیان فاطمه رمضانپور، آهنگری آمنه هزار خوانی زهرا، زاهدی ملیحه حاجی زاده مقدم و آقای محمد رضا میقانی بوده است ضمن اینکه زیبایی متن و صفحه‌آرایی آن مدیون حسن سلیقه سرکار خانم شهره خوری است زحمات لیتوگرافی چاپ و صحافی کتاب نیز برعهده جناب آقای امیر آبادی بوده که بر خود فرض میدانند از همه این بزرگواران صمیمانه تشکر کند. گمان پدید آمدن این اثر به همت مسؤولان گرانمایه سازمان اسناد و کتابخانه ملی جمهوری اسلامی بوده است و نگارنده امیدوار است خداوند آنان را در مسیر خدمت به فرهنگ و دانش ایران اسلامی مورد تأیید قرار دهد.

اللهم وفقنا لما تحب وترضی

غلامعلی منتظر

تهران- بهمن ماه یکهزار و سیصد و نود و یک خورشیدی

ص: 2

فصل اول : مبانی مدیریت و آرشيو وب

اشاره

ص: 3

چکیده فیدهای وب با فرمت آر.اس.اس یا مبتنی بر اتم(1). ایکس. ام. ال، اسناد توصیفی در حال توسعه ای هستند که کانون هاب) پویای وبگاهها مشخص میکنند و به مشترکان کمک میکنند تا با تازه ترین محتوای وبگاه مورد علاقه و روزآمد خود در ارتباط باشند در این مقاله نشان میدهم چگونه فیدهای وب میتوانند ابزاری مفید برای استخراج اطلاعات و تشخیص تغییر صفحه وب باشند. معمولاً صفحه های وب که با آیتمهای فید ارجاع می شوند عبارت اند از پستهای وب نوشت و یا مقاله های خبری، و داده هایی با ماهیتی پویا (سپس زودگذر) که به صورت موضعی در یک کانال فید خوشه بندی میشود ما کانالهای وب را پایش و از صفحه های وب مرتبط، متن و منابع مربوط به مقاله های وب را استخراج میکنیم نتیجه کار با بر حسب زمان و فراداده اضافی استخراج شده از فید غنی شده، در یک شیء داده ای محصور می شود شیء داده ای به شکل اطلاعات خاصی خواهد بود که فاقد هرگونه عناصر تمپلیت یا تبلیغات میباشد این عناصر بی ربط، که معمولاً boilerplate نامیده می شوند نه تنها از دید برنامه خزشگر وقت گیر و جاگیر هستند، بلکه مانع فرآیند تجزیه و تحلیل دادهها میشوند ما نخست با خزش فیدهای وب برای یک دوره زمانی و مشاهده جنبه های زمانی، آنها روی مجموعه ای از آنها نوعی بررسی الگوریتم مورد استفاده برای استخراج مقاله را ارائه میکنیم؛ الگوریتمی آماری کرده، سپس که از معانی فید(2) به طور اختصاصی تر شرح و عنوان آیتمهای فید به منظور شناسایی گره DOM در صفحه اچ تی ام ال که حاوی مقاله است، استفاده می. کند از اشیای داده ای ساخته شده با این شیوه میتوان برای مجموعه همپوشانی معنایی برای آرشیو و یا در زمینه یک خزش تدریجی استفاده کرد که آن را از طریق تشخیص تغییر در سطح شیء داده ای کارآمدتر میکند آزمایشهای انجام شده بر روی روش استخراج به منظور روایی رویکرد مورد نظر، با نتایج خوبی - حتی در مواردی که تکنیکهای دیگر شکست خورده بودند - انجام میشوند. در نهایت در مورد برنامه های مفید براساس استخراج و تغییر تشخیص شیء وبگاه بحث میکنیم. کلید واژه ها آرشیو کردن، وب شیء داده ای فید، وب پویایی صفحه های وب

ص: 4

Atom -1

Semantic -2

نوشته ماریلنا اویتا(1) و پیرسنلارت(2) | ترجمه لیلی سیفی(3)

1. مقدمه

آرشیو وب بایگانی کردن (تارنما) [15] به فرآیند جمع آوری مکرر محتوای بخش‌هایی از شبکه جهانی، وب به منظور حصول اطمینان از حفظ آن و اجازه دسترسی به اطلاعات - حتی پس از بین رفتن وب اطلاق می شود برنامه خزشگر آرشیو وب به دنبال همان مراحل اولیه به محض اینکه خزشگر موتور جست و جو شاخص‌هایی را برای صفحه های وب ایجاد میکند اجرا می.شود با این حال، خزشگر آرشیو زمانی که نسخه های جدید کشف میشوند، شاخصها را به نسخه های قدیمی انتقال نمی دهد، بلکه آنها را ذخیره و نسخه ها را به موقع ارجاع میدهد نتیجه نهایی مجموعه ای از صفحه های وب است که می توان به صورت ناپیوسته مرور کرد، و در شکل ایده آل، می توان به صورت موقتی و معنایی جست و جو کرد. خزش منسجم سایتهایی که در به روزرسانی عرضه اطلاعات خود بسیار سریع هستند، کار آسانی نیست گزینه مرسوم خزش تصویرهای لحظه ای است؛ اما خزش کل مجموعه در بازه زمانی دوردست و) به اندازه کافی مکرر از لحاظ استفاده از پهنای باند شبکه بسیار گران قیمت است و در پایان نیز هم برای برخی مناطق سایت کاهنده و برای بخشهای پویاتر سایت ناقص خواهد بود علاوه بر این چون

ص: 5

Marilena Oita -1

Pierre Senellart -2

3- استادیار دانشگاه بیرجند Leili.seifi@gmail.com

اجرای یک خزش تلفیقی زمان بر است و در عین حال ممکن است منبع تغییر کند مشکلاتی ناشی از عدم انسجام زمانی تصویر لحظه ای معین به وجود خواهد آمد.

هنگام تصمیم گیری برای خزش تدریجی آنچه که اندازه کافی مکرر است، باید مشخص شود که با سایتهایی انطباق بهتری دارد که در ساختارشان دارای بخشهای موقت ناهمگن هستند هر خزشگر تدریجی نسخه کامل سایت را در یک مرحله خزش میکند و این روند نوعی راه اندازی (1) است که به خزشگر اطلاع میدهد که کدام محتوای سایت اضافه یا به روزرسانی شده است، و خزشگر فقط محتوای اصلاح شده را خزش کرده و آن را به صورت یک ساختار داده ای دلتا (2) ذخیره می کند. مشکل انجام خزش تدریجی الزاماً تعیین این پویایی است چند وقت به چند وقت صفحه های وب جدید اضافه یا اینکه صفحه های موجود اصلاح میشوند و کارآمدی تشخیص تغییرات با علم به اینکه عوامل زیادی وجود دارند که میتوانند روی فرآیند تشخیص اثر منفی داشته باشند موضوع مورد توجه ما تکنیکهایی است که می تواند برای بهبود فرآیند تشخیص مورد استفاده قرار گیرد، در مورد خاص که در آن صفحه های که باید خزش شوند دارای آر.اس.اس یا فیدهای اتم الصافی هستند.

روش سنتی تشخیص تغییر بین دو نسخه پی در پی یک صفحه، وب مقایسه در محتوا با استفاده از معیارهای شباهت، است که با در هم سازی و امضاهای محتوا به منظور ویرایش فاصله ها برای انعطاف پذیری متفاوت است. تغییر به هر شیوه ای که شناسایی شود اهمیت آن با توجه به نوع محتوای یک صفحه وب مورد ارزیابی قرار نمی گیرد. با این حال درک اینکه آیا تغییرات مربوط به محتوای اصلی صفحه از نقطه نظر معنایی باشند بسیار مهم است؛ زیرا در برخی برنامه های کاربردی، تغییراتی که تنها روی قسمتهای [12 boilerplate] صفحه از قبیل، منوها تمپلیتهای ارائه یا تبلیغات اثر میگذارند ممکن است به راحتی نادیده گرفته شده است.

مطالعه اندکی در مورد فیدهای وب صورت گرفته است در حالی که پدیده ای به سرعت در حال تحول است. ما توجه خود را به این واقعیت معطوف می داریم که میتوان آنها را به عنوان ابزاری در تجزیه و تحلیل یک وبگاه قبل و در حین خزش استفاده کرد فیدهای وب علاوه بر اینکه راهی برای تبلیغ محتوا هستند جهت طبقه بندی منابع اطلاعاتی و نوع محتوا توسط موتورهای جست و جو استفاده می شوند. به طور خلاصه از طریق فیدهای، وب جنبه های مهم وبگاه پویا را میتوان در چوب یک خزش وب - به منظور آگاه تر ساختن آن از اطلاعاتی که تهیه میکند استخراج و بهره برداری کرد.

ماهیت فید

آموزنده است اطلاعات جدید و زمان ورود آنها، و انتشارشان توسط کانال را ترکیب می کند؛ توصیفی است توضیح میدهد چه نوع منابع جدیدی همراه با عنوان توضیح، و سایر عناصر تگ ممکن افزوده شود.

ص: 6

Trigger -1

Delta data Structure -2

هدف ما استخراج داده‌های ساختاری از صفحه‌های وب، با کمک فیدهاست. اساس و پایه رویکرد ما این است که یک آیتم از یک کانال [چگونه با یک شیء داده‌ای در یک صفحه وب انطباق میابد.

بنابراین، فراداده‌ای که در مورد آیتم مورد نظر در فید به دست می‌آید را میتوان برای تشخیص و استخراج شیء داده در صفحه وب استفاده کرد.

مفهوم شیء داده‌ای دارای تفسیرهای مختلفی در علوم رایانه است به منظور روشن ساختن اهمیت آن در زمینه مورد بررسی باید گفت که یک شیء داده‌ای نمونه‌ای از منابع ارجاعی فید است و فراداده‌هایی دارد که با عبارات «(1)».. [5] و «به(2)» [13] با هم مرتبط میشوند خود مفهوم شیء داده یک ترکیب است و از این نظر میتوان آن را به عنوان سند منطقی ما میگوییم معنایی در تقابل با «سند فرامتن»(صفحه اچ.تی.ام.ال) تلقی کرد.

حتی اگر شیء اغلب یک مقاله وب، باشد میتواند مدخل فرهنگ لغت، نظر، پیامی در یک نشست ویدئو، یک وضعیت، و هر نوع منابع دیگری باشد که به طور منحصر به فرد با آیت‌های فید وب مرتبط شده است رویکردهای مستقیم برای شناسایی محتوای اصلی یک صفحه وب، از جمله در نظر گرفتن عناصری که پس از تمیز کردن پایه مطالعه تراکم متن در مناطق خاصی از صفحه [12]، و یا حتی شناسایی برجسته‌ترین مناطق بصری صفحه، [29] با پیچیده تر شدن شیوه رمزگذاری و تکامل یافتن خود درک مستقیم مقاله، وب در حال منسوخ شدن هستند اشیای داده‌ای میتوانند مطابق با بخش‌های کوچکی از یک صفحه، وب ساده یا مرکب شامل چند رسانه ایها و حتی جاوا اسکریپت‌های خوبی باشند که باید مورد بهره برداری قرار گیرند. بسته به زمینه میتوانیم یک یا چند شیء داده چندگانه در هر صفحه داشته و تفاوت بین آنها را با استفاده از معناشناسی با قابلیت خوانده شدن توسط ماشین و یا با استفاده از فناوری هوشمند بشر انجام داد ما اهرم‌های معنایی هستیم که توسط فیدهای وب به صفحه‌های وب پویا آورده میشوند فیدها را - که در فرم ایکس ام ال مانند با عناصر استاندارد نوشته شده اند - می‌توان برای به تصرف در آوردن برخی جنبه‌های مهم اطلاعاتی استفاده کرد که می‌خواهیم استخراج کنیم با استفاده از تکنیک‌های کلاسیک بازیابی اطلاعات که از ساختار یک آیتم فید به دست می‌آوریم، توصیفگر معنایی شیء که به عنوان ورودی برای الگوریتم استخراج می‌شود، استفاده خواهد شد. در مرحله بعد، منطقه‌ای در یک صفحه وب را شناسایی میکنیم که شامل محتوای شیء داده بوده و آن را با استفاده از تجزیه و تحلیل تراکم معنایی استخراج می‌کند میتوان با داشتن محتوای استخراج شده در زمان و بخش مهمی از خواص آن پرس و جوهای پیچیده را از نقطه نظر زمانی و معنایی اجرا کرد.

در بخش بعدی برخی کارهای مرتبط را ارائه و در بخش 3 نتایج حاصل از مطالعاتمان بر روی فیدهای وب توصیف میکنیم که به منظور تعیین ارزش این فیدها در فرآیند آرشیو وب انجام شد. ما در بخش 4 شرح میدهیم که چگونه معناییهای مشخص شده توسط عناصر خاص موجود در یک فید را میتوان برای استخراج اشیای داده‌ای استفاده کرد. در بخش، آزمایش‌های مربوط به استخراج

ص: 7

and -1

to -2

اشیای داده ای نشان داده میشوند نتیجه گیری ما با ترسیم اهمیت اشیای داده ای استخراج شده در زمینه برنامه های کاربردی است.

2. کارهای مرتبط

آرشیو وب چه به عنوان یک ضرورت یا وظیفه دیده شود به تازگی از اهمیت زیادی برخوردار شده است که علت آن ماهیت قرار، وب و به ویژه ارزشی است که اطلاعات از دست رفته میتواند برای نسلهای آینده داشته باشد. آرشیو اینترنتی [11] یکی از مبتکران جنبش آرشیو وب است. بسیاری از بازیگران دیگر وجود دارد که به طور فعال خزش را به عنوان بخشی از رسالت حفظ میراث خود اجرا کرده و برای اینکه مجموعه آرشیو وب جهانی به صورت واحد تقارب، یابد تلاشهایی در شرف انجام است [28]. اگرچه فیدهای وب به طور معمول به عنوان انواع دیگری از اسناد وب توسط خزشگرهای آرشیو نمایه میشوند تلاش محدودی در بهره برداری از این ویژگیهای در روند آرشیو کردن صورت گرفته است آرشیو پرس که پروژه ای برای آرشیو کردن وب نوشت است [20] در حال توسعه نوعی پلاگین وردپرس است که پستها را با استفاده از فیدهای وب آرشیو می کند اشکال اصلی این است که تنها می توان محتوایی را گرفت که توسط فید آر.اس.اس قابلیت تحویل را داشته باشد. هر فید آر.اس.اس. در واقع میتواند پوشش کامل مقاله و فایل های رسانه ای را داشته باشد. اما این مورد بسیار نادر است زیرا هر فید اغلب فقط یک شیوه محتوای تبلیغاتی است در مقابل سرخ های فید را با هدف بهره برداری از اطلاعات واقعی مربوط تحت کنترل در می آوریم. علاوه بر این، خود را به یک چارچوب، وبلاگ نویسی یا وبلاگ نویسی به طور خاص محدود نمی کنیم.

هنگام مطالعه پویایی صفحه های وب دو دیدگاه وجود دارد این تغییر میتواند ناشی از انتشار محتوای جدید و یا ناشی از تغییراتی باشد که در صفحه های موجود رخ میدهد. برای جلوگیری از تجمع بیش از حد بر روی سرور وبگاههایی که محتوای به موقع را تهیه میکنند، [22]، رویکرد رایانه انطباقی پیشنهاد می شود. مدل ناهمگن پواسون (1) در فیدهای وب، به منظور یادگیری الگوهای ارسالی در بلاگها و پیش بینی بررسی مجدد بهینه برای محتوای جدید مورد استفاده قرار می گیرد در زمینه خزش افزایشی در حال حاضر مشکل مشابه این است که آمار خود را در جنبه های زمانی فیدها با هدف بیان ارزش خود در فرآیند یادگیری راهبرد انتشارات ارائه میدهد در زمینه خزش تدریجی نیز همان مشکل را داریم آمارهای ما در مورد جنبه های موقتی با هدف بیان ارزشهایشان در فرآیند یادگیری راهبرد انتشار انجام می شوند. به منظور مدل سازی رفتار وبگاهها در طول زمان و خزش مؤثرتر، آنها مطالعات دیگر درباره درک درست پویایی صفحه های وب متمرکز شده اند، [8، 16] میزان سرعت تغییرات محتوای وب و نوع ماهیت این تغییرات معانی ضمنی در مورد ساختار و همبستگی با موضوع صفحه ها برای انطباق خودکار خزش با آهنگ پیشبینی شده تغییرات مقدماتی وجود دارد [10] از مدل فضای برداری متنی

ص: 8

برای شناسایی الگوهای صفحه و برای آموزش فیلترهای کالمن(1) استفاده میشود در پایان تغییر رویدادی است که با پیش بینی مطابقت ندارد با این، حال فرضیه خطی سیستم و عدم تمامیت مدل فضای برداری اشکالهای ممکن را نشان میدهد. برای شناسایی و ذخیره سازی محتوای اصلاح شده یک صفحه وب خزشگر آرشیو کد منبع باز [23] Heritrix] با استفاده از عبارتهای غیر مستدل و منظم برای فیلتر کردن تغییرات بی ربط استفاده میشود در تلاش برای برآورد منصفانه آهنگ تغییر صفحه وب مدل رسمی تر پیش بینی توسط چو(2) و گارسیا-مولینا(3) [3] مورد مطالعه قرار گرفت که در آن، نویسنده بررسی می که آیا تغییرات یک صفحه وب از فرآیند پواسون همگن تبعیت میکند یا نه با این وصف، شناسایی آهنگ تغییر توسط هر دو رویکرد به عنوان چالشی مطرح میشود ما استخراج اشیای داده ای در داخل این موضوع را توصیه میکنیم هنگامی که دستیابی اشیای داده ای به موقع انجام می شود، یک معیار شباهت را میتوان در محتویات و یا بر روی خواص اشیای داده ای استفاده کرد. با این فرض که آنها دو نسخه پی در پی یک صفحه وب را با استفاده از چند روش برای تشخیص تغییر به دست آورده اند پهلوان(4) بن سعد(5)، و گانکارسکی(6) [19] روی تشخیص تغییراتی تمرکز میکنند که براساس یک نسخه قدیمی در نسخه جدید رخ داده است برای این منظور، الگوریتم ویس(7) [29] برای شناسایی بخشهای معنایی مرتبط یک صفحه وب استفاده میشود که به منظور تشخیص تغییرات ساختاری و محتوایی مقایسه می شوند. ابتکارهایی روی ظاهر بصری یک صفحه وب ایجاد میشوند تا محتوایی را با هم گروه بندی کنند که به نظر میرسد در صفحه از اهمیت یکسانی برخوردار باشند این ابتکارها مشکل را به طور جامع پوشش نمیدهند و الگوریتم محاسباتی گران قیمت است. رویکرد ما محدودتر است، چون نیاز داریم که از طریق فید عبور کنیم؛ در عین حال، می تواند مؤثرتر باشد: با شناسایی مناطق معنایی «مهم» در صفحه وب میتوانیم بر روی تغییراتی تمرکز کنیم که به تلاش کمتری وابسته هستند.

یک مقاله وب شناسایی شده با استفاده از عنوان و شرح آیتم فید باید از کد اچ.تی.ام.ال. صفحه وب مرتبط استخراج شود کار زیادی در استخراج بدون نظارت دادههای ساختاری از صفحه های وب انجام شده است؛ بسیاری از اینها مبتنی بر ام.دی. آر(8). [14]، اکس آلز(9) [2] یا دونده جاده [4] هستند. در واقع این روشها تلاش میکنند تا با استنتاج گرامر برای کد اچ تی ام ال به طور خودکار لفافهای تولید کند که حاوی اطلاعات مورد علاقه به طور کلی به، آن سوابق داده گفته میشود، باشد به شیوه ای که به دانش قیاسی در مورد صفحه های هدف و محتویاتشان وابسته نباشد معمولاً صفحه های مختلف که همان قالب را دارند به منظور مقایسه زوجی آنها و کشف الگوهای مشترک و قواعد کد گذاری مورد نیاز هستند این

ص: 9

Kalman -1

Cho -2

Garcia-Molina -3

Pehlivan -4

Ben Saad -5

Gancarski -6

VIPS -7

MDR -8

Ex Alg -9

قواعد یا از طریق بررسی شباهتها و تفاوت‌های بین صفحه‌ها [4] یا با ساخت کلاسهای هم ارز [2] تأکید میشوند. برخلاف کارهای پیشین ام.دی.آر. [14] ساختار درختی DOM صفحه اچ.تی.ام.ال. را در نظر گرفته و منطقه داده‌ها را با پیدا کردن گرهی کلی شناسایی میکند که شامل بیشترین تعداد فرزندان است که الگوهای مشابه را با توجه به اندازه گیری شباهت ارایه میدهند. حتی اگرچه ما توجه خود را روی سایتهای حاوی مقاله‌های ویی متمرکز میکنیم - الزاماً به شیوه‌ای که ما در یک صفحه وب میبینیم - متوالی نیستند اما ما با یک مشکل مشابه استخراج داده‌ها احتمالاً ساختاری با یک ماهیت پیچیده تر و منزوی تر از این رویکردها مواجه هستیم و همانطور که در بخش 4 توضیح داده خواهد شد، به نحوی از تکنیکهای مربوط استفاده می‌کنیم.

جدول 1. انواع فیدهای مجموعه‌ای

عکس

قواعد، یا از طریق بررسی شباهت‌ها و تفاوت‌های بین صفحه‌ها [۴]، یا با ساخت کلاس‌های هم‌ارز [۲] تأکید می‌شوند. برخلاف کارهای پیشین، ام‌دی‌آر [۱۴] ساختار درختی DOM صفحه اچ‌تی‌ام‌ال را در نظر گرفته، و منطقه داده‌ها را با پیدا کردن گرهی کلی شناسایی می‌کند که شامل بیشترین تعداد فرزندان است که الگوهای مشابه را با توجه به اندازه‌گیری شباهت ازایه می‌دهند. حتی اگر چه ما توجه خود را روی سایت‌های حاوی مقاله‌های وبی متمرکز می‌کنیم - الزاماً به شیوه‌ای که ما در یک صفحه وب می‌بینیم - توالی نیستند، اما ما با یک مشکل مشابه استخراج داده‌ها (احتمالاً ساختاری با یک ماهیت پیچیده‌تر و منزوی‌تر از این رویکردها) مواجه هستیم، و همانطور که در بخش ۴ توضیح داده خواهد شد، به نحوی از تکنیک‌های مربوط استفاده می‌کنیم.

جدول ۱. انواع فیدهای مجموعه‌ای

Type	Number	Proportion
Atom	21	6.1%
RDF	30	8.8%
RSS 0.91	1	0.2%
RSS 2.0	288	84.7%
Total	340	100.0%

دو رویکرد دیگر برای مشکل استخراج مقاله اصلی از یک صفحه وب به‌تازگی توسط خولچوتر^{۱۰}، فن خسار^{۱۱} و نجدی^{۱۲} (۱۲) و پاسترناک^{۱۳} و راث^{۱۴} [۱۸] پیشنهاد شده است. در حالی که [۱۲] از تراکم متن بر روی صفحه وب برای شناسایی مقاله استفاده می‌کند، [۱۸] از روش تقسیم‌بندی متن توالی برای رسیدن به نتیجه مشابه استفاده می‌کند. الگوریتم استخراج اشیای داده ما، مشابه به ام‌دی‌آر، از برخی ابتکارات در کد اچ‌تی‌ام‌ال، به‌منظور شناسایی منطقه‌ای که در آن می‌توان مقاله را یافت استفاده می‌کند، اما برخلاف تمام روش‌های برشمرده شده ما با استفاده از معنایی برگرفته شده از فید برای استخراج اشیای داده‌ای استفاده می‌کنیم. این امر به وضوح تنها برای صفحه‌هایی ممکن است که به یک فید وب پیوند شوند. در بخش ۵، ما نتایج را با نتایج الگوریتم [۱۲] Boilerpipe، که به‌عنوان پایه از آنچه که می‌توان بدون بهره‌گیری از اطلاعات فید وب انجام داد، مقایسه می‌کنیم. طبق بررسی‌های ما، هیچ کار قبلی، که اطلاعات معنایی را (که ممکن است از فید وب یا از هر منبع دیگر آمده باشد) جهت استخراج بخش مربوط به صفحه وب - به‌شیوه‌ای کلی و بدون نظارت - اهرم نمی‌شود.

10. Khoschutter
11. Fankhauser
12. Nejadi
13. Pasternack
14. Roth

دو رویکرد دیگر برای مشکل استخراج مقاله اصلی از یک صفحه وب به‌تازگی توسط خولچوتر⁽¹⁾، فن خسار⁽²⁾، و نجدی⁽³⁾ (12) و پاسترناک⁽⁴⁾ و راث⁽⁵⁾ [18] پیشنهاد شده است. در حالی که [12] از تراکم متن بر روی صفحه وب برای شناسایی مقاله استفاده میکند [18] از روش تقسیم بندی متن توالی برای رسیدن به نتیجه مشابه استفاده میکند. الگوریتم استخراج اشیای داده، ما مشابه به ام‌دی‌آر، از برخی ابتکارات در کد اچ تی ام ال به منظور شناسایی منطقه ای که در آن میتوان مقاله را یافت استفاده می کند، اما برخلاف تمام روشهای بر شمرده شده ما با استفاده از معنایی برگرفته شده از فید برای استخراج اشیای داده ای استفاده میکنیم این امر به وضوح تنها برای صفحه هایی ممکن است که به یک فید وب پیوند شوند در بخش 5 ما نتایج را با نتایج الگوریتم [12] Boilerie که به عنوان پایه از آنچه که میتوان

بدون بهره‌گیری از اطلاعات فید وب انجام داد مقایسه می‌کنیم طبق بررسیهای ما، هیچ کار قبلی که اطلاعات معنایی را (که ممکن است از فید وب یا از هر منبع دیگر آمده باشد) جهت استخراج بخش مربوط به صفحه وب - به شیوه ای کلی و بدون نظارت - اهرم نمی‌شود.

ص: 10

Khoschutter -1

Fankhauser -2

Nejdi -3

Pasternack -4

Roth -5

به منظور اثبات ارزش فیدها به عنوان ابزارهای تجزیه و تحلیل در فرآیند تشخیص تغییر وب طی یک دوره کمی بیش از یک ماه دو بار در روز تعداد 400 فید وب را همراه با تمام صفحه های وب مرتبط با آن خزش کرده ایم. نخست، چگونگی انتخاب این فیلها را توصیف و سپس برخی آمار جالب مجموعه داده هایمان را گزارش میکنیم.

3-1. فراهم آوری

مجموعه ای از فیدهای وب با عبور از بخش بزرگی از طریق یک موتور جست و جوی فید به نام جست و جوی آر.اس.اس.4 [21] (1) جمع آوری شد. این موتور جست و جو تعداد فیدها و صفحه های کانال مرتبط برای یک کلید واژه را بررسی میکند ما شیوه ای را که در آن نتایج به شکل سوابق برگردانده می، شد کنار گذاشتیم و تمام یو. آر. ال. (2) های تجزیه شده را جهت تجزیه و تحلیل بیشتر در فهرست فایل قرار دادیم.

کلید واژه های انتخابی برای ردیابی واسط جست و جو نام دامنه های زیر بود: هنر (3)، زیست شناسی (4)، محیط (5)، دارو (6)، علم (7)، و جهان (8). به منظور دریافت مترادفهای این کلمات برای (مثال مترادف هنر عکاسی (9) است) از شبکه ورد (10) استفاده کرده ایم و کیسه هایی از واژه های معرف زیر دامنه ساختیم.

این کیسه های واژه برای ردیابی خودکار واسط نام جست و جوی آر.اس.اس.4 و برای ساخت برای هر، دامنه فهرستی از یو آر ال های فید مورد استفاده بود [بار] معنایی واژه ها به ما امکان تمرکز جست و جو برای فیدها و شناسایی صفحه های وی را داد که به عنوان یک موضوع خاص تلقی می شدند. بازتاب در رسانه دامنه های مورد علاقه را میتوان از طریق چشمهای فید گرفتار شده به این شیوه - که معمولاً برخی از الگوها و ویژگیهای خاص را ارائه میدهد - مشاهده کرد.

هدف از این انتخاب نیمه خودکار به دست آوردن درک و بینش کلی نسبت به تنوع، فیدها، برحسب فرمتها الگوهای به روز رسانی و ساختارهای متنوع صفحه های وب متناظر بود. علاوه بر این، به منظور حصول اطمینان از پوشش سیستم عاملهای وب نوشت رایج و همچنین نتایج بیشتر خبرگرای بازگردانده شده توسط نام جست و جوی آر.اس.اس.4 تعدادی از سایتهای وب نوشت به صورت دستی از فهرست بهترین وب نوشت [24] انتخاب شدند. به این ترتیب، فهرستی از حدود 400 سایت به دست آمد که

ص: 11

Search4RSS -1

URL -2

Art -3

Biology -4

Environment -5

Medicine -6

Science -7

Universe -8

WordNet -9

Photography -10

به طور سیستماتیک خزش میشدند (این عدد پوشش انواع فیدهای وب است که بدون هیچ گونه زیر ساخت درگیر آرشیو مدیریت پذیر هستند). در پایان دوره، خزش متوجه شدیم که برخی از انواع فیدها هرگز به روزرسانی نشده برخی دیگر ناپدید شده و برخی را نتوانستیم تجزیه کنیم. با فیلتر کردن، آنها آرشیوی از 340 فید فعال و صفحه های مرتبط با آنها را به دست آوردیم.

فید وب به یک صفحه وب اولیه (کانال) است که معمولاً یا صفحه اصلی سایت یا هابی است که ما میتوانیم از آن پیوندهای اطلاعات ارائه شده به شکل مقاله های وب را پیدا کنیم. بقیه فید، آیتمهای فید فردی مربوط به مقاله های جدید و یا به روزرسانی شده را توصیف میکند برای هر دامنه برای هر سایت پیگیری شده فید و منابع مرتبط - به طور عمده صفحه کانال و صفحه های وبی - را ذخیره کرده ایم که با هر آیتمی اشاره شده بودند به منظور از دست ندادن آیتمهای جدیدی که بتوان در فید اضافه کرد، اجرای این خزشگرها دوبار در روز (علاوه بر اجرای خزشگرهای همان فیدها که روزانه توسط آرشیو وب اروپا اجرا میشود) انجام گرفت.

2-3- ویژگیهای فید وب برای تجزیه و تحلیل، فید از ادی (1) [7] - کتابخانه تجزیه فید برای جاوا (2) - استفاده کرده ایم، که بر اساس تجزیه مبتنی بر سکس (3) قادر به تجزیه حتی دنیای واقعی به شکل ایکس ام ال است. ادی، از فرمتهای آر.اس.اس استاندارد، اتم و آر.دی.اف (4) برای، فیدها پشتیبانی می. کند ساختار داده فید بازگردانده شده توسط این تجزیه را میتوان برای استخراج همه نوع اطلاعات مفید در مورد کانال و آیتمهای تشکیل دهنده مورد تفحص قرار داد به طور خاص جهت، کانال زبان شبیه، فراداده شعار شرح و عنوان و همچنین برای هر مورد دیگر به علاوه نویسنده و مقوله هایی که در آن مقاله طبقه بندی شده است.

نوع فید: اجازه دهید نخست نگاهی به انواع فرمتهای فید که در مجموعه داده هایمان برشمردیم، داشته باشیم همانطور که در جدول 1 نشان داده شده است بیشتر فیدها از گویش آر.اس.اس. 2/0 استفاده می میکنند در حالی که اقلیتی هم وجود دارد که از اتم یا آر.دی.اف استفاده میکنند از آر.اس.اس. 0/91، تنها یک بار در میان 340 فید استفاده شد و از آر.اس.اس. 1/0 هرگز استفاده نشد، که ممکن است منسوخ شدن این دو فرمت فید را نشان دهد با این حال اینکه این اعداد نیز به دلیل استفاده از جست و جوی آر.اس.اس. 4، به عنوان منبع اصلی ما برای فیدها دچار تورش شوند، کاملاً امکان پذیر است.

تعداد آیتمها: تعداد تعداد آیتمهای ارائه شده در یک فید مفروض را بررسی کردیم. در واقع هر چند به طور نظری این امکان وجود دارد برای فید به همه آیتمهای منتشر شده قبلی اشاره کند به ندرت برای محدود کردن اندازه فید وب حاصل استفاده میشود در واقع بسیاری از فیدها کوتاه شده، تنها جدیدترین آیتم k برای یک مفروض را ارائه میدهد در شکل 1 هیستوگرام تعداد آیتمها به ازای هر فید در مجموعه داده

ص: 12

Eddie -1

Java -2

SAX -3

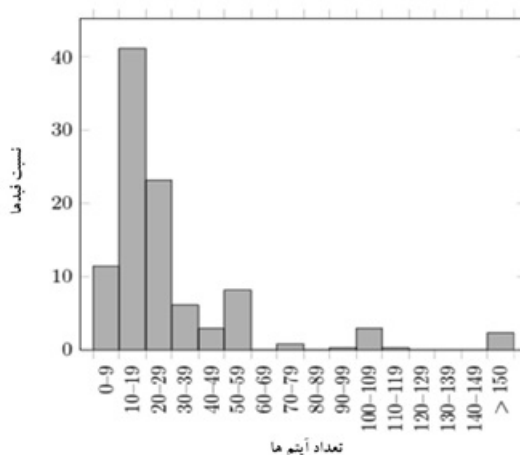
RDF -4

را نشان می‌دهیم. حدود 75 درصد فیدها، اطلاعاتی حدود کمتر از 30 آیتم را در یک زمان ارائه می‌دهند. سایر پیکهای مشاهده شده در شکل 1 با مقادیر جادویی « $K = 50$ و $K = 100$ توضیح داده می‌شوند. اگر فیدی تنها شامل 10 آیتم مرسوم ترین عدد باشد، بدان معنی است که با دو بار در روز خزش آن، ما توانسته ایم روزانه حداکثر 20 مقاله جدید را ذخیره کنیم همانطور که خواهیم دید تعداد خیلی از فیدها با یک فرکانس به روزرسانی بالاتر از آن وجود دارد، که به همین دلیل برخی از به روزرسانیها در خزش ما در واقع فیدها باید بیشتر از دو بار در روز خزش شوند از دست می‌روند.

اطلاعات موقتی: در خصوصیات آر.اس.اس به همین ترتیب در سایر فرمت‌های فید) اطلاعات

عکس

را نشان می‌دهیم. حدود ۷۵ درصد فیدها، اطلاعاتی حدود کمتر از ۳۰ آیتم را در یک زمان ارائه می‌دهند. سایر بیکهای مشاهده شده در شکل ۱ با مقادیر «جادویی» $k = 50$ و $k = 100$ توضیح داده می‌شوند. اگر فیدی تنها شامل ۱۰ آیتم (مرسوم‌ترین عدد) باشد، بدان معنی است که با دو بار در روز خزش آن، ما توانسته‌ایم روزانه حداکثر ۲۰ مقاله جدید را ذخیره کنیم. همانطور که خواهیم دید، تعداد قلیلی از فیدها با یک فرکانس به‌روزرسانی بالاتر از آن وجود دارد، که به همین دلیل، برخی از به‌روزرسانی‌ها در خزش ما (در واقع، فیدها باید بیشتر از دو بار در روز خزش شوند) از دست می‌روند.



شکل ۱. تعداد آیتم‌ها به‌ازای هر فید در مجموعه داده‌ها

اطلاعات موقتی: در خصوصیات آر.اس.اس (به همین ترتیب در سایر فرمت‌های فید)، اطلاعات زمانی را می‌توان از طریق عناصر `lastBuildDate`، `TTL`، و `updateFrequency` برای کانال، و `pubDate` و `lastModified` برای آیتم‌ها به‌دست آورد. از طریق آزمایش‌ها، ما مشاهده کرده‌ایم که اگرچه `pubDate` مؤلفه‌ای اختیاری است، اما در اکثریت قریب به اتفاق فیدها ارائه می‌شود. این در مورد انواع دیگر عناصر مرتبط به زمان (`timerelated`) مذکور صدق نمی‌کند، هر چند `lastBuildDate` را می‌توان به‌نحوی به‌عنوان تاریخ انتشار جدیدترین آیتم استنباط کرد. اهمیت این مشاهدات از آن روست که نشان می‌دهد که فیدها را می‌توان برای تعیین زمانی که داده جدیدی به یک کانال اضافه می‌شود، مورد استفاده قرار داد و در تشخیص تغییر کمک‌رسان است. با تجزیه و تحلیل یک فید برای یک دوره زمانی، می‌توانیم الگوها را در انتشار راهبرد، شناسایی و به‌طور خودکار خزش را با آن منطبق کنیم.

پازه به‌روزرسانی: تمام تاریخ انتشارات مربوط به آیتم‌های ظاهر شده که در طول دوره آزمایش را جمع‌آوری کرده‌ایم؛ از آنجا که هر آیتم یک تاریخ انتشار دارد، تعداد تاریخ‌های انتشار برابر است با آیتم‌ها. ما به طیف وسیعی از بازه‌های به‌روزرسانی بین دو انتشار، و همچنین نشانه‌های وجود راهبرد انتشار منظم

زمانی را می‌توان از طریق عناصر `lastBuildDate`، `TTL`، و `updateFrequency` برای کانال و `pubDate` و `lastModified` برای آیتم‌ها به‌دست آورد از طریق آزمایش‌ها ما مشاهده کرده ایم که اگر چه `pubDate` مؤلفه ای اختیاری است اما در اکثریت قریب به اتفاق فیدها ارائه میشود این در مورد انواع دیگر عناصر مرتبط به زمان (`timerelated`) مذکور صدق نمی کند هر چند `lastBuildDate` را میتوان به نحوی به عنوان تاریخ انتشار جدیدترین آیتم استنباط کرد اهمیت این مشاهدات از آن روست که نشان میدهد که فیدها را میتوان برای تعیین زمانی که داده جدیدی به یک کانال اضافه میشود مورد استفاده قرار داد و در تشخیص تغییر کمک رسان است. با تجزیه و تحلیل یک فید برای یک دوره زمانی، می توانیم الگوها را در انتشار راهبرد شناسایی و به طور خودکار خزش را با آن منطبق کنیم.

بازه به روزرسانی تمام تاریخ انتشارات مربوط به آیتمهای ظاهر شده که در طول دوره آزمایش را جمع آوری کرده ایم؛ از آنجا که هر آیتیم یک تاریخ انتشار دارد تعداد تاریخهای انتشار برابر است با آیتمها. ما به طیف وسیعی از بازه های به روزرسانی بین دو انتشار و همچنین نشانه های وجود راهبرد انتشار منظم

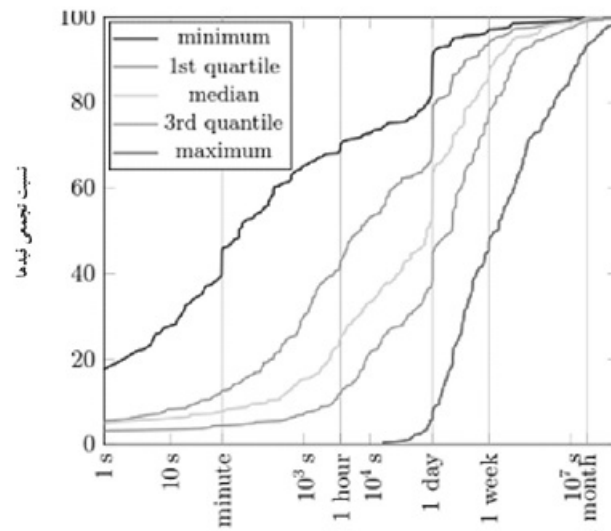
ص: 13

برای فید توجه زیادی میکنیم در شکل 2، متوسط بازه به روزرسانی بین دو انتشار هر فید را به صورت نمودار تجمعی (به رنگ سبز در وسط نشان داده ایم). توجه داشته باشید که محور x مقیاس لگاریتمی است. شکل 2 نشان میدهد که به عنوان مثال 20 درصد فیدها یک بازه به روزرسانی متوسط کمتر از یک ساعت و حدود 10 درصد فیدها یک بازه به روزرسانی متوسط دقیقاً یک روزه دارند که مربوط به فیدهایی است که به طور منظم و خودکار هر روز به روزرسانی میشوند. در سطح جهانی، مهم که توجه داشته باشید که هیچ فاصله به روزرسانی نمونه واری وجود ندارد و حتی بدون در نظر گرفتن موارد شدید میتواند به محدوده کمتر از یک ساعت تا به بیش از یک هفته برسد. همچنین شکل 2 مقادیر دیگر فاصله به روزرسانی هر فید را نشان میدهد که کمک میکند تنوع الگوهای به روزرسانی برای یک فید مفروض را بتوان دید به این ترتیب، حتی اگر چه 60 درصد فیدها دارای یک بازه به روزرسانی متوسط یک روز یا کمتر بودند، کمتر از حدود 10 درصد آنها همیشه حداقل یک به روزرسانی در هر روز، و بیش از 90 درصد آنها حداقل به روزرسانی روزانه در دوره مشاهده داشتهاند. در واقع، شکاف بزرگی بین مقادیر بازه به روزرسانی متوسط و مقادیر حداقل و حداکثر وجود دارد که باعث بروز مشکل پیش بینی به روزرسانی بعدی یک فید مفروض میشود.

نسبت تجمعی فیدها

عکس

برای فید توجه زیادی می‌کنیم. در شکل ۲، متوسط بازه به‌روزرسانی بین دو انتشار هر فید را به‌صورت نمودار تجمعی (به رنگ سبز، در وسط) نشان داده‌ایم. توجه داشته باشید که محور X مقیاس لگاریتمی است. شکل ۲ نشان می‌دهد که به‌عنوان مثال، ۲۰ درصد فیدها یک بازه به‌روزرسانی متوسط کمتر از یک ساعت، و حدود ۱۰ درصد فیدها یک بازه به‌روزرسانی متوسط دقیقاً یک روزه دارند، که مربوط به فیدهایی است که به‌طور منظم و خودکار، هر روز، به روز رسانی می‌شوند. در سطح جهانی، مهم است که توجه داشته باشید که هیچ فاصله به روز رسانی نمونه‌واری وجود ندارد، و حتی بدون در نظر گرفتن موارد شدید می‌تواند به محدوده کمتر از یک ساعت تا به بیش از یک هفته برسد. همچنین شکل ۲، مقادیر دیگر فاصله به‌روزرسانی هر فید را نشان می‌دهد، که کمک می‌کند تنوع الگوهای به‌روزرسانی برای یک فید مفروض را بتوان دید: به این ترتیب، حتی اگر چه ۶۰ درصد فیدها دارای یک بازه به‌روزرسانی متوسط یک روز یا کمتر بودند، کمتر از حدود ۱۰ درصد آنها همیشه حداقل یک به‌روزرسانی در هر روز، و بیش از ۹۰ درصد آنها حداقل به‌روزرسانی روزانه در دوره مشاهده داشته‌اند. در واقع، شکاف بزرگی بین مقادیر بازه به‌روزرسانی متوسط، و مقادیر حداقل و حداکثر وجود دارد، که باعث بروز مشکل پیش‌بینی به‌روزرسانی بعدی یک فید مفروض می‌شود.



شکل ۲. نسبت تجمعی فیدها با مقادیر یک چهارم مفروض فاصله به‌روزرسانی‌ها

در جدول ۲ برخی آمارهای دیگر در مورد فواصل به‌روزرسانی در سطح دامنه را نشان می‌دهیم. برای هر دامنه، میانگین فاصله زمانی به‌روزرسانی متوسط به‌عنوان انحراف استاندارد ادغام شده فواصل به‌روزرسانی داده‌ایم، که روشی است مبتنی بر آمار جهت خلاصه کردن انحراف از اتحاد مجموعه‌ای از

در جدول ۲ برخی آمارهای دیگر در مورد فواصل به‌روزرسانی در سطح دامنه را نشان می‌دهیم برای هر، دامنه میانگین فاصله زمانی به‌روزرسانی متوسط به‌عنوان انحراف استاندارد ادغام شده فواصل به‌روزرسانی داده‌ایم که روشی است مبتنی بر آمار جهت خلاصه کردن انحراف از اتحاد مجموعه‌ای از

اعداد. همان طور که دیده میشود تغییرات زیادی در میان دامنه ها وجود دارد که نشانه دیگری از عدم فاصله زمانی به روزرسانی معمولی است همچنین در اینجا انحراف استاندارد بسیار بالا را در برخی دامنه ها خاطر نشان میکنیم به خصوص یک دامنه مفروض میتواند وبگاههای با ماهیت بسیار متفاوت، اخبار، وب نوشتهها نوشته های و یکی و مانند آن را نشان دهند در دامنه، هنر به عنوان مثال مشاهده کرده ایم که سایتهای مختلفی هستند که مقاله های کوچک در مورد نقاشی یا عکس از رده 100 ورودی در هر روز چاپ میکنند بنابراین مفهوم آیتم متفاوت از یک مقاله تخصصی است که ممکن است حاوی متنهای (مثل آیتم خبری نسبت به یک مقاله ای باشد که به طور انحصاری از تصاویر و یا فیلم ها تشکیل شده است. هر چه رده وب نوشتههای محبوب دسته بندی همگن ساختاری داشته باشد، انحراف معقول تری از فواصل به روزرسانی را دارا خواهد بود.

در اینجا، مطالعه مجموعه دادهها را نتیجه گیری میکنیم که در یک اقدام خاص منعکس کننده وضعیت موجود و تنوع فیدهای وب به منظور گردش به سمت بحث و گفت و گوراجع به روش ما برای استخراج اشیای داده ای وب است.

4 - استخراج اشیای داده ای

ما در این بخش الگوریتم را برای پیدا کردن شیء داده ای در یک صفحه وب مفروض بحث میکنیم شیئی که مربوط به یک آیتم فید وب است.

4-1 جمع آوری اطلاعات معنایی

نخست از آیتمهای فیدی نوعی بافت معنایی ایجاد کردیم که جهت استخراج شیء متناظر عمل میکند آیتمهای فید سه مؤلفه اجباری دارند پیوند عنوان و شرح.

پیوند. پیوند به ما یو آر. ال. صفحه وبی را میدهد که در آن شیء داده ساکن است.

عنوان: آیتم فید باید دارای عنوان باشد، که معمولاً یک متن کوتاه به منظور توصیف محتوای یک مقاله است.

شرح بیشتر اوقات شرح شامل کل محتوای مقاله نمی شود بلکه تنها چند خط اول آن که در اچ. تی. ام. ال برای اهداف ارائه کدگذاری شده همراه با پیوندی در انتهای مانند ادامه مطلب...» صفحه اصلی وب در موارد دیگر شرح تنها شامل یک جمله می شود که به جای در برگرفتن چند خط اول آن، مقاله را خلاصه می. کند دقت هر چه باشد میتوانیم برخی معناهای قابل اعتماد در مورد مقاله را با بهره گیری از این شرح استخراج کنیم.

با بازیابی تمامی مطالب، متنی از عنوان و شرح شروع میکنیم کد اچ تی ام ال شرح علامت گذاری و فقط متن نگه داشته میشود توالی های منتج لغات به دو نوع از نهادهای معنایی تبدیل میشوند مفاهیم و ان-گرما(1).

برای به دست آوردن مفاهیم کلمات را لب خوانی (tokenize) و ریشه آن را قطع کرده، و براساس فراوانی lexemes مرتب کرده (به عنوان اندازه اهمیت در نظر گرفته شدند) و فقط آنها برجسته ها را نگه داشتیم تا، حدی مفهوم شبیه یک برچسب است چون اصطلاحی است که یک شیء داده ای را توصیف کند. در واقع کلماتی را که از مفاهیم می آیند میتوان به عنوان کلید واژه های جست و جو در یک صفحه، وب به منظور شناسایی منطقه شیء داده ای به کار برد اما میتوان توجه را به مناطق دیگری نیز معطوف داشت که سرشار از مفاهیم هستند مانند مناطقی که حاوی آرا یا دسته بندیهایی برای یک مقاله وب مفروض هستند.

به همین دلیل تمرکزمان را به آن - گرم بر می گردانیم آن گرم در زمینه ما معرف دنباله عناوین n است که از عنوان و، شرح همانگونه که ظاهر میشوند گرفته میشوند گزینه برای n ، مصالحه ای بین مثبت های کاذب و منفیهای کاذب در فرآیند استخراج بوده و بیشتر مورد بحث قرار خواهد گرفت.

2-4. استخراج

عکس

برای به دست آوردن مفاهیم، کلمات را لب‌خوانی (tokenize) و ریشه آن را قطع کرده، و براساس فراوانی lexemes مرتب کرده (به عنوان اندازه اهمیت در نظر گرفته شدند) و فقط آنها برجسته‌ها را نگه داشتیم. تا حدی، مفهوم شبیه یک برچسب است، چون اصطلاحی است که یک شیء داده‌ای را توصیف می‌کند. در واقع، کلماتی را که از مفاهیم می‌آیند می‌توان به عنوان کلید واژه‌های جست‌وجو در یک صفحه وب، به منظور شناسایی منطقه شیء داده‌ای به کار برد، اما می‌توان توجه را به مناطق دیگری نیز معطوف داشت که سرشار از مفاهیم هستند، مانند مناطقی که حاوی آرا یا دسته‌بندی‌هایی برای یک مقاله وب مفروض هستند.

به همیسن دلیل، تمرکزمان را به ان - گرم، برمی گردانیم. ان - گرم در زمینه ما معرف دنباله عناوین n است، که از عنوان و شرح، همانگونه که ظاهر می‌شوند، گرفته می‌شوند. گزینه برای n ، مصالحه‌ای بین مثبت‌های کاذب و منفی‌های کاذب در فرآیند استخراج بوده و بیشتر مورد بحث قرار خواهد گرفت.

جدول ۲. آمارهای فید به‌ازای هر دامنه

Domain	Number of feeds	Average mean update interval	Pooled standard derivation of update interval
Art	87	12 days, 14 hours, 12 min	82 days, 6 hours, 32 min
Biology	80	7 days, 13 min	8 days, 17 hours, 43 min
Blogs	29	15 hours, 35 min	8 hours, 39 min
Environment	7	19 hours, 49 min	4 days, 15 hours, 18 min
Medicine	8	3 days, 19 hours, 16 min	1 day, 22 hours, 43 min
Other	13	4 days, 16 hours, 48 min	4 days, 19 hours, 46 min
Science	112	22 days, 12 hours, 45 min	14 days, 21 hours, 35 min
Universe	4	4 hours, 44 min	7 hours, 5 min
Total	340	12 days, 15 hours, 17 min	37 days, 16 hours, 49 min

۴-۲. استخراج

در اینجا نوعی الگوریتم از پایین به بالا را نشان می‌دهیم، که با توجه به فید، آیت‌ها را شناسایی کرده و برای هر آیت فید مؤلفه لفاف بسته‌بندی شیء داده‌ای را با تطبیق دادن آن - گرم در برابر محتوای متنی گره‌های برگی استخراج شده از صفحه وب اچ.تی.ام.ال. می‌یابد. این الگوریتم در الگوریتم ۱ خلاصه شده است.

نخست مفهوم گره مفهومی را معرفی می‌کنیم:

تعریف ۱. گره مفهومی یک گره برگی (گره بدون فرزند) است که در مفهوم متنی اش شامل حداقل یک مفهوم (یا ان - گرم) از عنوان و شرح آیت باشد.

ما تمام گره‌های برگی صفحه را استخراج و برای هر یک، تراکم معنایی ایجاد می‌کنیم.

تعریف ۲. تراکم معنایی یک گره مفهومی به عنوان تعداد مفاهیم همسان (یا ان - گرم) تقسیم بر طول محتوای متنی گره‌های مربوط تعریف می‌شود.

ما گره‌های مفهومی را طبق نزدیک‌ترین جد(نیا) طبقه‌بندی می‌کنیم که یک مؤلفه در سطح بلوک است. یک غیرمستدل مؤثر در واقع گرفتن نزدیک‌ترین جد(نیا) است که یک مؤلفه div می‌باشد.

در اینجا نوعی الگوریتم از پایین به بالا را نشان می‌دهیم که با توجه به، فید آیت‌ها را شناسایی کرده و برای هر آیت فید مؤلفه لفاف بسته‌بندی شیء داده‌ای را با تطبیق دادن آن - گرم در برابر محتوای متنی گره‌های برگی استخراج شده از صفحه وب اچ.تی.ام.ال. می‌یابد. این الگوریتم در الگوریتم ۱ خلاصه شده است.

نخست مفهوم گره مفهومی را معرفی می‌کنیم:

تعریف ۱. گره مفهومی یک گره برگی (گره بدون فرزند) است که در مفهوم متنی اش شامل حداقل

یک مفهوم (یا آن - گرم) از عنوان و شرح آیتم باشد.

ما تمام گره های برگه را صفحه را استخراج و برای هر یک تراکم معنایی ایجاد میکنیم.

تعریف 2. تراکم معنایی یک گره مفهومی به عنوان تعداد مفاهیم همسان یا ان - گرم تقسیم بر طول

محتوای متنی گرههای مربوط تعریف میشود.

ما گره های مفهومی را طبق نزدیکترین جد (نیا) طبقه بندی میکنیم که یک مؤلفه در سطح بلوک است. یک غیر مستدل مؤثر در واقع گرفتن نزدیکترین جداست نیاهاست که یک مؤلفه `div` میباشد،

ص: 16

چون در آزمایشهای مان مشاهده کرده ایم که یک شیء داده تقریباً همیشه در یک مؤلفه div محدود است. پس از این تجزیه و تحلیل میتوانیم بگوییم که کدام یک گره های مفهومی هستند که همان جد را به اشتراک می گذارند فهرست اجداد به ما مناطق معنایی صفحه را میدهد که توسط گره های معنایی در سطح کد مدل شده اند.

تعریف 3. گره معنایی پایینترین جد مشترک سطح بلوک مجموعه ای از گره های مفهومی است. به منظور روشن شدن اینکه کدام یک از گره های معنایی نشان دهنده لفاف بسته بندی مقاله است، اندازه تراکم معنایی زیر را برای هر کدام محاسبه کرده و گرهی را در نظر میگیریم که بزرگترین مقدار را برای آن دارد.

تعریف 4. گره لفاف بسته بندی شیء داده ای گرهی معنایی است که حاوی بیشترین تعداد گره های متراکم مفهومی است.

در شرح قبلی اجازه استفاده از مفاهیم یا آن - گرمها برای پیدا کردن گره های مفهومی و محاسبه تراکم معنایی را داده ایم هنگام همسان کردن با شروطی که متناظر با مفاهیم است، تعداد مناطق معنایی افزایش خواهد یافت در حالی که همسان کردن با آن - گرمها به وضوح این تعداد را کاهش خواهد داد. دلیل وقوع این اتفاق این است که آن گرمها نسبت به مفاهیم نسبت به محتوای مقاله معنی دارتر هستند. اضافه بر این در بعضی موارد انتخاب آن - گرم بیش از حد محدود کننده است. این امر کاملاً به ندرت محدودکننده اتفاق میافتد، بیشتر زمانی که شرح بیشتر مقاله را با کلمات مختلف خلاصه می کند، به جای اینکه چند خط اول آن را ارائه دهد. در نتیجه به منظور تشخیص این نوع موارد که در آن - گرمها یک گزینه نیستند مفهوم ثبات گره معنایی را معرفی میکنیم.

تعریف 5. یک گره از لحاظ معنایی سازگار است اگر متن آن حاوی یک نسبت بزرگی از مفاهیم به

دست آمده از عنوان و شرح آیتیم باشد.

ما میگوییم یک نسبت بزرگی (در عمل $0/5$) از مفاهیم زیرا لازم نیست حضور همه را به منظور اثبات یک گره لفاف بسته بندی بررسی کرد. از سوی دیگر، زمانی که نامزد گره جد شامل نیمی از مفاهیم موجود نباشد ممکن است گمان بریم که آن لفاف بسته بندی مقاله نیست. اگر این اتفاق بیفتد میتوان نتیجه گرفت که آن گرمها به علت نقص مفروضات کسب معنایی مؤثر نبودند و در نتیجه مقدار π در آن گرمها را کاهش داده و روش گره لفاف بسته بندی را تکرار میکنیم در آزمایشهای انجام شده با $3 - \pi$ (هدف به دست آوردن تعداد مناسب آن - گرمهای قابل توجه است) شروع کردیم که بهترین نتیجه را داد و در صورت شکست، به طور مستقیم سعی در همسان کردن با مفاهیم داشتیم به طور کلی اشیای داده ای مقاله های خبری، پستهای وب نوشت دارای نظرات مرتبط هستند. ما میخواهیم بین نظرات یک مقاله و خود مقاله به دلایل زیر تمایز روشنی قائل شویم:

1- از نظر مفهومی اطلاعات مورد نظر در مقاله همان مواردی نیست که در موردش نظر داده شده است.

2- خزش مقاله باید از خزش آرا از هم تفکیک شود هر زمان که یک نظر اضافه شد، مقاله را باید

تغییر یافته در نظر گرفت یک خزش جدید لازم است و شیء حاصله میتواند در مقایسه با نسخه های قبلی آن خیلی اضافی باشد؛ در عوض این مقاله باید به آرای خود اشاره داشته، جداگانه ردیابی شده، و در به روزرسانی هماهنگ باشد.

معمولاً وب مستر مراقبت افزودن یو آر ال. در فیدی است که می تواند برای پیگیری آرا مورد استفاده قرار گیرد، اما وقتی که این مورد نیست میتوانیم منطقه آرا را با استفاده از غیر مستدلهای در الگوریتم شناسایی کنیم.

عکس

تغییر یافته در نظر گرفت، یک خزش جدید لازم است و شیء حاصله می تواند در مقایسه با نسخه های قبلی آن خیلی اضافی باشد؛ در عوض، این مقاله باید به آرای خود اشاره داشته، جداگانه ردیابی شده، و در به روزرسانی هماهنگ باشد.

۳- معمولاً، وب مستر مراقبت افزودن یو آر ال. در فیدی است که می تواند برای پیگیری آرا مورد استفاده قرار گیرد، اما وقتی که این مورد نیست، می توانیم منطقه آرا را با استفاده از غیرمستدلهای در الگوریتم شناسایی کنیم.

```

import re
import sys
import urllib
import urllib2
import urlparse
import hashlib
import time

def get_urls(urls):
    """
    Get all the URLs from a list of URLs.
    """
    new_urls = []
    for url in urls:
        new_urls.append(url)
    return new_urls

def get_content(url):
    """
    Get the content of a URL.
    """
    try:
        response = urllib2.urlopen(url)
        content = response.read()
        return content
    except:
        return None

def get_md5(content):
    """
    Get the MD5 hash of the content.
    """
    return hashlib.md5(content).hexdigest()

def main():
    """
    Main function.
    """
    urls = sys.argv[1:]
    new_urls = get_urls(urls)
    for url in new_urls:
        content = get_content(url)
        md5 = get_md5(content)
        # Do something with the MD5 hash

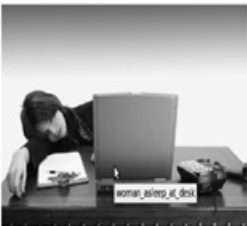
```

الگوریتم ۱۰ استخراج اشیا داده ای

study on how to study

one of the most delightful aspects of being a scientist is that you're always learning. Your colleagues teach you things. Your students teach you things. Journal articles teach you things. You sit quietly at your desk and figure things out. You're perennially a student. But how to be a better student?

On the morning of the New York Times, I read an article on "study habits". It was against the conventional wisdom (find a clean, neutral space, I hear down on a single topic), I was in favor of what might be called intellectual cross-training: remaining study environments, changing content, spacing study sessions, self-testing. The basic philosophy seems to be specialized:



شکل ۳. یک نمونه مقاله وبی و آیتم داده ای مرتبط با فید

-/http://feedproxy.google.com/r/Cosmic Variance Blog/3

/uatEVOIO0g

/http://blogs.discovermagazine.com

cosmicvariance/2010/09/07/a-study-on-how-to-study/comments

Wed, 08 Sep 2010 03:16:54 +0000

daniel

?/http://blogs.discovermagazine.com/cosmicvariance

p=5353

scientist is that you're always learning. Your colleagues teach you

things. Your students teach you things. Journal articles teach you things

You sit quietly at your desk and figure things out. You're perennially

a student. But how to be a better student? This morning the New

<[[[...]] York

One of the most delightful aspects

<[[<...of being a scientist is that you're always learning

/http://blogs.discovermagazine.com/cosmicvariance

/a-study-on-how-to-study/feed/2010/09/07

6

/http://blogs.discovermagazine.com/cosmicvariance

/a-study-on-how-to-study/2010/09/07

>>

به منظور شرح و توضیح بیشتر ما مثالی از یک صفحه وب را که قسمتی از آن در شکل 3 ارائه شده است در نظر خواهیم گرفت صفحه وب با استفاده از تمایز دهنده اچ تی ام ال پاک و خوب فرمت خواهد شد. این مرحله به منظور انتخاب معقول گرههای برگه درخت dom در صفحه لازم است (و برای

تجزیه و تحلیل تنها آنهایی نگهداری میشود که حداقل یک واژه معنایی را در بردارند).

مشاهده میکنیم که عنوان این مقاله در آیتم، فید و در دو خط اول متن وجود دارد که در شرح فید کد گذاری شده است. برچسب زمانی نیز وجود دارد که مربوط به تاریخ انتشار است. تگ مؤلفه ای نیست که به طور متداول ظاهر شود بنابراین، ما در الگوریتم خود در مورد مفید بودنش فرض ایجاد نمیکنیم.

عکس

۲۰ مدیریت منابع اطلاعاتی وب

تجزیه و تحلیل تنها آنهایی نگهداری می شود که حداقل یک واژه معنایی را در بردارند). مشاهده می کنیم که عنوان این مقاله در آیتم فید، و در دو خط اول متن وجود دارد، که در شرح فید کد گذاری شده است. برچسب زمانی نیز وجود دارد که مربوط به تاریخ انتشار است. تگ `<content:encoded>` مؤلفه ای نیست که به طور متداول ظاهر شود، بنابراین، ما در الگوریتم خود در مورد مفید بودنش فرض ایجاد نمی کنیم.

در شکل ۴، مشاهده می کنیم که عنوان مقاله همچنین در سمت راست منطبقه آرا وجود دارد، یعنی جایی که برخی از تازه ترین مقاله های وبگاه نیز ارائه شده است. عنوان تنها برای یک اشاره کافی نیست؛ بلکه به شرح آیتم هم نیاز داریم. برخی نمونه های تصادفی آن - گرمها از شرح عبارت اند از «جنبه های توضیحی»^۱ ($n=2$)، «دانشمندان بودن»^۲ ($n=3$)، «این صبح، جدید»^۳ ($n=4$)، به طور کلی آن - گرمها فرصت/احتمال کمتری برای خوشه بندی زیاد در سایر مناطق صفحه های وب دارند. هرچه توالی بزرگ تر باشد، به همان اندازه سرعت شناسایی مقاله بیشتر خواهد بود.



شکل ۴. توضیحی برای این واقعیت که عنوان به تنهایی برای شناسایی منطبقه مقاله وب کافی نیست.

1. Delightful Aspects
2. Being a Scientist
3. This morning the New

در شکل 4، مشاهده میکنیم که عنوان مقاله همچنین در سمت راست منطقه آرا وجود دارد، یعنی جایی که برخی از تازه ترین مقاله های وبگاه نیز ارائه شده است عنوان تنها برای یک اشاره کافی نیست بلکه به شرح آیتم هم نیاز داریم برخی نمونه های تصادفی آن - گرمها از شرح عبارت اند از «جنبه های توضیحی (1)» (n-2)، «دانشمند بودن (2)» (n-3)، «این صبح جدید (3)» (n-4). به طور کلی آن - گرمها / فرصت / احتمال کمتری برای خوشه بندی زیاد در سایر مناطق صفحه های وب دارند هر چه توالی بزرگتر باشد به همان اندازه سرعت شناسایی مقاله بیشتر خواهد بود.

شکل 4. توضیحی برای این واقعیت که عنوان به تنهایی برای شناسایی منطقه مقاله وب کافی نیست.

ص: 20

Delightful Aspects -1

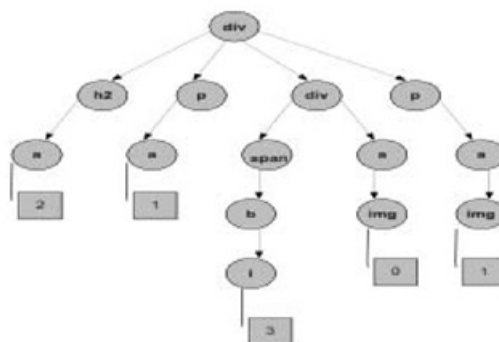
Being a Scientist -2

This morning the New -3

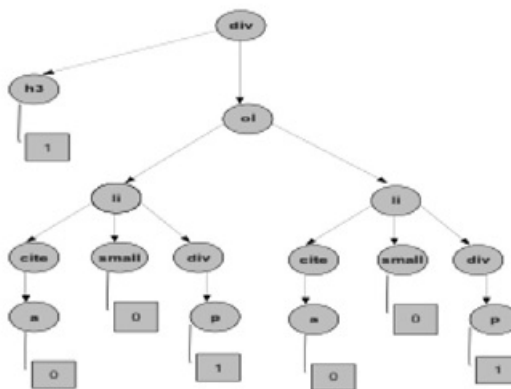
شکل 5. زیر شاخه DOM ساختگی (به منظور درک و توضیح بیشتر) متناظر با مقاله وب. مقدار الصاق شده به گره برگه معرف تعداد مفاهیم در برگرفته در متن اصلی است.

عکس

آرشیو اشیای داده ای با استفاده از فیدهای وب ۲۱



شکل 5. زیر شاخه DOM ساختگی (به منظور درک و توضیح بیشتر) متناظر با مقاله وب. مقدار الصاق شده به گره برگه معرف تعداد مفاهیم در برگرفته در متن اصلی است.



شکل 6: زیر شاخه ساده شده DOM متناظر با منطقه نظریه‌ها. متوجه شدیم که این منطقه به عنوان یک فهرست (مرتب شده با آن) ساماندهی شده و همچنین مفاهیم (یا آن - گرمها) مقاله را، در نسبت کوچک تر نشان می‌دهد.

شکل 6: زیر شاخه ساده شده DOM متناظر با منطقه نظریه‌ها متوجه شدیم که این منطقه به عنوان یک فهرست (مرتب شده با آن) ساماندهی شده و همچنین مفاهیم (یا آن - گرمها) مقاله را، در نسبت کوچک تر نشان می‌دهد.

برعکس وقتی شرح مقاله خیلی کوتاه باشد و یا فرمول بندی ایده وجود داشته باشد، استفاده از مفاهیم گزینه بهتری خواهد بود در این مثال، خاص الگوریتم زوج مقاله را فقط با این صبح جدید شناسایی می کند این اتفاق به دلیل اینکه این صبح جدید یک آن - گرم منحصر به فرد در صفحه وب است، رخ می دهد.

استفاده تنها از جنبه های توضیحی یا دانشمند بودن مقاله را شناسایی نخواهد کرد، چون آن گرمها

نیز در اولین نظر ظاهر خواهند شد (شکل 6، `(div/ol/li/div/p)`).

بنابراین نمیتوانیم فرض کنیم که توالی ها منحصر به فرد هستند و یا اینکه میتوانیم آنها را منحصر

به فرد کنیم. پس ما نمونه های مختلفی از عنوان و شرح را در نظر میگیریم.

برای توضیح اصل الگوریتم روند تطبیق آن - گرمهای دو گره معنایی توضیح داده شده در شکل 6 و 6 را در نظر میگیریم در عمل برای این صفحه وب الگوریتم سه گره معنایی ممکن را که در واقع ساختار پیچیده تری دارند بر می گرداند. همانطور که میبینیم گره های برگی در اولین زیر شاخه (شکل dom) (ه به لحاظ مفهومی غنی تر از دومی هستند (شکل 6))، بنابراین اولین گره را انتخاب میکنیم چون با توجه به اندازه تراکم معنایمان معنی دارترین است.

5. آزمایشها

به منظور اثبات روایی رویکرد مان برای استخراج اشیای داده ای به طور کامل سیستم را پیاده سازی کرده به منظور ارزیابی دقت، آن آزمایشهایی را انجام داده ایم.

آزمایشهایی با استفاده از فیدهای جمع آوری شده بر اساس پاسخهای موتور جست و جوی آر.اس. اس 4 همان مجموعه داده های ذکر شده در بخش 3 انجام شده است به یاد بسپارید که مجموعه داده ها، از لحاظ ساختار و نوع اشیای داده های خیلی متنوع بود برای هر فید ما کانال ساختار داده را بازیابی کرده و روش استخراج را برای تمام اجزای آیتم آن اعمال کرده ایم.

به عنوان اولین، آزمون تلاش کرده ایم تا منطقه صفحه وب مرتبط با عنوان آیتم را برگردانیم. با این حال، این روش به چند دلیل نتایج ضعیفی ارائه میکند تطبیق عنوان ممکن است به خاطر ویژگیهای کدگذاری به طور کامل امکان پذیر نباشد یا اینکه ممکن است در چندین محل مختلف در صفحه ظاهر شود (شکل 4) علاوه بر این با توجه به محل مؤلفه عنوان و سطوح انباشتی آن در بلوکهای کد اچ. تی. ام. ال.، شناسایی محدوده کل شیء داده ای کار آسانی نیست.

حال ما عملکرد الگوریتم را مقایسه میکنیم یعنی اینکه دقت اشیای داده ای استخراج شده با [12] Boilerpipe که پیشرفته ترین روش برای شناسایی محتوای اصلی یک صفحه در غیاب اطلاعات معنایی اضافی است. ما تأکید میکنیم که هر چه ما اطلاعات بیشتری نسبت به آنچه Boilerpipe بدان دسترسی دارد استفاده کنیم به دست آوردن دقت بیشتر از علاقه به این روش که کلی تر است - نمی کاهد.

مشاهده کردیم که نتیجه کارمان اغلب دقیقتر است چون تراکم متن در صفحه وب را در نظر نمی گیریم بلکه انسجام معنایی آن مطابق با آیتم فید را در نظر گرفتیم. مواردی وجود دارد که در آن گره

ممکن است حاوی مقدار زیادی متن باشد در حالی که ممکن است با توجه به اندازه گیری تراکم معنایی ما فاقد ارزش قضاوت شود.

علاوه بر این توجه داشته باشید که زمانی که صفحه وب حاوی مقاله های مختلف متوالی باشد، روش ما بین آنها تمایز قایل خواهد شد و مقاله خاص متناظر با یک آیتم را شناسایی خواهد کرد. در مقابل، Boilerpipe محتوای متنی تمام مقاله ها و یا تنها متراکم ترین نوع را بسته به مورد در بر خواهد گرفت.

جدول 3. نتایج آزمایشها

عکس

ممکن است حاوی مقدار زیادی متن باشد در حالی که ممکن است با توجه به اندازه گیری تراکم معنایی ما فاقد ارزش قضاوت شود.

علاوه بر این، توجه داشته باشید که زمانی که صفحه وب حاوی مقاله‌های مختلف متوالی باشد، روش ما بین آنها تمایز قایل خواهد شد و مقاله خاص متناظر با یک آیتم را شناسایی خواهد کرد. در مقابل، Boilerpipe محتوای متنی تمام مقاله‌ها و یا تنها متراکم ترین نوع را، بسته به مورد در بر خواهد گرفت.

روش	استخراج های درست	دقت (درصد)
Our technique	1038/1314	79.0%
Boilerpipe	821/1314	62.5%

جدول ۳. نتایج آزمایش‌ها

آزمایش‌ها برای ۶۰ سایت انتخابی تصادفی از مجموعه داده در دامنه هنر (اولین مورد خزش شده)، متناظر با تمام ۱۳۱۴ آیتم های فید انجام شد. ما به‌طور دستی مقاله وب را، نتیجه الگوریتم استخراج ما و نوع Boilerpipe [۱۲] بررسی کردیم. ما فقط نتایج متنی استخراجی را مقایسه کردیم چون خروجی رایگان قابل استفاده پیاده سازی Boilerpipe بود. روش ما در واقع استخراج محتوای کل منطقه شناسایی شده شامل پیوندها و تصاویر است.

نتایج عددی در جدول ۳ داده شده است. ما یک شیء را همانطور که به‌درستی استخراج شده زمانی که متن دقیقاً استاندارد طلایی است، در نظر می‌گیریم که با توضیح دستی صفحه وب به دست آمده است. انطباق‌های جزئی نادیده گرفته شدند. دقت این الگوریتم (در حدود ۷۹ درصد) در مقایسه با Boilerpipe (در حدود ۶۲ درصد) رضایت‌بخش است. در نهایت، توجه داشته باشید که زمانی که روش ما با شکست مواجه شد، الگوریتم در منطقه غنی‌تر مفهومی را شناسایی می‌کند که یکی از اشیای داده‌ای است، که هنوز هم مربوط به مقاله است، اگرچه توسط آن شناسایی نشده است.

۶. بحث در مورد کاربردها

ما این مقاله را با بحث در مورد تعداد کاربردها به پایان می‌رسانیم که می‌توان در فرآیند آرشیو وب تلفیق کرد، و نیز آنی که روش استخراج اشیای داده‌ای را استفاده می‌کند که ما پیشنهاد می‌کنیم. ماندگاری آرشیوهای وب. زمان عامل خیلی تأثیرگذاری برای تفسیر محتوای خزش شده است. درحالی‌که ممکن است داده‌ها دست نخورده باقی بمانند، روشی که ما آن را درک و ارائه می‌کنیم متفاوت است، که دلیل عمده آن واقعیت خود زبان، فرهنگ، و وسایل فناوری تکامل بیان است. یکی از جدی‌ترین مشکلات برشمرده شده در آرشیو وب زمانی است که فرمت داده‌های خزش شده منسوخ و یا به‌طور کلی استفاده نشود. راه‌حل‌های ارائه شده توسط نویسندگان [۹]، [۲۷] و [۲۶] شبیه‌ساز نرم‌افزاری یا سخت‌افزاری، انتقال محتوا، یا شامل یک پروکسی است که قابلیت‌های ترجمه فرمت را ترکیب خواهد

آزمایش‌ها برای 60 سایت انتخابی تصادفی از مجموعه داده در دامنه هنر (اولین مورد خزش شده)، متناظر با تمام 1314 آیتم‌های فید انجام شد ما به‌طور دستی مقاله وب را نتیجه الگوریتم استخراج ما و نوع [12] Boilerie بررسی کردیم ما فقط نتایج متنی استخراجی را مقایسه کردیم چون خروجی رایگان قابل استفاده پیاده سازی Boilerpipe بود روش ما در واقع استخراج محتوای کل منطقه شناسایی شده شامل پیوندها و تصاویر است.

نتایج عددی در جدول 3 داده شده است ما یک شیء را همانطور که به‌درستی استخراج شده زمانی که متن دقیقاً استاندارد طلایی است در نظر می‌گیریم که با توضیح دستی صفحه وب به دست آمده است. انطباق‌های جزئی نادیده گرفته شدند. دقت این الگوریتم در حدود 79

درصد در مقایسه با Boilerpipe در حدود 62 درصد رضایت بخش است در نهایت توجه داشته باشید که زمانی که روش ما با شکست مواجه شد الگوریتم در منطقه غنی تر مفهومی را شناسایی میکند که یکی از اشیای داده ای است که هنوز هم مربوط به مقاله است اگرچه توسط آن شناسایی نشده است.

6. بحث در مورد کاربردها

ما این مقاله را با بحث در مورد تعداد کاربردها به پایان میرسانیم که میتوان در فرآیند آرشیو وب تلفیق کرد، و نیز آنی که روش استخراج اشیاء داده ای را استفاده میکند که ما پیشنهاد میکنیم. ماندگاری آرشیوهای وب: زمان عامل خیلی تأثیرگذاری برای تفسیر محتوای خزش شده است. در حالی که ممکن است دادهها دست نخورده باقی بمانند روشی که ما آن را درک و ارائه میکنیم متفاوت است که دلیل عمده آن واقعیت خود، زبان فرهنگ و وسایل فناوری تکامل بیان است یکی از جدی ترین مشکلات بر شمرده شده در آرشیو وب زمانی است که فرمت دادههای خزش شده منسوخ و یا به طور کلی استفاده نشود راه حلهای ارائه شده توسط نویسندگان [9] [27] و [26] شبیه ساز نرم افزاری یا سخت افزاری، انتقال محتوا یا شامل یک پروکسی است که قابلیتهای ترجمه فرمت را ترکیب خواهد

کرد و این کار را به صورت پویا براساس درخواست کاربر و یا زمانی که نیاز تشخیص داده شود انجام خواهد داد. در حالی که این کارها در حال تلاش برای مبارزه با تکامل فناوری هستند این امکان را به وجود می آورد که روش مخالف را در بر بگیرد انطباق دادهها با فناوریهای موجود به منظور انجام این کار میتوان محفظه سازی اطلاعات مرتبط صفحه وب خزش شده را و ذخیره سازی شیء حاصل مستقل از فرمت کدگذاری اصلی تصور کرد. در این صورت این واقعیت که فناوری تکامل می یابد دیگر چیزی برای مقاومت در مقابل آن نخواهد بود به جای آن امکان ارائه دادههای موجود با روشهای جدید را با تطبیق محتوای واقعی با پیش تنظیمات ارتقا خواهد داد استخراج اشیای داده ای که ما در این مقاله نشان میدهم اولین گام در جهت ذخیره سازی اطلاعات آزاد از روش خاص کدگذاری است.

در حالی که میتوان استدلال کرد که برای آرشیو وب فرم اصلی صفحه وب اهمیت دارد این امر بیشتر با نیاز قرار دادن اطلاعات واقعی در مفهوم درست آن مرتبط است توجه داشته باشید که هدف ما تغییر روشهای موجود ذخیره سازی محتوای آرشیوها نیست بلکه برای مواردی که کاربرد دارد (داده های پویا با فیدهای مرتبط به آن) میتوان برخی کاربردهای همپوشانی جالب را برای کسانی که از این مجموعه استفاده میکنند یا برای خود آرشیویست ها ایجاد کرد استخراج اشیای داده ای در مفهوم آرشیو وب خیلی با تجزیه و تحلیل معنایی محتوا در زمان و امکان ارائه خدمات ارزش افزوده مرتبط است. در واقع آنچه که میتواند بسیار مفید باشد این است که قادر به اجرای پرس وجوهای پیچیده در محتوا باشد و تعاملات و امکانات را به مصرف کنندگان هدف اطلاعاتی که به عنوان آرشیو وب ارائه می شوند اضافه نماید.

بازسازی صفحه های وب از شکل 2 بیاد میآید که صفحه های وب مرتبط با فیدهای وب میتوانند بسیار پویا باشند و در مواردی با فواصل دقیقه ای روزآمد میشوند در این شرایط، تلاش برای جذب نسخه های پی در پی صفحه وب کانال متناظر غیر معقول به نظر میرسد با این حال از آنجا که فیدهای وب مرجعی از تعداد آیتمها را نگه میدارند و خوشبختانه تعداد نسبتاً زیادی برای چنین کانال پویایی، هنوز این امکان وجود دارد که در فید وب منظم خزش و آرشیو کرد کاربرد روش استخراج دادهها می تواند بازسازی صفحه وب در یک نقطه مفروض زمانی با استفاده از آیتمهای خزش شده فیدهای وب و ارجاع به مؤلفه های تمپلت، باشد بنابراین، به نسخه ای از صفحه وب اشاره می کند که در واقع به روش کلاسیک خزش نشده است.

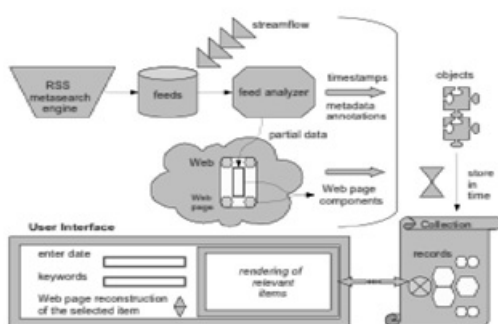
علاوه بر این، با استفاده از الگوریتم (احتمالاً برخی از فناوریهای هوشمند)، میتوانیم نه تنها گره DOM را - که شامل مقاله است - شناسایی کنیم بلکه سایر گره های معنایی صفحه را که شامل نظرها مقوله ها یا تگ هاست نیز شناسایی کنیم این مناطق میتوانند برای محتوای متنی و منابع به طور مستقل از تمپلت صفحه وب (که شامل باقی مانده های پس از استخراج است). استخراج و ذخیره شوند. رابطه بین اجزای به دست آمده را میتوان با استفاده از معناها و تجزیه و تحلیل استفاده شده برای این اجزا دوباره استنباط و ایجاد کرد ما میتوانستیم روش اطلاعات ارائه شده را با ترکیبی از اشیای (در ترکیب) دوباره اختراع کنیم تا با تنظیمهای کاربر انطباق داده شود.

شکل 7. کاربردهایی برای بازسازی صفحه های وب از اشیای داده ای

شکل 7 جریان بازسازی صفحه های وب از اشیای داده ای را همانطور که در [17] توصیف شده است نشان میدهد.

خزش و بهره برداری معنایی از آرشوها در نهایت ما میخواهیم به طور خلاصه دو برنامه کاربردی دیگر را یادآوری کنیم که قبلاً در این مقاله ذکر شده بودند. در مرحله نخست، این امکان وجود دارد که از روش استخراج شیء داده ای برای تشخیص تغییر استفاده کرد با استفاده از روش تجزیه و تحلیل زمانی روی فیدهای مشابه در همان موردی که در این مقاله ذکر شد میتوانیم راهبرد انتشار کانالی را که قسمت پویای وبگاه را ارائه میکند تعیین کنیم با این، آگاهی میتوان آهنگ خزش را با فرکانس تقریبی انطباق داد. علاوه بر این در مورد مقاله های خزش شده، قبلی می توانیم تشخیص بدهیم که نسخه جدیدی جهت خزش ظاهر شده یا نه این کار را میتوان با استفاده از الگوریتم استخراج اشیای داده ای روی مقاله های خزش شده و نسخه جدید محتمل آن صفحه وب (فعلی انجام داد با مقایسه نتایج اشیای داده ای اشاره کننده به همان یو آر. ال، میتوانیم ببینیم که آیا در متن یا منابع در دوره مداخله زمانی تغییر رخ داده است.

عکس



شکل ۷. کاربردهایی برای بازسازی صفحه‌های وب از اشیای داده‌ای

شکل ۷ جریان بازسازی صفحه‌های وب از اشیای داده‌ای را همانطور که در [۱۷] توصیف شده است نشان می‌دهد.

خزش و بهره‌برداری معنایی از آرشیوها. در نهایت، ما می‌خواهیم به‌طور خلاصه دو برنامه کاربردی دیگر را یادآوری کنیم که قبلاً در این مقاله ذکر شده بودند. در مرحله نخست، این امکان وجود دارد که از روش استخراج شیء داده‌ای برای تشخیص تغییر استفاده کرد. با استفاده از روش تجزیه و تحلیل زمانی روی فیدهای مشابه، در همان موردی که در این مقاله ذکر شد، می‌توانیم راهبرد انتشار کانالی را که قسمت پویای وبگاه را ارائه می‌کند، تعیین کنیم. با این آگاهی، می‌توان آهنگ خزش را با فرکانس تقریبی انطباق داد. علاوه بر این، در مورد مقاله‌های خزش شده قبلی، می‌توانیم تشخیص بدهیم که نسخه جدیدی جهت خزش ظاهر شده یا نه. این کار را می‌توان با استفاده از الگوریتم استخراج اشیای داده‌ای روی مقاله‌های خزش شده و نسخه جدید محتمل آن (صفحه وب فعلی) انجام داد. با مقایسه نتایج اشیای داده‌ای اشاره‌کننده به همان یو.آر.ال، می‌توانیم ببینیم که آیا در متن یا منابع در دوره مداخله زمانی تغییر رخ داده است. دوم، آرشیو وب حاوی اشیای داده‌ای (شاید علاوه بر صفحه‌های وب) را می‌توان توسط تحلیلگران به‌طور مؤثرتری نسبت به استفاده صرف آرشیو صفحه‌های وب استفاده کرد. به‌عنوان مثال، یک زبان‌شناس می‌تواند روی اصطلاحات جدیدی که در مقاله‌های روزنامه‌های ظاهر می‌شوند، بدون در نظر گرفتن اصطلاحاتی که در نظر ظاهر می‌شود، تمرکز کند. به‌طور کلی، هدف این است که معنای قابل بهره‌برداری و موقتی (در سطوح بهتر) به مجموعه‌ای از وب آرشیوهای رساتر و قابل انطباق با نیازهای کاربران اضافه شود. تشکر و قدردانی. ما از وب آرشیو اروپا (مخصوصاً جولین ماسان^۱، گابریل واسل^۲ و رادو پاپ^۳) به‌خاطر بحث و تبادل نظر در مورد عنوان این مقاله و کمک‌هایشان در به‌دست آوردن مجموعه داده‌های آزمایش‌هایمان سپاسگزار می‌کنیم.

1. Julien Masanès
2. Gabriel Vasile
3. Radu Pop

دوم، آرشیو وب حاوی اشیای داده‌ای شاید علاوه بر صفحه‌های وب را می‌توان توسط تحلیلگران به‌طور مؤثرتری نسبت به استفاده صرف آرشیو صفحه‌های وب استفاده کرد به‌عنوان مثال یک زبان‌شناس می‌تواند روی اصطلاحات جدیدی که در مقاله‌های روزنامه‌های ظاهر میشوند بدون در نظر گرفتن اصطلاحاتی که در نظر ظاهر میشود تمرکز کند به‌طور کلی هدف این است که معنای قابل بهره‌برداری و موقتی در سطوح بهتر به مجموعه‌ای از وب آرشیوهای رساتر و قابل انطباق با نیازهای کاربران اضافه شود.

تشکر و قدردانی ما از وب آرشیو اروپا (مخصوصاً جولین ماسان^(۱)، گابریل واسل^(۲) و رادو پاپ^(۳)) به‌خاطر بحث و تبادل نظر در مورد عنوان این مقاله و کمک‌هایشان در به‌دست آوردن مجموعه داده‌های آزمایش‌هایمان سپاسگزار می‌کنیم.

Julien Masanès -1

Gabriel Vasile -2

Radu Pop -3

- .1 E. Adar, J. Teevan, S. T. Dumais, and J. L. Elsas. The web changes everything .1
 ,Understanding the dynamics of web content. In Proc. WSDM, Barcelona, Spain
 .Feb.2009
- .2 A. Arasu and H. Garcia-Molina. Extracting structured data from Web pages. In Proc .2
 .SIGMOD, San Diego, USA, June 2003
- .3 J. Cho and H. Garcia-Molina. The evolution of the Web and implications for an incremental .3
 .crawler. In Proc. VLDB, Cairo, Egypt, Sept. 2000
- .4 V. Crescenzi, G. Mecca, and P. Merialdo. Roadrunner: Towards automatic data extraction .4
 .from large Websites. In Proc. VLDB, Roma, Italy, Sept. 2001
- .5 H. V. de Sompel, M. L. Nelson, C. Lagoze, and S. Warner. Resource harvesting within .5
 .the oai-pmh framework. In D-Lib Magazine, volume 10, number 12, Dec. 2004
- .6 P. Dmitriev, C. Lagoze, and B. Suchkov. As we may perceive: Inferring logical documents .6
 .from hypertext. In Proc. HT, Salzburg, Austria, Sept. 2005
- .7 .Eddie java feed parser. Website, 2010. <http://www.davidpashley.com/projects/eddie.html> .7
- .8 D. Fetterly, M. Manasse, M. Najork, and J. Wiener. A large-scale study of the evolution of .8
 .web pages. In Proc. WWW, Budapest, Hungary, May 2003
- .9 J. Hunter and S. Choudhury. Implementing preservation strategies for complex multimedia .9
 .objects. In Proc. ECDL, Trondheim, Norway, Aug. 2003
- .10 P. L. B. II, J. Johnson, U. P. Karadkar, R. Furuta, and F. Shipman. Application of kalman .10
 filters to identify unexpected change in blogs. In Proc. JCDL, Pittsburgh, USA, June

- .Internet Archive. Website, 2010. <http://web.archive.org/collections/web.html> .11
- C. Kholschutter, P. Fankhauser, and W. Nejdi. Boilerplate detection using shallow text .12
.features. In Proc. WSDM, New York, USA, Feb. 2010
- G. Knight and M. Pennock. Data without meaning: Establishing the significant properties .13
.of digital research. International Journal of Data Curation, 4(1), 2009
- ,B. Liu, R. Grossman, and Y. Zhai. Mining data records in Web pages. In Proc. KDD .14
.Washington, USA, Aug.2003
- .J. Masanès, editor. Web Archiving. Springer-Verlag, Heidelberg, Allemagne, 2006 .15
- A. Ntoulas, J. Cho, and C. Olston. What's new on the web? the evolution of the web from .16
.a search engine perspective. In Proc. WWW, New York, USA, May. 2004

- .M. Oita and P. Senellart. Archivage du contenu éphémère du Web à l'aide des flux Web .17
.In Proc. BDA, Toulouse, France, Oct. 2010. Conference without formal proceedings
(Demonstration)
- J. Pasternack and D. Roth. Extracting article text from the web with maximum .18
subsequence segmentation. In Proc. WWW, Madrid, Spain, Apr. 2009
- Z. Pehlivan, M. Ben Saad, and S. Gançarski. A novel Web archiving approach based on .19
visual pages analysis. In Proc. IAWW, Corfu, Greece, Sept. 2009
- M. Pennock and R. Davis. ArchivePress: A really simple solution to archiving blog .20
content. In Proc. iPRES, San Francisco, USA, 2009
- .Search4RSS feed search engine. Website, 2010. <http://www.search4rss.com> .21
- K. C. Sia, J. Cho, and H.-K. Cho. Efficient monitoring algorithm for fast news alerts. In .22
IEEE Trans. On Knowl. and Data Eng., vol. 19 nb.7, Piscataway, NJ, USA, 2007. IEEE
Educational Activities Department
- .K. Sigursson. Incremental crawling with Heritrix. In Proc. IAWW, Vienna, Austria, Sept .23
.2005
- .Blogger's choice awards. Website, 2010. <http://bloggerschoiceawards.com> .24
- M. Spaniol, D. Denev, A. Mazeika, and G. Weikum. Catch me if you can. Temporal .25
coherence of Web archives. In Proc. IAWW, Aarhus, Denmark, Sept. 2008
- ,S. Strodl, P. P. Beran, and A. Rauber. Migrating content in warc files. In Proc. IAWW .26
Corfu, Greece, Sept. 2009
- D. S. Swaney, F. McCown, and M. L. Nelson. Dynamic Web file format transformations .27

.with grace. In Proc. IAWW, Vienna, Austria, Sept. 2005

.H. van de Sompel, M. L. Nelson, R. Sanderson, L. L. Balakireva, S. Ainsworth, and H .28

,Shankar. Extracting structured data from Web pages. In Proc. SIGMOD to be modified

.San Diego, USA, June 2009

S. Yu, D. Cai, J.-R. Wen, and W.-Y. Ma. Improving pseudo-relevance feedback in web .29

,information retrieval using web page segmentation. In Proc. WWW.Budapest, Hungary

.May 2003

ص: 27

امروزه با افزایش اهمیت محتویات آرشیو وب یا حداقل بخشی از محتویات آن نگهداری منابع مهمی که برای مشورت و بررسی ضروری هستند به وظیفه ای حساس تبدیل شده است برای اطمینان از سازگاری آرشیو وب و نگهداری پیوسته، آن خزشگرها به طور متناوب نسخه‌های جدید اسناد را از وب بازیابی میکنند. در عین حال خزش به صفحات وب با تغییرات کم اهمیت مانند آگهیها که دائماً صفحه را روزآمد میکنند مکرراً اتفاق میافتد. بنابراین سیستمهای آرشیو وب زمان و فضا را برای شاخص گذاری و ذخیره نسخه‌های این صفحات کم اهمیت وب به هدر می دهند برای حل این مشکل و اطمینان از آرشیو مؤثر صفحات وب روش جدیدی را معرفی میکنیم که تغییرات مهم بین نسخه‌های اسناد آرشیو شده را تشخیص میدهد روش ما مفهوم نمایش دیداری صفحات را با مفهوم موجودیت در تشخیص تغییرات بین نسخه ها ترکیب می. کند روش مورد نظر شامل آرشیو ساختار دیداری صفحه وب است که به صورت فرم معنایی بلاکها ارائه میشود در این، مقاله الگوریتمی را برای تشخیص تغییرات ویژه در ساختار دیداری این صفحات ارائه میدهم همچنین روشی را برای ارزیابی اهمیت تغییرات مشخص شده معرفی میکنیم تجربیات به دست آمده از نتایج اسناد وب نشان میدهند که روش مورد نظر امیدوار کننده است.

کلیدواژه ها آرشیو وبگاه تشخیص تغییرات تجزیه و تحلیل دیداری صفحات

نوشته: میریام بن سعد (1) - استفان گانکار سکی (2) ازینب پهلوان (3)

ترجمه: مجیدرضا وحیدی (4)

1. انگیزه

با رشد سریع محتوای وب حفظ و نگهداری منابع پرکاربرد اطلاعات به وظیفه بسیار مهمی تبدیل شده است به همین دلیل انجام این وظیفه برای بسیاری از مؤسسه های ملی آرشیو در سراسر جهان اهمیت زیادی دارد. در عین حال، وب، بسیار پویا و در طول زمان در حال تکامل است (صفحات به طور دائم تغییر میکنند) در بیشتر موارد آرشیو کردن وب (5) به طور خودکار با استفاده از خزشگرهای وب اجرا می شود. خز شگرهای وب صفحات وب را که باید آرشیو شوند مشاهده و یک کپی نمونه (6) و / یا شاخص (7) از صفحات وب ایجاد میکنند برای بهروز نگهداری، آرشیو خزشگر باید به طور متناوب صفحات را بازبینی کند و آرشیو را با کپی جدید به روز رسانی کند به دلیل این که معمولاً خزشگر، منابع محدودی دارد پهنای باند، فضای ذخیره سازی و مانند آن و با توجه به حجم عظیم صفحاتی که باید آرشیو شوند،

ص: 29

Myriam Ben Saad -1

Stephane Gancarski -2

Zeyneb Pahlivan -3

4- دانشجوی کارشناسی ارشد مهندسی نرم افزار از سازمان اسناد و کتابخانه ملی ج.ا.ا.

Web crawlers -5

Snapshot -6

Index -7

ممکن است خزشگر نتواند در همه زمانها یک سایت را بازبینی کند و نسخه(1) جدید صفحه را بارگذاری کند. در حقیقت نگهداری آرشیو کل وب یا حتی قسمتی از آن که شامل همه نسخه های همه صفحه ها باشد امکان پذیر نیست بنابراین مسئله این طور بیان میشود که خزش برای بارگذاری مهمترین نسخه ها را چگونه باید بهینه کنیم تا ریزش اطلاعات مفید به حداقل برسد. البته این کار باید بدون دخالت مدیران وبگاهها انجام شود از این رو سیستم آرشیو باید رفتار سایت را تخمین بزند تا زمان و میزان تناوب(2) بازبینی صفحه را مشخص کند کارهای متعدد [3,4] بر روی تناوب تغییر متمرکز شده اند تا با کمک آنها بتوان خزشگرهای وب را بهبود بخشید. در عین حال ممکن است خزشگر با ذخیره کردن یک نسخه جدید صفحه با تغییرات کم اهمیت زمان و فضا را به هدر دهد مثالی از این مورد آگهی ها هستند که به طور دائم تغییر میکنند. بنابراین روش مؤثری مورد نیاز است که مشخص کند تغییرات بین نسخه ها دقیقاً چه موقع و چند وقت یکبار صورت می پذیرد. روشهایی که تا به حال مطرح شده اند فقط تناوب تغییرات را تخمین میزنند ولی اهمیت تغییرات را در نظر نمی گیرند اگر بتوانیم تناوب تغییرات مهم را با درستی بیشتری پیش بینی کنیم اثر بخشی سیستم آرشیو وب بهبود می یابد.

برای تخمین تناوب به روزرسانیها باید تغییرات بین نسخه های بازبایی شده اسناد مشخص شود. بسیاری از الگوریتمهای موجود [5,6] به طور ویژه برای مشخص کردن تغییرات بین اسناد نیمه ساخت یافته(3) (xml و html) طراحی شده اند در عین حال روشی وجود ندارد که تغییرات مرتبط / نامرتب را از اطلاعات پر استفاده بدون استفاده تمیز دهد کارهای قبلی [2] نشان میدهد که میتوان صفحه را به بخشهای(4) متعدد یا بلاک ها(5) تقسیم(6) کرد. معمولاً بلاکهای موجود در یک صفحه اهمیت متفاوتی دارند در حقیقت در صفحات وب نواحی مختلف بر حسب موقعیت، اندازه و محتوا دارای وزن اهمیت متفاوتی هستند. معمولاً اطلاعات مهمتر در مرکز صفحه قرار دارد. آگهی در بالای صفحه یا سمت چپ و حق نشر(7) در قسمت پانویس قرار دارد با تقسیم بندی صفحه باید به هر بلاک، یک میزان اهمیت نسبی داده شود. این کار برای نمونه با استفاده از الگوریتم [7] یا در کل با روش یادگیری ماشین به طور خود کار انجام میشود سپس میتوانیم اهمیت تغییرات را بین دو نسخه صفحه محاسبه کنیم این محاسبه بر مبنای دو مشخصه انجام میشود: 1) اهمیت نسبی بلاکها و 2) اهمیت نسبی عملیات (درج، حذف، به روزرسانی) که در این بلاکها انجام شده است و با مقایسه دو نسخه مشخص می شود.

در این تحقیق روشی را برای مشخص کردن تغییرات مهم بین نسخه ها برای آرشیو مؤثر وب پیشنهاد میکنیم این، روش روی یک انباره(8) برای مؤسسه ملی سمعی و بصری فرانسه (INA) به کار رفته است.

ص: 30

Version -1

Frequency -2

Semi-structured documents -3

Segments -4

Block -5

Partition -6

Copyright -7

Repository -8

یکی از وظایف INA ایجاد ذخایر قانونی(1) است که صفحات وب رادیو و تلویزیون فرانسه و صفحات مرتبط آن را نگهداری می‌کند. یکی از نیازهای این پروژه نگهداری جنبه دیداری صفحات است. بنابراین ایده ما به کارگیری تحلیل دیداری صفحه برای نسبت دادن اهمیت به بخشهای مختلف صفحه وب بر مبنای موقعیت نسبی آنهاست. مفاهیم تحلیل دیداری صفحه و اهمیت بخشهای صفحات وب جدید نیستند ولی تا آنجا که میدانیم برای آرشیو وب به صورت ترکیبی به کار نرفته اند.

ادامه مطالب این مقاله به این صورت است: بخش 2 معماری سیستم آرشیو وب را ارائه می‌دهد. بخش 3 مدل توسعه یافته تقسیم بندی دیداری صفحه برای صفحات وب HTML را توضیح می‌دهد. در بخش 4 الگوریتم تشخیص تغییر کافی برای محاسبه اختلاف بین دو نسخه صفحه بازسازی شده دیداری را ارائه می‌دهیم در بخش 5 روش ارزیابی اهمیت بلاکها تغییرات را ارائه می‌دهیم. در بخش 6 راهبرد به کار رفته برای زمانبندی مؤثر خزشگرهای وب و در بخش 7 نگاهی به همه مراحل روش مورد نظر داریم. در بخش 8 نیز کارهای آینده را مطرح می‌کنیم.

عکس

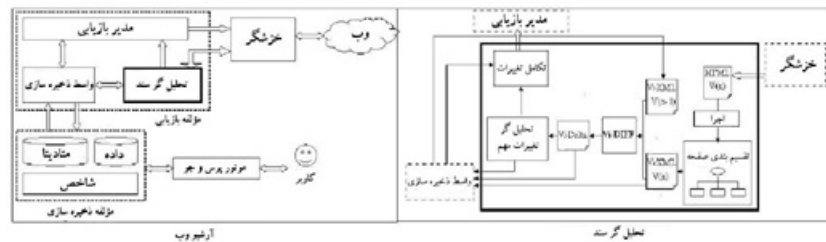
یکی از وظایف INA ایجاد ذخایر قانونی^۱ است که صفحات وب رادیو و تلویزیون فرانسه و صفحات مرتبط آن را نگهداری می‌کند. یکی از نیازهای این پروژه، نگهداری جنبه دیداری صفحات است. بنابراین ایده ما به‌کارگیری تحلیل دیداری صفحه برای نسبت دادن اهمیت به بخش‌های مختلف صفحه وب بر مبنای موقعیت نسبی آنهاست. مفاهیم تحلیل دیداری صفحه و اهمیت بخش‌های صفحات وب جدید نیستند، ولی تا آنجا که می‌دانیم برای آرشیو وب به‌صورت ترکیبی به‌کار نرفته‌اند.

ادامه مطالب این مقاله به این صورت است: بخش ۲ معماری سیستم آرشیو وب را ارائه می‌دهد. بخش ۳ مدل توسعه یافته تقسیم‌بندی دیداری صفحه برای صفحات وب HTML را توضیح می‌دهد. در بخش ۴ الگوریتم تشخیص تغییر کافی برای محاسبه اختلاف بین دو نسخه صفحه بازسازی شده دیداری را ارائه می‌دهیم. در بخش ۵ روش ارزیابی اهمیت بلاک‌ها/تغییرات را ارائه می‌دهیم. در بخش ۶ راهبرد به‌کار رفته برای زمان‌بندی مؤثر خزشگرهای وب و در بخش ۷ نگاهی به همه مراحل روش مورد نظر داریم. در بخش ۸ نیز کارهای آینده را مطرح می‌کنیم.

۲. معماری آرشیو وب

آرشیو وب شامل ۴ مؤلفه اصلی است: خزشگرهای وب، مؤلفه بازیابی^۲، مؤلفه ذخیره سازی^۳، و موتور پرس‌وجو^۴. شکل ۱ شمای کلی سیستم را ارائه می‌دهد.

خزشگر وب: خزشگرهای وب، با بارگذاری مکرر نسخه‌های جدید صفحات، مرور و بررسی وب را پوشش می‌دهند.



تصویر ۱. دید کلی از آرشیو وب

مؤلفه بازیابی: امکان نگهداری آرشیو را به‌صورت به‌روز فراهم می‌کند. که شامل ۳ ماژول است: (۱) مدیر بازیابی^۵، امکان بهینه‌سازی منابع تخصیص یافته را فراهم می‌کند تا اطلاعات کمتری از دست رود. همچنین، اسنادی را که باید به سرعت به روز رسانی شوند انتخاب می‌کند تا آرشیو در حد ممکن

- 1 - Legal deposit
- 2 - Freshness component
- 3 - Storage component
- 4 - Query engine
- 5 - Freshness Manager

۲. معماری آرشیو وب

آرشیو وب شامل ۴ مؤلفه اصلی است: خزشگرهای وب مؤلفه بازیابی^(۲)، مؤلفه ذخیره سازی^(۳)، و موتور پرس و جو^(۴). شکل ۱ شمای کلی سیستم را ارائه می‌دهد.

خزشگر: وب خزشگرهای وب با بارگذاری مکرر نسخه‌های جدید صفحات، مرور و بررسی وب

را پوشش میدهند .

تصویر 1 دید کلی از آرشیو وب

مؤلفه بازیابی امکان نگهداری آرشیو را به صورت به روز فراهم می‌کند که شامل 3 ماژول است (1)مدیر بازیابی(5)، امکان بهینه سازی منابع تخصیص یافته را فراهم میکند تا اطلاعات کمتری از دست رود. همچنین اسنادی را که باید به سرعت به روز رسانی شوند انتخاب میکند تا آرشیو در حد ممکن

ص: 31

Legal deposit -1

Freshness component -2

Storage component -3

Query engine -4

Freshness Manager -5

به روز نگهداری شود. (2) تحلیلگر سند (1) امکان تشخیص و تحلیل نسخه‌های صفحه وب بازیابی شده را فراهم می‌کند با توجه به اینکه تحلیلگر سند هسته اصلی روش ما را تشکیل می‌دهد، در اینجا جزئیات بیشتری در مورد آن ارائه می‌دهیم تحلیلگر سند شامل زیر ماژول‌های متعددی مطابق با مراحل مختلف تحلیل صفحه است که در شکل 1 نشان داده شده است. تحلیلگر سند برای به دست آوردن نسخه صفحه HTML خاصی که باید آرشیو شود با خزشگر در تعامل است. سپس برای بازیابی اطلاعات دیداری صفحه توسط موتور تفسیر (2) پردازش می‌شود. منفعت اصلی تفسیر فراهم کردن یک توضیح دیداری کامل و حقیقی از سند حتی با وجود اسکریپت‌های صفحه مانند جاوا اسکریپت است پس از آن، صفحه تفسیر شده تقسیم بندی میشود و ساختار طرح بندی دیداری (3) صفحه ساخته میشود الگوریتم [2] VIPS برای تقسیم بندی صفحه وب به بلاک‌های سلسله مراتبی معنایی (4) به کار میرود و قسمتهای کافی برای صفحه وب در نظر می‌گیرد. ما این الگوریتم را برای استخراج، پیوندها تصاویر و متنهای هر بلاک توسعه داده ایم الگوریتم توسعه یافته VIPS، یک سند Vi-XML را به عنوان خروجی تولید می‌کند که ساختار محتوای سلسله مراتبی صفحه را توضیح می‌دهد در پایان فرآیند تقسیم بندی الگوریتم تشخیص تغییر Vi-DIFF، از تغییرات ایجاد شده بین نسخه جدید Vi-XML تولید شده (Vn) و آخرین نسخه آرشیو شده (V(n-1)) توضیحی فراهم میکند تغییرات در یک فایل delta XML که Vi-Delta نام دارد ذخیره می‌شوند. این فایل اعمال اتفاق افتاده بین دو سند (درج، حذف و مانند آن) را توضیح می‌دهد. پس از آن، فایل Vi-Delta توسط زیر ماژول تحلیلگر تغییرات مهم (5) تحلیل میشود تا اهمیت تغییرات مشخص شده ارزیابی شود ماژول 3) تکامل تغییرات (6)، نتیجه این ارزیابی تغییر را برای بهینه سازی بیشتر منابع مدیریت شده توسط مؤلفه، بازیابی مورد استفاده قرار میدهد در پایان Vi-Delta و نسخه فعلی Vi-XML از طریق واسط ذخیره سازی (7) با اطلاعات اضافی مانند URL و تاریخ زمان خزش در پایگاه داده ذخیره می‌شوند واسط ذخیره سازی برای ذخیره کردن / شاخص کردن (8) نسخه های صفحه و فراداده آنها، که در طول تحلیل به دست آمده، است با مؤلفه ذخیره در تعامل است.

مؤلفه ذخیره سازی مؤلفه ذخیره شامل واحدهای ذخیره داده و متادیتاست. همچنین شامل شاخصی است که پرس و جو از آرشیو را تسهیل می‌کند.

موتور پرس و جو کاربران میتوانند از طریق موتور پرس و جو بین نسخه ها و نسخه های صفحه

آرشیو شده درخواستی حرکت (9) کنند.

ص: 32

Document Analyzer -1

rendering engine -2

visual page layout structure -3

semantic hierarchical blocks -4

Importance Changes Analyzer -5

Changes Evolution -6

Storage Interface -7

Index -8

Navigate -9

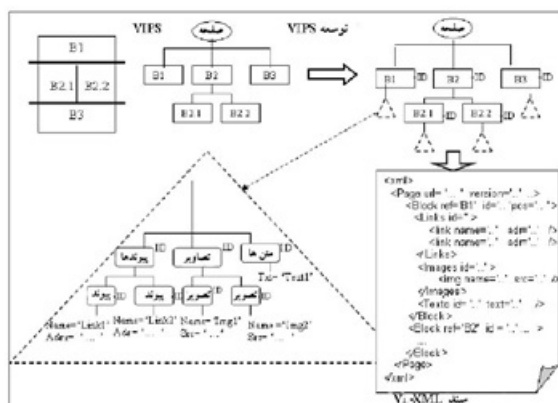
همانطور که قبلا ذکر شد ، [2] VIPS برای تقسیم بندی صفحه وب به بلاکهای معنایی تو در تو بر مبنای گره های مناسب در درخت HTML DOM در صفحه به کار میرود و جداکننده های (1) عمودی و افقی در صفحه را مشخص میکنند بر مبنای این جداکننده ها درخت معنایی صفحه و بی را که به بلاکهای متعدد تقسیم شده است تولید می. کند. ریشه درخت کل صفحه است. هر بلاک به عنوان گرهی از درخت در نظر گرفته میشود که در شکل 2 نشان داده شده است.

شکل 2. الگوریتم VIPS توسعه یافته

عکس

۳. تقسیم‌بندی دیداری صفحه

همانطور که قبلاً ذکر شد، [۲] VIPS برای تقسیم‌بندی صفحه وب به بلاک‌های معنایی تو در تو بر مبنای گره‌های مناسب در درخت HTML DOM در صفحه به کار می‌رود و جداکننده‌های عمودی و افقی در صفحه را مشخص می‌کند. بر مبنای این جداکننده‌ها، درخت معنایی صفحه وبی را که به بلاک‌های متعدد تقسیم شده است، تولید می‌کند. ریشه درخت، کل صفحه است. هر بلاک، به عنوان گرهی از درخت در نظر گرفته می‌شود که در شکل ۲ نشان داده شده است.



شکل ۲. الگوریتم VIPS توسعه یافته

برای تکمیل درخت معنایی کل صفحه، الگوریتم VIPS را با استخراج پیوندها، تصاویر، و متن برای هر بلاک توسعه دادیم. همانطور که در شکل ۳ نشان داده شده است، هر گره بلاک، دارای گره‌های فرزند اضافه‌ای است. گره‌های پیوندها، تصاویر، و متون هر کدام به ترتیب پیوندها، تصاویر، و متون هر بلاک را جمع‌آوری می‌کند. همه گره‌های صفحه، به‌طور منحصر به فرد، با یک خصوصیت ID شناسایی می‌شوند. این ID یک مقدار درهم‌سازی شده^۱ است که با استفاده از محتوای گره و محتوای گره‌های فرزند آن محاسبه می‌شود. اگر گره‌های همتا^۲ (گره‌های در موقعیت یکسان در دو نسخه متوالی) چند ID متفاوت داشته باشند، لزوماً محتویات آنها به‌روزرسانی شده است. گره‌های برگ، خصوصیات دیگری مانند نام و نشانی برای ابر پیوند دارند. الگوریتم VIPS توسعه یافته ما، یک سند Vi-XML را به‌عنوان خروجی تولید می‌کند. این سند، ساختار سلسله‌مراتبی صفحه وب را توضیح می‌دهد. که در شکل ۲ نشان داده شده است.

1 - Separators
2 - Hash value
3 - Matched nodes

برای تکمیل درخت معنایی کل صفحه الگوریتم VIPS را با استخراج پیوندها، تصاویر، و متن برای هر بلاک توسعه دادیم همانطور که در شکل ۳ نشان داده شده است هر گره، بلاک دارای گره‌های فرزند اضافه‌ای است گره‌های، پیوندها، تصاویر و متون هر کدام به ترتیب پیوندها تصاویر و متون هر بلاک را جمع‌آوری می‌کند. همه گره‌های صفحه، به‌طور منحصر به فرد، با یک خصوصیت ID شناسایی می‌شوند. این ID یک مقدار در هم‌سازی شده^(۲) است که با استفاده از محتوای گره و محتوای گره‌های فرزند آن محاسبه می‌شود. اگر گره‌های همتا^(۳) (گره‌های در موقعیت یکسان در دو نسخه متوالی) چند ID متفاوت داشته باشند، لزوماً محتویات آنها به‌روزرسانی شده است. گره‌های برگ، خصوصیات دیگری مانند نام و نشانی برای ابر پیوند دارند. الگوریتم VIPS توسعه یافته ما یک سند Vi-XML را به

عنوان خروجی تولید میکنند این سند ساختار سلسله مراتبی صفحه وب را توضیح میدهد که در شکل 2 نشان داده شده است.

ص: 33

Separators -1

Hash value -2

Matched nodes -3

الگوریتمهای تشخیص تغییر متفاوت [5, 6] برای اسناد xml منظور شده اند. با توجه به اینکه آنها الگوریتمهای عمومی هستند، به طور کامل نیازمندیهای ما را برآورده نمی کنند. ما میخواهیم بعضی ضوابط ویژه را برای مقایسه گرههای خصوصیت مانند تشخیص پیوندهای به روز رسانی شده در صورت تغییر نام یا نشانی یک پیوند که خصوصیات پیوند (هستند اضافه کنیم. همچنین می خواهیم یک متن به روز رسانی شده در دو بلاک همتا بر مبنای امتیاز شباهت فاصله متنی تعداد کلمات مشترک را مشخص کنیم با الگوریتمهای عمومی این گره ها به صورت حذف شده از نسخه قدیمی و افزوده شده به نسخه جدید در نظر گرفته می شوند. ویژگی دیگر روش ما این است که لازم است عناصر تغییر یافته درون یک بلاک و عناصر تغییر موقعیت داده از یک بلاک به بلاک دیگر را تشخیص دهیم. ولی عنصر تغییر موقعیت داده درون یک بلاک بدون اهمیت است؛ زیرا در این صورت نه اطلاعاتی به بلاک اضافه شده و نه اطلاعاتی از آن حذف شده است. همچنین میخواهیم تغییر ساختار صفحه در سطح بلاکها از نسخه ای به نسخه دیگر تشخیص داده شود بلاک حذف شده درج شده بنابراین دلایل الگوریتم تشخیص تغییرات به نام [1] Vi-DIFF را ارائه میدهیم این، الگوریتم اختلافات بین دو نسخه سند - Vi XML را محاسبه میکند و یک سند Vi-Delta تولید میکند که اختلاف (دلتهای) بین دو نسخه را ارائه میدهد Vi-DIFF مورد نظر امکان مشخص کردن دو نوع تغییر را فراهم میکند: تغییرات ساختاری و تغییرات محتوایی.

تغییرات ساختاری معمولاً ساختار سند XML را اصلاح میکنند در سطح بلاکها؛ در صورتی که تغییرات محتوایی، محتوای متنی را اصلاح میکنند در سطح پیوندها، تصاویر، و مانند آن. اگر فرض کنیم که تغییری در ساختار وجود ندارد پیچیدگی الگوریتم Vi-DIFF لگاریتمی - خطی ($n \log(n)$). است در اینجا n تعداد کل گرهها است. اگر تغییرات ساختاری وجود داشته باشد در بدترین حالت (حالتی که همه ساختار تغییر یابد) پیچیدگی الگوریتم از درجه دوم $O(n^2)$ است. لازم به ذکر است که n همیشه کوچک باقی میماند.

5. اهمیت تغییرات

بر مبنای Vi-Delta تولید شده توسط الگوریتم Vi-DIFF تابعی را در نظر میگیریم که اهمیت تغییرات مشخص شده را ارزیابی می کند. این، کار وظیفه زیر ماژول تحلیلگر تغییرات مهم (شکل 1) است این تابع، به به 3 پارامتر اصلی وابسته است:

اهمیت بلاک به روز رسانی شده معمولاً مهمترین اطلاعات در مرکز، و آگهی ها در بالای صفحه قرار دارند. سانگ و همکارانش، [7] نسبت دادن مقادیر اهمیت برای بلاکهای مختلف در صفحه وب به طور خودکار را پیشنهاد میدهند همچنین میتوانیم سایر پارامترها را برای ارزیابی اهمیت یک بلاک، با توجه به تاریخچه تغییرات این بلاک در نظر بگیریم. برای نمونه، بلاکی که با تناوب بیشتری تغییر میکند اهمیت کمتری دارد مطالعه بیشتر برای یافتن بهترین روش به منظور تخمین

اهمیت بلاکها ضروری است

اهمیت عملیات اهمیت عملیات به نوع عمل جابه جایی درج و از این قبیل) و عنصر تغییر یافته پیوند، تصویر، و مانند آن بستگی دارد برای نمونه ممکن است عملیات درج و حذف اهمیت بیشتری نسبت به جابه جایی داشته باشد، همچنین ممکن است درج تصویر از درج یک پیوند یا متن مهم تر باشد بنابراین برای انتخاب بهترین مقادیر پارامترها برای هر نوع عملیات، مطالعه روشهای یادگیری ماشین در برنامه ما قرار دارد.

میزان تغییرات در بلاک میزان عملیات تغییر درج و حذف و از این قبیل) روی داده در یک بلاک برای هر عنصر (پیوند) تصویر (و متن از V_i -Delta تولید شده استنباط میشود این میزان درصد عملیات تغییر مشخص شده برای هر بلاک تقسیم بر تعداد کل عناصر بلاک را مشخص می کند.

بر مبنای این پارامترها، تابع $E(V, V)$ را ارائه می دهیم که اهمیت تغییرات بین نسخه های V و V را که هر کدام از آنها از بلاکهای BK تشکیل شده است) تخمین میزند.

در این فرمول :

{ درج، حذف، به روزرسانی، جابه جایی } - OP -

{ پیوند، تصویر متن } - EI -

به ترتیب تعداد بلاکها تعداد نوع عملیات و تعداد نوع عناصر هستند = Na, NO, NB -

مقدار اهمیت X که میتواند یک بلاک یا یک عملیات تغییر را مشخص میکند = (x)

EI

عکس

اهمیت بلاکها ضروری است.

- اهمیت عملیات. اهمیت عملیات به نوع عمل (جابه‌جایی، درج، و از این قبیل) و عنصر تغییر یافته (پیوند، تصویر، و مانند آن) بستگی دارد. برای نمونه، ممکن است عملیات درج و حذف اهمیت بیشتری نسبت به جابه‌جایی داشته باشد. همچنین، ممکن است درج تصویر از درج یک پیوند یا متن مهم‌تر باشد. بنابراین، برای انتخاب بهترین مقادیر پارامترها برای هر نوع عملیات، مطالعه روش‌های یادگیری ماشین در برنامه ما قرار دارد.
- میزان تغییرات در بلاک. میزان عملیات تغییر (درج و حذف، و از این قبیل) روی داده در یک بلاک برای هر عنصر (پیوند، تصویر، و متن) از Vi -Delta تولید شده استنباط می‌شود. این میزان، درصد عملیات تغییر مشخص شده برای هر بلاک تقسیم بر تعداد کل عناصر بلاک را مشخص می‌کند. بر مبنای این پارامترها، تابع $E(v1, v2)$ را ارائه می‌دهیم که اهمیت تغییرات بین نسخه‌های $v1$ و $v2$ را (که هر کدام از آنها از بلاک‌های Bk_i تشکیل شده است) تخمین می‌زند.

$$E = \sum_{i=1}^{N_{Bk}} I(Bk_i) * \left[\frac{1}{N_{Op}} \sum_{j=1}^{N_{Op}} I(Op_j) * \frac{1}{N_{El}} \sum_{k=1}^{N_{El}} \frac{N(Op_j, El_k)}{N(El_k, Bk_i)} \right]$$

- $Op_j = \{ \text{درج، حذف، به‌روزرسانی، جابه‌جایی} \}$

- $El_k = \{ \text{پیوند، تصویر، متن} \}$

- N_{El}, N_{Op}, N_{Bk} = تعداد نوع عناصر هستند

$I(x)$ = مقدار اهمیت X که می‌تواند یک بلاک یا یک عملیات تغییر را مشخص می‌کند

- $N(Op_j, El_k)$ = تعداد عملیات تغییر j را که روی عنصر k می‌دهد، مشخص می‌کند

- $N(El_k, Bk_i)$ = تعداد کل عناصر k در بلاک i را مشخص می‌کند

برای نرمال سازی نتیجه تابع $E()$ ، محدودیت زیر را روی اهمیت بلاکها اضافه می‌کنیم:

$$\sum_{i=1}^{N_{Bk}} I(Bk_i) = 1; 0 \leq I(Op) \leq 1$$

تابع $E()$ با ضرب درصد تغییرات، برای هر عمل (Op_j) و بلاک Bk_i ، در اهمیت عملیات $I(Op_j)$ و بلاکها $I(Bk_i)$ محاسبه می‌شود و مقدار نرمالی بین ۰ و ۱ بر می‌گردد.

۶. زمان‌بندی خزش وب

یکی از اهداف روش ما، بارگذاری مهم‌ترین نسخه‌ها بر مبنای زمان‌بندی خزش است. مهم‌ترین وظیفه زمان‌بند عبارت است از خزش به مهم‌ترین نسخه اسناد بر مبنای تاریخچه تغییرات. زمان‌بند، فهرستی از اسناد مرتب شده توسط یک تابع ضرورت بازیابی^۱ را مدیریت می‌کند. این تابع، برای هر صفحه

1 - Freshness urgency function

تعداد عملیات تغییر Z را که روی عنصر R می‌دهد مشخص می‌کند $EI = (NOP) -$

تعداد کل عناصر در بلاک i را مشخص می‌کند $(NEBK) =$

برای نرمال سازی نتیجه تابع $E()$ ، محدودیت زیر را روی اهمیت بلاکها اضافه می‌کنیم

تابع $E()$ با ضرب درصد تغییرات برای هر عمل (Op) و بلاک BK در اهمیت عملیات (OP) و

بلاکها (BK) محاسبه میشود و مقدار نرمالی بین 0 و 1 بر می گرداند.

6. زمان بندی خزش وب

یکی از اهداف روش ما بارگذاری مهمترین نسخه ها بر مبنای زمان بندی خزش است. مهمترین وظیفه زمان بند عبارت است از خزش به مهمترین نسخه اسناد بر مبنای تاریخچه تغییرات زمان بند، فهرستی از اسناد مرتب شده توسط یک تابع ضرورت بازیابی (1) را مدیریت می. کند این تابع برای هر صفحه

ص: 35

Freshness urgency function -1

میزان ضروری بودن بازیابی آن در یک زمان داده شده را مشخص میکند این تابع اهمیت تغییرات را به حساب می آورد تخمین زده شده توسط تابع $E()$ که بین نسخه اصلی و آخرین نسخه آرشیو شده روی داده است. تابع ضرورت بازیابی برای زمانبند را به صورت زیر تعریف میکنیم:

که در آن :

اولویت صفحه $p=$

متوسط اهمیت تغییرات تخمین زده شده بین دو نسخه اسناد $AvgE =$

زمان آخرین نسخه بازیابی شده $DATE-$

زمان اولین نسخه اسناد $DATE-$

این تابع به موارد زیر بستگی دارد :

1 - اولویت صفحه،

2- اهمیت تغییرات نسخه‌های آرشیو شده قبلی و

3- زمان آخرین بازیابی صفحه

برای تولید به روزرسانی‌های ایجاد شده با اهمیت تغییرات متفاوت روی اسناد شبیه سازی شده، نوعی شبیه ساز تولید کرده ایم. سپس راهبرد زمان بندی با استفاده از تابع ضرورت توسعه داده میشود و با دو سیاست خزش موجود مقایسه می شود.

([Round Robin t, Cho [4)

نتایج اولیه امیدوار کننده هستند. در واقع راهبرد به کارگیری تابع ضرورت نسبت به سایر سیاستهای موجود، بازیابی نسخه های مهم تر را امکان پذیر می کند.

عکس

میزان ضروری بودن بازیابی آن در یک زمان داده شده را مشخص می‌کند. این تابع، اهمیت تغییرات را به حساب می‌آورد (تخمین زده شده توسط تابع $E()$) که بین نسخه اصلی و آخرین نسخه آرشیو شده روی داده است. تابع ضرورت بازیابی برای زمان‌بند را به صورت زیر تعریف می‌کنیم:

$$U(doc, date) = p * \frac{AvgE}{date_{lastUpd} - date_{OrigDoc}} * (date - date_{lastUpd})$$

که در آن :

p = اولویت صفحه

$AvgE$ = متوسط اهمیت تغییرات تخمین زده شده بین دو نسخه اسناد

$date_{lastUpd}$ = زمان آخرین نسخه بازیابی شده

$date_{OrigDoc}$ = زمان اولین نسخه اسناد

این تابع به موارد زیر بستگی دارد :

- ۱- اولویت صفحه،
 - ۲- اهمیت تغییرات نسخه های آرشیو شده قبلی، و
 - ۳- زمان آخرین بازیابی صفحه.
- برای تولید به روزرسانی های ایجاد شده (با اهمیت تغییرات متفاوت) روی اسناد شبیه سازی شده، نوعی شبیه ساز تولید کرده ایم. سپس راهبرد زمان بندی با استفاده از تابع ضرورت توسعه داده می شود و با دو سیاست خزش موجود مقایسه می شود.

(Round Robin †, Cho [‡] †)

نتایج اولیه، امیدوار کننده هستند. در واقع، راهبرد به کارگیری تابع ضرورت نسبت به سایر سیاست های موجود، بازیابی نسخه های مهم تر را امکان پذیر می کند.

۷. جمع بندی

در این قسمت، همه مراحل روش را بیان می کنیم: تقسیم بندی دیداری تشخیص تغییرات و ارزیابی اهمیت تغییرات. این موارد بر روی صفحات مختلف وب HTML رادیو و تلویزیون اعمال می شود.

• تقسیم بندی دیداری

ابتدا نشان می دهیم که صفحات وب، چگونه به بلاک های معنایی دیداری تقسیم بندی می شوند و چگونه پیوندها، تصاویر، و متن ها برای هر بلاک استخراج می شوند. یک سند Vi-XML، به عنوان خروجی تولید می شود که ساختار دیداری سلسله مراتبی صفحه را نشان می دهد. نتایج تجربی (از لحاظ زمان اجرا و اندازه خروجی) تقسیم بندی دیداری ارائه می شوند.

• تشخیص تغییرات

در این مرحله، نشان می دهیم که الگوریتم Vi-DIFF، تغییرات بین نسخه های مختلف اسناد را

7. جمع بندی

در این قسمت همه مراحل روش را بیان می کنیم تقسیم بندی دیداری تشخیص تغییرات و ارزیابی اهمیت تغییرات این موارد بر روی صفحات مختلف وب HTML رادیو و تلویزیون اعمال میشود

• تقسیم بندی دیداری

ابتدا نشان می‌دهیم که صفحات، وب چگونه به بلاک‌های معنایی دیداری تقسیم بندی می‌شوند و چگونه پیوندها، تصاویر و متن‌ها برای هر بلاک استخراج میشوند یک سند ViXML، به عنوان خروجی تولید میشود که ساختار دیداری سلسله مراتبی صفحه را نشان می‌دهد. نتایج تجربی (از لحاظ زمان اجرا و اندازه خروجی) تقسیم بندی دیداری ارائه میشوند.

● تشخیص تغییرات

در این مرحله نشان می‌دهیم که الگوریتم Vi-DIFF تغییرات بین نسخه‌های مختلف اسناد را

ص: 36

مشخص می. کند آزمایشها روی نسخه های گوناگون اسناد Vi-XML با تغییرات متفاوتی اعمال می شوند: تغییرات محتوا و تغییرات ساختاری.

Vi-Delta خروجی عملیات تغییر روی داده بین دو نسخه از اسناد را توضیح می دهد که با موقعیت تغییرات در سند اصلی HTML به صورت شهودی در آمدهاند پس از آن کارایی Vi-DIFF از لحاظ زمان اجرا ارائه میشود.

اهمیت تغییرات

بر مبنای Vi-Delta تولید شده، اهمیت تغییرات مشخص شده را با استفاده از تابع $E()$ را که، در بخش ه توضیح داده شد، ارزیابی میکنیم از طریق شبیه سازی نشان میدهیم که چگونه این تابع را میتوان برای زمان بندی خزشگرها به کار برد. همانطور که در بخش 6 مطرح شد زمان بند با استفاده از تابع ضرورت، ضروری ترین سندی که باید بازیابی شود را انتخاب می. کند. مطالعه مقایسههای راهبردهای زمان بندی شبیه سازی شده ارائه میشود.

8. کارهای آینده

کارهای آینده، مرتبط به مدیر بازیابی و تخمین اهمیت است. ما در حال حاضر بهترین روش یادگیری ماشین را برای به دست آوردن خودکار اهمیت نسبی بلاکها و اهمیت عملیات تغییر جست و جو می کنیم. کار، دیگر مشخص کردن، انتقال جداسازی و پیوستن بلاکها به عنوان تغییرات ساختاری است. همچنین قصد داریم روشمان را به منظور آشکار سازی و تحلیل تغییرات بین دو نسخه از سایت به جای دو صفحه گسترش دهیم.

منابع

M. Ben Saad, S. Gançarski, and Z. Pehlivan. A Novel Web Archiving Approach based on Visual Pages [1]
.Analysis

.In IAWW '09: 9th International Web Archiving Workshop, Corfu, Greece, 2009

D. Cai, S. Yu, J.-R. Wen, and W.-Y. Ma. VIPS: a Vision-based Page Segmentation Algorithm. [2]
,Technical report

.Microsoft Research, 2003

J. Cho and H. Garcia-Molina. The Evolution of the Web and Implications for an Incremental Crawler. In [3]
VLDB

.Proceedings of the 26th International Conference on Very Large Data Bases, 2000 :00"

J. Cho and H. Garcia-Molina. Estimating frequency of change. ACM Trans. Interet Technol., 3(3), [4]

G. Cobena, S. Abiteboul, and A. Marian. Detecting changes in XML documents. In ICDE '02: [5]
Proceedings of

.18th International Conference on Data Engineering, 2002

R. La-Fontaine. A Delta Format for XML: Identifying Changes in XML Files and Representing the [6]
Changes in

.XML. In XML Europe, 2001

R. Song, H. Liu, J.-R. Wen, and W.-Y. Ma. Learning block importance models for web pages. In WWW [7]
:'04

.Proceedings of the 13th international conference on World Wide Web, 2004

آرشیو داران وب برای گردآوری منابع ویدئویی وب بیش از همیشه با ابزارها و پروتکل‌های غیراستاندارد میزبانی میشوند این ، مقاله به وضعیت فعلی فناوری در این حوزه میپردازد. در این نوشتار، توجه به تجربیات چندین ساله گردآوری محتوای وب ویدئو به ذکر نمونه های مفصلی میپردازیم که به درک مسائل و راه حل‌های محتوای ویدئویی در وب کمک میکنند ، همچنین به معرفی چارچوب معماری که برای مقیاس بندی گردآوری محتوای ویدئویی وب می پردازیم که بخشی از پروژه تحقیقاتی اتحادیه اروپا موسوم به [LIWA\(1\)](#) است.

ص: 38

نوشته: رادو پاپ، (1) گابریل واسیلی، (2) ژولین ماسانه (3)

ترجمه: فروزان رضائی نیا (4)

مقدمه

ویدئو، امروزه به عنوان بخش مهمی از وب شناخته شده است فناوری گسترش یافته به ویدئو کمک میکند که همواره بر نیاز صنعت رسانه مسلط باشد به ویژه زمانی که دسترسی مستقیم به فایلها توسط کاربران منع شده باشد. به عنوان مفهوم، دیگر وظیفه جمع آوری مطالب آرشیو وب، بسیار سخت است و به توسعه رویکردها و ابزارهای خاص نیاز دارد. هدف این مقاله بررسی مشکلات اصلی آرشیو منابع ویدئویی وب است. براساس تجربیات به دست آمده در بنیاد آرشیو اروپا که در سالهای اخیر بر طیف متنوعی از این نوع محتوا در وب کار کرده است انواع حالتها در سال گذشته مشکلات مختلفی در ضبط بارگذاری و گردآوری ویدئو بوده که به دو دسته اصلی تقسیم و با استفاده از برخی مثالهای مناسب چندین راه حل فنی شرح داده می شود.

دسته نخست شامل وبگاههایی است که با استفاده از پروتکل استاندارد HTTP به ارائه محتوای

ویدئویی مبادرت میکنند.

ص: 39

Radu Pop -1

Gabriel Vasile -2

Julien Masanes -3

4- کارشناس سازمان اسناد و کتابخانه ملی جمهوری اسلامی ایران

مشکل آنها شامل ناشی از فنون مختلفی است که از پیوندهای مبهم برای فایل ویدئویی استفاده میکنند (برای مثال 2 یا 3 مسیره‌ی مجدد(1) یا پرش(2)). نمونه گویای این دسته مانند نماینده وبگاه YouTube است.

دسته دوم مشکلات در وبگاههایی نمودار میشوند که با استفاده از پروتکل‌های انتقال داده‌های غیر از حمل و نقل از HTTP نشان داده شده است از بین پروتکل‌های مختلف پخش زنده ای(3) که در حال حاضر مورد استفاده قرار میگیرند جدیدترین و مشکل ترین آنها یعنی پروتکل RTMP را برگزیدیم. ذکر این نکته مهم است که فناوریهای استفاده شده برای میزبانی منابع ویدئویی در وب بسیار سریع دستخوش تحول میشود و به احتمال زیاد جزئیات مواردی که در این مقاله ارائه میشوند به سرعت در شرف تغییر است همچنین در این نوع نیز به احتمال زیاد تغییر سریع در جزئیات آنها ارائه شده است. با این حال قصد این است که جزئیات به اندازه کافی ذکر، شوند تا درک منطق این فنون میسر گردد. حتی اگر این جزئیات به سرعت در شرف کهنه شدن باشند نتیجه آنچه که امروز به یک بازی موش و گربه تبدیل شده آن است که ابزارهای ما هم در سطح گردآوری و هم در سطح دسترسی، باید به طور مداوم بهبود یافته و روزآمد شود.

در بخش دوم مقاله، طرحی برای معماری آرشیو منابع ویدئویی وب پیشنهاد میکنیم که قادر به انطباق سریع و همچنین مقیاس پذیری براساس تفکیک فرآیندهای مربوط به بارگذاری منابع ویدئویی از خزشگر است. این طرح، باعث بهبود کارآمدی هم از لحاظ مقیاس پذیری و هم از لحاظ انعطاف پذیری میشود چرا که ابزارهای بارگذاری را که توسط متخصصان چندرسانه ای توسعه یافته اند، آسانتر میتوان ادغام و روزآمد کرد بالاخره با استفاده از خزشگرها به صورت غیر همزمان بهتر میتوان کنترل خطا و مدیریت فرآیند را پشتیبانی کرد.

2 - نمونه هایی برای بارگذاری ویدئو

پیچیدگیهای مختلفی در بارگذاری منابع ویدئویی وجود دارد که در دو نمونه زیر با تفصیل بیشتری بیان شده اند.

2-1- تغییر مسیر HTTP به YouTube.com

هر ویدئوی YouTube منحصراً با یک شناسه در هم سازی مشخص شده است (یک رشته از 11 کاراکتر) و عموماً میتواند در یک صفحه HTML با نشانی شبیه به زیر دسترس پذیر باشد

<http://www.youtube.com/watch?v=uniqueID>

احتمالاً سختترین چالش در برداشت ویدئوهای YouTube روزآمد کردن مکرر مکانیسمهایی است

که YouTube برای دسترسی ویدئوها در اختیار قرار میدهد و عبارت است از تلاش مستمر برای

ص: 40

Hop -2

Streaming Protocols -3

مخفی کردن نشانی مستقیم فایل‌های ویدئویی.

با استفاده از روش‌های گردآوری، کلاسیک خزشگر باید 5 مرحله متمایز از پیوندهای مستقیم

یا غیر مستقیم را به منظور دسترسی به مطالب ویدئویی طی کند الگوهای کلی از نشانیهای حد واسط در جدول 1 خلاصه شده است.

جدول . الگوهای YouTube برای دسترسی به فایل‌های ویدئویی

اگر شناسگر یک ویدئو (foobar) و آدرس صفحه (1 ouTube)، باشد خزشگر ابتدا URL فلش پلیر (2) صفحه را کشف می‌کند با عبور از فهرست پارامترها به، پلیر خزشگر ممکن است عبارت پرس وجوی HTML را برای شناسایی به کار ببرد که برای درخواست یک محتوای ویدئویی ارسال می‌شود قبل از به دست آوردن URL فایل ویدئویی (5) خزشگر باید یک واسطه(1) تغییر مسیر را دنبال کند که حاوی توکن های(2) رمزگذاری شده باشد مانند نشانیهای IP از میزبان و برچسب زمان(3) در پرس وجو.

عکس

مخفی کردن نشانی مستقیم فایل‌های ویدئویی.

با استفاده از روش‌های گردآوری کلاسیک، خزشگر باید ۵ مرحله متمایز از پیوندهای مستقیم یا غیرمستقیم را به منظور دسترسی به مطالب ویدئویی طی کند. الگوهای کلی از نشانی‌های حد واسط در جدول ۱ خلاصه شده است.

جدول ۱. الگوهای YouTube برای دسترسی به فایل‌های ویدئویی

#	URL	mime-type	HTTP response
1	http://www.youtube.com/watch?v=foobar	text / html	200
2	http://s.ytimg.com/yt/swf/watch-vfl118818.swf	application /x-shockwave-flash	200
3	http://www.youtube.com/get_video?video_id=foobar&t=...	text/html	204 (OK no content)
4	http://v1.lscache4.c.youtube.com/videoplayback?ip=0.0.0.0&sparams=id...	Redirect	302
5	http://v1.cache2.c.youtube.com/videoplayback?ip=0.0.0.0&sparams=id...	video / x-lv	200

اگر شناسگر یک ویدئو (foobar) و آدرس صفحه YouTube (#۱) باشد، خزشگر ابتدا URL فلش پلیر (#۲) صفحه را کشف می‌کند. با عبور از فهرست پارامترها به پلیر، خزشگر ممکن است عبارت پرس‌وجوی HTML (#۵) را برای شناسایی به‌کار ببرد، که برای درخواست یک محتوای ویدئویی ارسال می‌شود، قبل از به دست آوردن URL فایل ویدئویی (#۵)، خزشگر باید یک واسطه^۱ تغییر مسیر (#۵) را دنبال کند. که حاوی توکن‌های^۲ رمزگذاری شده باشد، مانند نشانی‌های IP از میزبان و برچسب زمان^۳ در پرس‌وجو.

URL ویدئو (#۵) در پارامترهای فلش، که هنگام بارگذاری صفحه (#۱) به صورت دینامیک تولید شده‌اند یکپارچه سازی می‌شوند.

مسئله‌ای که خزشگر با آن مواجه است، شناسایی صحیح این URL است، چراکه کاراکترهای مختلفی از آن حذف شده‌اند و پارامترهای مختلف دیگری توسط شیء فلش به آن الصاق شده‌اند.

به‌عنوان نمونه، در نسخه فعلی صفحه‌های YouTube، پارامتر «flashvars» رشته‌ای^۴ به طول ۳۳۷۴ کاراکتر است، در حالی که، URL آن می‌تواند ۳۹۲ کاراکتر داشته باشد. در تجزیه و تحلیل دقیق‌تر از

1 Intermediary.
2. Tokens
3. Time-stamp
4. String

URL ویدئو (5) در پارامترهای فلش که هنگام بارگذاری صفحه به صورت دینامیک تولید

شده اند یکپارچه سازی میشوند.

مسئله ای که خزشگر با آن مواجه است شناسایی صحیح این URL است چراکه کاراکترهای مختلفی از آن حذف شده اند و پارامترهای مختلف دیگری توسط شیء فلش به آن الصاق شده اند.

به عنوان نمونه در نسخه فعلی صفحه های YouTube پارامتر «flashvars» رشته ای (4) به طول 6374 کاراکتر است در حالی که URL آن میتواند 392 کاراکتر داشته باشد در تجزیه و تحلیل دقیق تر از

ص: 41

Intermediary -1

Tokens -2

Time-stamp -3

String -4

صفحه (1) URL ویدئو (5) در قطعه (1) پردازنده جاوا (2) شناسایی میشود که پارامتر فلش را ایجاد میکند در این نسخه صفحه های YouTube متنی که باید تجزیه شود حاوی رشته «PLAYER_CONFIG» و به دنبال آن دارای فهرستی از URLها به صورت تصادفی است. در بین دو کارکتر «[]» میتوان URL ویدئو را برداشت، کرد که توکنهای محاسبه شده را براساس نشانی IP و برچسب زمان در خود دارد.

همان طور که ملاحظه میشود مسیری که برای برداشت فایل های ویدئویی باید طی کرد، بسیار پیچیده است و URLهای منابع ویدئویی با تغییر مسیرهای مجدد و اضافه شدن توکنهای موقتی پیچیده شده اند.

در خصوص پارامترهای خزشگر میتوان خزش را برای 5 سطح در صفحه YouTube تنظیم کرد.

علاوه بر این URLها در و 5 به زیر دامنه های مختلف باز میشوند یا ارجاع میدهند، بنابراین باید با صراحت به دامنه های خزشگر اضافه شوند.

بالاخره مشکل بزرگی که به علت URLهای تغییر مسیر داده به وجود می آید، به ابزارهای دسترسی آرشیو مربوط است. حتی اگر هر فایل ویدئویی در YouTube یک URL متفاوت به صورت پویا در هر بارگذاری تولید می شود. به همین دلیل، باید به نمایه URLها برای ردیابی بین صفحه اصلی (1) و URL فایل ویدئو (5) اضافه شود.

در عمل، دوروش برای ضبط فایل های ویدئویی در YouTube وجود دارد: (1) ضبط ویدئویی برخط، (2) فنون بارگذاری غیر برخط فایل های ویدئویی

با توجه به فرآیند، خزش در روش برخط فایل های ویدئو در زمان خزش بارگذاری می شوند و این کار با استفاده از یک پردازنده اضافی به خزشگر Heritrix صورت می پذیرد عنوان مثال اسکریپت BeanShell که توسط آدام تیلور (3) نوشته شده، همه URLهای واسطه (در 5-2) را به فهرست خزش (4) میفرستد چرا که اغلب خارج از دامنه مطلوب برای خزش هستند (.s.ytimg.com, v.cache2.c.youtube.com). فایل های ویدئویی به فایل های WARC اضافه میشوند که توسط خزشگر ایجاد شده.

در روش غیر برخط فایل های ویدئویی پس از مرحله پردازش بر اساس URL صفحه های YouTube می شوند. در این روش از یک بارگذاری کننده بیرونی استفاده می کنند، که یک نمونه آن بارگذاری کننده توسعه یافته توسط ریکاردو گارسیا گونزالس (5) است که محتوای ویدئویی را بارگذاری و به فایل های flv منتقل می کند. کاربرد ابزار (6) WARC، فایل های flv به فایل های مجزای WARC وابسته بندی میکنند.

هر دوروش نیاز به ایجاد یک پیوند از URL اصلی ویدئو دارند همان گونه که در صفحه وب ظاهر میشود و نام فایل یا URL جدیدی که به برای محتوای ویدئو باز میشود چون فایل ویدئویی با تعقیب

JavaScript -2

<http://webarchive.jira.com/wiki/display/Heritrix/BeanShell+Script+For+Downloading+Video> -3

Frontier -4

<http://bitbucket.org/rg3/youtube-dl/wiki/Home> -5

<http://code.google.com/p/warc-tools> -6

URL های واسطه در آرشیو ذخیره میشود).

مزیت روش برخط این است که بارگذاری فایل‌های ویدئویی توسط خود خزشگر انجام میشود و

هیچ ابزار بیرونی برای پایش و همزمانی وجود ندارد. علاوه بر این همه سرآیندهای HTTP در آرشیو در تعامل با سرورهای YouTube ذخیره می‌شوند. اشکال این روش این است که URL نهایی محتوای ویدئویی (5) دیگر شناسگر اولیه ویدئو را در خود ندارد (1) مدیریت ردیابی شناسگر ویدئویی دشوار است؛ همچنین آرشیوها به همه URL های پرشی (1) آلوده میشوند (این URL ها دیگر معتبر نیستند چرا که توکنهای بارگذاری اعتبار موقت دارند).

از سوی دیگر رویکرد غیر برخط در پایش و مدیریت بارگذاری کننده‌های بیرونی انعطاف پذیرترند به عنوان مثال کنترل خطا) فایل‌های ویدئویی نام شناسگرهای اولیه خود را حفظ میکنند و URL های پرشی در آرشیو ذخیره نمیشوند لازم است یک سرآیند HTTP ساختگی (2) (برای هر فایل flv) در فایل‌های WARC وارد شود زیرا بارگذاری کننده بیرونی پاسخ سرور را حفظ نمیکند.

RTMP streaming on SWR.de -2-2

-1-2-2- مروری بر پروتکل‌های داده در جریان

Streaming، که با استانداردهای کارگروه مهندسی اینترنت (IETF) تطابق دارد اجازه میدهد تا سرور تبادلات را کنترل کند و برای نگهداری موجودیها در وضعیت بهینه سازی شده است. لازم نیست کاربران فایل‌های عظیم را بارگذاری کنند و این رویکرد به خصوص برای اطلاع رسانی و پخش زنده مناسب است.

در واقع پروتکل داده در جریان از دو نوع پروتکل داده در جریان به سرعت استفاده میکند: (3) [RTP 3550] برای ارسال بسته های داده رسانه ای و (4) RTSP [RFC2326] برای کنترل اطلاعات RTP از UDP استفاده میکند که بسته های گمشده را دوباره منتقل نمیکند، بنابراین حمل بر این است که همه طرفها پذیرفته اند که هنگام انتقال بعضی از بسته ها ناپدید شوند این بدان معنی است که کاربران باید موقرانه عدم دریافت همه داده مربوط به یک ویدئو و یا محتوای شنیداری را بپذیرند و خودشان مدیریت کنند. این نسبت به رویکرد مبتنی بر TCI/IP ترجیح داده میشود که گرفتن بسته های گم شده ممکن است مجبور باشد دفعات نامعینی تلاش کند و بنابراین زمان نامعینی نیز وقت لازم خواهد بود. RTSP نوعی پروتکل کنترل شبکه ای برای استفاده در صنعت سرگرمی و سیستمهای ارتباطی به منظور کنترل سرورهای رسانه streaming است. این، پروتکل برای ایجاد و کنترل تراکنش رسانه ای بین کاربران است کاربران سرورهای رسانه ای فرمانهای VCR مانند پخش و توقف برای تسهیل کنترل زمان بلادرنگ پخش فایل‌های رسانه از سرورها را ارسال میکنند.

ص: 43

Redirect URLs -1

Fake -2

Real-Time Streaming Protocol -3

Real - Time Transport Protocol -4

پروتکل RTSP به HTTP شباهت دارد با این تفاوت که RTSP در خواسته‌های جدیدی را اضافه می‌کند. HTTP فاقد هویت است. حال آنکه RTSP کاملاً دارای هویت است شناسه تراکنش برای ردگیری در مواقع لزوم مورد استفاده قرار می‌گیرد از این رو هیچ گونه ارتباط دائمی مورد نیاز نیست. پیامهای RTSP از کاربر به سرور فرستاده می‌شود هر چند استثنائاتی وجود دارد و گاه سرور به کاربر پیام می‌فرستد.

سرویس پیام‌رسانی چندرسانه‌ای (1)(MMS)

یک استاندارد ارتباطات از راه دور برای ارسال پیام به وسیله ابزارهای چند رسانه‌ای (عکس، صدا، تصویر و متن است. MMS توسعه یافته استاندارد SMS است که پیامهای طولانی تر را با استفاده از WAP برای نمایش محتوای پیامها را امکان پذیر می‌کند. پیامهای MMS در یک روشی تقریباً مشابه SMS تحویل داده می‌شوند اما محتوای چندرسانه‌ای ابتدا کدگذاری و در یک پیامی متنی به شیوه‌ای شبیه به ارسال ایمیل MIME درج می‌شود.

RTMP یک پروتکل دارای حق مالکیت است که توسط شرکت سیستمهای Adobe برای جریان داده‌های صوتی، تصویری و اطلاعات موجود در اینترنت بین فلش پلیر و سرور گسترش یافته است. این پروتکل برای تضمین تحویل جریان داده صوتی و تصویری در عین حفظ ارسال حجم بیشتری از، اطلاعات داده و تصاویر را از هم مجزا می‌کند اندازه قطعات میتواند به صورت پویا به وسیله کاربر و سرور تعیین شود، و حتی این امکان را میتوان پشت صفحه در صورت تمایل غیر فعال کرد.

شرکت سازنده Adobe مشخصات پروتکل RTMP را در 15 ژوئن 2009 (2) در دسترس عموم قرار داد اما این به نظر میرسد در این مشخصات بسیاری از جزئیات مربوط به پیاده سازی پروتکل ارائه نشده است.

2-2-2 - ضبط داده تصویری در جریان

بر اساس یک خزش در اینترنت که توسط European Archive برای ضبط تصاویر در حال پخش انجام شده است جزئیات فنی ای را که از وبگاه SWR.de گرفته شده اند ارائه می‌کنیم.

صفحه آرایی صفحه‌های ویدئو

ساختار نمایشی صفحه وبی که ویدئو نمایش میدهد از الگویی مشابه که مثال قبلی در YouTube.com پیروی می‌کند صفحه HTML شامل پانلهای اصلی از جمله پخش ویدئو (یعنی JW فلش ویدئو پلیر با کد منبع باز از ویدئو Long Tail) است. محتوای هر صفحه HTML به صورت پویا به وسیله سرور ایجاد و در عناصر خاص HTML ذخیره میشود مانند:

ص: 44


```
<<type=>>application/x-shockwave-flash
```

```
<<id=>>player46
```

```
/data=http://www.swr.de/static
```

```
<"jwplayer/player46.swf
```

...

فلش پلیر در یک عنصر `<object>` جاسازی شده و از طریق سرور وب بارگذاری میشود همه پارامترهایی که برای فلش پلیر مورد نیاز هستند به وسیله JavaScript آماده و (با استفاده از پارامترهای فلش به شیء فلش منتقل میشوند از اینجا همه تعاملات با سرور داده در جریان به طور مستقیم به وسیله فلش پلیر انجام میشود

ص: 45

بارگذاری فیلم در حال پخش

در این مثال، ویژه URL فایل تصویری به صورت واضح در پرده جاوا نوشته شده است. از دید خزشگر، این نمونه مناسبی است چون برداشت کننده (1) پرده جاوا قادر به شناسایی و برداشت URL فایل‌های تصویری است اما این URL برای خزشگر URL معتبری نیست زیرا برنامه های پروتکل HTTP/HTTPS را پشتیبانی نمیکنند.

از این رو URL مربوط به RTMP به عنوان یک URL غیر معتبر در رخدادهای خطای خزشگر ارائه گزارش داده میشود برای بارگیری بهتر فایل‌های تصویری (URL) `rtmp://.../foobar-video.flv` از رخدادهای خطا خارج و به بارگذاری کنندگان بیرونی واگذار میشود (مانند FLVStreamer). بارگذاری کننده RTMP محتوای تصاویر را در فایل‌های flv که در فایل WARC بسته بندی شده اند کپی می کند.

عکس



شکل ۱- وبگاه زنده - ویدئوی ATMP

بارگذاری فیلم در حال پخش

در این مثال ویژه، URL فایل تصویری به صورت واضح در پرده‌جاوا نوشته شده است. از دید خزشگر، این نمونه مناسبی است چون برداشت‌کنندهٔ پرده‌جاوا قادر به شناسایی و برداشت URL فایل‌های تصویری است. اما این URL برای خزشگر، URL معتبری نیست، زیرا برنامه‌های پروتکل HTTP / HTTPS را پشتیبانی نمی‌کند.

از این رو، URL مربوط به RTMP به‌عنوان یک URL غیرمعتبر در رخدادهای خطای خزشگر ارائه گزارش داده می‌شود. برای بارگیری بهتر فایل‌های تصویری (rtmp://.../foobar-video.flv)، URLها از رخدادهای خطا خارج و به بارگذاری‌کنندگان بیرونی واگذار می‌شود (مانند FLVStreamer). بارگذاری‌کننده RTMP، محتوای تصاویر را در فایل‌های flv که در فایل WARC بسته‌بندی شده‌اند، کپی می‌کند.

دسترسی به آرشیو محتوای ویدئویی

از نقطه نظر دسترسی، تفاوت اساسی بین صفحه وب زنده و صفحه وب آرشیو شده، استفاده از پروتکل انتقال برای تحویل تصاویر است. قرار دادن سرور برای داده جریان در زیرساخت آرشیو، مستلزم تلاش فراوان و راه‌حلی پرهزینه برای توسعه و نگهداری از آن است. همچنین، سرورهای مختلفی برای پروتکل‌های مختلف مورد نیاز خواهد بود. به‌طور کلی، دسترسی به آرشیو محتوای ویدئویی از طریق پروتکل HTTP و فایل‌های flv، به‌طور مستقیم از فایل‌های WARC و صفحه‌های HTML و دیگر

1. Extractor

دسترسی به آرشیو محتوای ویدئویی

از نقطه نظر دسترسی تفاوت اساسی بین صفحه وب زنده و صفحه وب آرشیو شده، استفاده از پروتکل انتقال برای تحویل تصاویر است. قرار دادن سرور برای داده جریان در زیرساخت آرشیو مستلزم تلاش فراوان و راه‌حلی پرهزینه برای توسعه و نگهداری از آن است. همچنین سرورهای مختلفی برای پروتکل‌های مختلف مورد نیاز خواهد بود به‌طور کلی دسترسی به آرشیو محتوای ویدئویی از طریق پروتکل HTTP و فایل‌های flv، به‌طور مستقیم از فایل‌های WARC و صفحه‌های HTML و دیگر

منابع ایجاد می شود. در حال حاضر دسترسی به کارکردهای خاص جریان داده وجود ندارد (مانند fastforward یا بارگیری حجیم) اما دسترسی پایه به محتوا قطعی است و آن را میتوان بدون هزینه اضافی برای نمونه های مختلف تنظیم کرد.

ایراد مهم در جایگزینی داده در جریان با پروتکل HTTP ساده در ویدئوهای بزرگ قابل مشاهده است و اگر هیچ گونه تنظیمی صورت نگیرد اجراکننده باید تمام فایل های flv را قبل از اجرای تصاویر بارگذاری کند ما ابزارهای دسترسی را در EA برای کمتر کردن این مشکل را به وسیله بارگذاری زیادی از فایل های ویدئویی از آرشیو تنظیم کردیم. حجم این تصاویر بهینه سازی می شود، اما باید تعادلی بین دسترسی سریع برای اجرای ویدئوها و پیچیدگی روش های buffering انجام شود.

در زمان دسترسی، تنظیم اصلی که باید صورت گیرد جایگزینی فلش پلیر اصلی است. چون فایل ویدئویی آرشیو شده دیگر در جریان نیست اجراکننده اولیه در صفحه نمیتواند برای نسخه آرشیو شده مورد استفاده قرار گیرد؛ از این رو قالب HTML باید به تناسب روزآمد شود:

```
flashvars=file=http://collection.europarch
```

```
.ive.org/swr/...7.1.AEVA/rtmp://fcondemand
```

```
swr.de/at/e/foobar-video.flv
```

```
src=http://collections.europarchive.org/me
```

```
dia/player.swf
```

```
type=application/x-shockwave-flash
```

```
<id=video645568_plr
```

در مقایسه با صفحه، زنده در یک نسخه آرشیوی عناصر script و <object را با عنصر خاص

که حاوی اجراکننده خاص آرشیو است و URL تصاویر آرشیوی جایگزین میکنیم:

"http://collections.europarchive.org/media/player.swf"

http://collection.europarchive.org/swr/20100601084708/rtmp://.../foobar-video.flv

ص: 47

توجه داشته باشیم که URL فایل ویدئویی آرشیوی به HTTP URL که به فایل flv در آرشیو باز

می شود تبدیل شده اند.

جایگزینی عناصر HTML با شیوه های خاصی که در کد دسترسی در سرور پیاده سازی صورت می گیرد. فایل های WARC همیشه نسخه اصلی صفحه ها را نگه میدارد پیدا کردن یک الگوی مشترک برای شناسایی صحیح و جایگزینی عناصر حاوی اجراکننده هنوز چالش برانگیز است، چرا که ساختار و ویژگیهای یک عنصر <object> ممکن است از یک وبگاه به وبگاه دیگر متفاوت باشد.

در این بخش ما به گرفتن نمونه فیلمهای RTMP، پرداختیم زیرا فرآیند بارگذاری مستلزم پروتکل خاص داده در جریان است. با این حال فنون جایگزینی اجراکننده که در کد دسترسی پیاده می شوند در مورد فیلمهای بارگذاری شده در پروتکل HTTP به همان شیوه انجام میگیرد.

عکس



شکل ۲

توجه داشته باشیم که URL فایل ویدئویی آرشیوی به HTTP URL که به فایل flv در آرشیو باز می شود تبدیل شده اند.

جایگزینی عناصر HTML، با شیوه های خاصی که در کد دسترسی در سرور پیاده سازی صورت می گیرد. فایل های WARC همیشه نسخه اصلی صفحه ها را نگه می دارد. پیدا کردن یک الگوی مشترک برای شناسایی صحیح و جایگزینی عناصر حاوی اجراکننده هنوز چالش برانگیزست، چرا که ساختار و ویژگی های یک عنصر <object> ممکن است از یک وبگاه به وبگاه دیگر متفاوت باشد.

در این بخش ما به گرفتن نمونه فیلم های RTMP پرداختیم، زیرا فرآیند بارگذاری مستلزم پروتکل خاص داده در جریان است. با این حال، فنون جایگزینی اجراکننده که در کد دسترسی پیاده می شوند در مورد فیلم های بارگذاری شده در پروتکل HTTP به همان شیوه انجام می گیرد.

۳- ضبط ویدئو با استفاده از بارگذاری کننده خارجی

به عنوان جزئی از فناوری های جدید در پروژه LiWA¹، ماژول خاصی برای افزایش قابلیت های ضبط توسط خزشگر با توجه به انواع محتوای چندرسانه ای طراحی شد [در مورد نخستین تلاش هایی که در این مورد در EA انجام شد نگاه کنید به Baly 2006]. نسخه فعلی Heritrix مبتنی بر پروتکل HTTP/HTTPS است و نمی تواند پروتکل های دیگر را که در چندرسانه ای ها (نظیر داده در جریان) مورد استفاده قرار می گیرند، اجرا کند.

ماژول Liwa Rich Media Capture² بازیابی محتوای چندرسانه ای را به یک برنامه کاربردی خارجی

1. <http://www.liwa-project.eu/>

2. The LiWA Rich Media Capture module, <http://code.google.com/p/liwatechnologies/source/browse/rich-media-capture>

3- ضبط ویدئو با استفاده از بارگذاری کننده خارجی

به عنوان جزئی از فناوری های جدید در پروژه LiWA (1)، ماژول خاصی برای افزایش قابلیت های ضبط توسط خزشگر با توجه به انواع محتوای چندرسانه ای طراحی شد در مورد نخستین تلاش هایی که [در این مورد در EA انجام شد نگاه کنید به Baly 2006]. نسخه فعلی Heritrix مبتنی بر پروتکل HTTP/HTTPS است و نمیتواند پروتکل های دیگر را که در چندرسانه ایها (نظیر داده در جریان) مورد استفاده قرار میگیرند اجرا کند.

/http://www.liwa-project.eu -1

The LiWA Rich Media Capturemodulehttp://code.google.com/p/liwatechnologies/source/browse/rich- - 2
media-capture

واگذار می کند (مانند 1) MPlayer یا (2) FLVStreamer که قادر است طیف وسیعی از پروتکل‌های انتقال را مدیریت کند.

این ماژول به عنوان پلاگین (3) خارجی برای Heritrix ساخته شده است. در این رویکرد، شناسایی و بازیابی داده‌های در جریان کاملاً از هم جدا شده اند و امکان استفاده از ابزارهای کارآمدتر برای تجزیه و تحلیل محتوای ویدئویی و شنیداری فراهم شده است.

-3-1- معماری

ماژول از چندین جزء فرعی ساخته شده است که از طریق پیام ارتباط برقرار میکنند ما از پروتکل ارتباط استاندارد استفاده میکنیم که (4) AMQP نامیده میشود.

تلفیق ضبط رسانه های غنی (5) با خزنگر در تصویر 3 نشان داده شده و گردش کار پیامها را به صورتی که ذکر میشود میتوان خلاصه کرد پلاگین متصل به هریتریکس URL منابع داده‌های در جریان را شناسایی و برای هر کدام از آنها یک پیام AMQP ایجاد میکند این پیام به یک سرور پیام مرکزی منتقل میشود و نقش سرور پیام این است که هریتریکس را از بارگذاری کننده های خوشه بندی شده داده در جریان مجزا کند یعنی همان ابزارهای بارگذاری کننده خارجی سرور پیام یوآرال را در صف (6) میکند و زمانی که یکی از بارگذاری کننده ها در دسترس باشد URL بعدی را برای پردازش میفرستد.

در معماری ماژول سه ماژول فرعی را شناسایی میکنیم:

- ماژول کنترل اولیه مسئول دسترسی به سرور، پیام آغاز کارهای جدید متوقف کردن آنها و ارسال هشدار؛
- ماژول دوم که برای شناسایی و بارگذاری داده‌های در جریان استفاده میشود (از یک ابزار خارجی مانند MPlayer استفاده میشود)
- ماژول سوم که داده در جریان بارگذاری شده را در فرمتی که به وسیله ابزارهای دسترسی قابل شناسایی باشند دوباره بسته بندی میکند (ضبط کننده WARL)

ص: 49

1- <http://www.mplayerhq.hu>

2- <http://savannah.nongnu.org/projects/flvstreamer>

3- Plugin

4- Advanced Message Queuing Protocol (AMQP), <http://www.amqp.org/confluence/display/AMQP/Advanced+Message+Queuing+Protocol>

5- Rich Media

6- Queue

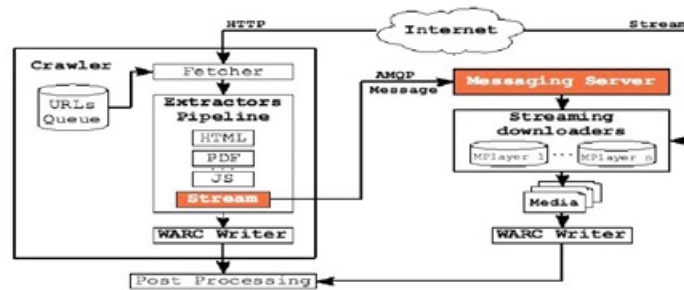
تصویر 3 - ماژول ضبط داده در جریان در تعامل با خزشگر

بارگذاری کننده، فیلم به سرور پیام دهی متصل میشود تا URL فیلم جدیدی را برای ضبط کردن درخواست کند به محض دریافت URL جدید یک تحلیل اولیه به منظور شناسایی بعضی پارامترها انجام می شود از جمله نوع و طول زمان جریان داده، البته اگر فیلم زنده باشد ممکن است زمان ثابت و قابل پیکر بندی در نظر گرفت.

بعد از یک شناسایی، موفق بارگذاری واقعی شروع میشود ماژول کنترل برنامه کاری ای را ایجاد میکند که همراه با تمهیدات حفاظتی به MPlayer منتقل میشود تا مدت زمان بارگذاری از تخمین اولیه طولانی تر نشود بعد از ضبط موفقیت آمیز آخرین مرحله بسته بندی فیلمهای ضبط شده در فایل WARC است که بعد از آن به مرحله ذخیره سازی نهایی هدایت میشود.

2-3- بهینه سازی(1)

عکس



تصویر ۳- مازول ضبط داده در جریان در تعامل با خزشگر

بارگذاری کننده فیلم، به سرور پیام دهی متصل می شود تا URL فیلم جدیدی را برای ضبط کردن درخواست کند به محض دریافت URL جدید، یک تحلیل اولیه به منظور شناسایی بعضی پارامترها انجام می شود از جمله نوع و طول زمان جریان داده. البته، اگر فیلم زنده باشد ممکن است زمان ثابت و قابل پیکربندی در نظر گرفت.

بعد از یک شناسایی موفق، بارگذاری واقعی شروع می شود. مازول کنترل، برنامه کاری ای را ایجاد می کند که همراه با تمهیدات حفاظتی به MPlayer منتقل می شود تا مدت زمان بارگذاری از تخمین اولیه طولانی تر نشود. بعد از ضبط موفقیت آمیز، آخرین مرحله بسته بندی فیلم های ضبط شده در فایل WARC است که بعد از آن به مرحله ذخیره سازی نهایی هدایت می شود.

۳-۲- بهینه سازی^۱

مسائل اصلی که از آزمایش های اولیه مطرح می شوند، همزمان کردن خزشگر و مازول ضبط خارجی است. در مورد مجموعه ای از فیلم های پر حجم در وبگاه، بارگذاری پی در پی فیلم قطعاً از فرآیند جست و جوی صفحه های متن طولانی تر خواهد بود. و خزشگر باید منتظر باشد تا مازول بیرونی بارگذاری فیلم را تمام کند. با افزایش بارگذاری کننده ها می توان سرعت بارگذاری تصاویر را افزایش داد. از سوی دیگر، اگر بخواهیم فرآیند به صورت موازی انجام شود.

راه حل دیگر برای کنترل بارگذاری کننده های فیلم این است که مازول ضبط فیلم را از خزشگر جدا کنیم و آن را در مرحله بعد از پردازش به کار بیندازیم. این به معنای جایگزینی پلاگین خزشگر با یک ثبت کننده ورود به سیستم و یک مدیر مستقل برای بارگذاری کننده های فیلم است. از مزایای این روش (که برای نمونه در EA استفاده شده است) عبارتند از:

- نگرشی همه جانبه از تعداد کل URL های فیلم

1. Optimizations

مسائل اصلی که از آزمایش های اولیه مطرح می شوند، همزمان کردن خزشگر و مازول ضبط خارجی است. در مورد مجموعه ای از فیلم های پر حجم در وبگاه بارگذاری پی در پی فیلم قطعاً از فرآیند جست و جوی صفحه های متن طولانی تر خواهد بود و خزشگر باید منتظر باشد تا مازول بیرونی بارگذاری فیلم را تمام کند. با افزایش بارگذاری کننده ها می توان سرعت بارگذاری تصاویر را افزایش داد. از سوی دیگر، اگر بخواهیم فرآیند به صورت موازی انجام شود.

راه حل دیگر برای کنترل بارگذاری کننده های فیلم این است که مازول ضبط فیلم را از خزشگر جدا کنیم و آن را در مرحله بعد از پردازش به کار بیندازیم. این به معنای جایگزینی پلاگین خزشگر با یک ثبت کننده ورود به سیستم و یک مدیر مستقل برای بارگذاری کننده های

فیلم است.

از مزایای این روش (که برای نمونه در EA استفاده شده است) عبارت اند از:

- نگرشی همه جانبه از تعداد کل URL های فیلم

ص: 50

• مدیریت بهتر منابع تعداد بارگذاری کننده های فیلم که مشترکاً از پهنای باند استفاده می کنند). اشکال اصلی این روش ناهماهنگی است که ممکن است بین زمان خزش از وبگاه و ضبط فیلم در مرحله بعد از پردازش ظاهر شود.

• برخی محتویات فیلم ممکن است ناپدید شود (با فاصله یک یا دو روز)

• بارگذاری فیلم در انتظار پایان یافتن خزش متوقف شود.

بنابراین باید تعادلی هنگام مدیریت بارگذاری فیلم برقرار کرد کوتاه کردن زمان بارگذاری کامل کنترل اشتباهات (برای فیلمهایی که سرورهایشان کند است) و بهینه سازی پهنای باند که توسط بارگذاری کننده های مختلف استفاده میشود.

4- نتیجه گیری و چشم انداز

همانگونه که ملاحظه شد ارائه راه حلی کلی برای پرداختن به تمام وبگاهها که محتوای ویدئویی دارند مشکل است. براساس روش ارائه شده در این مقاله، باید برای هر مورد خاص تکنیک گردآوری خاصی به کار گرفته شود. مهندسی خزش برای تنظیم با ابزارها بستگی به پیچیدگی وبگاه دارد.

کارهایی که در این زمینه میتوان انجام داد در سه گروه جای می گیرد:

• مقیاس بندی ضبط ویدئو به احتمال زیاد با جدا کردن آن از خزش و مدیریت بهتر خطاهای متعدد و وقفه هایی که سرورهای ویدئویی به طور کلی و سرورهای جریان به طور خاص ایجاد میکنند.

• بهبود کشف خودکار از پیوندهای سردرگم و تعقیب آنها با قواعد خاص (allowing off-domains) کشف منابع در قالبهایی غیر از Html و....).

• توسعه روشی عام برای دسترسی و ارائه این مستلزم تشخیص ویژگی نمایش دهنده ها به طور خودکار، است به گونه ای که آنها را به شیوه ای عام جایگزین کند دسترسی بهتر را مدیریت و گزینه هایی برای فایل های بزرگ و غیره در اختیار قرار دهد.

5- تقدیر و تشکر

بخشی از هزینه های این کار توسط کمیسیون اروپا زیر نظر LIWA تامین اعتبار شده است.

منابع

Baly, N., Sauvin, F. (2006). Archiving Streaming Media on the Web, Proof of Concept and First Results. International Web Archiving Workshop (IWAW 06), Alicante, Spain

توسعه خدمات اطلاعاتی و محتوایی در وب و از دیگر سو توسعه فناوریهای ارتباطی و معرفی زیر ساختهای نرم افزاری سخت افزاری و مدل‌های جدید مدیریتی در سطوح مختلف سازمانها باعث افزایش روزافزون نیاز به ایجاد زیرساختهای لازم جهت پشتیبانی از ذینفعان مختلف این حوزه شده است نیاز به توسعه دسترسی در سطح، وب درخواست جستجوهای، پیشرفته وابستگی به خدمات معنایی تقاضای بهبود رابطها و محیطهای تعاملی و مبتنی به مدل برای کاربران رشد رویکردهای بهینه سازی روشهای گردآوری و ذخیره سازی اطلاعات در وب از جمله مؤلفه‌هایی است که توجه به زیر ساختهای متناسب و نوین در این حوزه را به سرعت مطرح مینماید گسترده شدن نوع خدمات اختصاصی در سازمانها که منجر به گسترده شدن زنجیره های ارزش شده است و همچنین توسعه سازمانهای مجازی و افزوده شدن حجم تعاملات الکترونیکی درون و برون سازمانهای مختلف باعث شده است که قابلیت تعامل پذیری در سطوح مختلف، مدیریتی فرایندها، خدمات و دادگان به عنوان یکی از حوزه های اصلی پژوهش مطرح شده و از جنبه های مختلف مورد بررسی قرار گیرد نگاه سلسله مراتبی فوق به این تعاملات علاوه بر تحقق خدمات مشترک یکپارچگی لازم در تمام فضای زنجیره ارزش را تا فیزیکی ترین تعاملات داده ای و نرم افزارهای کاربردی ایجاد مینماید در حوزه تعاملات نرم افزاری به عنوان یکی از زیر ساختهای فنی مهم در سازمانها پشتیبانی از نیازمندیهای سیستمی نظیر پروتکل‌های ارتباطی واسط‌های ارتباطی دسترسی به دادگان گونه های، اطلاعاتی معناشناسی پارامترها ایجاد قابلیت‌های تحرک و انجام وظایف کاربردی به عنوان اولویت اصلی شناخته شده است. یکی از شیوه های تحقق زیر ساخت نرم افزاری مناسب برای تعاملات سازمانی بکارگیری عامل‌های هوشمند است.

عامل‌های هوشمند، سیستم‌های نرم افزاری خود، مختار واکنش گرا پیش فعال و با توانایی برقراری ارتباط هستند که می توانند به عنوان اجزاء کلیدی خدمات اطلاعاتی و محتوایی مبتنی بر وب در سامانه های مختلف مورد استفاده قرار گیرند. بسیاری از خدمات محتوایی تحت وب از جمله فرایندهای مدیریتی توزیع شده اشتراک گذاری محتوا و یکپارچه سازی خدمات و محتوا از مخازن و سامانه‌های مدیریت محتوای مختلف با استفاده از عاملها محقق میشوند و علاوه بر صرفه جویی در زمان و هزینه انعطاف لازم جهت پشتیبانی از رخدادهای غیر منتظره را فراهم می آورند به کارگیری عامل‌های نرم افزاری همانگونه که در محیط درون سازمان زیر ساخت مناسبی برای انجام مؤثر و کارای تعاملات ایجاد مینماید در فضای برون سازمانی نیز تأثیرات زیادی خواهد داشت به دلیل تفاوت در دامنه کاری و فناوریهای مورد استفاده در هر سازمان لازم است که روشهایی که برای تبادل و به اشتراک گذاری داده میان آنها استفاده میشود مستقل از هر نوع فناوری خاصی بوده و امکان تولید و به اشتراک گذاری اطلاعات سازگار و با معنی را داشته باشند برای تبادل مؤثر دادهها باید درک مشترکی از داده میان سازمانها وجود داشته باشد که این امر با استفاده از روشهای هستان شناسی محقق می.شود پیشرفتهای وب 2 و امکان به اشتراک گذاری و درک داده ها توسط موتورهای جستجو از جمله دستاوردهای استفاده از روشهای هستان شناسی است. از مجموع مباحث فوق میتوان نتیجه گرفت که استفاده از عاملها نه تنها در فرایندهای درون سازمان منجر به افزایش کارایی و اثربخشی میشوند؛ بلکه در افزایش کیفیت رابطه با نهاد‌های برون سازمانی از قبیل تأمین کنندگان، مشتریان، رقبا یا شرکای احتمالی تأثیر بسزایی دارند و استفاده از معماری مبتنی بر عامل میتواند ضمن ایجاد صرفه اقتصادی برای سازمانها به عنوان مزیت رقابتی ماندگار محسوب شود. در این مقاله ضمن معرفی کاربرد عامل‌های نرم افزاری در سازمانها، به معرفی چارچوبی مبتنی بر عاملها برای یکپارچه کردن تعاملات اطلاعاتی در آنها پرداخته میشود همچنین در انتها با استفاده از یک استاندارد شناخته شده در کنسرسیوم جهانی، توسعه همگرایی و پذیرش خدمات تحت وب و کسب و کار الکترونیکی (1) OASIS، نشان داده خواهد شد که چگونه میتوان چارچوب پیشنهادی را برای پشتیبانی تعامل پذیری در سیستم‌های مدیریت محتوا و مخازن اطلاعاتی استفاده کرد و ضمن بهره مندی از تعامل گسترده میان منابع و خدمات مختلف در سطح وب، از قابلیت‌های عامل‌های هوشمند برای بهبود این گونه خدمات استفاده کرد.

1. مقدمه

استفاده از فناوری اطلاعات در نهادهای مختلف منجر به رخداد تغییرات بسیاری در درون سازمانها و در طول زنجیره ارزش شده است. اگرچه این فناوری مزیت‌های بسیار زیادی برای سازمانها به همراه دارد، ولی مسائلی را در بر میگیرد که میتواند مدیریت ارشد را با مشکلاتی روبرو کند تغییرات زیر در نتیجه ورود فناوری اطلاعات و تحولات فضای سازمانها در زنجیره ارزش رخ داده است:

- تغییر تمرکز مدیریت زنجیره ارزش از تولیدات خاص به درخواستهای کاربر و همزمان سازی

فعاليتها

- به دست آوردن دیدگاههای جدید راجع به شبکه گسترده شرکا و روابط میان آنها
- افزایش عدم قطعیت در زنجیره ارزش و به طبع در اجزای کوچکتر زنجیره و سازمانها
- افزایش اطلاعات و دادههای ورودی و خروجی سازمان و پیچیده شدن نگهداری و کنترل این داده ها

ص: 53

1- عضو هیئت علمی پژوهشگاه ارتباطات و فناوری اطلاعات kharrat@itrc.ac.ir

2- دانشجوی کارشناسی ارشد مهندسی فناوری اطلاعات دانشگاه تهران m.mosharraf@ut.ac.ir

3- استادیار دانشکده مهندسی برق و کامپیوتر دانشگاه تهران ftaghiyar@ut.ac.ir

● افزایش سیستمهای اطلاعاتی سازمانها توزیع شدن این سیستمها در شبکه و تغییر مدیریت آنها از مرکزی به توزیع شده

تبدیل شدن امکان بازیابی اطلاعات و استفاده مجدد از امکانات موجود به عنوان یک ضرورت

● انجام تعاملات سازمانها از طریق شبکههای اینترنتی و به صورت مجازی

● ضروری شدن توانایی پیکر بندی پویا و عملکرد سریع در برابر تغییرات

● نیاز به ایجاد یکپارچه سازی هم شکلی و اتحاد میان سازمانها و سامانه های موجود برای ارائه خدمات با توجه به نکات اشاره شده لازم است راهکارهای گوناگونی برای تحقق فرایندهای سازمانی و روند مدیریت جدید در این حوزه مد نظر قرار بگیرد؛ به طوریکه علاوه بر ایجاد زنجیره ارزش مورد نظر تعاملات لازم در این خصوص را نیز پشتیبانی کند.

ظرفیتهای و قابلیتهای تعامل میان سازمانها و سامانه ها حوزه جدیدی را تحت عنوان تعامل پذیری (1) وارد ادبیات سازمانی نموده است. مفهوم تعامل پذیری در بردارنده موارد و نکات متنوعی است که تعریف آن را با پیچیدگیهای خاصی روبرو می کند. جمع بندی تعاریف معتبر و در عین حال متفاوتی که محققان و دست اندرکاران حوزه تعامل پذیری ارائه کرده اند میتواند تعریف ساده و در عین حال نسبتاً کامل ذیل را برای تعامل پذیری قبول کرد تعامل پذیری عبارت است از توانایی دو سیستم به منظور شناخت و استفاده از کارکردهای یکدیگر از دیدگاه فناوری رایانه تعامل پذیری نشان دهنده توانایی دو سیستم رایانههای ناهمگن برای کارکرد مشترک و همچنین برای دادن دسترسی به منابع یکدیگر به صورت متقابل و دو طرفه است. در زمینه سازمانهای شبکه ای نیز تعامل پذیری اشاره به توانایی تعاملات (تبادل اطلاعات و خدمات میان سیستمهای سازمانی دارد در صورتی که تعاملات میان دو طرف حداقل در سه سطح داده، خدمات و فرایندها و در زمینه کسب و کاری مشخص صورت پذیرد، آنگاه تعامل پذیری حساس و مهم جلوه میکند (David et al. 2008).

به طور کلی، تعامل پذیری به معنای همزیستی خودمختاری و محیط متحد است؛ در حالی که یکپارچگی بیشتر اشاره به مفاهیمی نظیر هماهنگی وابستگی و متحدالشکل شدن دارد. مقایسه رایانه های به هم پیوسته قوی (2) و پیوسته ضعیف (3) میتواند به درک بهتر تفاوت این دو مفهوم کمک کند. پیوستگی قوی یعنی ارتباط اجزا به گونه ای است که جدا شدن آنها غیر ممکن بوده و به هم وابستگی دارند. این تعریف قابل مقایسه با مفهوم یکپارچگی است. در مقابل، پیوستگی ضعیف بیانگر گونه ای از ارتباط است که اجزا میتوانند در عین اتصال به یکدیگر با حفظ منطق عملیاتی خاص خود با یکدیگر تبادل خدماتی برقرار نمایند این گونه از ارتباط از درجه به هم پیوستگی قابل قیاس با تعامل پذیری است. یکپارچگی شرط لازم و کافی برای تعامل پذیری است ولی رابطه معکوس آن برقرار نیست.

دیدگاه دیگری در این خصوص که توسط ISO 14258 بیان شده عنوان میکند که دو سیستم در صورتی یکپارچه در نظر گرفته میشوند که قالبی استاندارد و تفصیلی برای تمامی اجزا مؤلفها وجود داشته

ص: 54

Interoperability -1

Tightly Coupled -2

Loosely Coupled -3

باشد (ISO 14258, 1999) رویکرد همشکل سازی در تعامل پذیری منوط به وجود ساختاری فراسطحی و مشترک (ارائه ابزاری برای برقراری تعامل معنایی) در تمامی مدل‌های مؤلفه‌ها است. از سوی دیگر در محیطی که لازم است مدل‌ها به صورت پویا و مستمر با یکدیگر تطابق یابند، رویکرد تعامل پذیری از ضروریات آن محیط قلمداد می‌شود لذا اگر ضرورت فوق در محیطی نباشد با توجه به اینکه شرایط اولیه تحقق تعامل پذیری پیچیدگی بالایی دارد تعهد به تعامل پذیری توصیه نمی‌شود (David et al. 2008).

1. تعامل پذیری می‌تواند در حوزه‌های مختلف سازمانی مورد توجه قرار گیرد که این حوزه‌ها عبارتند از: 1. تعامل پذیری داده این مفهوم اشاره به توانایی سازگار کردن مدل‌های داده‌ای متفاوت و زبان‌های جستجو دارد تعامل پذیری داده مشتمل بر یافتن و به اشتراک گذاری اطلاعات بدست آمده از پایه‌های غیریکنواخت که توسط سیستم‌های کاربردی و مدیریتی بانک‌های اطلاعاتی مختلف استخراج می‌شوند، است.

2. تعامل پذیری خدمات اشاره به شناسایی و سازگار کردن کارکردهای بسیاری از خدماتی که مجزا طراحی و اجرا شده‌اند. دارد واژه خدمات تنها به برنامه‌های کاربردی مبتنی بر رایانه محدود نشده و کارکردهای شرکتها و سازمانهای شبکه‌ای را نیز شامل می‌شود.

3. تعامل پذیری فرایندها هدف از تعامل پذیری فرایندها، هماهنگ کردن فرایندهای متعدد سازمانی برای کار کردن با یکدیگر است یک فرایند سلسله مراتب خدمات را با توجه به برخی نیازهای مشخص سازمان تعریف می‌کند همچنین بررسی چگونگی اتصال فرایندهای داخلی دو سازمان به منظور به وجود آوردن فرایندی مشترک ضروری است.

4. تعامل پذیری کسب و کار اشاره به کارکرد هماهنگ و یکنواخت در سطح سازمانی و پرهیز از سبک‌های مختلف و متفاوت تصمیم‌گیری روش‌های مختلف، کاری قانونی و رویکردهای تجاری به منظور توسعه کسب و کار مشترک و هماهنگ در داخل سازمان و با سازمانهای دیگر دارد (David 2008, et al).

بر اساس استاندارد ISO 14258 سه رویکرد اساسی و پایه برای مرتبط کردن موجودیتهای (سیستمها) به یکدیگر به منظور برقراری تعامل پذیری وجود دارد (ISO 12581999):

1. رویکرد ایجاد یکپارچگی (1): در این رویکرد یک قالب مشترک برای تمامی مدل‌ها وجود دارد این قالب مشترک، لزوماً یک استاندارد نیست بلکه شیوه‌ای است که باید توسط تمامی ذی‌نفعان در امور ساخت مدل‌ها و سیستمها مورد تأیید قرار گیرد.

2. رویکرد ایجاد یک شکلی (2): رویکردی است مشترک که تنها در فراسطح (3) وجود دارد. فرامدل (4) موجودیتی قابل اجرا نیست؛ بلکه وسیله‌ای برای تبادل معنایی است که تبدیل و تفسیر مدل‌ها را امکان‌پذیر می‌کند.

ص: 55

Integrated Approach -1

Unified Approach -2

Meta-level -3

Meta Model -4

3. رویکرد ایجاد اتحاد(1): در این رویکرد قالب مشترکی وجود ندارد برای برقراری تعامل پذیری هر یک از ذی نفعان باید خود را با شرایط محیطی تعامل پذیری تطبیق دهد استفاده از این رویکرد این نکته را در بردارد که هیچ یک از ذی نفعان نمیتواند مدلها، زبانها و روشهای کاری خودش را به دیگران تحمیل کند این مطلب نشان از آن دارد که باید برای ترسیم مفاهیم در سطح معنایی، یک هستان شناسی مشترک ایجاد گردد (Berre et al. 2004).

هر یک از سه رویکرد فوق امکان تعامل پذیری میان سیستمهای سازمانی را فراهم میکنند اما رویکرد ایجاد اتحاد، نسبت به رویکردهای دیگر از مقبولیت بیشتری برای ایجاد تعامل پذیری برخوردار است.

به طور کلی میتوان گفت که انتخاب رویکرد مناسب بستگی زیادی به نوع نیازها و زمینه مورد نظر دارد. اگر هدف از تعامل پذیری ادغام سازمانها باشد استفاده از رویکرد ایجاد یکپارچگی مناسب است. در چنین موردی تنها نیاز به ایجاد یک قالب مشترک برای طرفین بوده و همه مدلهای بر اساس این قالب مشترک میبایست ساخته و تفسیر شوند. اگر نیاز به تعامل پذیری برای ایجاد همکاری بلندمدت باشد، رویکرد ایجاد یک شکلی رویکردی مناسب است در این حالت یک فرا مدل مشترک برای تمامی طرفها تعریف شده و سپس براساس آن امکان همترازی معنایی و ایجاد نگاشت میان مدلهای مختلف فراهم میشود در صورتی که نیاز به تعامل پذیری برای ایجاد همکاری در پروژه های کوتاه مدت باشد میتوان از رویکرد ایجاد تعهد و همکاری استفاده کرد همچون سازمانهای مجازی در این حالت طرفها باید برای کسب توافق به صورت پویا با هم سازگاری لازم را ایجاد کنند.

تحقق راهکارهای ذکر شده نیازمند زیر ساختهای گوناگونی از جمله زیر ساختهای فنی است. در همین راستا سامانه های الکترونیکی و نرم افزارهای تحلیلی - مدیریتی، امکانات بسیاری را برای سازمانها فراهم کرده اند که هر یک وجوهی از نیازهای مربوطه را محقق مینمایند. استفاده از هستان شناسی در رویکرد ایجاد اتحاد به سازماندهی دانش می انجامد و همزمان قالب مناسبی را پدید می آورد که از آن رهگذر چگونگی استفاده از دانش قابل درک میشود. علاوه بر این استفاده از هستان شناسی نه تنها باعث شناسایی عناصر دانش می شود، بلکه به شناسایی و اصلاح ناهمخوانیهای واژه ای میان واحدهای مختلف کمک میکند (Navarretta et al. 2006). شباهت و تناظر ویژگیهایی چون پیمانه ای بودن(2)، توزیع شدگی تغییر پذیری ساختار یافتگی و پیچیدگی در سازمانها با همین خصوصیات در عاملهای هوشمند آنها را به گزینه مناسبی برای پشتیبانی از تغییرات و حل مشکلات سازمان تبدیل کرده است.

به دلیل اهمیت اقتصاد دانش محور استفاده از عاملهای هوشمند در سازمانها با معرفی مدلهایی در سازمانهای تجاری شروع شد. در سال 1997 مدلی مبتنی بر عاملها برای نشان دادن اهمیت یکپارچگی در سازمانهای تجاری توسط Chu و همکارانش پیشنهاد شد (1997). اگرچه در فرایندی که توسط آنها مطرح شده بود، روند یک سازمان به صورت جامع و از طراحی تا تولید پوشانده نشده و بیشترین تمرکز بر معماری نرم افزاری عاملها قرار گرفته بود ولی به عنوان یکی از نخستین گامهای برداشته شده

در این موضوع مورد توجه است در همین سال Shen و همکارانش مدلی برای زیر ساختهای مبتنی بر، عامل به نام (1) ABMEI ارائه دادند که به عنوان شبکه ای از واسطها بین زیر بخشهای سازمان قرار گرفته و به حل مشکل واگرایی آنها میپرداخت (1997) این سیستم به دلیل مشکلاتی که در مذاکرات عاملها وجود داشت با موفقیت زیادی همراه نبود در ادامه این روند دیدگاه مشابه دیگری توسط Maturana و همکارانش مطرح و سیستمی با نام Metamorph پیشنهاد شد (1999). Metamorph نیز مانند ABMEI از عاملها برای برقراری ارتباط استفاده میکرد در سال 2000 این دیدگاه با گسترش دامنه کار از طراحی تا اجرا ادامه پیدا کرده و Metamorph II پا به عرصه نهاد (Shen et al. 2000) در سیستم جدید نیز کاستیهایی مانند عدم پشتیبانی از بخشهای کیفیت و تحقیق و توسعه نتوانست انتظارات پیش بینی شده را به خوبی محقق کند به این منظور محققان سامانه ای را برای اعمال کنترل توزیع شده به نام ManAge پیشنهاد کردند (Heikkila et al.) عملکرد این سیستم مبتنی بر چهار عامل اصلی عملگر کنترل کننده اجراکننده و ناظر بنا شده بود. مشکل اصلی این مدل نیز قرار گرفتن تمرکز اصلی کار بر ساختار درونی خود عاملها به جای فرایندهای سازمان بود.

زنجیره ارزش یکپارچه مبتنی بر سیستمهای چند عاملی با تمرکز بر تولید، برنامه ریزی، کنترل و حرکت هدف محور توسط Frey و همکارانش ارائه شد (2003). بعد از آن قالب دیگری مبتنی بر سیستمهای چند عاملی برای انجام یکپارچه فرایند برنامه ریزی و زمان بندی تولید در سازمانها مطرح شد (Lim et al. 2004). هدف این، سیستم ایجاد یکپارچگی به منظور افزایش انعطاف و پویایی سازمان در مقابل رقیبان بود Feng و همکارانش به ارائه مدل دیگری به منظور پشتیبانی از برنامه ریزی، پیش بینی و کنترل در سازمان پرداختند (2005) در این مدل عاملها به پایگاه دانش پایگاه داده سازمان سیستم کنترل و طراحی کامپیوتری دسترسی داشته و از اطلاعات آنها استفاده میکردند.

برای ایجاد یکپارچگی در سازمان قالبی مفهومی به نام MIBIS توسط Kishore و همکارانش ارائه شد (2006). در این مدل سیستمهای چند عاملی به عنوان ابزاری برای شبیه سازی یکپارچگی در سیستمهای اطلاعاتی - تجاری مورد استفاده قرار میگرفتند عاملها در این مدل با استفاده از هشت کلمه کلیدی عامل، نقش هدف تعامل کار، منابع اطلاعات و دانش ورودیها را درک کرده و با برقراری یکپارچگی در فرایندهای مربوط به این سیستمها زیر ساخت بهینه ای را برای سازمانهای تجاری فراهم میکردند به منظور پوشش فرایندهای طراحی محصولات و خدمات ارزیابی سازمان برنامه ریزی زمان بندی و مدیریت تولید بلادرنگ، عاملها در سامانه ای به نام MMA (2) مطرح شدند. (Mahesh et al, 2007 MMA) به عنوان کنترل کننده مرکزی با ارسال پیام به دیگر عاملها باعث سازماندهی آنها و حفظ اطلاعات سیستم میشد. در ادامه تحقیقات این دامنه عاملهای چند، رفتار، در زنجیره ارزش و با استراتژیهای برنامه ریزی مختلف توسط Forget و همکارانش پیشنهاد شدند (2008). تمرکز این مدل بر عاملهای برنامه ریزی، تولید شامل، تحویل، نگهداری انبار و حمل و نقل قرار گرفته بود.

ص: 57

1- Agent-Based Manufacturing Enterprise Infrastructure

2- Manufacturing Management Agent

در ادامه این مقاله ابتدا عاملهای هوشمند و ساختار استنتاجی آنها را معرفی خواهیم کرد. در بخش سوم به بیان کاربرد اصلی عاملهای هوشمند در مدیریت اطلاعات سازمانها پرداخته و عملکرد آنها از منظر تعاملات درون سازمانی و برون سازمانی را بررسی میکنیم در بخش چهارم، چارچوبی مبتنی بر عاملهای هوشمند برای یکپارچگی سازمانها و زنجیره ارزش معرفی شده و از جنبه های مختلفی بررسی میشود. در ادامه نیز ملاحظات لازم جهت انطباق چارچوب پیشنهادی برای پشتیبانی تعامل پذیری در سیستمهای مدیریت محتوا و مخازن اطلاعاتی بیان میگردد در بخش پایانی نیز جمع بندی و نتیجه گیری ارائه شده است.

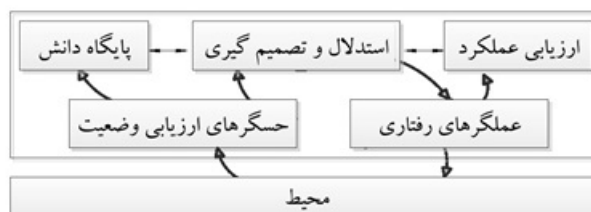
2. عاملهای هوشمند

عامل هوشمند سیستم نرم افزاری خود مختار با تواناییهای اجتماعی واکنش گرا و پیش فعال است تواناییهای تعاملی این عاملها و قدرت پشتیبانی آنها در برخورد با تغییرات پیش بینی نشده، آنها را از دیگر نرم افزارها متمایز کرده است. عاملها به عنوان جزئی از سیستمهای پیچیده و برای نیازمندیهای، توزیعی تعاملی موازی، سازی مقیاس، پذیری هماهنگی و همکاری بین اجزا استفاده می شوند (2005) (Cao et al. 2009; Zhang et al.) به همین دلیل در دیدگاههایی که بر مبنای عاملهای هوشمند نرم افزاری ارائه میشوند راه حلها به صورت توزیع شده در نظر گرفته شده و استانداردهایی برای برقراری تعامل میان آنها تعریف میشود که قدرت مذاکره و ایجاد فهم مشترک را برای آنها ممکن میکند. شبیه سازی اجتماع انسانها در سیستمهای چند عاملی تعریف سلسله مراتب برای آنها و تنظیم نظارت در ساختارشان، و استفاده از مکانیزمهای حراجیها در شکستن کارها و تقسیم خرده کارها بین عاملها به طوریکه با افزایش کیفیت باعث بهبود روند اجرای کارها نیز شود از دیگر ویژگیهای این عاملها هستند. به کارگیری عاملها به عنوان زیر ساخت نرم افزاری در تمام کاربردهای مبتنی بر فناوری توانسته است بستری یکپارچه برای تعامل و هماهنگی همه اجزای آنها فراهم آورده و بهبود چشمگیری در سامانه های مرتبط با آنها و کیفیت ارائه خدمات ایجاد کند. علاوه بر این سیستمهای چند عاملی به عنوان ابزاری برای پشتیبانی و بهینه سازی پایگاههای اطلاعاتی سیستمها نیز میتوانند استفاده شوند. در این دیدگاه عاملها میتوانند به عنوان کاربردهای انحصاری و ناهمگن که در فرایند تصمیم سازی شرکت می کنند مطرح شوند در ادامه ساختار استنتاجی یک عامل هوشمند نشان داده شده است. همانطور که شکل 1 نشان میدهد یک عامل میتواند با دریافت دادههای محیط و پردازشهایی که روی این دادهها انجام میدهد، قوانین استنتاجی خود و در نتیجه عملکردش را در طول زمان بهبود بخشد (Oztemel et al. 2009)

شکل 1: ساختار استنتاجی عامل هوشمند

استفاده از عاملها به عنوان قدرت اجرایی در تعامل با کاربران یک سیستم الکترونیکی و نیز در ارتباط سیستمهای الکترونیکی با همدیگر از ویژگی تعامل پذیری عاملها و عملکرد اجتماعی آنها ناشی شده است. استفاده از بسترهای مختلف برای پیاده سازی عاملها و ارائه واسط کاربری متناسب با کاربر باعث شده است کاربران سیستمهای الکترونیکی هنگام استفاده با انگیزه و اشتیاق بیشتری عمل کرده و اعتماد بیشتری نیز به امنیت تعاملات داشته باشند.

عکس



شکل ۱: ساختار استنتاجی عامل هوشمند

استفاده از عامل ها به عنوان قدرت اجرایی در تعامل با کاربران یک سیستم الکترونیکی و نیز در ارتباط سیستم های الکترونیکی با همدیگر، از ویژگی تعامل پذیری عامل ها و عملکرد اجتماعی آنها ناشی شده است. استفاده از بسترهای مختلف برای پیاده سازی عامل ها و ارائه واسط کاربری متناسب با کاربر، باعث شده است کاربران سیستم های الکترونیکی هنگام استفاده، با انگیزه و اشتیاق بیشتری عمل کرده و اعتماد بیشتری نیز به امنیت تعاملات داشته باشند.

۳. کاربرد عامل های هوشمند در مدیریت اطلاعات سازمان ها

یکپارچگی در یک سازمان و در بعد وسیع تر در زنجیره ارزش به همبستگی بین بخش های مختلف زنجیره که از دریافت نیازمندی ها شروع شده و تا تحویل آنها ادامه می یابد اشاره می کند. این همبستگی در تمام قسمت ها، شامل عناصر پیش پردازش، بخش ورودی، پردازش، خروجی، سنجش و تعامل با محیط وجود دارد. در صورت وجود یکپارچگی در سیستم های اطلاعاتی سازمان، خدمات ترتیبی هر بخش طوری زمان بندی می شوند که به محض پایان کار بخش پیشین، آغاز شده و بلافاصله پس از تمام شدن (به موقع و یا پیش از موعد) آنها بخش بعدی شروع به کار کند. مشکلی که در این باب وجود دارد، این است که منابع لازم برای این بخش ها به طور قطعی مشخص نیست. نمود اصلی این مشکل در درخواست های ارائه شده به سازمان (نه در فرایندهای معمول سازمان) است. راهکار ارائه شده برای حل این مشکل، تقسیم بلادرنگ درخواست ارائه شده به چند زیر مسئله و تخمین زدن منابع برای هر کدام است. نکته قابل توجه در اینجا این است که باید بین زیر مسئله ها، راه حل ها و منابع ارائه شده برای هر کدام هماهنگی وجود داشته باشد (Wang et al. ۲۰۰۸). به همین جهت ضمن انجام کارها به صورت موازی، لازم است بخش های مستقل نیز تعاملات نسبی با هم داشته و بین آنها زیر ساخت یکپارچه ای وجود داشته باشد.

۳-۱. درون سازمانی

با قرار گرفتن زیر ساخت سازمان ها بر سیستم های کامپیوتری و نیز در نظر گرفتن این موضوع که با گسترده شدن ابعاد سازمان و گسترش حجم داده ها، این سیستم ها نیز باید از قابلیت توزیع شدگی

۳. کاربرد عامل های هوشمند در مدیریت اطلاعات سازمانها

یکپارچگی در یک سازمان و در بعد وسیع تر در زنجیره ارزش به همبستگی بین بخش های مختلف زنجیره که از دریافت نیازمندیها شروع شده و تا تحویل آنها ادامه می یابد اشاره می کند. این همبستگی در تمام قسمت ها شامل عناصر پیش پردازش بخش ورودی پردازش، خروجی، سنجش و تعامل با محیط وجود دارد. در صورت وجود یکپارچگی در سیستم های اطلاعاتی سازمان، خدمات ترتیبی هر بخش طوری زمان بندی میشوند که به محض پایان کار بخش پیشین آغاز شده و بلافاصله پس از تمام شدن به موقع و یا پیش از موعد آنها بخش بعدی شروع به کار کند مشکلی که در این باب وجود دارد، این است که منابع لازم برای این بخشها به طور قطعی مشخص نیست. نمود

اصلی این مشکل در درخواستهای ارائه شده به سازمان نه) در فرایندهای معمول (سازمان) است راهکار ارائه شده برای حل این مشکل تقسیم بلادرنگ درخواست ارائه شده به چند زیر مسئله و تخمین زدن منابع برای هر کدام است. نکته قابل توجه در اینجا این است که باید بین زیر مسئلهها، راه حلها و منابع ارائه شده برای هر کدام هماهنگی وجود داشته باشد (Wang et al. (2008) به همین جهت ضمن انجام کارها به صورت موازی لازم است بخشهای مستقل نیز تعاملات نسبی با هم داشته و بین آنها زیر ساخت یکپارچه ای وجود داشته باشد.

1-3 درون سازمانی

با قرار گرفتن زیر ساخت سازمانها بر سیستمهای کامپیوتری و نیز در نظر گرفتن این موضوع که با گسترده شدن ابعاد سازمان و گسترش حجم داده ها این سیستمها نیز باید از قابلیت توزیع شدگی

ص: 59

برخوردار باشند ضرورت این نکته که نه تنها عناصر یک سامانه باید به طور مستقل به خوبی مدیریت شوند، بلکه روابط بین آنها نیز از اهمیت ویژه ای برخوردار خواهد شد دو چندان میشود (Sheory 2006). مقیاس پذیری سازمانها و ایجاد قابلیت توزیع شدگی در زیر ساختهای آن، با پشتیبانی از مسائلی چون همزمانی ناسازگاری و سربار اطلاعاتی ممکن خواهد بود با استفاده از عاملهای هوشمند می توان مدیریت توزیع شده را در چهار گام جمع آوری داده ها مشخص کردن پارامترهای وابستگی، انجام تحلیلها به صورت توزیعی تجمیع و همکاری پیاده سازی کرد (Wang et al. 2009).

عاملهای هوشمند میتوانند در سطوح مختلف زیر ساختهای یک سازمان را تشکیل دهند بسته به سطحی که عاملها در آن قرار میگیرند و وظایفی که در آن به عهده دارند قابلیتهای آنها و در نتیجه پیچیدگیهای آنها تفاوتهای اساسی خواهد داشت. معمولاً زیر ساختهای فراهم شده به وسیله عاملها هم سطح نبوده و دارای سلسله مراتب معنی داری است (Pawlewski et al, 2009 Wang et al. 2009). این سلسله مراتب حداقل سه لایه زیر را در بر خواهد گرفت:

● نظارت

● تحلیل و کاربرد

● داده

بر حسب شرایط و اندازه سازمان پیچیدگیهای امنیتی و تسهیلاتی که در دسترسها فراهم می کند، لایه های دیگری نیز میتوانند به این ساختار افزوده شده و یا لایه های معرفی شده خود به چندین زیر لایه تقسیم شوند.

معمولاً پایین ترین سطح به پایگاههای اطلاعاتی و ساختارهای داده ای سازمان می پردازد. عاملهای مربوط به این سطح ضمن مدیریت دسترسی به ساختارهای اطلاعاتی سازمان به عنوان واسطی برای تبدیل اطلاعات و ایجاد هماهنگی و سازگاری در اطلاعات توزیع شده عمل می کنند (Wang et al. 2006).

لایه تحلیل گر معمولاً در بخش اصلی این ساختار قرار گرفته و میتواند به عنوان یک عنصر نرم افزاری یک پارچه یا مجموعه ای از عناصر توزیع شده عمل کند این لایه دادههای دریافت شده از لایه پایین تر را بازیابی کرده و ضمن انجام پالایش آن دادهها تحلیلهای مرتبط را بر آنها انجام خواهد داد. نتایج تحلیل میتواند به عنوان گزارشی برای به روز کردن سیاستهای برنامه ریزی و فعالیتهای سازمان به بخش مدیریت ارسال شود. این لایه میتواند بنا به اقتضا از تعدادی عامل کاربردی نیز تشکیل شود که بخش فیزیکی کار را انجام دهند (Yang et al. 2010)

لایه نظارت از ابزارهایی برای نمایش روند کارها شامل توصیف فرایندهای تجاری، قوانین و منطق آنها تشکیل شده است کاربردهای مناسب برنامه ریزی و زمان بندی برای لایه های پایینتر در این لایه مشخص می شود. علاوه بر این روند کار با دیگر سیستمهای تجاری عملگرهای انسانی، مشتریها و ماشینها و دیگر کاربردها از وظایف این لایه خواهد بود (Wang et al. 2006).

اگرچه عاملها به عنوان نرم افزارهایی که دارای قابلیتهای اجتماعی هستند، برای برقراری تعامل و ارتباط بین بخشهای مستقل یا توزیع شده معرفی شدند؛ ولی باید این نکته در نظر گرفته شود که استفاده

از عاملها تنها به همین مورد محدود نمی شود نکته ای که در ارتباط عاملها باید در نظر گرفته شود این است که این ارتباطات میتوانند از طریق شبکه و به صورت ناهمگام نیز انجام شود و همین مورد میتواند استقلال عملکردی آنها را افزایش دهد.

3-2 بین سازمانی

موفقیت یک سازمان در برقراری سریع ارتباط با دیگر سازمانها به طرز اشتراک گذاری اطلاعات وابستگی مستقیم دارد به این دلیل طراحی و ساخت ابزارهای ارتباطی و نحوه تبادل دادهها اهمیت بسیار زیادی خواهد داشت فناوریهای ارتباطی و اطلاعاتی میتوانند با فراهم آوردن مواردی نظیر مدیریت شبکه مدیریت و نظارت بر اجتماع، شرکا پیکر بندی سازمان مجازی و کنترل مشارکتی زمان هزینه و کیفیت از تعاملات مجازی سازمانها پشتیبانی کنند.

وب سرویسها به عنوان راه حلی برای سیستمهای کاربردی سازمانی و به دلیل قابلیت پیکر بندی بالا و پشتیبانی از عملکرد سازمانها برای رسیدن به سطوح بالاتری از توزیع شدگی مورد توجه قرار گرفتند. اگر چه این نرم افزارها به عنوان پارادایم مناسب برای کاربردهای با حجم بزرگ مانند زنجیره ارزش و معماریهای مبتنی بر سرویس به عنوان شیوه ای پیشرو در معماریهای سازمانی مطرح شده اند، با این حال نتوانسته اند انتظارات سازمانها در تعاملات برون سازمانی را محقق کنند دلایلی چون توصیف سخت فرایندهای سازمانی برای وب سرویسها کشف زمان بر ساختار معنایی نابالغ تجمیع سخت و نبود تضمین برای امنیت دادههای سازمانها از جمله دلایل عدم موفقیت وب سرویسها در برقراری ارتباط بین سازمانی محسوب میشوند (Shen et al. (2007).

با در نظر گرفتن قابلیتهایی که اجتماع عاملهای هوشمند در توزیع شدگی برای سازمانها فراهم کنند و نیز ویژگیهای واکنش گرایی و پیش فعالی، عاملها ترکیب آنها با وب سرویسها می تواند تقایص آنها را جبران کند عاملهای هوشمند با تحلیل نیازمندها ساخت پویای سرویس، مذاکره با سازمانهایی که از سوی وب سرویسها کشف شده اند و ایجاد قرار داد در صورت حاصل شدن توافق به برقراری ارتباطهای بین سازمانی کمک میکنند.

در فرایند مذاکره که به صورت بین سازمانی انجام میشود علاوه بر اینکه ممکن است دو عامل مذاکره کننده دارای درک مشترکی از مفاهیم نبوده و دچار مشکل شوند در مذاکرات کاربر انسانی با عاملها نیز ممکن است بحرانهایی رخ دهد به همین دلیل برای استفاده از عاملها در سطوح پایین تعاملات داشتن زبان مشترک برای برقراری ارتباط میان آنها کفایت میکند؛ ولی در سطوح بالا وجود هستان شناسی مشترک برای فراهم کردن درک مشترکی از مفاهیم لازم است.

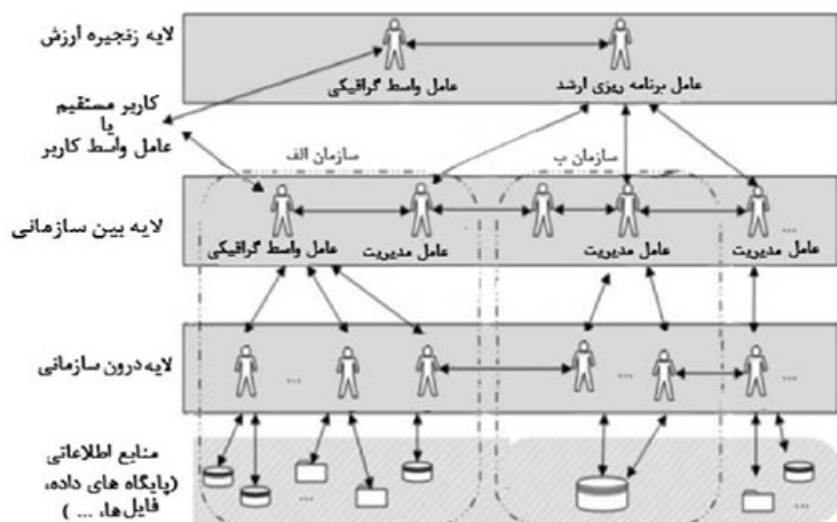
4. پیشنهاد چارچوبی مبتنی بر عامل در یکپارچه سازی مدیریت اطلاعات سازمانها ایجاد یکپارچگی و امکان تعامل پذیری در سازمانها نیازمند فراهم کردن چارچوبی منطقی برای جریان یافتن اطلاعات در میان بخشهای مختلف و نیز تصمیم گیریهای مشارکتی است به این منظور لازم

است ضمن ارائه زیر ساخت پیوسته در تمامی بخشهای سازمان امنیت اطلاعات و ساختارهای نظارتی نیز در نظر گرفته شود استفاده از شیوه های مطرح در معماری نرم افزار مانند معماریهای چند لایه، می تواند برای این منظور مفید واقع شود. در این معماریها ضمن فراهم کردن ساختار سلسله مراتبی و اعمال کنترل از بالا به پایین امنیت اطلاعات نیز از پایین به بالا تأمین می شود.

عکس

است ضمن ارائه زیر ساخت پیوسته در تمامی بخشهای سازمان، امنیت اطلاعات و ساختارهای نظارتی نیز در نظر گرفته شود. استفاده از شیوه های مطرح در معماری نرم افزار مانند معماریهای چند لایه، می تواند برای این منظور مفید واقع شود. در این معماریها ضمن فراهم کردن ساختار سلسله مراتبی و اعمال کنترل از بالا به پایین، امنیت اطلاعات نیز از پایین به بالا تأمین می شود.

با در نظر گرفتن این موضوع که تعاملات سازمانها علاوه بر اطلاعات، لایه فرایند را نیز در بر می گیرد، معماریهای سلسله مراتبی نمی تواند جوابگوی نیاز آنها باشد. از طرفی برای برقراری امنیت در تعاملات سازمان و اعمال دیدگاههای نظارتی، این معماریهای سلسله مراتبی مفید هستند. با افزودن تعاملات در سطوح فرایندی به ساختار سلسله مراتبی، می توان بر مشکل مطرح شده غلبه کرد. شکل ۲ چارچوب چهار لایه پیشنهاد شده مبتنی بر عاملهای هوشمند در سازمانها را نشان می دهد. در چارچوب ارائه شده سه لایه پایین در سطح سازمان و لایه چهارم به منظور اعمال مدیریت در سطح زنجیره ارزش خواهد بود. زنجیره ارزش شبکه ای از نهادهایی است که در تولید و توزیع مواد اولیه با ارزش افزوده و تحویل نهایی آن به مشتری همکاری دارند. عملکرد مؤثر زنجیره ارزش را می توان در رساندن به موقع و با کیفیت خدمات، اطلاعات یا محصولات تولیدی به مشتریها اندازه گیری نمود. به همین دلیل زنجیره ارزش باید بتواند در فرایندهای برنامه ریزی، زمان بندی و کنترل با پیشامدهای غیر قطعی داخلی و خارجی مقابله کند. مدیریت زنجیره ارزش یک رویکرد یکپارچه سازی برای برنامه ریزی و کنترل مواد و اطلاعات است که از تأمین کنندگان تا مشتریان جریان دارد.



شکل ۲: چارچوب مبتنی بر عامل در یکپارچه کردن مدیریت اطلاعات سازمانها

با در نظر گرفتن این موضوع که تعاملات سازمانها علاوه بر اطلاعات لایه فرایند را نیز در بر می گیرد، معماریهای سلسله مراتبی نمیتواند

جوابگوی نیاز آنها باشد. از طرفی برای برقراری امنیت در تعاملات سازمان و اعمال دیدگاههای نظارتی این معماریهای سلسله مراتبی مفید هستند با افزودن تعاملات در سطوح فرایندی به ساختار سلسله مراتبی میتوان بر مشکل مطرح شده غلبه کرد. شکل 2 چارچوب چهار لایه پیشنهاد شده مبتنی بر عاملهای هوشمند در سازمانها را نشان میدهد در چارچوب ارائه شده سه لایه پایین در سطح سازمان و لایه چهارم به منظور اعمال مدیریت در سطح زنجیره ارزش خواهد بود زنجیره ارزش شبکه ای از نهادهایی است که در تولید و توزیع مواد اولیه با ارزش افزوده و تحویل نهایی آن به مشتری همکاری دارند عملکرد مؤثر زنجیره ارزش را میتوان در رساندن به موقع و با کیفیت خدمات اطلاعات یا محصولات تولیدی به مشتریها اندازه گیری نمود به همین دلیل زنجیره ارزش باید بتواند در فرایندهای برنامه ریزی زمان بندی و کنترل با پیشامدهای غیر قطعی داخلی و خارجی مقابله کند مدیریت زنجیره ارزش یک رویکرد یکپارچه سازی برای برنامه ریزی و کنترل مواد و اطلاعات است که از تأمین کنندگان تا مشتریان جریان دارد.

شکل 2: چارچوب مبتنی بر عامل در یکپارچه کردن مدیریت اطلاعات سازمانها

به دلیل اهمیتی که منابع اطلاعاتی و داده‌های سازمان دارند این بخش در پایین‌ترین لایه چارچوب پیشنهادی قرار گرفته و دسترسی به آن به وسیله عامل‌های کنترل‌کننده محدود می‌شود. بخش‌های مختلف سازمان به فراخور نیاز به منابع اطلاعاتی دسترسی داشته و در لایه بعدی قرار می‌گیرند. بخش مدیریت در بالاترین قسمت قرار گرفته و بر تعاملات تمام بخش‌های درونی سازمان نظارت دارد. به همین ترتیب مدیریت زنجیره ارزش بر کل سازمانهایی که در زنجیره قرار می‌گیرند نظارت خواهد داشت. لایه‌های مطرح شده در این چارچوب به این شرح هستند:

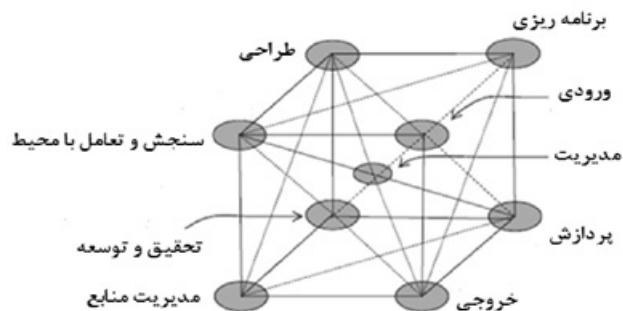
● لایه زنجیره ارزش: وظیفه این لایه ایجاد یکپارچگی در طول زنجیره ارزش و برقراری ارتباط با کاربر نهایی است. با توجه به درخواست کاربر و موقعیت سازمانها این لایه وظایف سطوح پایینتر را مشخص کرده و به آنها اعلام میکند عامل برنامه ریزی ارشد که وظیفه اصلی برنامه ریزی را بر عهده دارد در طول زنجیره ارزش یکتا است. این عامل با جمع‌آوری داده‌های سازمانها در طول زنجیره و اتصال این داده‌ها در مورد زمان بندی تولید، توزیع‌کننده‌ها و تأمین‌کننده‌ها تصمیم‌های لازم را کند. عامل برنامه ریزی ارشد مسئول برنامه ریزیهای کلان برای عامل‌های مدیریت است.

عکس

به دلیل اهمیتی که منابع اطلاعاتی و داده‌های سازمان دارند، این بخش در پایین‌ترین لایه چارچوب پیشنهادی قرار گرفته و دسترسی به آن به وسیله عامل‌های کنترل‌کننده محدود می‌شود. بخش‌های مختلف سازمان، به فراخور نیاز به منابع اطلاعاتی دسترسی داشته و در لایه بعدی قرار می‌گیرند. بخش مدیریت در بالاترین قسمت قرار گرفته و بر تعاملات تمام بخش‌های درونی سازمان نظارت دارد. به همین ترتیب مدیریت زنجیره ارزش بر کل سازمان‌هایی که در زنجیره قرار می‌گیرند نظارت خواهد داشت.

لایه‌های مطرح شده در این چارچوب به این شرح هستند:

- لایه زنجیره ارزش: وظیفه این لایه ایجاد یکپارچگی در طول زنجیره ارزش و برقراری ارتباط با کاربر نهایی است. با توجه به درخواست کاربر و موقعیت سازمان‌ها، این لایه وظایف سطوح پایین‌تر را مشخص کرده و به آنها اعلام می‌کند. عامل برنامه ریزی ارشد که وظیفه اصلی برنامه ریزی را بر عهده دارد در طول زنجیره ارزش یکتا است. این عامل با جمع‌آوری داده‌های سازمان‌ها در طول زنجیره و اتصال این داده‌ها، در مورد زمان بندی تولید، توزیع کنندگان و تأمین کنندگان تصمیم‌های لازم را اتخاذ می‌کند. عامل برنامه ریزی ارشد مسئول برنامه ریزی‌های کلان برای عامل‌های مدیریت است.
- لایه مدیریت سازمان: وظیفه این لایه ایجاد یکپارچگی در سازمان، تصمیم‌گیری‌های مدیریتی، برنامه ریزی و نظارت بر لایه‌های درونی سازمان است. عامل مدیریت که عاملی یکتا در سازمان است؛ بر روابط سایر عامل‌ها در درون سازمان نظارت کرده و مسئول اجرای طرح‌های محول شده از سوی عامل برنامه ریزی ارشد است. این عامل می‌تواند تا زمانی که در طرح اصلی و برنامه‌های استراتژیک سازمان تناقضی پیش نیامده است، با تدوین برنامه‌های محلی، زمان بندی سایر بخش‌های سازمان را مشخص کند.
- لایه درونی سازمان: این لایه از بخش‌های مختلفی تشکیل شده که هر کدام وظایف مشخصی دارند. بر حسب نوع و پیچیدگی سازمان، تعداد عامل‌های این لایه و روابط آنها متغیر است. در حالت کلی این لایه از هشت عامل تشکیل می‌شود. شکل ۳ عامل‌های این لایه را که نماینده بخش مرتبط با نامشان هستند را نشان می‌دهد.



شکل ۳: عامل‌های لایه درون سازمان و روابط بین آنها

• لایه مدیریت سازمان: وظیفه این لایه ایجاد یکپارچگی در سازمان تصمیم‌گیری‌های مدیریتی، برنامه ریزی و نظارت بر لایه‌های درونی سازمان است. عامل مدیریت که عاملی یکتا در سازمان است؛ بر روابط سایر عامل‌ها در درون سازمان نظارت کرده و مسئول اجرای طرح‌های محول شده از سوی عامل برنامه ریزی ارشد است. این عامل می‌تواند تا زمانی که در طرح اصلی و برنامه‌های استراتژیک سازمان تناقضی پیش نیامده است، با تدوین برنامه‌های محلی زمان بندی سایر بخش‌های سازمان را مشخص کند.

• لایه درونی سازمان: این لایه از بخش‌های مختلفی تشکیل شده که هر کدام وظایف مشخصی دارند. بر حسب نوع و پیچیدگی سازمان تعداد عامل‌های این لایه و روابط آنها متغیر است. در حالت کلی این لایه از هشت عامل تشکیل می‌شود. شکل ۳ عامل‌های این لایه را که

نماینده بخش مرتبط با نامشان هستند را نشان میدهد.

شکل 3: عاملهای لایه درون سازمان و روابط بین آنها

ص: 63

رأسهای مکعب فوق نشان دهنده عاملهای درونی سازمان هستند. هر کدام از گره های نشان داده شده می توانند شامل زیر لایه ها و عاملهای دیگری بوده و خودشان مدیریت بخش را بر عهده داشته باشند. پیوندهای رسم، شده ارتباط و وابستگی اطلاعاتی بخشهای مورد نظر را نشان میدهد. این مجموعه از عاملها علاوه بر اینکه خود ساختار دانش مستقلمی دارند با دسترسی به پایگاههای اطلاعاتی مجاز، سازمان به انجام وظیفه میپردازند گره نشان داده شده در قطر مکعب نشان دهنده عامل مدیریتی سازمان است که در لایه دوم ساختار پیشنهادی قرار دارد.

● لایه اطلاعاتی سازمان: این لایه شامل سیستمهای اطلاعاتی سازمان، اعم از پایگاه داده اطلاعات کاربران و نیز اطلاعات و اسناد سازمان خواهد بود عاملهای مختلف با رعایت قوانین دسترسی مجاز به استفاده از این اطلاعات خواهند بود عاملهای مختلفی مسئولیت مدیریت این سیستمهای اطلاعاتی را بر عهده دارند وظیفه این عاملها کنترل دسترسیهای مجاز به این اطلاعات و نیز ایجاد هماهنگی در ورودی و خروجی پایگاههای داده است. در صورتی که پایگاههای داده توزیع شده، باشند عامل مرتبط با آن وظیفه برقراری یکپارچگی بین اطلاعات آنها را بر عهده خواهد داشت.

در مدل مطرح شده عامل برنامه ریزی ارشد مسئول شکست وظیفه محول شده به زیر وظایف و اعلام آن به عاملهای مدیریت است عاملهای مدیریت مبادرت به تعیین وظایف بخشهای مختلف سازمان در راستای وظیفه انتسابی از سوی عامل برنامه ریزی ارشد نموده و زمان بندیها را به این بخش اعلام می کنند.

برقراری ارتباطهای بین سازمانی به وسیله عاملهای لایه سوم ممکن می شود. عامل واسط مسئول برقراری ارتباط در این لایه خواهد بود. این عامل نیز با دریافت سیاستهای مذاکره از سوی عامل مدیریت با همتای خود در سازمان دیگر و یا با کاربر انسانی مذاکره میکند به منظور ایجاد درک صحیح از تعاملات انجام شده لازم است عاملهای مذاکره کننده در سازمانهای مختلف هستان شناسی یکسان و ساختار ارتباطی منطبق بر استانداردهای تعریف شده داشته باشد.

همانند لایه سوم برقراری ارتباط توسط بخشهای مختلف یک سازمان با بخشهای داخلی سازمانهای دیگر توسط عاملهای واسط در لایه دوم ممکن می شود. تعیین سیاستهای ارتباطی عامل واسط، بر عهده عامل مسئول بخش مرتبط با او است نکته قابل توجه این است که قوانین تعیین شده توسط عاملهای بخشهای مختلف در راستای سیاستهای تعیین شده توسط عامل مدیریت است به دلیل نزدیکی این سطح به منابع اطلاعاتی و دسترسی به جزئیات برنامه های سازمان، تعاملات انجام گرفته در این لایه با کنترل و اعمال قوانین امنیتی بیشتری صورت می گیرد. تعاملات در سطح مدیریت ارشد در این چارچوب بر عهده عاملهای دو لایه مختلف قرار داده شده است. با توجه به اینکه دیدگاه زنجیره ارزش نگاهی راهبردی است و ممکن است سازمان به دلیل پویایی و چابکی در سطح راهبرد با تغییراتی در این سطح مواجه شود ملاحظات درون و برون سازمانی به عاملهای مختلف در لایه های متفاوت واگذار شده است. در چارچوب پیشنهادی در این مقاله آن بخش از رویکردها و تصمیمات حوزه مدیریت ارشد که بر اساس راهبردها و مقررات موجود محقق می شود به

عاملهای لایه مدیریت سازمان سپرده شده است عاملهای لایه زنجیره ارزش تعاملات بیرونی زنجیره (عمدتاً جریان بالایی و پایینی) را با توجه مجموعه قوانین و اهداف دراز مدت محقق مینمایند و مرجعی برای یکپارچگی میان عملکرد عاملها در لایه های پایین تر خواهند بود.

چارچوب مطرح شده ضمن فراهم آوردن زیر ساخت یکپارچه برای سازمان و زنجیره ارزش، قابلیت برقراری ارتباط با دیگر سازمانها را نیز فراهم می کند. تنها شرط موفقیت تعاملات برای هر دو سازمان وجود هستان شناسی مشترک و استفاده از استانداردهای تعاملی است. در این چارچوب با برقراری یکپارچگی در سازمان علاوه بر ایجاد دسترسی به اطلاعات برای بخشهای مختلف، با حفظ سطوح دسترسی، امنیت اطلاعات نیز تأمین می شود. ساختار لایه ای ضمن اعمال قوانین مدیریتی و کنترل بر تعاملات لایه های پایینتر جهت تصمیم گیرها را به سرعت از بالاترین سطح در زنجیره ارزش به پایین ترین سطح سازمان هدایت میکند.

1-4. مثالی از به کارگیری چارچوب پیشنهادی در سامانه تعامل پذیری مدیریت محتوا

همانطور که در مقدمه ذکر شد، امروزه یکی از دغدغه های خدمات محتوایی امکان برقراری ارتباط میان سازمانها و ذینفعانی است که از گونه های مختلف مدیریت محتوا و مخازن اطلاعاتی بهره میبرند سامانه تعامل پذیری مدیریت محتوا (1) (CMIS) استاندارد بازی است که در یک سطح انتزاع بالا جهت ایجاد ظرفیت ارتباط و تعامل مناسب میان بنگاههای مختلف محتوایی ایجاد شده است. این استاندارد مورد قبول و پشتیبانی کنسرسیوم استانداردهای وب OASIS نیز میباشد (OASIS Committee 2010). هدف اصلی از ارائه این استاندارد ارائه ویژگی اختصاصی برای سامانه های مدیریت محتوا یا ارائه یک معماری جامع برای ارتباط انواع مخازن اطلاعات و مانند آنها نیست بلکه هدف ایجاد ظرفیتی است که سامانههای مختلف مدیریت محتوا و مخازن اطلاعات در سطح وب بتوانند ضمن تعامل از امکانات و ظرفیتهای یکدیگر استفاده نمایند این استاندارد از یک هسته مدل داده جهت تعریف هسته های اطلاعاتی موجود و یک مجموعه خدمات پایه تشکیل شده است.

در اینجا به منظور بررسی قابلیت های چارچوب پیشنهادی نحوه انطباق آن در بافتار ذخیره سازی و ارائه خدمات محتوایی بررسی شده است. با عنایت به این واقعیت که موجودیتهایی که توسط CMIS مدیریت می شوند در قالب انواع اشیاء بیان میشوند برای انطباق چارچوب پیشنهادی در این مقاله و CMIS مدل اشیاء این استاندارد مدنظر قرار گرفته است. بدین منظور نحوه انطباق به شرح زیر خواهد بود:

مدل سازی اشیاء سیاستهای (2) مدیریتی در لایه زنجیره ارزش سیاست گذاری و کنترل بر اشیاء اصلی توسط عاملهای این لایه انجام میشود. این اشیاء میتوانند بطور توزیع شده در سازمانهای مختلف قرار داشته و از طریق Object Identity تعریف شده در سرویسهای CMIS و مبتنی بر دانش سیاستگذاری موجود در عاملهای این، لایه مورد مدیریت قرار گیرند.

ص: 65

مدل سازی اشیاء ارتباطی (1) در لایه مدیریت سازمان این اشیاء در استاندارد CMIS با هدف ایجاد ارتباط میان سایر اشیاء تعریف شده‌اند و سرویسهای مربوط به آنها به امر یکپارچگی و ارتباط میان سازمانها کمک میکنند عاملهای این لایه میتوانند ضمن ایجاد یکپارچگی در سیاستهای، مدیریتی امکان مدیریت بر سایر عاملها در لایه های دیگر را هم بر عهده گیرند. ارجاع در خواست مستندات و سایر اشیاء دیجیتال از سازمانهای مختلف براساس دانش این عامل مسیریابی و رهگیری میشود.

مدل سازی اشیاء پوشه (2) در لایه درونی سازمان اشیاء پوشه در استاندارد CMIS با هدف ارائه اقدامات و تصمیمات منطقی مورد نظر در خصوص اشیاء دیجیتال تعریف شده‌اند این تصمیمات منطقی میتوانند به انواع پردازشها برنامه ریزی تصمیم گیری سنجش و رویکردهای دیگر مدیریتی بر اشیاء دیجیتال اختصاص . یابند همان طور که در شکل نشان داده شده است عاملهای این لایه به خوبی میتوانند پشتیبان ملزومات این اشیاء باشند.

مدل سازی اشیاء مستندات و سرمایه های دیجیتال (3) در لایه اطلاعاتی سازمان اشیاء دیجیتال مورد پشتیبانی در این استاندارد با مشخصات و خصیصه های مختلفی شناسایی میشوند. نکته حائز اهمیت امکان وجود این اشیاء در سامانه های مختلف مدیریت محتوا و مخازن گوناگون دیجیتال است که در سازمانها و بنگاههای مختلف و بطور توزیع شده وجود دارند عاملهای این لایه میتوانند به شکل خود مختار و هوشمند مراحل پیش پردازش جمع آوری، ذخیره سازی و سایر پردازشهای پسین را به سامانه مدیریت یکپارچه محتوا اعمال نمایند به عنوان مثال "Get Repositories Information" سرویسی است که توسط عاملهای این لایه جهت شناسایی اطلاعاتی نظیر نام ارائه کننده سرویس، محتوایی نام محتوا، نسخه محتوا مکانهای ذخیره سازی محتوا و مانند آن را در اختیار لایه درونی قرار میدهد تا برنامه ریزی لازم در خصوص آن انجام شود.

5. نتیجه گیری

گسترش منابع مختلف توسعه محتوا و خدمات محتوایی تحت وب امروزه در وضعیتی قرار گرفته است که مرز میان اجزای مختلف زنجیره ارائه خدمات کمی غیر شفاف شده است به طوریکه ضمن از بین رفتن فاصله میان مصرف کننده ها توزیع کنندهها تأمین کننده ها و تولیدکنندگان؛ حوزه های مرتبط با این اجزا و به خصوص در سطح عملکرد خدمات تحت وب گسترش یافته و به ازای سطوح دسترسی ذینفعان قابل بهره برداری شده است. حجم زیاد تنوع و ساختار متفاوت اطلاعات سازمانها و اطلاعات موجود در اینترنت باعث شده است که فرایند مدیریت بازاریابی و استخراج آنها اهمیت ویژه ای پیدا کند. از طرف دیگر افزایش سیستمهای اطلاعاتی سازمانها و توزیع شدن بخشهای مختلف زنجیره ارزش

ص: 66

Relationship Object -1

Folder Object -2

Documents and Digital Assets -3

باعث پیچیده و پویا شدن روابط شرکا و عدم قطعیت در زنجیره ارزش شده است. اگرچه سیستمهای نرم افزاری با تغییر و تسریع فرایندها به بالابردن بازدهی زنجیره ارزش کمک کرده اند، با این حال پویایی این زنجیره و گردش اطلاعات مختلف پیش بینی نشده باعث ایجاد مشکلاتی در این امر شده است. عامل های هوشمند با ایجاد یکپارچگی در زیر ساختهای سازمانها و زنجیره ارزش امکان انجام بهینه مدیریت برنامه ریزی، زمان بندی کنترل و مدیریت عدم قطعیتها را فراهم آورده و بازدهی آنها را بالا میبرند با تعریف هسته شناسی مشترک برای عاملهای تعاملی سازمانها انجام ارتباطهای بین سازمانی نیز به وسیله عاملهای هوشمند ممکن شده و با بهبود کیفیت انعطاف و پویایی در تعاملات خارجی، سازمان بازدهی و تأثیر آن بالا رفته و به مزیت ماندگار منجر می شود.

در این مقاله چارچوبی مبتنی بر عاملهای هوشمند برای مدیریت اطلاعات سازمانها پیشنهاد شده است. این چارچوب با در نظر گرفتن معماریهای مطرح نرم افزار ساختارهای سازمانی و چگونگی تعاملات بخشهای مختلف سازمان؛ یکپارچگی و پویایی را در سطح زنجیره ارزش برقرار کرده و تعاملات سازمان با شرکا را نیز محقق می کند چارچوب پیشنهادی میتواند ضمن حفظ سطوح دسترسی برای بخشهای مختلف سازمان سازگاری هماهنگی و امنیت اطلاعات را نیز تأمین کند. عامل برنامه ریزی ارشد در زنجیره تأمین با اشراف بر سازمانهای لایه، پایین ضمن دریافت درخواستهای ارائه شده آن را به بخشهای مختلف تقسیم کرده و با کنترل زمانبندی باعث کم شدن زمان تحویل درخواست و افزایش کیفیت آن خواهد شد. ویژگی دیگر چارچوب پیشنهادی، پشتیبانی و ظرفیت سازی برای تعاملات سازمانی است. وجود عاملها در هر لایه و همچنین رویکرد پیشنهادی در این مقاله سبب شده است که تعاملات سازمانی در سطوح مختلف و با حفظ ملاحظات خاص درون هر سازمان و با ایجاد استقلال لازم به انجام برسند این قابلیت سبب میشود که چارچوب پیشنهادی به گسترش زنجیره ارزشی که متشکل از سازمانهای مختلف است بیانجامد.

منابع

Berre, A. et al. 2004. State-of-the art for interoperability architecture approaches, Model driven and dynamic, federated enterprise interoperability architectures and interoperability for .non-functional aspects. Information Society Technology

Cao, L., V. Gorodetsky, and P.A. Mitkas. 2009. Agent mining: The synergy of agents and data .mining. Intelligent Systems, 24:64-72

.Chu, B. et al. 1997. Towards intelligent integrated manufacturing planning-execution .International Journal of Advanced Manufacturing Systems, 1:77-83

David, C., G. Doumeingts, and F. Vernadat. 2008. Architecture for enterprise integration and .interoperability: Past, present and future. Computers in Industry, 59:647-659

Feng, S.C., K.A. Stouffer, and K.K. Jurrens. 2005. Manufacturing planning and predictive

- ,process model integration using software agents. *Advanced Engineering Informatics*, 19:135–142
- Forget, P., S. D'Amours, and J.M. Frayret. 2008. Multi-behavior agent model for planning in supply chains: An application to the lumber industry. *Robotics and Computer-Integrated Manufacturing*, 24:664-679
- Frey, D., T. Stockheim, P.O. Woelk, and R. Zimmermann 2003. Integrated multi-agent based supply chain management. In *Twelfth IEEE International Workshops on Enabling Technologies: Infrastructure for Collaborative Enterprises (WETICE'03)*, Washington
- Heikkila, T., M. Kollingbaum, P. Valckenaers, and G. Bluemink 2001. An agent architecture for manufacturing control: manage. *Computers in Industry*, 46:315-331
- .ISO 14258.1999. Industrial automation systems-concepts and rules for enterprise models .ISO TC184/SC5/WG1, April 14, 1999
- :Kishore, R., H.Zhang, and R.Ramesh.2006. Enterprise integration using the agent paradigm Foundations of multi-agent-based integrative business information systems. *Decision Support Systems*, 42:48–78
- Lim, M.K., D.Z. Zhang. 2004. An integrated agent-based approach for responsive control of manufacturing resources. *Computers Industrial Engineering*, 46:221-232
- Mahesh, M., S.K. Ong, A.Nee, J.Fuh, and Y.F. Zhang. 2007. Towards a generic distributed and collaborative digital manufacturing. *Robotics and Computer Integrated Manufacturing*, 23:267–275
- Maturana, F., W.Shen, and D.H. Norrie.1999. MetaMorph: An adaptive agent-based

,architecture for intelligent manufacturing. International Journal of Production Research

.37:2159–2173

-Navarretta, C., B.S.Pedersen, and D.H. Hansen. 2006. Language technology in knowledge

.organization systems. New Review of Hypermedia and Multimedia, 12:29–49

OASIS Committee. 2010.The CMIS v1.0 OASIS standard specification. Retrieved from

<http://docs.oasis-open.org/cmisis/CMIS/v1.0/os/cmisis-spec-v1.0.pdf>, [Accessed 18 Jan

2013]

Oztemel, E., E.K.Tekez. 2009.A general framework of a Reference Model for Intelligent

Integrated Manufacturing Systems (REMIMS). Engineering Applications of Artificial

.Intelligence, 22:855–864

.Shehory, O. 2006. The role of agents in enterprise system management: A position paper

ص: 68

,Fifth International Joint Conference on Autonomous Agents and Multi Agent Systems

.Hokkaido

Shen, W., Q. Hao, S.Wang, Y.Li, and H.Ghenniwa.2007.An agent-based service-oriented integration architecture for collaborative intelligent manufacturing. Robotics and .Computer-Integrated Manufacturing, 23:315-325

Shen, W., F.Maturana, and D.H. Norrie.2000. MetaMorph II: An agent-based architecture for ,distributed intelligent design and manufacturing. Journal of Intelligent Manufacturing .11:237-251

Shen, W., D. Xue, and D.H. Norrie. 1997. An agent-based manufacturing enterprise infrastructure for distributed integrated intelligent manufacturing systems. In Third -International Conference on the Practical Application of Intelligent Agents and Multi .agents, London

Wang, M., J. Liu, H.Wang, W.K. Cheung, and X. Xie.2008. On-demand e-supply chain ,integration: A multi-agent constraint-based approach. Expert Systems with Applications .34:2683-2692

Pawlewski, P., P.Golinska, M. Fertsch, J.A.Trujillo, and Z.J. Pasek. 2008. Multiagent approach for supply Chain integration by distributed production planning, scheduling and control system. International Symposium on Distributed Computing and Artificial .(Intelligence (DCAI

Wang, G., J. Zheng, H. Wu, and Y.Tang. 2009. Research of the enterprise application ,integration platform based on multi-agent. Fifth International Joint Conference on INC

-Wang, S., W. Shen, and Q. Hao. 2006. An agent-based web service workflow model for inter
.enterprise collaboration. *Expert Systems with Applications*, 31:787-799

Yang, L., X. Gen, and X. Cao. 2010. Application research of enterprise information integration
based on intelligent agent on ciagent platform. *Second International Workshop on
.Education Technology and Computer Science*, Wuhan

Zhang, C., Z. Zhang, and L. Cao. 2007. Agents and data mining: Mutual enhancement by
integration. *2nd International Conference on Autonomous intelligent systems: Agents
.and Data Mining*, Berlin

شبکه جهان گستر وب (WWW) شبکه ای از محتویات و ساختارهای پیوندهای تو در توست که دائماً در حال تکامل و تغییر است از این رو، یک آرشیویست هرگز ممکن نیست که از قوام و صحت محتویاتی که تاکنون جمع کرده با چیزی که بعداً نیاز پیدا خواهد کرد اطمینان حاصل کند بنابراین، پرسشهایی در مورد تشخیص و اندازه گیری این محتویات و در نهایت در مورد درک نقص در انسجام به وجود خواهد آمد به این منظور برخی راهبردهای نموداری ارائه شده‌اند که ممکن است در سطوح مختلف به کار روند کار با آخرین تغییرات زمانی صحیح، براساس فراداده های استخراج شده از خزشگرها یا از فایل‌های WARC. برای کمک به آرشیویست جهت درک ذات این نارساییها این مقاله به بررسی روشهایی برای نمایش رفتار تغییرات و انسجام آرشیوها می پردازد.

مارک اسپانیول (2) | آرتوراس مازیکا (3) | دیمیتار دنوو (4) | اگرهارد ویکوم (5)

ترجمه: عبدالله حسینیان (6)

مقدمه

اگر میتونی منو بگیر عنوان فیلمی براساس یک داستان واقعی درباره یک کلاهبردار معروف به نام فرانک ابا گنیل (7) است. شخصی که در نقش یک خلبان، دکتر و وکیل ظاهر میشود. این فیلم که سختی های دستگیر کردن یک کلاهبردار را در دنیای واقعی شرح میدهد می تواند با شکل مشابهی در آرشیو وب مقایسه شود.

دنیای جهانگستر وب (www)، میلیونها کاربر را قادر میسازد تا محتویات روی وب را تألیف، تغییر، یا حتی حذف کنند. درست مثل تعقیب یک کلاهبردار حفظ و جمع آوری این دادهها نیز کاری جزئی نیست و میتواند شامل موضوعهای کیفیت داده نیز بشود برای مثال یک سیستم مدیریت محتوا (cms) را در نظر بگیرد که وبگاه مؤسسه ای تحقیقاتی را نگهداری می کند. هر گاه دو محقق به طور مشترک

ص: 71

Defects in Web Archiving Catch me if you can": Visual Analysis of Coherence -1

Marc Spaniol -2

Arturas Mazeika -3

Dimitar Denev -4

Gerhard Weikum -5

6- کارشناس کامپیوتر سازمان اسناد و کتابخانه ملی

Frank Abagnale -7

مقاله ای را منتشر، کنند CMS به صورت خودکار مرجعی برای اتصال مقاله ها در صفحه اصلی دو محقق ایجاد میکند در این حین خزش ممکن است یکی از این صفحه ها را قبل از روز آمدسازی و صفحه دیگر را بعد از آن ملاقات کند. در این صورت آرشیو این صفحه ممکن است غیر منسجم پایان پذیرد. از این رو آرشیویست ممکن نیست که از قوام اطلاعاتی که اکنون جمع آوری کرده برای درخواستهای بعدی اطمینان حاصل. کند به این ترتیب خزش سایت باید برای جلوگیری از بارگذاریهای صفحه بی مورد بین درخواستهای HTTP مکث های قابل توجهی داشته باشد.

در نتیجه ثبات یک وبگاه عظیم ممکن است ساعتها یا حتی روزها به طول انجامد. تغییرات در

خلال این دوره زمانی به طور موقت غیر قابل دسترسی است.

1 - انسجام در واژه نامه آکسفورد چنین تعریف میشود: «عمل یا حقیقت به هم چسباندن» یا «نظم اتصال چندین قسمت یا بخش به منظور ایجاد تمامیت آنها با یکدیگر». در نتیجه زمانی نقص انسجام وجود دارد که بعضی از عناصر به شرایط انسجام حمله ور شوند در مورد آرشیو وب، انسجام، دارای بعد زمانی است محتویات بازه زمانی X یا زمان بین X و Y .

چیزی که به عنوان یک نیاز ساده ظاهر می شود به صورت پیچیده توسعه پیدا میکند و در نهایت غیر ممکن شدن عمل آرشیو را به همراه دارد چون نشر دهندگان نمیتوانند از همه وبگاه خود کپی تهیه کنند، که به صورت بخش بخش است و لزوماً همه با هم کار نمیکند این مسئله ضمانت کیفیت خزش را محدود میکند شکل 1 اشکال انسجام را در آرشیو وب به تصویر میکشد در این مورد نقص انسجام در مواقعی یک صفحه به صفحه دیگر رجوع میکنند که آن صفحه قبلاً در نسخه اخیر از اعتبار افتاده است، اتفاق می افتد. در این مورد سند های آرشیو شده در سمت چپ در ارتباط با صفحه ورودی غیر منسجم هستند که با چارچوب قرمز نمایش داده شده اند (با زمان رجوع 2007/2/17). به هر حال، پیوندها از صفحه ورودی به صفحه های سمت راست که در تاریخ 2007/2/19 آرشیو شده است - منسجم هستند (که توسط قاب بنر نمایش داده شده اند) چون هر دو صفحه از 2007/2/17 معتبر است و تغییری نکرده اند. اما تشخیص چیزی که به عنوان غیر منسجم برای انسان آسان است برای یک ماشین بسیار سخت خواهد بود. یک رایانه ممکن است تنها سطح محدودی از جنبه های زمانی یک صفحه را تفسیر کند.

با وجود این به دست آوردن آخرین تاریخ اصلاح به عنوان یک نقطه زمانی، ما را قادر می سازد در مورد نقص انسجام بین دو نمونه از یک سند تصمیم بگیریم به این منظور چندین تکنیک در سطوح مختلف جزئیات را برای تعیین نقطه زمان ویرایش محتویات معرفی خواهیم کرد.

تحقیقات بر روی تحلیل تصویری از اشکال انسجام در آرشیو وب در درجه تغییراتی که یا در زمان خزش یا در میان رشته ای از خزشهای سایت رخ خواهد داد کمک بسیاری خواهد کرد. بنابراین، دقت و قدرت تغییر خزش قابل سنجش است؛ دقیقاً همانند تعقیب یک کلاهبردار و ما قادر نیستیم که از تغییر شکل یک وبسایت جلوگیری کنیم اما قادر هستیم این تغییرات را مشخص و راهبرد خزش را تنظیم تا در آینده تا حد ممکن منسجم باشد برای فهم بهتر از نقص در انسجام و تغییرات در آرشیو وب ما تحلیل دادهها را با استفاده از چهار روش تصویری پیشنهاد میکنیم:

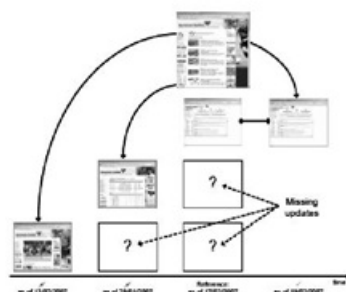
1) تصویر سازی تغییرات در درخت Spanning خزش با 2; visone) تصویر سازی Scatterplot در تحلیل محتوای تغییر؛ 3) منطقه مکانی تصویر سازی سری زمانی تغییرات؛ و 4) تصویر سازی Scatterplot سری زمانی تغییرات.

شکل 1. اشکال انسجام در آرشیو وب

تحلیل تصویری اجازه میدهد که نقص انسجام در ساختار و محتوا را بررسی کنید سطح تغییرات را به دست آورید و الگوی تغییرات را در میان تک تک صفحهها در توالی آرشیو وب کشف کنید. UKGOV, MPI از تحلیل تصویری استفاده میکنند UKGOV, MPI، شامل 120 خزش هفتگی از هفت سایت دولتی در انگلستان است. MPI شامل خزش روزانه از سازمانهای اساسی برای سایتهای اطلاع رسانی است.

عکس

(۱) تصویرسازی تغییرات در درخت Spanning خزش با visone ۲) تصویر سازی Scatterplot در تحلیل محتوای تغییر؛ (۳) منطبقه مکانی تصویرسازی سری زمانی تغییرات؛ و (۴) تصویر سازی Scatterplot سری زمانی تغییرات.



شکل ۱. اشکال انسجام در آرشیو وب

تحلیل تصویری اجازه می‌دهد که نقص انسجام در ساختار و محتوا را بررسی کنید، سطح تغییرات را به دست آورید، و الگوی تغییرات را در میان تک تک صفحه‌ها در توالی آرشیو وب کشف کنید. UKGOV, MPI از تحلیل تصویری استفاده می‌کند. شامل ۱۲۰ خزش هفتگی از هفت سایت دولتی در انگلستان است. MPI شامل خزش روزانه از سازمان‌های اساسی برای سایت‌های اطلاع رسانی است.

این مقاله، براساس روش زیر سازمان‌دهی شده است: ابتدا کارهای مرتبط انجام شده را بازبینی می‌کنیم. در بخش ۳ به بررسی تغییرات می‌پردازیم و مفهوم نقص انسجام را معرفی خواهیم کرد. در بخش ۴ در مورد موضوع جمع‌آوری استخراج و آماده سازی داده‌ها برای تحلیل نقص انسجام صحبت خواهیم کرد. بخش ۴ از تحلیل‌های تصویری را در مورد نقص انسجام معرفی می‌کند در نهایت، در بخش ۶ به نتیجه‌گیری کارهای آینده پرداخته خواهد شد.

۲- فعالیت‌های مرتبط

جامع‌ترین دید در مورد آرشیو وب توسط ماسانه ارائه شده است. این مقاله جنبه‌های مختلف در آرشیو وب را پوشش می‌دهد. همچنین فرآیند دسترسی‌های پشت سر هم را بررسی می‌کند. موضوع نقص انسجام نیز معرفی می‌شود، اما تنها بعضی از فرآیندهای کاوشی در مورد اینکه چگونه عمل آرشیو را توسط خزشگر اندازه‌گیری کنیم و ارتقا دهیم مورد بررسی قرار گرفته است. در این مورد، مفهوم انسجام زمانی به صورت جزئی تر بررسی شده است، اما ذات نقص انسجام در دنیای واقعی را بررسی نکرده است.

این مقاله براساس روش زیر سازمان‌دهی شده است ابتدا کارهای مرتبط انجام شده را بازبینی می‌کنیم. در بخش 3 به بررسی تغییرات می‌پردازیم و مفهوم نقص انسجام را معرفی خواهیم کرد. در بخش 4 در مورد موضوع جمع‌آوری استخراج و آماده سازی داده‌ها برای تحلیل نقص انسجام صحبت خواهیم کرد. بخش 4 از تحلیل‌های تصویری را در مورد نقص انسجام معرفی میکند در نهایت در بخش 6 به نتیجه‌گیری کارهای آینده پرداخته خواهد شد.

جامع ترین دید در مورد آرشیو وب توسط ماسانه ارائه شده است. این مقاله جنبه های مختلف در آرشیو وب را پوشش میدهد. همچنین فرآیند دسترس‌یهای پشت سر هم را بررسی میکند. موضوع نقص در انسجام نیز معرفی میشود اما تنها بعضی از فرآیندهای کاوشی در مورد اینکه چگونه عمل آرشیو را توسط خزشگر اندازه گیری کنیم و ارتقا دهیم مورد بررسی قرار گرفته است در این مورد مفهوم انسجام زمانی به صورت جزئی تر بررسی شده است اما ذات نقص انسجام در دنیای واقعی را بررسی نکرده است.

ص: 73

مک، کاون سیاستهای خزش را که میتواند برای یافتن سایتها استفاده شود ارزیابی کرده است (13,14). نتیجه این تحقیق نشان میدهد که نوع محتویاتی که قرار است دوباره برگردانده شود مهم است و چگونه یک وبگاه دوباره ساخته میشود، نویسنده ها همچنین تفاوت بین سایت اصلی و سایتی که دوباره ساخته شده است را بررسی میکنند به هر حال آنها راهبردهای ارزیابی های خودکار را برای منابع مبتنی متن پیشنهاد نمی کنند.

در این مقاله نویسندگان به کشف عمر محتویات صفحه ها و منظور از ساخت تاریخچه صفحه ها به صورت پویا در درخواست مشتری می پردازند به طور مشابه Nuneset dl تلاش میکند که سندهای وب را توسط تحلیل همسایه هایش تاریخ گذاری کند (17). هر دو مقاله، مکمل روشهای سنتی بر پایه سرایند HTTP headers HTTP یا تنها روی محتویات فراداده ها هستند به هر حال هیچ یک از دو روش تأثیر تغییرات را روی تمام انسجام سایت اندازه گیری نمیکند و شامل هیچ تفسیر تصویری از آنها نیستند.

سایر مقاله های مرتبط اغلب روی چیدمان خزشها برای سرعت و تأثیر بیشتر اندیسهای وب تمرکز میکنند. ساینکو، (1) تغییرات روی وبگاه را تحلیل میکند و اینکه چگونه باید اندیس گذاری شوند. موضوع اینکه چگونه عمل خزش را مؤثرتر انجام دهیم توسط چو و همکارانش (5) ارائه شده است.

آنها دلایلی را که نشان میدهند که طراحی یک خزش خوب تا چه اندازه مهم است (برای مثال، ترتیب و توالی نشانی جهت دیده شدن) و نوعی الگوریتم برای به دست آوردن صفحه های مرتبط ارائه می دهند. در مطالعات بعدی چو موهینا و گارسیا توسعه یک خزشگر افزایشی را شرح میدهند (2) آنها به ارتقای مجموعه بی تجربگی ها توسط صفحه های جدید در یک محدوده زمانی کمک می کنند. در همین مطالعات آنها روی سیاستهای تازه کردن مؤثر صفحهها میباشد (3). آنها یک فرآیند پواسون براساس مدل تغییرات منابع داده معرفی کردند در تحقیقی دیگر آنها توالی تغییرات را روی داده های برخط تخمین میزنند به این، منظور چندین تخمین زننده توالی را با هدف ارتقای خزشگرهای وب و کشهای وب بررسی میکنند در همین مسیر میتوان به تحقیقات اولستون و پندی (18) اشاره کرد که هدفشان تهیه یک جدول زمانبندی Recrawl بر اساس اطلاعات طول عمر به منظور به دست آوردن یک بازده خوب است ایپروتیس و همکارانش (10) تحلیل بقا را برای بررسی اطلاعات طول عمر به کار میبرند. آنها یک جدول زمانبندی به روزرسانی براساس رگرسیون طول عمر متناسب اختراع کردند. تان و همکارانش (20)، از نمونه برداری برای بررسی و پیش بینی به روزرسانی صفحه ها استفاده میکنند. آنها ویژگیهای انعکاسی ساختار پیوندها ساختار، و محتویات صفحه های وب را تعیین میکنند. راهبردهای بارگذاری انطباقی آنها بر اساس بررسی مجموعه صفحه هاست تحقیق دیگری درباره راهبردهای خزشگر توسط ناتورک و وینر (20) ارائه شده است آنها کشف کردند که پهنای بارگذاری صفحه ها در ابتدا بهتر است اما میانگین کیفیت آنها به مرور زمان کاهش مییابد بنابراین آنها اکیداً جست و جوی - breadth first را به منظور افزایش روی خط بودن برای فراخوانی صفحه ها مهم پیشنهاد کرده اند. تحلیل و فهم

نقص انسجام کاملاً متفاوت و مشکل تر است. ما تغییرات و نقص انسجام را به طور مناسب به تصویر میکشیم و برای تعیین صفحه ها و زیر گرافهای وب خزشهایی که باید در آینده تنظیم شوند، به شما کمک خواهیم کرد.

3- بررسی تغییرات و نقص در انسجام

انسجام یک نوع کاراکتر کیفیتی داده است. به عبارت دیگر در تنظیمات عمومی مجموعه ای از آیتمهای اطلاعاتی هیچ تضادی با محدودیتهای از پیش تعیین شده ندارند در سیستمهای پایگاه داده ای رایج زیر سیستم مدیریت تعاملات (1) مطمئن میشود که کیفیتهای الزامی داده بدون مشکل است. در سیستمهای توزیعی تک تک اجزا باید با یکدیگر همکاری کنند و از الگوریتمهای خاص برای اطمینان از این الزامات استفاده نمایند.

انسجام موضوع پیچیده ای در آرشیو وب است تولید کننده محتویات (ناشر) ممکن است اطلاعاتی را پست کند که با سیستم تضاد دارند برای مثال یک صفحه وب از یک مسابقه فوتبال به عکسی از مسابقه دیگری اشاره کند تهیه کنندگان محتویات وبگاهها تواناییهای محدود شده ای دارند و تمایل به همکاری دارند و منطق در cmsها نیز متفاوت است (صفحه به سرعت میتواند بروز شود در حالی که دیگران ممکن است در تغییرات تأخیر ایجاد کنند محتویات یکی ممکن است به صورت دینامیک ایجاد شود و دیگری خیر).

در این مقاله نقص در انسجام را از دیدگاه زمانی پیگیری میکنیم این مسئله، بخشی از زمان را برای آرشیو تمام وبگاه میگیرد در حالی که اگر آرشیو شامل نسخه صفحه هایی باشد که میتوانند نقطه ای از زمان دیده شوند یا بارگذاری شوند، ما میگوییم که آرشیو بدون نقص در انسجام است. یا اگر بهتر بخواهیم بگوییم اگر یکی از صفحه ها در خلال خزش تغییر کند هیچ ضمانتی وجود ندارد که آن صفحه بتواند در زمان دیگری دیده شود و ما در این مواقع میگوییم که نقص در انسجام وجود دارد.

به منظور اثبات نقص در انسجام دو نمونه از محتویات هم تاریخ و هم محتوا را بررسی می کنیم. به طور رسمی برای ثبت زمان یک صفحه وب از آخرین ویرایش سرانند HTTP استفاده می شود، که متأسفانه غیر قابل اعتماد است (11.6) (Cf) به همین دلیل از روش تاریخ گذاری دیگری بهره می گیریم که تاریخ گذاری معنایی محتوا نام دارد این تکنیک ممکن است یک روش تاریخ گذاری کلی باشد برای نمونه، اولویت تاریخ توسط آخرین ویرایش در پاورقی صفحه وب قرار بگیرد) یا به صورت مجموعه ای از تاریخ تنها روی تک تک آیتمهای صفحه قرار بگیرد (مثل، داستانهای خبری، پستهای بلاگها، کامنتها).

به هر حال استخراج زمان به صورت معنایی مستلزم یک برنامه کاوشی است که در مواردی که به عدم قطعیت در مورد زمان میرسیم به کار برده شود. در نهایت پرهزینه ترین اما صددرصد مطمئن روش مقایسه صفحه ها با نسخه قبلی بارگذاری شده آن است به دلیل پرداخت هزینه ها و دلایل، مؤثر یک

ص: 75

روش چند مرحله ای قوی را دنبال میکنیم:

1. کنترل تاریخ ضمیمه (Time Stamp HTTP) اگر ارائه شود و قابل قبول باشد در این مرحله متوقف می شود؛

2. کنترل تاریخ ضمیمه محتویات: اگر تاریخ ارائه شده مطمئن و قابل قبول باشد در این مرحله متوقف می شود؛

3. مقایسه مجموعه صفحهها با مجموعه قبلی که بارگذاری شده است؛

3. حذف تفاوت های کم اهمیت

تنها مجموعه متنهای محتوا یا متنهای مفید محتوا؛

مقایسه توزیع آن - گرام؛ و

محاسبه مقصد ویرایش از نسخه قبلی.

بر اساس این تکنیکهای تاریخ گذاری قادر هستیم که راهبردهای ارتقای انسجام را گسترش دهیم که به ما اجازه میدهند اطلاعات وابسته به زمان را با چندین خزش یا چندین آرشیو وفق دهیم.

4- استخراج و تهیه داده

این بخش در مورد جمعآوری، استخراج و تهیه داده برای تحلیل نقص انسجام و کشف تغییرات در آرشیو وب صحبت میکند.

الگوهای تصویری و تحلیل نقص انسجام نیازهای مختلفی را برای ورود داده ها در نظر می گیرد. در ساده ترین نوع یک تحلیل ممکن است به صفحه های آرشیو شده یک سایت نیاز پیدا کنیم، در حالی که تحلیلهای وسیع تر ممکن است تغییرات پویا را هم برای صفحه ها محتویات و هم پیوندهای (ساختار) یک سایت بررسی کند.

در این بخش ما یک راهنمایی که یک شمای بانک اطلاعاتی چگونه باید باشد به شما میدهم (بخش 4-1) و اینکه چگونه دادهها را با Standard SQL وارد یا پاک کنیم.

1-4- شمای بانک اطلاعاتی

به طور خاص شمای بانک اطلاعاتی شامل صفحهها (Cft_Pages در تصویر 10) و پیوندهای (Cft Links در تصویر 10) مرتبط با هم است. اطلاعات صفحه های مرتبط با یک صفحه در یک وبگاه شامل نشانی اندازه حالت کد و آخرین زمان ویرایش است. به علاوه URL را با (Cf.t_urls URL-id) کدگذاری میکنیم به این ترتیب، سریع تر؛ به طور مؤثر انتخابهایی از صفحه های مختلف با شماره سایت خاص و Crawl_id خواهیم داشت و میتوانیم چک کنیم که آیا صفحه در دو خزش پشت سر هم تغییر کرده است یا خیر سپس باید به طور مؤثر به پاک کردن دادههای تکراری پردازیم. محتویات صفحههایی که تغییر کرده اند در صفت content برای مقایسه با خزش قبلی ذخیره میشود (Cfvs_Page_idSection).

همچنین اطلاعات پیوندها در جدول t_link ذخیره میشود. هر دو صفت From_url_id و To_url_id تمامی پیوندهایی را که از یک صفحه به صفحه ای دیگر وجود دارد برای خزش تعیین میکنند برای جست و جوی ترتیبی وب ارشیو Parent_Page_id در جدول T_Page قابل دسترسی است.

دو درخت چند بعدی B روی صفت های Crawl_id, Url_id, Site_id از جدول Tpage- و (From_Site_id, To_Site_id, Crawl_id, From_url_id, To_url_id, Visited_TimeStamp) برای T_Link و همچنین یک صفت درخت B برای کلیدهای اصلی وجود دارد. این اصل ساده اما مؤثر در سازماندهی اطلاعات باعث واکنشی بسیار سریع دادهها و گزینه ها برای دریافت خزش و شماره (سایت در محاسبه نقص انسجام میشود (بخش 5)

2-4- ورود اطلاعات از فایل WARC ورود اطلاعات از ARC و WARC فایلها اصولاً شامل دو وظیفه است:

1) بارگذاری کردن داده ها در بانک اطلاعاتی و (2) حذف اطلاعات تکراری از بانک (ARC) و در حالت موفقتر آن (WAR) استاندارد کاربردی در آرشیو کردن وبها هستند. آنها برای ذخیره صفحههای آرشیو شده فرادادهای صفحه هایی مثل (نشانی زمان بارگذاری) و کنترل صفحه وب استفاده می. شوند متأسفانه فرمت ARC و WARC اطلاعات مربوط به پیوندهای وبگاهها را پشتیبانی نمی کنند ما این اطلاعات را از DAT فایلها به دست میآوریم یا به طور راحت تر آنها را توسط (HERITRIX) در حین آرشیو کردن و استخراج نشانها ایجاد میکنیم. اگر ARC و WARC قابل دسترسی باشند ساختار پیوند بین صفحه ها میتواند با کمک استخراج نشانی توسط Heritrix از صفحه های HTML آرشیو شده دوباره ایجاد شود.

داده های وب آرشیوها قبل از تحلیل نقص انسجام نیاز به تمیز شدن دارند. رایج ترین مشکل در اینجا بارگذاری چند باره یک صفحه / URL است این، عمل به چند دلیل اتفاق میافتد بعضی از صفحهها به دلیل سیاستهای بارگذاری وبگاه چند دفعه بارگذاری شده اند مثل (Robot.txt)؛ زمانی که بعضی از نیازها در زمان بارگذاری قابل دسترسی نیستند یا به دلیل اینکه آرشیویست برای ارتقای کیفیت پوشش یا کیفیت سایت ممکن یک صفحه را چندین بار بارگذاری کند حتی اگر بخواهیم دلایل بیشتری هم میتوانیم بیاوریم برای مثال صفحهها میتوانند به خاطر فرمت نشانی خود چندین بار بارگذاری شود (اگر) وب سرور حروف بزرگ و کوچک در نشانی تشخیص ندهد به خصوص که طراحان صفحه ها تمایل دارند هم از حروف کوچک و هم از حروف بزرگ برای نام گذاری فایلها و مسیرها استفاده کنند) به همین دلیل، قسمت اعظمی از سایت میتواند چندبار خزش شود. حذف اطلاعات تکراری به دلیل تغییر سندها و پیچیده شدن تحلیل تاریخچه این تغییرات برای صفحه ضروری است.

کد SQL برای حذف این داده های تکراری در ضمیمه موجود است این الگوریتم چند زمان بالاتر را در همه گروههای صفحه ها که دارای نشانی یکسان هستند تعیین میکنند (Cf.Line).

حذف دادههای تکراری از فرمت نشانی نیز به طور مشابه انجام میشود همه نشانها باید به صورت

کوچک و گروه شده توسط نشانیهای یکسان و نشانیهایی با تاریخ جدیدتر گرفته شود. به هر حال، این راه نیز به دلیل مقایسه رشته ها هزینه بر است. به جای آن ما گروهی روی Id های Url ها (Cf AppendixLine 1-AIN Listing) تأسیس میکنیم و کوچکترین مقدار هر گروه را محاسبه می کنیم (Cf.Line 10-33).

3-4- جمع آوری دادهها با Heritrix

جمع آوری دادهها توسط Heritrix

در به دست آوردن داده هایی که باید در بانک به طور مستقیم از خزشگر ذخیره شود به جای فرآیند وقت گیر فایل های WARC استفاده میشود حتی اطلاعاتی از جمله مسیر برای بعضی از صفحه ها به طور مستقیم از خزشگر استخراج میشود اما نیاز به دوباره سازیهای پیچیده فایل های WARC دارند به علاوه ما نوعی مکانیسم Crawl-Revisit را به منظور کاهش زمان برای تحلیل انسجام توسعه دادیم به صورت تکنیکی انسجام زمانی ما به نسخه ویرایش شده از خزشگر Heritrix بانک اطلاعاتی همراه آن و یک تحلیل و یک محیط تصویری تقسیم میشود. در درون بانک اطلاعاتی فرادادهها و دادههای استخراج شده توسط Heritrix ذخیره میشود و به علاوه، از مکانیسم Crawl-Recrawl نوعی راهبرد مؤثر برای انجام دوباره خزش استفاده می شود و اجازه می دهد که محتویات را بعد از کامل شدن عمل خزش دوباره تست کنیم.

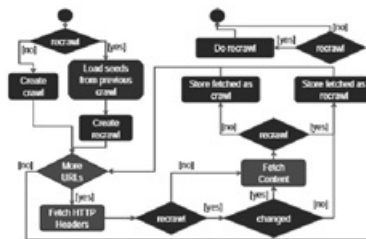
عکس

کوچک و گروه شده توسط نشانی‌های یکسان و نشانی‌هایی با تاریخ جدیدتر گرفته شود. به هر حال، این راه نیز به دلیل مقایسه رشته‌ها هزینه بر است. به جای آن، ما گروهی روی Idهای Urlها (Cf. Appendix Line 1-11N Listing) تأسیس می‌کنیم و کوچک‌ترین مقدار هر گروه را محاسبه می‌کنیم (Cf. Line 10-33).

۳-۴- جمع‌آوری داده‌ها با Heritrix

جمع‌آوری داده‌ها توسط Heritrix در به دست آوردن داده‌هایی که باید در بانک به‌طور مستقیم از خزشگر ذخیره شود به جای فرآیند وقت‌گیر فایل‌های WARC استفاده می‌شود. حتی اطلاعاتی از جمله مسیر برای بعضی از صفحه‌ها به‌طور مستقیم از خزشگر استخراج می‌شود، اما نیاز به دوباره‌سازی‌های پیچیده فایل‌های WARC دارند. به علاوه، ما نوعی مکانیسم Crawl-Revisit را، به منظور کاهش زمان برای تحلیل انسجام، توسعه دادیم. به صورت تکنیکی انسجام زمانی ما به نسخه ویرایش شده از خزشگر Heritrix، بانک اطلاعاتی همراه آن و یک تحلیل و یک محیط تصویری تقسیم می‌شود. در درون بانک اطلاعاتی فراداده‌ها و داده‌های استخراج شده توسط Heritrix ذخیره می‌شود و به علاوه، از مکانیسم Crawl-Recrawl نوعی راهبرد مؤثر برای انجام دوباره خزش استفاده می‌شود و اجازه می‌دهد که محتویات را بعد از کامل شدن عمل خزش دوباره تست کنیم.

به این منظور، Getهای شرطی را به کار می‌بریم که از محتویات Etagها استفاده می‌کند. در نتیجه، اعتبارسنجی با کاهش پهنای باند و بارگذاری سرور به‌طور همزمان بسیار سریع‌تر خواهد شد. بنابراین، تمام خزش‌ها از Crawl-Recrawlهای مجزا ساخته شده‌اند. البته خزش‌های دلخواه می‌توانند به عنوان محصولی از جفت Crawl-Recrawl ترکیب شوند که همان خزش تعریف می‌شود. شکل ۲ فلوچارت قسمت‌های اصلی انسجام زمانی در Heritrix را نشان می‌دهد. عناصر سبز رنگ شامل عناصری هستند که در مقایسه با Heritrix Crawler استاندارد تغییری نکرده‌اند. عناصر آبی رنگ روش خزشگرهای حاضر را که با راهبرد recrawl ما سازگار هستند را نمایش می‌دهد. و در نهایت واحدهای قرمز رنگ یک مرحله اضافه‌ای را که برای شروع Recrawlها نیاز است نمایش می‌دهد.



شکل ۲. فلوچارت قسمت‌های اصلی در انسجام زمانی در هر هیتریکس

به این منظور Getهای شرطی را به کار می‌بریم که از محتویات Etagها استفاده می‌کند. در نتیجه، اعتبارسنجی با کاهش پهنای باند و بارگذاری سرور به‌طور همزمان بسیار سریع‌تر خواهد شد. بنابراین، تمام خزش‌ها از Crawl-Recrawlهای مجزا ساخته شده‌اند البته خزش‌های دلخواه می‌توانند به عنوان محصولی از جفت Crawl-Recrawl ترکیب شوند که همان خزش تعریف می‌شود. شکل ۲ فلوچارت قسمت‌های اصلی انسجام زمانی در Heritrix را نشان می‌دهد عناصر سبز رنگ شامل عناصری هستند که در مقایسه با Heritrix Crawler استاندارد تغییری نکرده‌اند عناصر آبی رنگ روش خزشگرهای حاضر را که با راهبرد recrawl ما سازگار هستند را نمایش می‌دهد و در نهایت واحدهای قرمز رنگ یک مرحله اضافه‌ای را که برای شروع Recrawlها نیاز است نمایش می‌دهد.

شکل 2 فلوجارت قسمت‌های اصلی در انسجام زمانی در هر هیتریکس

ص: 78

تحلیل و محیط تصویری انسجام زمانی ما به عنوان وسیله ای برای اندازه گیری کیفیت یک جفت Crawl-Recrawl یا هر جفتی از دو خزش استفاده میشود. به این منظور آمار دادههای هر خزش (از جمله تعداد نقصهای رخ داده و ذخیره شده) از بانک اطلاعاتی مرتبط بعد از اجرای خزش محاسبه می شود. در این فرآیند سایت در یک درخت پوشای مشتق شده از مسیرهای خزش درون سایت نمایش داده میشود.

عکس

تحلیل انسجام و مصورسازی در آرشیو وب ۷۹

تحلیل و محیط تصویری انسجام زمانی ما به عنوان وسیله ای برای اندازه گیری کیفیت یک جفت Crawl-Recrawl یا هر جفتی از دو خزش استفاده می شود. به این منظور آمار دادههای هر خزش (از جمله تعداد نقصهای رخ داده و ذخیره شده) از بانک اطلاعاتی مرتبط بعد از اجرای خزش محاسبه می شود. در این فرآیند، سایت در یک درخت پوشای مشتق شده از مسیرهای خزش درون سایت نمایش داده می شود.

تصورسازی درخت پوشا، بینشی در مورد موقعیت و ذات تغییرات در محتویات وب در مقایسه با خزش قبلی به ما می دهد. به هر حال درختهای پوشا معمولاً بسیار بزرگ هستند و برای بسیاری از ابزارهای تصورسازی غیرعملی می باشند. برای برطرف کردن این مشکل تمرکز روی نقصهای درخت را فشرده می کنیم و تنها قسمت های مطرح را به صورت تصویر در می آوریم (cf. algorithm).

```
input : CrawlTree tree
begin
  collapseNode (tree.root)
  drawNode (tree.root)
end
```

Algorithm 1: processCrawlTree

الگوریتم ۱. فرآیند درخت خزش

در نخستین گام، محتویات درخت پوشا را تحلیل و آنها را طبق حالتشان تقسیم بندی می کنیم: سبز اگر بدون تغییر مانده باشند، زرد در مواردی که فقط عدم انسجام از نوع متنی باشد، قرمز در عدم انسجام از نوع ساختار (پیوندها) و در نهایت سیاه برای همزمانی که محتویات در خزش بعدی فراموش شده یا از بین رفته باشند استفاده می شود. در نهایت، راهبرد درهم کردن^۱ را به کار می بریم (cf. Algorithm^۳). گره های فشرده شده رنگ شده را به عنوان گره های بررسی شده می کشیم. به علاوه برای هر زیر درختی که فشرده شده است یک گره پایه به اندازه تعداد گره هایی که به این درخت وصل بوده است، رسم می کنیم.

```
input : Node node
begin
  node.collapsing=true
  if hasLinkChange (node) then
    node.color=red
  else if hasContentChange (node) then
    node.collapsing=false
    node.color=yellow
  else node.collapsing=false
  else node.color=green
  forall children of node do
    collapseNode (child)
    if child.collapsing=false then node.collapsing=false
    else node.collapsingSize=child.collapsingSize+1
  end
end
```

Algorithm 2: collapseNode

الگوریتم ۲. کنگره درهمکرد

1. Spanning tree
2. Collapsing

تصویر سازی درخت پوشا(1)، بینشی در مورد موقعیت و ذات تغییرات در محتویات وب در مقایسه با خزش قبلی به ما میدهد به هر حال درختهای پوشا معمولا بسیار بزرگ هستند و برای بسیاری از ابزارهای تصویر سازی غیر عملی میباشند برای برطرف کردن این مشکل تمرکز روی نقصهای درخت را فشرده میکنیم و تنها قسمتهای مطرح را به صورت تصویر در می آوریم الگوریتم 1. فرآیند درخت خزش

در نخستین گام محتویات درخت پوشا را تحلیل و آنها را طبق حالتشان تقسیم بندی میکنیم: سبز اگر بدون تغییر مانده باشند زرد در مواردی که فقط عدم انسجام از نوع متنی، باشد، قرمز در عدم انسجام از نوع ساختار (پیوندها) و در نهایت سیاه برای همزمانی که محتویات در خزش بعدی فراموش شده یا از بین رفته باشد استفاده میشود در نهایت راهبرد درهم کردن(2) را به کار میبریم (cf.Algorithm).

گرههای فشرده شده رنگ شده را به عنوان گره های بررسی شده میکشیم به علاوه برای هر زیر درختی که فشرده شده است یک گره پایه به اندازه تعداد گرههایی که به این درخت وصل بوده است رسم میکنیم.

الگوریتم 2. کنگره در همکرد

ص: 79

snanningtree -1

collapsing -2

برای یک نمایش گرافیکی قبلاً درخت محاسبه شده و در یک فایل graphML ذخیره شده است

(1) cf.Listing فایل graphML بر پایه استانداردهای XML و یک درخت برای گرافهاست این، فایل برای شرح تمام محاسبات قبلی سازگار است و در بسیاری از نرم افزارهای مرتبط با گرافها به کار می رود.

الگوریتم 3. گره خزش

عکس

برای یک نمایش گرافیکی، قبلاً درخت محاسبه شده و در یک فایل graphML ذخیره شده است (۱ cf. Listing). فایل graphML بر پایه استانداردهای XML و یک درخت برای گرافهاست. این فایل، برای شرح تمام محاسبات قبلی سازگار است و در بسیاری از نرم‌افزارهای مرتبط با گرافها به کار می‌رود.

```
<?xml version="1.0" encoding="UTF-8" standalone="no"?>
<graphml xmlns="http://graphml.graphdrawing.org/xmls/graphml" xmlns:xsi="http://www.w3.org/2001/XMLSchema-instance" xmlns:y="http://www.yworks.com/xml/graphml" xsi:schemaLocation="http://graphml.graphdrawing.org/xmls/graphml http://www.yworks.com/xml/schemas/graphml/1.0/ygraphml.xsd">
...
<graph edgedefault="directed" id="G229">
  <node id="http://www.mpi-inf.mpg.de/index.html">
    <data key="do">
      <y:ShapeNode>
        <y:Geometry width="10.003" height="10.003"/>
        <y:Fill color="#00FF00" transparent="false"/>
        <y:Shape type="ellipse"/>
      </y:ShapeNode>
    </data>
    <data key="di">http://www.mpi-inf.mpg.de/index.html O6</data>
  </node>
  <edge source="http://www.mpi-inf.mpg.de/index.html" target="docwww.mpi-inf.mpg.de"/>
</graph>
</graphml>
```

Listing 1: Coherence defect graphML-file (excerpt)

الگوریتم ۳. گره خزش

۵ - تحلیل نقص انسجام

تحلیل نقص انسجام، کیفیت یک خزش یا Crawl-recrawl، بین دو خزش یا یک سری از خزشها را اندازه‌گیری می‌کند. به این منظور، روشی را برای تولید آمارهای صحیح و تصویری توسعه داده‌ایم. برای مثال تعداد نقص‌های رخ داده و ذخیره شده توسط انواع کاستی‌ها.

۵-۱- تحلیل تغییر محتویات و ساختار با visone

همان‌طور که در بخش ۴-۳ شرح داده شد، گسترش Heritrix به ما اجازه می‌دهد تا فرآیند خزش را با داده‌های آماری دنبال کنیم و این داده‌ها را برای grapgML ارسال نماییم. با به کار بردن graphML نرم‌افزارهای مرتبط قادر خواهند بود درخت پوشا و ظاهر نقص در انسجام را نمایش دهند. این تصویرسازی‌ها به‌عنوان وسیله‌ای اضافی برای آمارهای خودکار جهت کشف مشکلی که در حین ثبت رخ می‌دهد در نظر گرفته می‌شود. قسمت اصلی این برنامه، تحلیل با کیفیت از خزش یک وبگاه است. شکل ۳ تصویر ساده‌ای از یک خزش از نشانی mpi-inf.mpg.de با نرم‌افزار ویژن را به تصویر کشیده است. بسته به اندازه گره‌ها، شکل و رنگ آنها کاربر یک دید کلی از فرآیندهای موفق و شکست خورده این جست‌وجو به‌دست خواهد آورد. به‌طور خاص اندازه یک گره مبنایی برای تعداد محتویات منسجم سایت (هر چه بزرگ‌تر انسجام در آن بخش بیشتر) در آن زیر درخت است. در همین شکل، رنگ یک گره حالت انسجام آنرا نمایش می‌دهد. جدی‌ترین نقص از دست دادن محتویات است که به رنگ مشکی نمایش داده شده. در نهایت شکل نت‌ها جنس نقص‌ها را نمایش می‌دهد مثلاً دایره

۵- تحلیل نقص انسجام

تحلیل نقص انسجام کیفیت یک خزش یا Crawl-recrawl بین دو خزش یا یک سری از خزشها را اندازه‌گیری میکند به این منظور روشی را برای تولید آمارهای صحیح و تصویری توسعه داده‌ایم. برای مثال تعداد نقصهای رخ داده و ذخیره شده توسط انواع کاستیها.

۱۵- تحلیل تغییر محتویات و ساختار با visone همان‌طور که در بخش ۴-۳ شرح داده شد گسترش Heritrix به ما اجازه می‌دهد تا فرآیند خزش را با داده‌های آماری دنبال کنیم و این داده‌ها را برای grapgML ارسال نماییم. با به کار بردن graphML نرم افزارهای

مرتبط قادر خواهند بود درخت پوشا و ظاهر نقص در انسجام را نمایش دهند. این تصویر سازه‌ها به عنوان وسیله‌ای اضافی برای آمارهای خودکار جهت کشف مشکلی که در حین ثبت رخ می‌دهد در نظر گرفته می‌شود. قسمت اصلی این برنامه تحلیل با کیفیت از خزش یک وبگاه است. شکل 3 تصویر ساده‌ای از یک خزش از نشانی mpi-inf.mpg.de با نرم افزار ویژن را به تصویر کشیده است. بسته به اندازه گره‌ها شکل و رنگ آنها کاربر یک دید کلی از فرآیندهای موفق و شکست خورده این جست و جو به دست خواهد آورد به طور خاص اندازه یک گره مبنایی برای تعداد محتویات منسجم سایت هر چه بزرگتر انسجام در آن بخش بیشتر در آن زیر درخت است. در همین شکل، رنگ یک گره حالت انسجام آن را نمایش می‌دهد جدی‌ترین نقص از دست دادن محتویات است که به رنگ مشکی نمایش داده شده. در نهایت شکل تنها جنس نقصها را نمایش می‌دهد مثلاً دایره

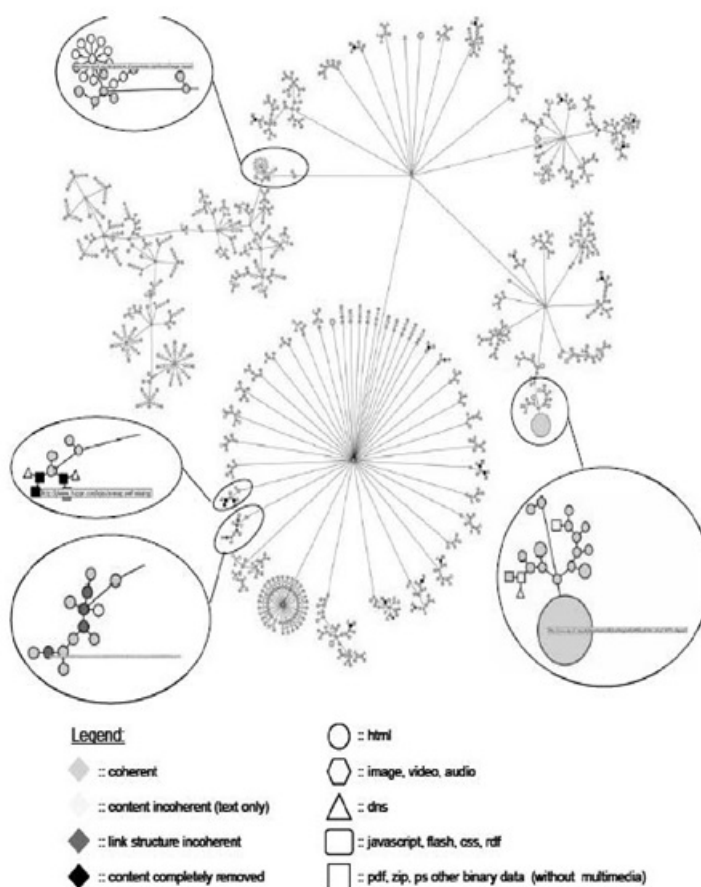
(HTML content) چند وجهی (محتویات چندرسانه‌ای) چار گوشه با گوشه های گرد فلش و مشابه آن مربع PDF و دیگر باینریها.

شکل 3

عکس

تحلیل انسجام و مصورسازی در آرشیو وب ۸۱

(HTML content) چند وجهی (محتویات چندرسانه‌ای) چار گوشه با گوشه های گرد فلش و مشابه آن، مربع PDF و دیگر باینریها.



شکل ۳

یک روش برای تحلیل یک جفت recrawl-Crawl تحلیل آن براساس زمان است. ایده بینایی این روش، به منظور دنبال کردن نقص در انسجام در میان چند خزش و برای تعیین محتویاتی است که کمتر تغییر می کنند و احتمالاً منسجم تر هستند. شکل ۴ تصویری از نقص انسجام شش خزش پشت سر هم روی نشانی dmoz.org/news است. هر جفت خزش های یک نقص انسجام درست شبیه مثال قبلی انجام می شود فقط برخلاف مورد قبل، اکنون ما خزش ها را به جای مقایسه با recrawl-Crawl

یک روش برای تحلیل یک جفت recrawl-Crawl تحلیل آن براساس زمان است. ایده بینایی این روش، به منظور دنبال کردن نقص در

انسجام در میان چند خزش و برای تعیین محتویاتی اس-----ت که کمتر تغییر میکنند و احتمالاً منسجم تر هستند شکل 4 تصویری از نقص انسجام شش خزش پشت سر هم روی نشانی dmoz.org/news است هر جفت خزشهای یک نقص انسجام درست شبیه مثال قبلی انجام میشود فقط برخلاف مورد قبل اکنون ما خزشها را به جای مقایسه با `recrawl-Crawl` با

ص: 81

خودشان مقایسه می‌کنند در انتقالهای هر دو تا از این جفتها همه گره ها مخفی میشوند، آنهایی که در تحلیل دچار نقص انسجام هستند ناپدید میشوند و در مقابل آن محتویاتی نمایش داده می‌شود که کاراکترهای منسجم یکسانی بین دو خزش حفظ کرده اند و گره هایی که جدید ظاهر میشوند اطراف آنها قرار می‌گیرند نکته جالبی که در این مثال دیده میشود این است که یک هسته محکم از یک زیر درخت بزرگ منسجم و محتویات غیر منسجم آن در اینجا وجود دارد.

شکل 4

5-2- تحلیل تغییرات محتوا

عکس

خودشان مقایسه می کنند. در انتقال های هر دو تا از این جفت ها همه گره ها مخفی می شوند، آنهایی که در تحلیل دچار نقص انسجام هستند ناپدید می شوند و در مقابل آن محتویاتی نمایش داده می شود که کاراکترهای منسجم یکسانی بین دو خزش حفظ کرده اند و گره هایی که جدید ظاهر می شوند اطراف آنها قرار می گیرند نکته جالبی که در این مثال دیده می شود این است که یک هسته محکم از یک زیر درخت بزرگ منسجم و محتویات غیر منسجم آن در اینجا وجود دارد.

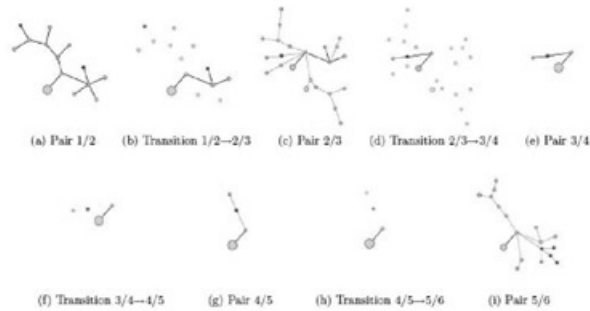


Figure 4: Tracing of coherence defects in crawl-recrawl pairs of the moz.org/news subdomain over time

شکل ۴.

۲-۵- تحلیل تغییرات محتوا

نمودارهای پراکندگی دوبعدی و سه بعدی می توانند برای تصویرسازی مکان و تحلیل تغییرات محتوا به کار روند. شکل ۵ و ۶ یک نمودار پراکندگی را برای سایت های *sabre* و *royal-navy* از بانک اطلاعاتی UKGOV نمایش می دهد، در اینجا اندازه و نشانی هر صفحه به روی محورهای *x,y,z* از سه بعد یک مکعب قرار گرفته اند. در همین حال، رنگ ها تغییراتی را که صورت گرفته اند نمایش می دهند (صفحه های جدید به رنگ آبی، صفحه هایی که تغییر کرده اند قرمز و صفاتی که تغییر نکرده اند به رنگ مشکی در آمده اند). آرشیویست باید الگویی از صفحه هایی که اضافه شده و تغییر کرده اند را بیابد. برای مثال از شکل ۵(a) یک نفر می تواند ببیند که چند تغییر در فایل های HTML در خروجی و زیرشاخه های متنی (cf. نطقه قرمز در شکل) و یک زیر مسیر جدید از فایل ها در آرشیو وب اضافه شده است (نقاط آبی در بالای تصویر). تغییرات صفحه ها (چهار نقطه قرمز) وابستگی های بین صفحه ها را نمایش می دهد. اگر صفحه ای در مسیر خروجی تغییر کند، صفحه منطبق با آن در آن مسیر نیز تغییر خواهد کرد. صفحه هایی که تازه اضافه شده اند نشان می دهند که ضوابط ساختاری سایت متحمل تغییرات می شوند. اگرچه محتویات سایت خیلی تغییر نکرده باشند.

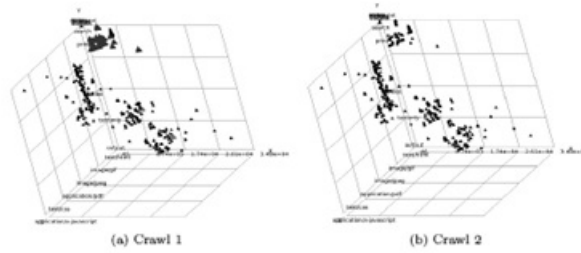
همچنین، الگوهای عکس های اضافه شده برای سایت *royal-navy* در شکل ۶(a)-۶(b) مشابه است. در حالی که، الگوی صفحه ها HTML اضافه شده و تغییر یافته اندکی متفاوت است که تغییر

نمودارهای پراکندگی دوبعدی و سه بعدی می توانند برای تصویرسازی مکان و تحلیل تغییرات محتوا به کار روند شکل ۵ و ۶ یک نمودار پراکندگی را برای سایت های *sabre* و *royal navy* از بانک اطلاعاتی UKGOV نمایش می دهد در اینجا اندازه و نشانی هر صفحه به روی محورهای *X,Y,Z* از سه بعد یک مکعب قرار گرفته اند. در همین حال رنگها تغییراتی را که صورت گرفته اند نمایش می دهند صفحه های جدید به رنگ آبی، صفحه هایی که تغییر کرده اند قرمز و صفاتی که تغییر نکرده اند به رنگ مشکی در آمده اند آرشیویست باید الگویی از صفحه هایی که اضافه شده و تغییر کرده اند را بیابد. برای مثال از شکل ۵(a) یک نفر می تواند ببیند که چند تغییر در فایل های HTML در خروجی و زیرشاخه های متنی (cf. نطقه قرمز در شکل) و یک زیر مسیر جدید از فایلها در آرشیو وب اضافه شده است نقاط آبی در

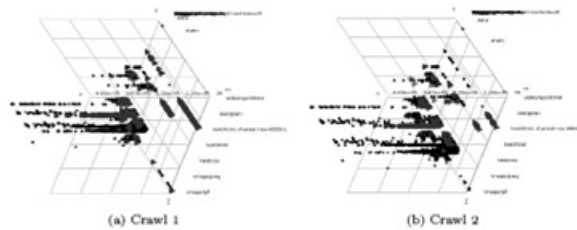
بالای تصویر) تغییرات صفحه‌ها (چهار نقطه قرمز) وابستگی‌های بین صفحه‌ها را نمایش می‌دهد. اگر صفحه‌ای در مسیر خروجی تغییر کند صفحه منطبق با آن در آن مسیر نیز تغییر خواهد کرد. صفحه - هایی که تازه اضافه شده اند نشان می‌دهند که ضوابط ساختاری سایت متحمل تغییرات می‌شوند. اگرچه محتویات سایت خیلی تغییر نکرده باشند.

همچنین الگوهای عکسهای اضافه شده برای سایت royal-navy در شکل 6(b)-6(a) مشابه است در حالی که الگوی صفحه‌ها HTML اضافه شده و تغییر یافته اندکی متفاوت است که تغییر

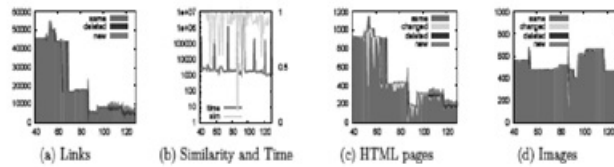
ساختار وبگاه را نشان می‌دهد.



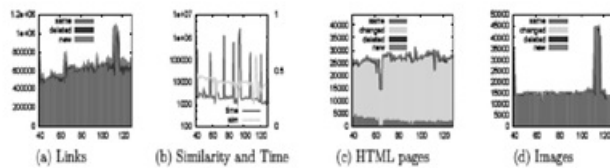
شکل ۵. تحلیل توالی خزش www.sabre.m.d.uk.sits



شکل ۶. تحلیل توالی خزش در www.royal-navy.mod.uk.site



شکل ۷. تحلیل توالی خزش در www.sabre.mod.us.site



شکل ۸. تحلیل توالی خزش در www.royal-navy.mod.uk.site

شکل ۵. تحلیل توالی خزش www.sabre.m.d.uk.sits

شکل ۶. تحلیل توالی خزش در www.royal-navy.mod.uk.site

شکل 7 . تحلیل توالی خزش در www.sabre.mod.us.site

شکل 8 . تحلیل توالی خزش در www.royal-navy.mod.uk.site

ص: 83

مجموعه نقشه های محیطی (شکل های 7 و 8) میتواند برای گرفتن یک دید کلی از درصد تغییرات در خزش آرشیو وبگاهها همان طور که `Crawl_id` افزایش مییابد استفاده میشود محور `X` ها شماره خزش و محور لها تعداد صفحه هایی که تغییر کرده اند/ نکرده اند/ اضافه شده اند حذف شده اند یا درصد زمان بارگذاری را نمایش میدهد. اصولاً شکل های مجزایی از ساختار پیوندها کشیده میشود (گراف ساختار سایت شکل (a)7-8(a) و صفحه های HTML شکل (c)7(a)) عکسهای شکل (d) و (d) تغییرات محتوا و به موازات آن کلیه اشکال نیز برای زمان بارگذاری مشابه هستند (شکل (b)8(b)-7(b)) بایگانی کننده باید الگوهایی را پیدا کند که باعث تغییرات قابل توجه در نقطه هایی از زمان میشود. برای مثال، یک شخص میتواند که تغییرات قابل توجهی در وبگاه در `Crawl52` و `93` ملاحظه کند. زمان خزش پیشنهاد میدهد که در `Crawl52` وبگاه متحمل تغییرات قابل ملاحظه ای میشود. به هر حال در `Crawl93` کیفیت آرشیو کاهش مییابد و عمل آرشیو لازم به رعایت مسائل خاص می شود.

محاسبه گراف محیطی میتواند در SQL هم بیان شود و توسط بهینه ساز query بهینه گردد (4) `cf.Listing` در نتیجه همه این اعمال پیچیدگی به دست آمده (`nlogn`) یا کمی بهتر خواهد بود. الگوریتم چندین سایت را از `Crawl xxx` و خزش قبلی آن `10 Lines`. با استفاده از یک `outer join` برای اتصال خزشها انتخاب می. کند فایل های که در یک یا دیگر خزشها هستند اما در هر دو وجود ندارند صفحه های جدید یا حذف شده میباشند در حالی که `tuples` که در نتیجه باقیمانده اند یا تغییر کرده اند و یا بدون تغییر مانده اند (10-15 `tuples of Lines`) اضافه شده حذف شده، تغییر یا تغییر نکرده همگی گروه میشوند (1-110 `cf.Lines`)

4-5- الگوهای تغییرات در صفحه های یک سایت

در شکل 9 محور `Y` ها صفحه ها و محور `X` ها شمار خزش را نمایش میدهد و نقطه تقاطع آنها نشان میدهد که اگر تغییری در صفحه وب در آن خزش وجود داشت با خزش قبلی مقایسه شود. صفحه ها بر روی محور `Y` ها براساس رفتارهای تغییر آنها که مشابه اند مرتب شده اند. شکل به آرشیویست وب برای پیدا کردن و تحلیل صفحه هایی که تغییرات مشابهی دارند کمک میکند الگوها و نقص در انسجام را کشف کند، تصویر به طور واضح صفحه های وبگاه را در بلوکهای مجزایی جدا میکند (مستطیلهای در شکل و الگوهای متفاوتی از تغییرات و نقص در انسجام را تعیین می کنند).

۸۵ تحلیل اسجام و مصورسازی در آرشیو وب

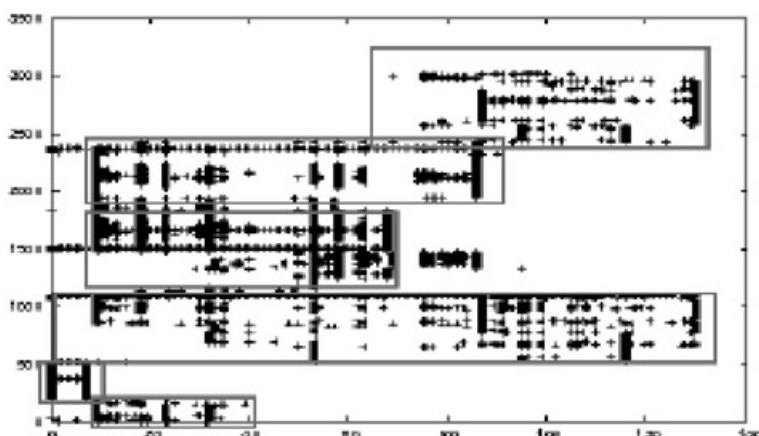


Figure 9: Scatterplot of lines for www.sabre.mod.

شکل ۹. الگوی تغییرات صفحه های یک سایت

```

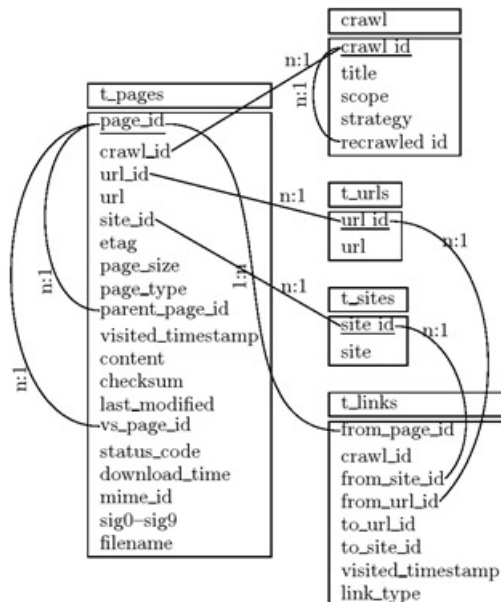
1 create table t_pages_dd as
2 select t_pages.* from (
3   select crawl_id, url_id, max(visited_timestamp)
4     as latest_timestamp
5   from t_pages
6   group by crawl_id, url_id) as x, t_pages
7 where t_pages.crawl_id = x.crawl_id and
8        t_pages.url_id = x.url_id and
9        t_pages.visited_timestamp
10       = x.latest_timestamp
11
12 create table t_links_dd as
13 select t_links.*
14 from t_pages_dd, t_links
15 where t_pages_dd.url_id = t_links.from_url_id and
16        t_pages_dd.visited_timestamp
17       = t_links.visited_timestamp
    
```

عکس

```

1 create table clean_mapping as
2 select dirty_url.url_id as dirty_url_id,
3       clean_url.url_id as clean_url_id
4 from (
5       select min(url_id) as url_id, lower(url) as url
6       from t_urls group by lower(url)
7 ) as clean_url, t_urls as dirty_url
8 where clean_url.url = lower(dirty_url.url);
9
10 create table lower_url_dd as
11 select crawl_id,
12        clean_mapping.clean_url_id as url_id,
13        lower(min(url)) as url, site_id as site_id,
14        min(etag) as etag,
15        min(page_size) as page_size,
16        min(page_type) as page_type,
17        min(visited_timestamp) as visited_timestamp,
18        min(checksum) as checksum,
19        min(last_modified) as last_modified,
20        min(status_code) as status_code,
21        min(download_time) as download_time,
22        min(sig0) as sig0, min(sig1) as sig1,
23        min(sig2) as sig2, min(sig3) as sig3,
24        min(sig4) as sig4, min(sig5) as sig5,
25        min(sig6) as sig6, min(sig7) as sig7,
26        min(sig8) as sig8, min(sig9) as sig9,
27        min(mime_id) as mime_id,
28        min(filename) as filename
29 from t_pages_dd, clean_mapping
30 where t_pages_dd.url_id
31       = clean_mapping.dirty_url_id
32 group by t_pages_dd.crawl_id, t_pages_dd.site_id,
33         clean_mapping.clean_url_id
    
```

فهرست ۳. SQL به کار رفته برای پاکسازی URL های پایین تر



نمودار ۱۰ شمای DB

6. درسهای آموخته شده و کارهای آینده از نظر یک آرشیویست آرشیو کردن مطلوب، وب جلوگیری از تغییرات محتویات در حین عمل خزش است. البته این یک توهم و عملاً نشدنی است در نتیجه ممکن است یک نفر هرگز مطمئن نشود که محتویات که تا کنون جمع کرده است هنوز با محتویاتی که بعداً جمع خواهد شد منطبق است. به هر حال انسجام در آرشیو وب یک موضوع کلیدی برای انسجام خزش جهت دادههای رقومی، در یک حالت قابل تکثیر و تفسیر است به این منظور ما گستره ای از Heritrix که با ارتباطات صحیح و همچنین تاریخ نامناسب محتویات سازگار است را توسعه داده ایم. به علاوه ما قادر هستیم شکل انسجام را مؤثرتر بدون توجه به تکیه بر وب سرور تحلیل کنیم به همین ترتیب توسعه تحلیل و تصویرسازی ویژگیها در کمک به مهندسان خزش برای درک بهتر ذات نقص در انسجام درون و بین وبگاهها و سازگاری راهبردهای Crawling برای خزشهای آینده مفید است در نتیجه مقاله به افزایش انسجام آرشیو هم کمک خواهد کرد.

در حالی که اکنون نقص در انسجام به ما در درک عدم تطابق به صورت سیستمی تر کمک میکند، تحقیقات آینده نیازمند یک بینش مولد و تولید کننده است. به علاوه، تحقیقات در حال پیشرفت به عمل خزش ناتمام و افزایش پوشش آرشیو کمک میکند حتی ترکیب بخشی از Recrawl در ترکیب با یک بخش افزایشی خزش ممکن است جذاب و مؤثر شود. به علاوه نتیجه به دست آمده از تحلیل خزشهای حقیقی برای ایجاد محیطهای شبیه سازی پیشرفته مفید خواهد بود همچنین قادر خواهیم بود رفتارهای تغییرات را در دنیای واقعی وبگاه در یک محیط شبیه سازی شده مشاهده کنیم.

سپاسگزاری

این کار توسط هفتمین برنامه FrameworkIST از E توسط تمرکز کوچکی یا متوسط پروژه های تحقیقی (STREP) روی وب آرشیوهای زنده (LiWA) با شماره 2162670 حمایت شده است. همچنین ما از همکارانمان برای مباحث امید بخششان متشکریم.

منابع

[1] Brian E. Brewington and George Cybenko. Keeping up with the changing web. Computer

52(5):33, May 2000.

[2] Junghoo Cho and Hector Garcia-Molina. The evolution of the web and implications for

an incremental crawler. In VLDB 400: Proceedings of the 26th International Conference

on Very Large Data Bases, pages 200-209, San Francisco, CA, USA, 2000. Morgan

Kaufmann Publishers Inc

[3] Junghoo Cho and Hector Garcia-Molina. Effective page refresh policies for web crawlers

ACM Transactions on Database Systems, 28(4), 2003

- .Junghoo Cho and Hector Garcia-Molina. Estimating frequency of change. ACM Trans [4]
Inter. Tech., 3(3):256{290, August 2003
- Junghoo Cho, Hector Garcia-Molina, and Lawrence Page. Efficient crawling through url [5]
ordering. In WWW7: Proceedings of the seventh international conference on World Wide
Web 7, pages 161{172, Amsterdam, The Netherlands, The Netherlands, 1998. Elsevier
Science Publishers B. V
- L. Clausen. Concerning etags and timestamps. In A. Rauber J. Masanés, editor, 4th [6]
International Web Archiving Workshop (IWAW'04), 2004
- Dimitar Denev, Arturas Mazeika, Marc Spaniol, and Gerhard Weikum. Share: Framework [7]
for qualityconscious web archiving. In VLDB '09: Proceedings of the 35th international
conference on Very Large Data Bases. VLDB Endowment, 2009
- .International Internet Preservation Consortium. Arc ia, internet archive arc le format [8]
<http://www.digitalpreservation.gov/formats/fdd/fdd000235.shtml>
- //:International Internet Preservation Consortium. Warc, web archive_le format. http [9]
www.digitalpreservation.gov/formats/fdd/fdd000236.shtml
- Panagiotis G. Ipeirotis, Alexandros Ntoulas, Junghoo Cho, and Luis Gravano. Modeling [10]
and managing changes in text databases. ACM Trans. Database Syst., 32(3):14, 2007
- Adam Jatowt, Yukiko Kawai, and Katsumi Tanaka. Detecting age of page content. In [11]
WIDM, pages 137{144, 2007
- .Julien Masanés. Web Archiving. Springer, New York, Inc., Secaucus, NJ, 2006 [12]
- Frank McCown and Michael L. Nelson. Evaluation of crawling policies for a web [13]

.repository crawler. In Hypertext, pages 157{168, 2006

Frank McCown, Joan A. Smith, and Michael L. Nelson. Lazy preservation: reconstructing [14]

.websites by crawling the crawlers. In WIDM, pages 67{74, 2006

G. Mohr, M. Kimpton, M. Stack, and I. Ranitovic. Introduction to heritrix, an archival [15]

.quality web crawler. In 4th International Web Archiving Workshop (IWAW'04), 2004

Marc Najork and Janet L. Wiener. Breadth-first search crawling yields high-quality [16]

.pages. In In Proc. 10th International World Wide Web Conference, pages 114{118, 2001

Sergio Nunes, Cristina Ribeiro, and Gabriel David. Using neighbors to date web [17]

.documents. In WIDM, pages 129{136, 2007

Christopher Olston and Sandeep Pandey. Recrawl scheduling based on information [18]

longevity. In WWW '08: Proceeding of the 17th international conference on World Wide

.Web, pages 437{446. ACM, 2008

M. Spaniol, D. Denev, A. Mazeika, P. Senellart, and G. Weikum. Data Quality in Web [19]

.Archiving. In Proceedings of WICOW, Madrid, Spain, April 20, 2009, pages 19 { 26

.ACM Press, 2009

Qingzhao Tan, Ziming Zhuang, Prasenjit Mitra, and C. Lee Giles. E_ciently detecting [20]

.webpage updates u In ICWE, pages 285{300, 2007

ص: 89

بایگانی وب پنهان سخت تر از بایگانی وب سطحی است. روش اصلی گردآوری محتوای وب بر پایه یافتن مسیر است هر صفحه در وهله اول باید به وسیله خزشگر پیدا شود تا بتوان آن را واکنشی و بایگانی کرد تاکنون روش مناسبی برای بایگانی وب پنهان پیش بینی نشده است. این امر نیازمند پیشرفتهایی برای حفاظت از وب پنهان از طریق راههای ساده و تکامل فنی وب است دو دلیل وجود دارد که بایگانی وب پنهان نباید مورد غفلت واقع شود نخست وب گستره وسیعی دارد و دارای منابع ارزشمندی است که بسیاری از مؤسسه های میراث فرهنگی به آن علاقه مندند. اینکه احتمالاً وب با معماریهایی از اطلاعات تکامل می یابد که در برابر شیوههای سنتی خزشگر مقاومت می کنند این مقاله به بررسی ویژگیهای وب، پنهان مسائل بایگانی وب، پنهان مسیرهای بایگانی کردن و فناوریهای بایگانی وب پنهان می پردازد.

نوشته ژولین ماسانه [\(2\)](#) | ترجمه افسانه تیموری خانی [\(3\)](#)

مقدمه

همان طور که در فصلهای قبلی ملاحظه کردیم روش اصلی گردآوری محتوای وب بر پایه یافتن مسیر است با توجه به اینکه پیمان نامه HTTP قابلیت تهیه سیاهه کامل را، ندارد هر صفحه در وهله اول باید به وسیله خزشگر پیدا شود تا بتوان آن را واکنشی و بایگانی کرد در فصل یک دیدیم که میدانیم که خزشگر محدودیت زمانی قابل توجهی را برای پردازش گردآوری کامل ارائه میکند اما لازم است که حداقل یک مسیر برای بایگانی هر مدرک وجود داشته باشد که البته این امر همیشه بعید به نظر میرسد در واقع بخش عظیمی از وب به همین دلیل توسط ابزارهای خودکار قابل دسترس نیست. این بخش برای اولین بار، در 1994، توسط جیل اچ الزوورث [\(4\)](#)، وب نامرئی نامیده شد (برگمن 2001 [\(5\)](#))؛ زیرا بخشی از وب است که توسط موتورهای کاوش خزشگر نمایه نمی شود. بعدها، پیشنهاد شد که این وب را وب عمیق نامگذاری کنند - در مقابل وب سطحی یا وب قابل نمایه سازی عمومی (پی آی دبلیو) [\(6\)](#) (لارنس [\(7\)](#))

ص: 91

Archiving the Hidden Web: in Masanes, Julien (ed.), Web Archiving. Berlin Heidelberg New - 1
York:Springer.pp.115-128

Julien Masane's 2-بایگانی وب اروپا julien@iwaw.net

3- دکترای کتابداری و اطلاع رسانی و کارشناس سازمان اسناد و کتابخانه ملی ایران

Jill H. Ellsworth 4-

Bergman 5-

(Publicly indexible web (PIW 6-

Lawrence 7-

و گیلز (1999) (1) - زیرا خزشگرها به راحتی میتوانند به آن دسترسی داشته باشند.

ما در اینجا از این دو اصطلاح استفاده نمیکنیم زیرا سبب ابهام بیشتر می شود. نخستین اصطلاح یعنی وب «نامرئی میتواند به این فکر منتهی شود که مشکل اصلی نمایش یا ارائه صفحه هاست؛ در صورتی که مشکل اصلی دسترسی ابزارهای خودکار است اصطلاح دوم وب عمیق ممکن است با عمق منابع در ساختار سلسله مراتبی ابر متنی وبگاهها اشتباه گرفته شود. بنابراین ترجیح داده ایم برای تعیین بخشی از وب که خزشگرها قادر به رسیدن به آن نیستند از اصطلاح وب پنهان مخفی» استفاده کنیم. به یاد داشته باشیم که حدود آن از نظر فنی تعریف شده است و به هیچ نوع خاصی از محتوا و تجربیات مشترک انسانی در محدودیتهای ناویری اشاره نمیکند و دیگر اینکه این تعریف کاملاً فنی است و مرز این دسته (گروه) ممکن است با تکامل فناوری تغییر یابد.

برای مثال، سایتهایی با توابع پیمایش رمزگذاری شده در فلش FLASH ساخته شده اند که قبل

از اینکه ماکرومدیا Macromedia نوعی SDK قادر به استخراج پیوندها از کد توابع را منتشر کند، به عنوان وب پنهان در نظر گرفته میشدند این سایتهای به خزشگر امکان میدهند تا از طریق وبگاههای کدگذاری شده فلش راه خود را پیدا کنند بنابراین این نوع وبگاهها نمیتوانند بیش از این به عنوان وب پنهان در نظر گرفته شوند.

اگر چه بیشتر توجه این فصل و بقیه فصلها به محدودیتهای خزشگرها معطوف است. شایان ذکر است که خزشگرها تمایل دارند صفحه های را کشف کنند که به ندرت توسط انسان دیده میشوند. در واقع مطالعات انجام شده در این زمینه به وسیله بوف خواد (2) و ووی نوت (3) (2003)، از طریق استفاده از تراکنش های خدمات وب اینریا (4) 2002 نشان داد تعداد قابل توجهی از صفحه های را که خزشگرها پیدا میکنند کاربران از دست میدهند.

پیدا کردن حداقل یک مسیر به اسناد

صفحه های متصل به وبگاه ها با هم یک مسیر وب را برای دسترسی به اسناد تشکیل می دهند. هر صفحه میتواند به یک یا چندین صفحه پیوند داده شود و خزشگرها نخستین صفحه ای را که به آن میرسند میگیرند به این ترتیب صفحههای پیوند داده نشده از وبگاههای دیگر حتی اگر در همان سرور ساکن باشند، شامل آن نمیشوند این مسیر وب میتواند به هر سندی گسترش یابد، به طوریکه از طریق مسیره های مجازی که با پرس وجو ساخته میشوند تولید یا قابل دسترس شوند به طور مثال صفحه های، پویا که از طریق محتوای ذخیره شده در پایگاه داده ایجاد میشود تنها به صورت مجازی وجود دارند که به عنوان بخشی از یک مسیر وب در نظر گرفته میشوند نکته مهم این است که چگونه مسیر اشیا ساخته می شود و آیا خزشگرها میتوانند آن را پیدا کنند یا نه.

ص: 92

Giles -1

Boufkhad -2

Viennot -3

INRIA -4

در اینجا میتوان دو مورد کلی را مشخص کرد نخست اینکه مجموعه هایی از مسیرهای با ارزش از پیش تعریف شده و محدود وجود دارد؛ که در واقع نمونه ای از ابر پیوندها و منوهاست. خزشگرها میتوانند چنین مسیرهایی را با استفاده از زبانی که برای رمزگذاری آنها استفاده شده استخراج، تفسیر، و دنبال کنند.

دومین مورد که برای خزشگرها هم بسیار مشکل است نیاز به تعامل صریح کاربر دارد (منظور بیشتر از یکبار کلیک کردن است. برای خزشگر تقلید مشکل است مانند زمانی که کاربران مجبورند از طریق فرمت HTML وارد پرس و جو شوند تا بتوانند به اسناد خاصی مثل یک، تصویر یک مقاله و یا یک صفحه پویا دسترسی یابند این نوع وب پنهان - که وب پنهان ساختاری نامیده می شود- شامل مقدار زیادی محتواست که از آن برای انتشار در انبار اسناد بزرگ وب استفاده شده است، که هم ساختار یافته است (پایگاه داده) و هم غیر ساختار یافته مثل مجموعه ای از تصاویر مقالات علمی موسیقی و مانند آن (استوری(1) و جنکه 1999 (2)) راه مناسب برای انتشار این قبیل مجموعه های بزرگ، محتوا استفاده از پیوندهای اختصاصی نیست؛ بلکه بهتر است از دسترسی پویا از طریق دروازه های پایگاه دادهها استفاده شود که شامل اطلاعات توصیفی برای هر مورد نیز میباشد. امروزه سایتهای زیادی از این نوع معماری اطلاعات شکل، پایگاه داده و گردآوری استفاده میکنند که بخشی از وب پنهان محسوب می شوند. این نوع معماری را دروازه مستند می نامیم زیرا ورود به فضای بزرگ اطلاعاتی را از طریق تعاملات (جست و جو) فراهم میکند (ماسانه(3) 2002). در برخی موارد رابط جایگزین مرورگر به گونه ای ارائه شده است که خزشگر هنوز به محتوای چنین انبارهایی دسترسی دارد در این صورت آنها در نوع اول قرار میگیرند.

در اینجا انواع تعامل کاربر را با جزئیات بیشتر نقد میکنیم تا ببینیم چه مشکلی برای خزشگرها

ایجاد میشود.

انواع تعامل با کاربر

لدسچر(4) و گوپتا(5) (1999)، نوعی مدل منبع تعاملی پیشنهاد دادند این مدل شامل چهار نوع عنصر ورودی است:

- ابر پیوندها روشی کلاسیک برای تهیه ورودی کاربر که محدود به یک نفر است؛

- منوها برای انتخاب زیر مجموعه ای از مقادارها از مجموعه ای از پیش تعریف شده ؛

- فرمها پیوندهای پویا با ویژگیهای ورودی متعدد (چندگانه)؛ و

- نوع چهارم که میتواند تنها بر اساس این مدل با تعامل صریح کاربر حمایت شود و عناصر بدون لفاف نامیده میشوند تصاویر نقشهها تعامل گرافیکی کاربر بر اساس جاوا و مانند آن).

ص: 93

Storey -1

Jahnke -2

Masanè s -3

Ludäscher -4

Gupta -5

مثال آن در دنیای واقعی توزیع عناصر ورودیهای مختلف در وبگاههای تصاویر غنی پزشکی است. (فرانکوویچ (1) و پروکوش (2) را ببینید (2001) هر یک از این عناصر ورودی موردی متفاوتی برای خزشگر است. دو مورد اول (فرایوندها و منوها تنها مشکلات تفسیر را برای خزشگرها ایجاد میکنند و با هم چیزی را میسازند که ما آن را عناصر ورودی تعریف شده می نامیم.

عناصر ورودی ارزش باز (معمولاً فرمت HTML) که در مقابل مورد دوم هستند؛ فضای امکان به یک یا چند مسیر تعریف شده را کاهش میدهد آنها در زمینههای مختلف استفاده می شود و ما در مورد آنها بیشتر بحث خواهیم کرد. نوع چهارم، گرچه نسبت به انواع دیگر کمتر معمول است؛ دسترسی به آن برای خزشگرها بسیار سختتر است. برای آگاهی از جزئیات نوع شناسی مشکلاتی که خزشگرها با آن مواجه شده اند دو گزارش کنسرسیوم بین المللی حفاظت اینترنت را ملاحظه کنید (بویکو (3) 2004؛ ماریل و همکاران 2004) در بخش زیر موارد اصلی و مباحثی در مورد مشکلات به وجود آمده برای خزشگرها را مطرح خواهیم کرد.

تعریف مقدار عناصر ورودی

پیوندهای ساده HTML یکی از رایج ترین این نوع است که مسیر راحتی را برای جست و جو فراهم میکند. پیوندهای نسبی، ممکن است بعضی از موارد تفسیر را افزایش دهند آنها بر اساس قوانین شرح داده شده RFC 1808 و RFC 2369 شبیه نشانیهای یونیکس هستند مکان نسبی، از راهنمای جاری شروع میشود و با استفاده از اسلش (/) پایین می رود و از طریق راهنمای والد با استفاده از دو نقطه (...) بالا می رود یک مسیر نسبی که با اسلش شروع میشود به این معنی است که راهنمای ریشه میزبان است مشکلات در اثر استفاده بد پدیدآرندگان صفحه ها یا مدیریت سیستم محتوا به وجود می آید. به طور مثال شامل نقاط، اضافی راهنماهای فوق عددی، والد و یا دیگر ترکیبات عجیب و غریب می شود. در صورت، وجود تگ، اصلی در بخش رأس یک صفحه HTML یک مکان پیشفرض جایگزین برای پیوندهای نسبی تهیه میکند که تمامی، پیوندها به جای شروع از راهنمای جاری، از آن شروع شوند. از پیوندهای جاوا اسکریپت به خوبی دیگر زبانهای برنامه نویسی برای ایجاد پیوندهای خاص برای منوها و ناوبری طومارنمایی (پیمایش) تقویم، پویا و مانند آن استفاده میشود از آنجاکه هر پیوند نتیجه ترکیبی از دستور متغیر یا ورودی تعامل کاربر است؛ در بعضی موارد تفسیر پیوندها بدون اجرای اسکریپت تقریباً غیر ممکن خواهد بود.

خزشگرها می توانند تفسیر مبتنی بر قواعد (4) و یا هر ترکیب ممکن از مسیر یا نام فایل را برای پیدا کردن اسکریپت (5) داشته باشند در هر مورد موفقیت کامل تضمین شده نیست راه حل جایگزین - که هنوز

ص: 94

Frankewitsch -1

Prokosch -2

Boyko -3

4- (Roche 2006) (This is the case for instance of the open-source web copier HTTrack (see Chap. 3

Heritrix, the open source archive-quality crawler developed jointly by the Internet Archive and the nordic

(libraries in the IIPC, implements this approach (Mohr et al. 2004

آزمایش نشده - تفسیر با اجرای کدها به جای تجزیه است که به معنی اجرای مرورگرها شبیه سازی زمینه ها و تعامل کاربر است. بعضی از انواع، پیوندها مثل پیوندهای فریم و پیوندهای تصویر ممکن است مشکلاتی برای تفسیر به وجود آورند؛ ولی به طور کلی میتوانند با موفقیت دنبال شوند.

عناصر ورودی مقدار باز

در همه موارد، قبلی با وجود مشکلات ایجاد شده دامنه محدودی از مسیرهای ممکن توسط کد تعریف شده اند که باعث میشود تفاوت زیادی با نوع دوم به وجود آید که در آن میتوان یک مجموعه بی نهایت و تعریف نشده از مقادیر را به پیوندها اختصاص داد.

مکانیسم اصلی برای این مورد فرمت HTML است آنها کاربران را قادر میسازند که یک مقدار دلخواه را به سرور عبور دهند مقادیر ورودی به طور مثال میتوانند برای پرس و جوی اسناد مورد استفاده قرار گیرند که ناشی از تولید مجموعه ای از پیوندها به صفحه ها یا اسناد است. وارد کردن یک پرس و جو حاوی اطلاعات نویسنده و عنوان در سیاهه میتواند یک سیاهه انتشار همراه با پیوند به هر یک از آنها ایجاد کند این پیوندها از یک پایگاه داده ایجاد و در نتایج صفحه HTML جاسازی شده اند. اگر هیچ پیوند دیگری به این اسناد - به عنوان مثال از طریق سیاهه مرتب شده نشریات بر اساس حروف الفبا وجود نداشته باشد خزشگرها تنها از طریق مسیر مجازی ایجاد شده توسط رابط پرس و جو میتوانند به آنها دسترسی پیدا کنند.

این نوع معماری اطلاعات به دروازه مستند معروف است (مازانه 2002) که بسیار معمول است و باید از استفاده دیگر فرمها متمایز باشد (شکل 1) فرمها به طور عمده برای جمع آوری ورودی کاربر مانند ورود به سیستم و ورود اطلاعات تماس و یا باز خورد ارسال نظرات جعبه جست و جوی عمومی و مانند آن مورد استفاده قرار میگیرند کوپ(1) و همکارانش در سال 2003 دریافتند که حدود 50 درصد از فرمهای HTML رابط جست و جو هستند و لاگ(2) و همکارانش در سال 2002 در نمونه خود یافتند که 95 درصد از فرمها از جمله جعبه جست و جوی عمومی فرمهای ناخواسته بودند.

اما حتی اگر بسیاری از فرمهای موجود بر روی وب برای مقاصد دیگر استفاده شوند، آنهایی که باقی مانده اند نقطه ورود به فضای اطلاعات بزرگی هستند که وب پنهان نمایش میدهد. این دو مطالعه تلاش کرده اند تا آن را مشخص کنند.

ویژگیهای وب پنهان

اولین مطالعه در سال 2000 توسط برگمن برگمن (2001) با استفاده از تحلیل همپوشانی بین جفت موتورهای کاوش به منظور برآورد تعداد وبگاههای پنهان صورت گرفت. آنها مشخص کردند که طیف وسیعی بین 43000-96000 وبگاه پنهان براساس حضور فرم وجود دارد. متأسفانه، فیلترهای مورد

ص: 95

استفاده مستند نشده اند؛ بنابراین به سختی میتوان درباره ارزش این نتایج و مقایسه آنها با دیگران قضاوت کرد. آنها همچنین به تحلیل 60 موتور کاوش بزرگ پرداختند و اندازه آنها را 550 بیلیون صفحه برآورد کردند که 550 برابر بزرگتر از وب سطحی در آن زمان است.

فرض اصلی در پشت این مطالعه این است که هر مورد در پایگاه داده با یک صفحه ایجاد شده

مرتبط است و اندازه آن با اچ تی ام ال. تخمین زده می شود.

این، اصل شامل تمامی HTML و اطلاعات مربوط به کد (HTML) به اضافه محتوای متن

، استاندارد تصاویر تعبیه شده منحصر به فرد و اطلاعات استاندارد سرآیند HTTP (پیمان نامه انتقال ابرمتن) است. استفاده از این پیمان نامه استاندارد اجازه میدهد تا مقایسه دقیقی بین وب سطحی و عمیق صورت گیرد (برگمن 2001).

برای مثال ورودی پایگاه دادههای آب و هوایی ملی ایالات متحده (بزرگترین مثال در نمونه خود)

با یک صفحه 13 کیلوبایتی مطابقت میکند.

واقعیت این است که این بانک اطلاعاتی و پایگاه داده NASA EOSDIS تقریباً 80 درصد از کل نمونه ها را ارائه میدهد و نشان میدهد که این مطالعه نگاهی سو گرفته به وب پنهان نسبت به محتوای تکراری و غیر مستند دارد اگر چه در وهله اول به نظر میآید این مطالعه با هدف توجه به اهمیت و غنای این بخش از وب صورت گرفته است.

مطالعه اخیر توسط چانگ (1) و همکارانش (2004)، جزئیات بیشتر و ویژگیهای مستندی از وب پنهان را آشکار میکند. این مطالعه تمایز بین پایگاه داده ساختار یافته یعنی پایگاه داده مرتبط با مقادیر زوج کلیدی و محتوای بدون ساختار (متن، عکس، شنیداری، دیداری) را مشخص مینماید که در این کتاب محتوای دروازه مستند نامیده میشوند در مطالعه چانگ و همکارانش همانند بسیاری از مطالعات انجام شده بر روی وب، پنهان توجه اصلی به نوع اول پایگاه داده ساختار یافته است که توسط پژوهشگران جامعه پایگاه داده که علاقه مند به یکپارچه سازی دادههای وب هستند - ساخته شده اند. اگر چه یافته های بسیاری مثل یافته ما بیشتر از منظر محتوا گرا مورد نظر هستند.

آنها مطالعات خرد و کلانی انجام داده اند مطالعه کلان بر روی یک میلیون نشانی تصادفی تولید IP صورت گرفت که برای پیدا کردن سرور HTTP مورد آزمایش قرار گرفتند 2260 وبگاه پیدا شد و مورد خزش قرار گرفت؛ 126 مورد وبگاه پنهان شناسایی شد که شامل 190 پایگاه داده است. به گزارش وب، جهانی این به معنی 307000 سایت حاوی 102000 دروازه مستند و 348000 پایگاه دادههای ساختاریافته است. این نکته اهمیت وبگاههای پنهان را نشان میدهد. همچنین، چانگ و همکارانش توزیع پایگاه داده وب را در عمق شناسایی کردند که نشان میدهد 91/6 درصد از آنها در داخل عمق 3 یافت شدند. در مطالعه خرد، 441 منبع را با جزئیات بیشتر مورد بررسی قرار دادند. آنها در آغاز نشان دادند که در بسیاری از موارد مسیر ناوبری جایگزینی برای رسیدن به محتوا وجود دارد که در آن، محتوای نمونه برای موتورهای کاوش واقعا پنهان نیست. در واقع پنهان بودن وبگاه بستگی به دامنه دارد.

آنها طرح رابط پرس وجو و تعدادی از ویژگیهای دقیقتر را مورد مطالعه قرار دادند که کوچکترین اندازه طرح 1 بزرگترین 18، و حد وسط 6 است. همچنین به مطالعه واژگان طرح پرداختند و پنج ویژگی اصلی طرح عنوان کلید واژه ها، «قیمت»، «ساخت» و «هنرمند» را نشان دادند.

گفته چانگ و همکارانش تشویق به پردازش خودکار میتواند باعث ایجاد نظم و قاعده شود ظاهراً بررسیهای ما پدیدههای دوگانه را نشان میدهد که با هم سیاهه و ویژگیهای منحصر به فرد مرز

وب عمیق را مشخص میکنند:

نخست به عنوان یک چالش منابع در ونخطی عملاً نامحدود هستند؛ حتی برای یک حوزه خاص مورد علاقه منابع جایگزین بی شماری پدیده تکثیر منابع وجود دارد. بنابراین، یکپارچه سازی در مقیاس بزرگ چالشی واقعی است.

دوم به عنوان یک فرصت مطرح است با این حال زمانی که منابع در حال تکثیر هستند، در مجموع، پیچیدگی آنها برای نشان دادن زیر بنای ساختار تمایل به هماهنگی دارد به طور خاص ما این ساختار هماهنگ را در ویژگیهای واژگان و الگوهای پرس وجو در منابع وب مشاهده کردیم این نوع تجمیع واژگانی در موقعیت و اندازههای همگرا دسته بندی میشود (چانگ و دیگران 2004).

بایگانی وب پنهان سرویس گیرنده

همان طور که در بخش قبلی نشان داده شده است نقطه ورودی وب پنهان (فرم) قواعد مربوط را نشان میدهد که برای استخراج خودکار محتوا با لفافه مورد استفاده قرار می گیرد و گاهی عوامل وب پنهان نامیده می شوند (راگهاوان(1) و گارسیا-مولینا 2001(2)؛ لاگ(3) و دیگران 2002؛ برای مقدمات کلی درباره این عنوان هرست(4)، 1998؛ و آدامز(5)، 2001 را ملاحظه بفرمایید).

نقش این عوامل شامل شناسایی فرمهای HTML آموزش پر کردن آنها شناسایی و واکنشی نتایج محتوای این فرآیند میتواند برای ارائه یک رابط جست و جوی یکپارچه اجرا شود (برای مثال این عناوین را بررسی کنید، فلورسکو(6) و دیگران 1998 مثال برای اجرای خدمات جست وجو بورگمن 200 و نمونه پیشرفتهای اخیر در این حوزه هی(7) و دیگران 2005) معمولاً-ردیابی در حالت طبیعی خزشگرها از طریق تجزیه و تحلیل صفحه های حاوی فرمهای HTML صورت میگیرد برای از بین بردن فرمهای نامطلوب از روش اکتشافی (فناوری در هوش مصنوعی) استفاده میشود (صفحه ورود به سایت و یا ارتباط با صفحه اطلاعات جعبه جست و جوی عمومی و مانند آن). سپس عوامل زمینه ای پرس وجو و برچسبها را استخراج میکنند و سعی میکنند آنها را با برچسبهای شناخته شده مقایسه

ص: 97

Raghavan -1

Garcia- Molina -2

Lage -3

Hearst -4

Adams -5

Florescu -6

He -7

کنند و گاهی اوقات برای ارزیابی، موضوع آنها را با واژگان شناخته شده بررسی میکنند (گراوانو(1) و دیگران 2003).

در نهایت فرمها به صورت خودکار پر شده و در نتیجه صفحه ها یا اسناد ذخیره میشوند. روش اکتشافی (فناوری در هوش مصنوعی) که توسط لاگ و همکارانش (2002) استفاده شده است فرمهای ناخواسته را فیلتر میکند و اجازه میدهد تا فرمها با کمترین عناصر از بین برده شوند و عناصر نوع HTML با هر رمز عبوری ساخته شوند کوپ و همکارانش (2003) درخت تصمیم گیری شفافی را برای کشف رابط جست و جو بر روی وب شکل (2) بر اساس تجزیه فرمهای HTML با توجه به ویژگیهای ذره ای (اتمی) مانند HTTP، دامنه کنترل، متن کنترل رمز عبور و نظیر آن پیشنهاد کردند که براساس مجموعه آموزش و یادگیری الگوریتم ساخته شده است.

پرکردن فرمها، به طور خودکار مستلزم درک زمینه است. لاگ و همکارانش (2002)، فرض کردند که برچسبها معمولاً در گوشه سمت چپ و یا بالایی زمینه فرم قرار داده میشود. جانگ(2) و همکارانش (جانگ و همکارانش 2004 فرض کردند که نظم و قاعده ای یا الگوهای طراحی خاصی در میان فرم پرس و جوی وب وجود دارد که با هم نوعی زبان دیداری قابل تجزیه را تشکیل میدهند آنها نوعی ابزار تبدیل پرس و جوی فرم HTML به مجموعه ای از نشانه ها را به وجود آوردند که هر یک نشان دهنده یک عنصر دیدنی ذره ای در چارچوبی دو بعدی است. آنها اخذ الگوهایی مانند قرابت همجواری و رابطه معنایی میان اصطلاحات را پیشنهاد کردند. با استفاده از الگوریتم تجزیه آنها دقت و پوشش 80/0 را برای تشخیص فرمها به دست آوردند.

هنگامی که این کار انجام شد لازم است اصطلاحات به صورت خودکار، از طریق فرمها، جهت ایجاد پاسخ ارائه شوند مشکلات دیگری در هنگام تکمیل فرم به طور خودکار وجود دارند(نگاه کنید لیدل(3) و دیگران 2002) ورود به فرمها میتواند محدود به فیلدهای متنی کلید رادیویی(4) جعبه بازبینی(5) سیاههها، و مانند آن و یا هر نوع فایل پیوست کدگذاری شده MIME باشد درخواست منطقی می تواند به اشکال مختلف با اطلاعات دولتی گرفته شده از روی سرور تقسیم گردد کوکیها فیلدهای پنهان مقادیر کدگذاری شده به URL پایه برخی فرمها قبل از ارسال فرم برای تغییر فیلد بر اسکرپتها تکیه میکنند (بازبینی گسترده(6)، دیگر اعتبارات، فیلد و محاسبه خودکار بعضی از فیلدها).

اگر واژگان اصطلاحات ارائه شده درست تعریف شده باشند فرصت مناسبی برای واکنشی محتوا به وجود میآید این مسئله زمانی مهم است که دامنه این واژگان محدود باشد (مانند کد پستی، تاریخ، و مانند آن) در واقع یکی از این زمینه ها کافی است به عنوان مثال دروازه مستند به متون فیلسوفان فرانسوی که یک ورودی تاریخ را ارائه میدهد میتواند با تمام تاریخها از سال 1100 - 2005 بدون هیچ

ص: 98

Gravano -1

Zhang -2

Liddle -3

radio buttons -4

check boxes -5

range checking -6

پیشفرضی در مورد نام نویسنده یا عنوان نوشته ها پرس و جو شود در 905 پرس و جو فرد می تواند مطمئن باشد که میتواند تمام متون را به دست آورد.

محدودیت این روش زمانی است که دامنه پرس و جو برای تمامی فیلدهای بیش از حد باز یا تعریف نشده به صورت نظام مند بررسی شود. در این موارد امکان استفاده از یک رویکرد دیگر امکان پذیر است که از یک پرس و جو برای استخراج اصطلاحات پرس و جوهای جدید سوء استفاده میکند که سپس ارسال و تکرار خواهد شد روشهای ارائه شده توسط کالان(1) و کانل(2)(2001) نمونه برداری پرس و جو محور نامیده میشود و با موفقیت استفاده شده است (آژیچتین(3) و همکارانش 2003، و (بربوزا(4) و فریره(5) 2004) ندولاس(6) و همکارانش (2005) پیشنهاد کرده اند که از الگوریتم تطابق برای انتخاب بهترین و مهمترین کلمات - کلماتی که به بسیاری از سندها مربوط اند - استفاده شود آنها نشان میدهند که تنها با 83 پرس و جو تقریباً 80 درصد از 14 میلیون اسناد ذخیره شده در پاپ مد(7) را میتوان بارگذاری کرد.

همکاری سرور خزشگر

آشکار کردن دروازه مستند خزشگر

زمانی که همکاری با پدید آورندگان وبگاه امکان پذیر است این امکان وجود دارد تا محتوای وبگاه پنهان برای خزشگر به منظور بایگانی شدن آشکار شود. در اینجا اغلب برای بایگانی کردن از روشهای ارائه شده برای موتورهای کاوش استفاده میشود آنها شامل ایجاد سیاهه کاملی از اسناد و یا فراهم آوردن امکان دسترسی به یک خدمت هستند که میتوانند به صورت خودکار توسط خزشگر مورد پرس و جو قرار گیرند آنها طیف وسیعی از صفحه های تولید شده هستند که اغلب برای کاربران انسانی پنهان اند و به همه اشیای وبگاه حتی پیمان نامه های پرس و جو اختصاص دارند پیمان نامه متعددی به جز چند مورد از جمله پیمان نامه ..ای.آی. از اواخر دهه 90 با موفقیتهای محدودی پیشنهاد شده اند. ما مهمترین آنها را در زیر مرور میکنیم.

1 - صفحه های پیوندهای پنهان

ساده ترین راه برای فعال کردن خزشگر به منظور دریافت محتوای، پنهان ایجاد سیاهه پیوندهاست که به همه اسناد وبگاه پنهان اشاره دارد این امر به ویژه برای دروازه مستند مناسب است و میتواند به طور کامل از دید کاربران عادی با پنهان کردن این صفحه ها از ناوبری، طبیعی پنهان بماند؛ به عنوان مثال با تعبیه پیوندهای پنهان در صفحه اصلی.

ص: 99

Callan -1

Connell -2

Agichtein -3

Barbosa -4

Freire -5

Ntoulas -6

اینکار برای مثال توسط کتابخانه ملی فرانسه و کتابخانه ملی استرالیا از طریق آشکار کردن مجموعه برای موتورهای کاوش انجام شد:

سیاهه جداگانه ای از URLها برای هر یک از مجموعه های دیجیتال کتابخانه یعنی عکسها، نقشه ها، صفحه های موسیقی نسخه های خطی کتابها و پیوندها ساخته شده است. هر مجموعه شامل هزاران مورد است که در یک سری از صفحه های وب سیاهه شده اند هر کدام شامل 100 پیوند هستند که اقلام (موارد) مجموعه را تفکیک میکنند این صفحه با استفاده از دستور ربات `anoindex, follow` موتورهای کاوش را مستقیم هدایت میکند تا با دنبال کردن پیوند به محتوا برسد؛ اما صفحه ها را سیاهه نمی کند. سیاهه URL خود به صورت پویا تولید میشود و با محتوای جدید به سیاهه اقلام جدید دیجیتالی شده به صورت خودکار اضافه میگردد و توسط اینترنت قابل دسترس میشود (بوستن:1) (2005)

برای مؤثر بودن طرح پیوند پایدار باید در جای خود قرار گیرد و با تمام اسناد در پیوند باشد. براند من(2) و همکارانش (2000) محاسبه کردند که چنین مکانیزمی میتواند تا 80 درصد در پهنای باند تبادل خزشگر سرویس دهنده صرفه جویی کند. از سال 1990، طرحهای پیشنهادی متعددی برای استاندارد کردن این نوع مکانیسمهای رسمی ارائه شده اند یکی از آنها که در حال حاضر مورد استفاده قرار می گیرد استاندارد باز RSS است:

RSS مخفف خلاصه سایت RDF سایت مختصر غنی یا «پیوند واقعاً ساده» است. در ابتدا، تصور می شد که این توسعه RDF برای خدمات مای نت اسکپ(3) است، اما از آن امروزه در وب نوشتهها یا سایتهای خبری برای ارائه سیاهه کوتاه از آخرین اخبار و / یا روز آمد سازی سایتهای، به طور وگسترده استفاده میشود چنین استاندارد را میتوان برای تولید دورههای فایل RDF حاوی URL تاریخ آخرین تغییرات در تمام صفحه های سایت مورد استفاده قرار داد پس از آن، خزشگر می تواند برای اولین بار فایل را بررسی و از آن برای خزش در سایت استفاده کند یا آن را با سیاهه های موجود خود مقایسه کند اگر صفحه های سایت قبلاً خزش شده اند تنها صفحه های تغییر یافته را واکنشی کند کاستیلو(4) (2004، ص 109) این نوع پیاده سازی را ارائه کرده است.

چنین مکانیسمی می تواند برای سیاهه و آشکار کردن هر نوع صفحه ای که به وبگاههای دیگر (و نه به صفحه های وب پنهان) پیوند داده شده اند مورد استفاده قرار گیرد.

1 - 2 سطح پیمان نامه

گامی به جلو را میتوان با پیمان نامه های ارتباطی اختصاصی برداشت مانند پیمان نامه برداشت طرح ابر داده بایگانی آزاد(5) (-OAL MHP) که توسط لاگوز(6) و ون دی سامپل(7) در سال 2001 پیشنهاد شد و با استفاده

ص: 100

Boston -1

Brandman -2

MyNetscape -3

Castillo -4

Open Archive Initiative Metadata Harvesting Protocol -5

Lagoze -6

Van de Sompel -7

از ترکیب نحوی (1) XML به آشکارسازی ابر داده بر روی HTTP می پردازد. پیاده سازی در سطح پیمان نامه، باعث ارتباط واقعی از طریق درخواست و پاسخ سرور میشود و از این طریق احتمالات افزایش می - یابد به عنوان مثال امکان پرس و جواز سند به کمک تاریخ و نوع اسناد فراهم می شود. خزشگر می تواند به طور مستقیم با سرور OAI ارتباط برقرار کند تا سیاهه ای از اسناد مرتبط با ابر داده (فایل) را به دست آورد. اگر انباره سازگار OAI دسترسی نامحدود به اسناد را فراهم کند این امکان به وجود می آید که آنها را واکنشی و با ابر داده خود در بایگانی ذخیره کند برخی خدمات دروازه واسطه نیز برای خزشگرهایی اجرا شده است که قادر به استفاده از پیماننامه OAI برای ایجاد صفحه های پیوند به تمامی اسناد از سرور OAI نیستند (لیو(2) و دیگران 2002)

حتی اگر این نوع مکانیسم همکاری در جای خود استفاده شود باید آن را در بیشتر موارد به عنوان روش مکمل جمع آوری محتوا برای بایگانی سازمان در نظر گرفت همانطور که سایتهای تولید کننده معمولاً از آنها برای کاهش بار وارده بر سرور خود - با هدف قرار دادن خزشگر موتور کاوش برای بازدید صفحه های روزآمد شده - استفاده میکنند در حالی که نمایه سازی خزشگر بر روی سطح سایتهای باقی میماند و نتایج را روز آمد میکند بایگانی خزشگر نیازمند یکپارچگی و تکامل است.

دروازه بایگانی مستند

عکس

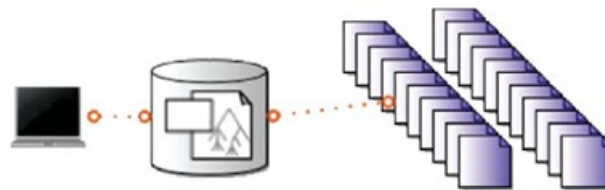
از ترکیب نحوی XML^۱ به آشکارسازی ابر داده بر روی HTTP می پردازد. پیاده سازی در سطح پیمان نامه، باعث ارتباط واقعی از طریق درخواست و پاسخ سرور می شود، و از این طریق احتمالات افزایش می یابد. به عنوان مثال، امکان پرس و جو از سند به کمک تاریخ و نوع اسناد فراهم می شود. خزشگر می تواند به طور مستقیم با سرور OAI ارتباط برقرار کند تا سیاهه ای از اسناد مرتبط با ابر داده (فایل) را به دست آورد. اگر انباره سازگار OAI دسترسی نامحدود به اسناد را فراهم کند این امکان به وجود می آید که آنها را واکنشی و با ابر داده خود در بایگانی ذخیره کند. برخی خدمات دروازه واسطه نیز برای خزشگرهایی اجرا شده است که قادر به استفاده از پیمان نامه OAI برای ایجاد صفحه های پیوند به تمامی اسناد از سرور OAI نیستند (لیو^۲ و دیگران ۲۰۰۲).

حتی اگر این نوع مکانیسم همکاری در جای خود استفاده شود باید آن را در بیشتر موارد به عنوان روش مکمل جمع آوری محتوا برای بایگانی سازمان در نظر گرفت، همانطور که سایت های تولید کننده معمولاً از آنها برای کاهش بار وارده بر سرور خود - با هدف قرار دادن خزشگر موتور کاوش برای بازدید صفحه های روزآمد شده - استفاده می کنند. در حالی که نمایه سازی خزشگر بر روی سطح سایت ها باقی می ماند و نتایج را روز آمد می کند. بایگانی خزشگر نیازمند یکپارچگی و تکامل است.

دروازه بایگانی مستند

در برخی موارد، روش های قبلی را نمی توان به کار برد، زیرا به از دست دادن غنا و ساختار ابر داده منجر می شوند، به خصوص در مورد اسنادی که نمی توانند یک پیوند ساده را ترسیم کنند. به عنوان مثال، مجموعه ای از تصاویر علمی را تصور کنید: بایگانی تمام تصاویر بدون ابر داده می تواند بی فایده باشد. در این موارد، ابر داده مرتبط با تصاویر به اندازه خود تصاویر مهم هستند.

گزینه جایگزین برای جلوگیری از اجرای پیمان نامه جدید توسط سرور، استخراج مستقیم ابر داده از بانک اطلاعاتی و بایگانی آن همراه با اسناد در یک فرمت آزاد است (شکل ۳). کتابخانه ملی فرانسه در سال ۲۰۰۲، این روش را با موفقیت در چند وبگاه پنهان به کار برده است. این روش همکاری تولید کننده را می طلبد و در مقایسه با پیاده سازی یک سرویس جدید خواهان کمتری دارد، زیرا نیاز به استخراج دارد. مشکل اصلی ناشی از ناهمگونی نظام های پایگاه داده، طرح پایگاه داده، و طرح نگاشت (ترسیم) اشیا است.



شکل ۳

1. fv'lkXML
2. Liu

در برخی موارد روشهای قبلی را نمیتوان به کار برد، زیرا به از دست دادن غنا و ساختار ابر داده منجر می شوند به خصوص در مورد اسنادی که نمیتوانند یک پیوند ساده را ترسیم کنند. به عنوان مثال، مجموعه ای از تصاویر علمی را تصور کنید بایگانی تمام تصاویر بدون ابر داده میتواند بی فایده باشد. در این موارد ابر داده مرتبط با تصاویر به اندازه خود تصاویر مهم هستند.

گزینه جایگزین برای جلوگیری از اجرای پیمان نامه جدید توسط سرور استخراج مستقیم ابر داده از بانک اطلاعاتی و بایگانی آن همراه با اسناد در یک فرمت آزاد است (شکل 3) کتابخانه ملی فرانسه در سال 2002، این روش را با موفقیت در چند وبگاه پنهان به کار برده است. این روش همکاری تولید کننده را می طلبد و در مقایسه با پیاده سازی یک سرویس جدید خواهان کمتری دارد زیرا نیاز به استخراج

دارد. مشکل اصلی ناشی از ناهمگونی نظامهای پایگاه داده طرح پایگاه داده و طرح نگاشت (ترسیم) اشیاست.

ص: 101

fv'lkXML-1

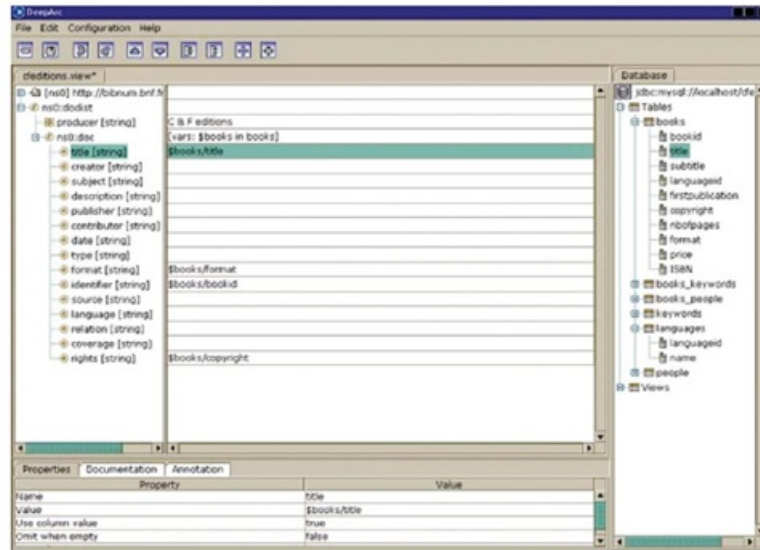
Liu-2

به منظور کاهش در دو مورد، اول کتابخانه ملی فرانسه نوعی استخراج کننده پایگاه داده گرافیکی منبع آزاد را توسعه داده است که میتواند در چندین نظام پایگاه داده اجرا شود. (1) سپس، ابزار می تواند محتوای پایگاه داده را به اسناد XML مطابق با طرح انتخاب شده توسط بایگانی صادر کند. تولید کننده، که به خوبی طرح داخلی خود را می شناسد می تواند نقشه گرافیکی آن را با هدف طرح ارائه شده توسط بایگانی ترسیم کند (شکل 4).

شکل 4

عکس

به منظور کاهش در دو مورد اول، کتابخانه ملی فرانسه نوعی استخراج کننده پایگاه داده گرافیکی منبع آزاد را توسعه داده است که می‌تواند در چندین نظام پایگاه داده اجرا شود^۱. سپس، ابزار می‌تواند محتوای پایگاه داده را به اسناد XML، مطابق با طرح انتخاب شده توسط بایگانی صادر کند. تولید کننده، که به خوبی طرح داخلی خود را می‌شناسد، می‌تواند نقشه گرافیکی آن را با هدف طرح ارائه شده توسط بایگانی ترسیم کند (شکل ۴).



شکل ۴

به‌عنوان مثال، اگر تولید کننده «AUT» را به عنوان نام فیلد در پایگاه داده نویسنده به‌کار برده باشد می‌تواند با کمک ویزارد این دو را به‌صورت گرافیکی رسم کند. ویزارد، همچنین، تولید کننده را قادر می‌سازد تا اطلاعات خاصی را که به‌دلیل حفظ حریم خصوصی نیاز به بایگانی شدن ندارند، فیلتر کند. تمام این مراحل را می‌توان در مدت زمان محدود انجام داد و به تولید خروجی نسخه XML ساختار پایگاه داده پرداخت که می‌تواند با اسناد شرح داده شده توسط این ابر داده صادر و حفظ شود. مشکل اصلی این روش نامگذاری و پیوند با طرح و چگونگی ترجمه در محیط بایگانی است. شناسایی اشیا در پایگاه داده می‌تواند برای پیوند با داده در محیط بایگانی پیچیده و دشوار باشد. اگر از نوعی ساختار راهنما برای سازماندهی اشیا استفاده شده باشد و یا یک اسکریپت در محیط اصلی راه را

1. DeepArc, an opensource database extractor. <http://bibnum.bnf.fr/downloads/deeparc/> (last visited May 2006).

به عنوان مثال، اگر تولید کننده «AUT» را به عنوان نام فیلد در پایگاه داده نویسنده به‌کار برده باشد می‌تواند با کمک ویزارد این دو را به صورت گرافیکی رسم کند، ویزارد، همچنین تولید کننده را قادر می‌سازد تا اطلاعات خاصی را که به دلیل حفظ حریم خصوصی نیاز به بایگانی شدن ندارند فیلتر کند. تمام این مراحل را می‌توان در مدت زمان محدود انجام داد و به تولید خروجی نسخه XML ساختار پایگاه داده پرداخت که می‌تواند با اسناد شرح داده شده توسط این ابر داده صادر و حفظ شود.

مشکل اصلی این روش نامگذاری و پیوند با طرح و چگونگی ترجمه در محیط بایگانی است. شناسایی اشیا در پایگاه داده می‌تواند برای پیوند با داده در محیط بایگانی پیچیده و دشوار باشد اگر از نوعی ساختار راهنما برای سازماندهی اشیا استفاده شده باشد و یا یک اسکریپت

DeepArc, an opensource database extractor.<http://bibnum.bnf.fr/downloads/deeparc/> (last visited May - 1
.2006).

برای اشیا ایجاد کند بایگانی باید طرح پیوند خود را بسازد که با ساختار و طرح نامگذاری سازگار است. از آنجا که مکانیسم پیوند اصلی را میتوان تعریف کرد روش دیگری برای به کارگیری وجود ندارد؛ روش مورد به مورد برای ایجاد یک بایگانی کارآمد امری ضروری است.

این بایگانی باید فرمت HTML خود را برای پرس و جوی XML ابر داده بایگانی شده ایجاد کند و به مجموعه ای از اشیا پیوند داده شود.

توجه داشته باشید که ابر داده را میتوان در بایگانی سیستمهای سنتی مدیریت پایگاه داده رابطه ای به راحتی تزریق کرد. نکته مهم این است که نسخه XML از ابر داده اصلی در پایگاه داده به منظور اطمینان از اینکه در آینده قابل خواندن خواهد بود باقی میماند و محافظت خواهد شد.

نتیجه گیری

همانطور که ملاحظه کردیم بایگانی وب پنهان سختتر از بایگانی وب سطحی است با اینکه برخی از روشها به موفقیت رسیده اند تا زمان نوشتن این مطلب هیچ یک از آنها را نمی توان به عنوان روش مناسب در نظر گرفت این امر نیازمند پیشرفتهایی برای حفاظت از وب پنهان از طریق راههای ساده است و البته به تکامل فنی وب نیز بستگی دارد اما حداقل دو دلیل وجود دارد که این بخش از وب مورد غفلت واقع نشود نخست اینکه وب گستره وسیعی دارد و دارای منابع ارزشمندی است که بسیاری از مؤسسه های میراث فرهنگی به آن علاقه مندند دوم، اینکه احتمالاً وب با معماریهایی از اطلاعات تکامل می یابد که در برابر شیوههای سنتی خزشگر مقاومت می کنند نخستین عامل ضد تعادلی فشاری است که موتور کاوش در سایت قرار میدهد برای نمایه شدن باید خزش شوند. اما اگر هماهنگیهای مستقیم و به روزرسانی دو طرفه وجود داشته باشد ممکن است تغییر کند که دلیل خوبی برای ادامه کار و دیدگاهی محافظت گرا می باشد.

منابع

Adams, K. C. (2001). The Web as Database: New Extraction Technologies and Content

Management. Online, March

Agichtein, E., Ipeirotis, P. G., Gravano, L. (2003). Modeling Query-Based Access to Text

Databases

Barbosa, L. Freire, J. (2004). Siphoning Hidden-Web Data through Keyword-Based

Interfaces. Paper presented at the SBBD

Bergman, M. I. K. (2001). The Deep Web: Surfacing Hidden Value. The Journal of Electronic

(Publishing, 7(1

Boston, T. (2005). Exposing the deep web to increase access to library collections. Paper

,presented at the AusWeb05. The Twelfth Australasian World Wide Web Conference

Queensland, Australia

Boufkhad, Y. Viennot, L. (2003). The Observable Web. RR Boyko, A. (2004). Test Bed

Taxonomy. IIPC Reports, 16

Brandman, O., Cho, J., Garcia-Molina, H., Shivakumar, N. (2000). Crawler-Friendly Web

Servers. SIGMETRICS Performance Evaluation Review, 28(2), 9-14

Callan, J. Connell, M. (2001). Query-based sampling of text databases. ACM Transactions

on Information Systems 19(2), 97-130

Castillo, C. (2004). Effective Web Crawling. University of Chile

Chang, K. C.-C., He, B., Li, C., Patel, M., Zhang, Z. (2004). Structured databases on the

web: observations and implications. SIGMOD Records, 33(3), 61-70

Cope, J., Craswell, N., Hawking, D. (2003). Automated discovery of search interfaces

on the web. Paper presented at the Proceedings of the Fourteenth Australasian Database

Conference on Database Technologies 2003

Florescu, D., Levy, A., Mendelzon, A. (1998). Database techniques for the World-Wide

Web: A survey. SIGMOD Records, 27, 59-74

.Frankewitsch, T. Prokosch, U. (2001). Navigation in medical Internet image databases

Medical Informatics and the Internet in Medicine, 26(1), 1-15 5 Archiving the Hidden

Web 129

Gravano, L., Ipeirotis, P. G., Sahami, M. (2003). QProber: A System for Automatic

(Classification of Hidden-Web Databases. ACM Transactions on Information Systems, 21(1

He, H., Meng, W., Yu, C., Wu, Z. (2005). WISE-Integrator: a system for extracting and integrating complex web search interfaces of the deep web. Trondheim, Norway

Hearst, M. (1998). Information Integration. IEEE Intelligent Systems, 13(5), 12-24

/HTTrack. <http://www.httrack.com>

Lage, J. P., Silva, A. S. D., Golgher, P. B., Laender, A. H. F. (2002). Collecting hidden Web

pages for data extraction. Paper presented at the Proceedings of the fourth international

workshop on Web information and data management low-barrier interoperability

framework. Roanoke, Virginia, United States

,Lawrence, S. Giles, C. L. (1999). Accessibility of Information on the Web. Nature, 400

107-109

Liddle, W. S., Yau, S. H., Embley, D. W. (2002). On the Automatic Extraction of Data from

ص: 104

,.the Hidden Web. Springer, Berlin Heidelberg New York Liu, X., Maly, K., Zubair, M
Nelson, M. (2002). DP9 – an OAI gateway service for Web crawlers. Paper presented at
the Second ACM/IEEE Joint Conference on Digital Libraries mation Mediation. Paper
presented at the Intl. Workshop on the World–Wide Web and Conceptual Modeling
WWWCM'99), Paris)

Marill, J., Boyko, A., Ashenfelder, M. (2004). Web Harvesting Survey, 10

Masanè s, J. (2002). Archiving the deep web. Paper presented at the 2nd International
Workshop on Web Archives (IWA W'02), Roma, Italy

Mohr, G., Kimpton, M., Stack, M., Ranitovic, I. (2004). Introduction to Heritrix, an
archival quality web crawler. Paper presented at the 4th International
Web Archiving Workshop (IWA W'04), Bath, UK

Ntoulas, A., Zerfos, P., Cho, J. (2005). Downloading textual hidden web content through
keyword queries. Denver, CO, USA sented at the Proceedings of the 27th International
Conference on Very Large Data Bases

.Roche, X. (2006). Copying web sites. In J. Masanè s (Ed.), Web Archiving
Springer, Berlin Heidelberg New York integration of data and its representation. Paper
,presented at the 1st International Workshop on Web Site Evolution (WSE'99), Atlanta
USA

-Zhang, Z., He, B., Chang, K. C.-C. (2004). Understanding Web query interfaces: Best
effort parsing with hidden syntax. Paper presented at the Proceedings of the 2004 ACM
SIGMOD International Conference on Management of Data

,Lagoze, C. Van de Sompel, H. (2001). The open archives initiative: building a Ludäscher

.B. Gupta, A. (1999). Modeling Interactive Web Sources for Infor-Raghavan, S

,Garcia-Molina, H. (2001). Crawling the Hidden Web. Paper pre- Storey, M.-A. Jahnke

J. H. (1999). Web site evolution – Towards a flexible

ص: 105

پژوهش حاضر پژوهشی مفهومی است که با هدف تبیین تاثیر بستر نحوی بر میانکنش پذیری استانداردهای فراداده ای انجام شد جامعه پژوهش عبارت بودند از استانداردهای فراداده ای مارک 21 در بستر زبان نشانه گذاری گسترش پذیر مارک ایکس ام ال استاندارد انتقال و کدگذاری فراداده ها (متس)، "طرح فراداده ای توصیف شیء (مودس) طرح فراداده ای توصیف مستند مدس"، طرح فراداده ای هسته دوبلین (دی) سی ام آی فراداده برای نگهداری اشیای

آی، دیجیتالی (پریمیس)، فراداده فنی برای اشیای دیجیتالی متنی تکست ام. دی. فراداده برای مدیریت تصاویر در بستر ایکس ام ال (میکس)؛ و گردآوری داده ها با روش کتابخانه ای صورت گرفته است. در بخش نخست انواع و ابزارهای میانکنش پذیری استانداردهای فراداده ای توصیف گردید در بخش دیگر با مبنا قرار دادن استاندارد انتقال و کدگذاری فرادادهها (متس) به عنوان استاندارد هسته، نحوه تعامل استانداردهای فراداده ای مورد مطالعه با یکدیگر و با استاندارد متس با رویکرد تحلیلی-سیستمی مورد بررسی قرار گرفته و الگوهایی متناسب با هر یک ترسیم شده اند تحلیل صورت گرفته بیانگر آن است که استفاده از بستر نحوی مناسب در میانکنش پذیری استانداردهای فراداده ای نقش به سزایی ایفا میکند و به یکپارچه سازی درونی و برونی نظامهای اطلاعاتی می انجامد.

کلیدواژه ها بستر، نحوی استانداردهای فراداده ای میانکنش پذیری یکپارچه سازی نظامهای اطلاعاتی

مفهوم میانکنش پذیری بر قابلیت تعامل و کار متقابل میان چند نظام اطلاعاتی با هدف تبادل داده ها و خدمات دلالت دارد اجرای فرایند میانکنش پذیری در راستای یکپارچه سازی درونی و برونی نظام اطلاعاتی با اجزای درونی خود و دیگر نظامهای اطلاعاتی صورت میگیرد و منجر به ارزش افزوده برای نظامهای موجود در فرایند میشود این تعامل در دو سطح نحوی و معنایی رخ میدهد در سطح نحوی، تبادل داده ها بر اساس قالبهای مشترک و یا استفاده از پروتکل های ارتباطی، و در سطح معنایی تفسیر دادههای مبادله شده به صورت معنادار به منظور تولید نتایج مفید همخوان با نیازها و سطح شناختی کاربران مد نظر است از آنجا که صفات و ویژگیهای هر شیء محتوایی (ورودی) در قالب استانداردها و طرحهای فراداده ای به صورت معنادار توصیف شده (پردازش)، و در قالب محصولی جدید به نام پیشینه های فراداده ای بازنمون می گردند فراداده نیز یک نظام اطلاعاتی به شمار میآید. بنابراین همانند دیگر نظامهای اطلاعاتی نیاز به تعامل میان نظامهای فراداده ای برای

نیل به اهداف فرایند میانکنش پذیری بدیهی مینماید، و به "میانکنش پذیری فراداده ای(1)" که نوعی میانکنش پذیری معنایی، است، شهرت یافته است.

به عبارت دیگر با توجه به حجم فراوان اشیای محتوایی منتشر شده در هر یک از حوزه های دانش، بشری و تنوع خدماتی که به کمک پیشرفتهای حوزه فناوریهای اطلاعاتی و ارتباطی برای ارائه این اشیاء امکان پذیر گردیده است و نیز پشتیبانی هر یک از استانداردهای فراداده ای از کارکرد هایی، خاص بهره مندی از طیفی از استانداردهای فراداده ای برای مدیریت اشیای محتوایی و خدمات ارائه شده در نظامهای اطلاعاتی مورد نیاز است و میانکنش پذیری این استانداردها به منظور یکپارچه سازی اجزا و فرایندهای نظام اطلاعاتی ضروری است در سالهای اخیر تلاشهای نظری و کاربردی گسترده ای برای انجام و تسهیل فرایند میانکنش پذیری نظامهای اطلاعاتی به ویژه نظامهای فراداده ای صورت گرفته است. برگزاری همایشهای متعدد با دامنه بین المللی برای تقویت ادبیات این موضوع(2)، و طراحی پروفایلهای کاربردی(3)، عناصر ارتباطی(4)، جداول یا گذرگاههای تطبیقی(5) و نظیر آنها برای اجرای فرایند میان کنش پذیری نشان دهنده اهمیت این موضوع است.

از سوی دیگر فراداده برای بازنمون خود نیاز به بستر نحوی(6) دارد، یعنی ماشین - خوان و ماشین - فهم شدن فراداده منوط به استفاده از بستر نحوی مناسب است. با توجه به این مهم مسئله ای که در اینجا مطرح میگردد آن است که آیا بستر نحوی بر میانکنش پذیری استانداردهای فراداده ای موثر است؟ آیا تعامل میان استانداردهای فراداده ای در محیط این بستر صورت میگیرد؟ انتخاب بسترهای نحوی گوناگون فرایند میانکنش پذیری را تغییر خواهد داد؟ و در پایان این بستر میتواند زمینه را برای مدیریت بهینه فراداده ها و به پیروی از آن اشیای محتوایی به عنوان هدف اصلی یکپارچگی نظام اطلاعاتی فراهم نماید؟

روش شناسی

این پژوهش یک پژوهش مفهومی است که به تبیین تاثیر بستر نحوی بر میانکنش پذیری استانداردهای فراداده ای می پردازد. جامعه پژوهش را استانداردهای فراداده ای قالب " فرادادهای مارک 21 در بستر زبان نشانه گذاری گسترش پذیر (مارک ایکس ام ال)(7) استاندارد انتقال و کدگذاری فرادادهها متس(8)

ص: 108

1- Metadata interoperability

2- از جمله همایش- <http://www.digcur-education.org/eng/Events/Metadata-Harmonization-Bridging-Languages-of-Description>

3- Application profiles

4- Linking Devices

5- Crosswalks or Mapping table

6- بستر نحوی عبارت است از مجموعه ای از، قواعد دستورالعملها و نشانهها برای اعمال ساختاری خاص بر روی محتوای متنی به منظور ذخیره سازی فهم و انجام پردازشهای خاص توسط ماشین (رایانه)

7- (Machine-readable Cataloguing in XML (MARCXML

طرح فراداده‌های توصیف شیء (مودس)(1)، طرح فراداده‌های توصیف مستند (مدس)(2)، طرح فراداده‌ای هسته‌دوبلین (دی سی ام آی)(3) فراداده برای نگهداری اشیای دیجیتال (پریمیس)(4)، فراداده فنی برای اشیای دیجیتال (متی (تکست. ام دی)(5) و فراداده برای مدیریت تصاویر در بستر ایکس ام ال (میکس)(6) تشکیل می‌دهند. در بخش نخست، پژوهش انواع ابزارهای مورد استفاده برای فرایند میانکنش پذیری با تاکید بر نوع میانکنش پذیری که میان استانداردهای فراداده‌ای ایجاد مینمایند، توصیف میشوند تبیین تاثیر بستر نحوی بر میانکنش پذیری استانداردهای فراداده‌ای بخش بعدی و اصلی پژوهش است در این بخش استاندارد "انتقال و کدگذاری فراداده‌ها متس به دلیل قابلیت مدیریت فراداده‌ها و امکان جاسازی دیگر استانداردهای فراداده‌ای درون آن به عنوان استاندارد هسته مد نظر قرار گرفته، و تعامل دیگر استانداردها با یکدیگر و با این استاندارد بر اساس بخشهای هفتگانه آن با رویکرد تحلیلی بررسی میشود همچنین عناصر ارتباطی مورد استفاده برای برقراری تعامل میان استانداردها بر پایه توصیه‌های ارائه شده از سوی استانداردهای مورد مطالعه تعیین می‌شوند برای گردآوری داده‌ها از روش کتابخانهای اسنادی) استفاده شده است و الگوهای تعاملی ارائه شده بر مبنای رویکرد تحلیلی سیستمی طراحی گردیده اند.

انواع و ابزارهای میانکنش پذیری استانداردهای فراداده‌ای

همان طور که پیشتر اشاره شد میانکنش پذیری میان نظامهای اطلاعاتی موجب یکپارچه سازی درونی و برونی آنها شده ارزشهای افزوده فراوانی برای این نظامها به ارمغان می آورد. میانکنش پذیری فراداده‌ای عبارت است از توانایی، نظامها، خدمات و سازمانها در تعامل با یکدیگر، تبادل دادهها، از داده‌های مبادله شده بدون نیاز به تلاشی خاص بر روی نظام مبدأ این فرایند در سه سطح انجام می شود نخست سطح فرامها که در آن سطح عناصر فراداده‌ای مد نظر قرار می گیرند، و از محیط فنی شبکه‌ای سخت افزاری و نرم افزاری مستقل هستند محصول این سطح از فرایند، مجموعه‌ای از عناصر استخراج شده، گذرگاههای تطبیقی پروفایلهای کاربردی ثبت‌های فراداده‌ای(7) هستند. سطح دیگر، به میانکنش پذیری پیشنهادی فراداده‌ای اختصاص دارد در این سطح یکپارچه سازی پیشنهادی فراداده‌ای از طریق همخوانی عناصر از بعد معناشناختی صورت می‌گیرد پیشنهادی تبدیل شده و تولید پیشنهادی جدید با ترکیب ارزشهای عناصر پیشنهادی موجود خروجی سطح پیشنهادی به شمار می آیند و در سطح دو دیگر که به سطح مخازن اطلاعاتی موسوم است رشته‌های ارزشهای برخی عناصر خاص با گردآوری فراداده‌ها از نظامهای مختلف و یکپارچه نمودن آنها همایند میشوند این سطح امکان جستجوی یکپارچه

ص: 109

(Metadata Object Description Schema (MODS -1

(Metadata Authority Description Schema (MADS -2

(Dublin Core Metadata Initiative (DCMI -3

(PREservation Metadata: Implementation Strategy (PREMIS -4

(Technical Metadata for Text (TextMD -5

(Metadata for Images in XML (MIX -6

Metadata registries -7

میان چند نظام اطلاعاتی را فراهم مینماید (معارف (1) و یحیی (2)، 2009؛ هیروید (3)، 2011).

برای اجرای فرایند میانکنش پذیری فراداده ای ابزارهای گوناگونی طراحی شده است. پروفایلهای کاربردی عناصر، ارتباطی گذرگاهها یا جداول تطبیقی و بستر، نحوی ابزارهای فرایند میانکنش پذیری فراداده ای محسوب میشوند. پروفایلهای کاربردی مجموعه عناصر فراداده ای (استخراج شده از یک یا چند استاندارد فراداده ای) خط مشی، ها تجربیات برتر و رهنمودهایی که به منظور کاربردهای خاص (محلی) تعریف شده است، یا اعلان ضوابطی که یک سازمان یک منبع اطلاعاتی، یک برنامه کاربردی، یا جامعه استفاده کنندگان در به کارگیری فراداده هایشان استفاده می کنند (طرح فراداده ای هسته دویلین 2013) و میانکنش پذیری سطح فرامها را پشتیبانی میکنند عناصر ارتباطی به صفات یا خصایص اشیای محتوایی مانند موضوع پدیدآورنده، ناشر و مانند آن گفته می شود که ارتباط میان چند شیء محتوایی را برقرار مینمایند و به میانکنش پذیری سطح پیشینه ها و نیز مخازن سازمانی می انجامد جداول یا گذرگاههای تطبیقی به جداولی اطلاق میشود که عناصر معادل در بیش از یک استاندارد فراداده ای نشان میدهند، و همانند پروفایلهای کاربردی موجب میانکنش پذیری در سطح فرامها میشوند.

و اما بستر، نحوی میانکنش پذیری استانداردهای فراداده ای در سطح فرامها انجام میشود. هر استاندارد فراداده ای دارای فرامایی ویژه است که میزان سازگاری پیشینههای تولید شده بر مبنای آن استاندارد را اعتبار سنجی مینماید استانداردهای فراداده ای مجموعه ای از عناصر مرتبط و ساختارمند از لحاظ معناشناختی هستند که برای پشتیبانی از کارکردهای خاص و متناسب با نیازهای جامعه کاربران خود طراحی شدهاند (فتاحی و طاهری، 1388) این استانداردها برای پیاده سازی پیشینههای مبتنی بر خود یک یا چند قالب ذخیره سازی و نمایش دادهها را به عنوان بستر نحوی بر میگزینند. طیف گسترده ای از قالبهای ذخیره سازی وجود دارد که برخی مبتنی بر پایگاه داده ها (4) و برخی مبتنی بر فایل (5) هستند.

مهمترین این قالبها، زبانهای نشانه گذاری (اس. جی. ام. ال اچ تی ام ال.، و ایکس. ام. ال.)، قالب مدارک قابل انتقال پی دی اف با استفاده از چارچوب توصیف منبع (6)، قالب متن تکست)، و قالب بومی نظامهای مدیریت پایگاههای داده ای دی. بی. ام. اس میباشند طاهری، (1391). هر یک از این قالبها قابلیتهای خاصی برای ذخیره سازی و نمایش دادهها دارند و بر پایه مقاصد خاصی تولید شده اند بنابراین انتخاب آنها از سوی استانداردهای فراداده ای میبایست سازگار با کارکردهای خاص آنها باشد. افزون بر این، به دلیل آن که نظامهای اطلاعاتی برای مدیریت محتوا و خدمات خود از چند استاندارد فراداده ای به طور همزمان استفاده میکنند تعامل میان این استانداردها در جهت نیل به اهداف نظام ضروری است. از این رو این ویژگی نیز در انتخاب بستر نحوی حائز اهمیت است.

ص: 110

Maarof -1

Yahya -2

Hirwade -3

Database-based format -4

File-based format -5

(Resource Description Framework (RDF -6

زبان نشانه گذاری فرامتن (اچ تی ام ال) قالبی برای توصیف ساختار صفحات وب به منظور نمایش آنهاست از مهمترین قابلیت‌های این زبان امکان استفاده از فناوری فرایوند، و ذخیره داده های چند رسانه ای است اما در طراحی این قالب انتقال دادهها مورد اقبال نبوده است به همین دلیل تعداد برچسبها و فرابر چسبهای آن محدود و از پیش تعریف شده هستند و نمیتوان آنها را گسترش داد. توصیف دادههای ذخیره شده در این قالب به ویژگیهای نرم افزاری نظام اطلاعاتی که از اچ تی ام ال بهره میبرد وابسته است. این ویژگی تعامل استانداردهای فراداده ای که پیاده سازی پیشینههای خود در بستر این زبان را توصیه میکنند محدود مینماید (کنسرسیوم وب جهانی، 2012). کنسرسیوم قالب مدارک قابل انتقال (پی دی اف) برای بازنمون اشیای محتوایی به صورت مستقل از سخت افزار نرم افزار و نظام عامل طراحی شده است. هنگامی که حفظ ویژگیهای صفحه آرای یک شیء دیجیتالی ذخیره شده در قالب الکترونیکی دیگر و یا یک شیء آنالوگ مطرح باشد، از قالب پی دی. اف. استفاده میشود (ویکی پدیا 2013) بنابراین یکی از بهترین قالبها برای تهیه نسخه چاپی از شیء دیجیتالی است. اگر چه این قالب به دلیل حفظ ویژگیهای صفحه آرای مستقل از پلت فرم میباشد اما دادههای ذخیره شده در آن به صورت معنادار توصیف نمی شوند و صرفاً تصویری از شیء تبدیل شده به قالب پی دی. اف. قلمداد میشوند به عبارت دیگر در انتقال دادهها از نظامی به نظام دیگر ساختار دادههای ذخیره شده در آن چندان قابلیت پردازشی ندارد و هدف اصلی آن همانند قالب اچ تی ام ال نمایش داده هاست از آنجا که در پیشینههای فراداده توصیف معنادار عناصر و روابط میان آنها از اهمیت فراوانی برخوردار است این قالب چندان مورد توجه بافت فراداده ای واقع نشد.

قالب متن، (تکست) برای ذخیره سازی دادهها بدون استفاده از نشانه ها یا اعمال ساختاری خاص طراحی شده است. دادهها ذخیره شده در این قالب به دلیل عدم وجود نشانه های اضافی در آن حجم بسیاری کمی را اشغال می. کنند در برخی از مواقع با افزودن نشانههایی به دادههای ذخیره شده در این قالب میتوان پردازشهای خاصی بر روی آن اعمال نمود. مهمترین ضعف این قالب در فرایند میانکنش پذیری ساختارمند نبودن و عدم توصیف دادههای ذخیره شده در آن است.

قالب بومی نظامهای مدیریت پایگاههای داده ای (دی بی ام اس) هر نظام مدیریت پایگاه داده ای از یک قالب محلی و بومی متناسب با ویژگیها و قابلیت‌های فنی خود سود میبرد. این قالب بر اساس اهداف کارکردها و همچنین لحاظ منحصر به فرد بودن نظام طراحی شده است. دادههای ذخیره شده در قالب بومی یک دی بی. ام اس در یک نظام مدیریت پایگاه داده ای دیگر قابل پردازش نیستند و حتماً بر روی آنها فرایند تبدیل به قالب جدید صورت گیرد. این ویژگی، با توجه به این که نظامهای اطلاعاتی از پلت فرمی خاص و در نتیجه نظامهای مدیریت پایگاه دادهای متفاوت استفاده میکنند باعث عدم استفاده از قالبهای بومی در فرایند تعامل میان نظام شده است.

زبان نشانه گذاری گسترش پذیر (ایکس ام ال) یک قالب ساده مبتنی بر متن است که به عنوان استاندارد بین المللی برای بازنمون دادههای ساختارمند نظیر اشیای محتوایی و تبادل و اشتراک دادهها گسترش یافته است (بری و دیگران 2008) دادههایی که در قالب ایکس ام ال نشانه گذاری می شوند.

به داده هایی ساختارمند، تبدیل و اشیای محتوایی خود-توصیف (1) بوجود میآورند. این ویژگی موجب استقلال اشیای محتوایی مبتنی بر ایکس ام ال از هر پلت فرمی، شده تبادل آنها را میان نظامهای ناهمگن (2) ممکن، و بنابراین میانکنش پذیری نظامهای اطلاعاتی را باعث می گردد. ایکس. ام. ال. همانند اچ تی ام ال یک مجموعه ثابت از برچسبها نیست با استفاده از این استاندارد، کاربران میتوانند برچسبهای مورد نیاز خود را، تعریف و در محیطهای اطلاعاتی دیگر استفاده کنند قابلیتهای منحصر به فرد این زبان گرایش طراحان استانداردهای فراداده ای به استفاده از این زبان به عنوان بستر نحوی پیشینههای فراداده ای را در پی داشته است افزون بر آن که پیاده سازی برخی از استانداردها مانند مارک 21 که تا پیش از این در قالب زبانهای نشانه گذاری ممکن نبود در قالب زبانهای نشانه گذاری با بکارگیری ایکس ام ال مهیا شده است (طاهری، 1387 ب؛ کین (3)، 2000؛ گیجی (4) و کلی (5)، 2006)، ساختارمند بودن و خود توصیف بودن آن میانکنش پذیری نظامها و استانداردهای فراداده ای را تسهیل نموده است (طاهری، 1391)

در ادامه مقاله و در قالب طراحی الگوهایی، چند تاثیر بستر نحوی بر فرایند میانکنش پذیری استانداردهای فراداده ای که امکان استفاده از طیفی از استانداردها در یک نظام اطلاعاتی به صورت همزمان را توجیه میکند با رویکرد تحلیلی سیستمی مورد بررسی قرار میگیرد.

تبیین تاثیر بستر نحوی بر میانکنش پذیری استانداردهای فراداده ای

در این بخش از، مقاله با استفاده از استاندارد انتقال و کدگذاری فرادادهها (متس) به عنوان استاندارد هسته تعامل دیگر استانداردهای فراداده ای با این استاندارد و در هنگام لزوم تعامل دیگر استانداردها با یکدیگر، با ارائه چند الگو نشان داده میشود دلیل انتخاب استاندارد متس به عنوان استاندارد، هسته کارکرد اصلی آن یعنی مدیریت فرادادههاست (طاهری 1387 الف) متس همانند بسته ای عمل مینماید که میتواند دیگر استانداردهای فراداده ای با کارکردهای گوناگون را در بر گرفته به مدیریت یکپارچه اشیای محتوایی پردازد. (6) استاندارد متس دارای هفت بخش است هر یک از این بخشها دارای کارکردی خاصی هستند. برخی برای جاسازی طرحهای فراداده ای و برخی برای مدیریت محتوا طراحی شده اند. در عین تمامی این بخشها قابلیت تعامل با یکدیگر دارند و میانکنش پذیری این بخشها نیز بر اهمیت متس می افزاید.

ص: 112

Self-description -1

Heterogeneous or Disparate systems -2

Qin -3

Gigee -4

Kelly -5

6- برای اطلاعات بیشتر در مورد کارکردها و دیگر ویژگیهای استانداردهای فراداده ای به این دو منبع مراجعه کنید سید مهدی طاهری 1387 طراحی یک کتابخانه دیجیتالی استاندارد در مجموعه مقالات نخستین همایش کتابخانه های دیجیتالی به کوشش شرکت پارس آذرخش: تهران سبزان سید رحمت الله فتاحی سید مهدی طاهری 1388 فهرست نویسی رایانه ای، مفاهیم شیوهها و کاربرد نرم افزارهای رایانه ای در سازماندهی اطلاعات با همکاری فرشته ناقدی احمدی تهران کتابدار

این بخشها به ترتیب عبارتند از بخش سرپیشینه(1)، بخش فراداده های توصیفی بخش فراداده های مدیریتی، بخش مربوط به فایلها بخش نقشه های ساختاری بخش پیوندهای ساختاری و بخش رفتارهای شیء (طاهری 1387 الف؛ دفتر استانداردهای مارک و توسعه، شبکه C2013). تعامل هر یک از استانداردها توسط عناصر ارتباطی، و در قالب بخشهای هفتگانه متس صورت می گیرد. هر پیشینه فراداده ای به دو روش با پیشینه متس تعامل ایجاد میکند نخست روش درونی که پیشینه یاد شده درون پیشینه متس به دو صورت داده های کدگذاری شده با ایکس ام ال (توسط برچسب) و داده های مبتنی بر کدهای دودویی یا متن خام توسط برچسب درج (جاسازی)(2) میشود و دیگر تهیه پیوند توسط یو. آر. آی. یک پیشینه یا دیگر شناسگرها (پی یو آر ال ای. آر. کی، و دی. ا. آی.) از درون عنصر مرتبط متس به پیشینه فراداده ای مبتنی بر استاندارد فراداده ای دیگر لازم به ذکر است امکان درج پیشینه های بیش از استاندارد در هر بخش از پیشینه های متس وجود دارد. پیشینه های تولید شده بر مبنای هر استاندارد فراداده ای دارای یک عنصر ریشه هستند این عنصر نقش ارتباطی را برای ارتباط با پیشینه های متس ایفا می کند در ادامه شیوه تعامل هر یک از استانداردها و عناصر ارتباطی آنها مورد تحلیل قرار میگیرد.

قالب فراداده های مارک 21 در بستر زبان نشانه گذاری گسترش پذیر (مارک ایکس ام ال).

این قالب توسط دفتر استانداردهای مارک و توسعه شبکه(3) کتابخانه کنگره آمریکا به منظور پیاده سازی داده های مارک در بستر ایکس ام ال طراحی شده است. انعطاف پذیری و گسترش پذیری این چارچوب امکان پاسخگویی به نیازهای گوناگون و خاص کاربران را میسر ساخته است (دفتر استانداردهای مارک و توسعه شبکه، 2013a). وجود عناصر متعدد باعث شده قالب مارک از چند کارکرد به صورت موثر پشتیبانی نماید دو کارکرد مدیریت و توصیف کارکردهای اصلی قالب مارک هستند. در ذیل نحوه میانکنش پذیری قالب مارک 21 در بستر زبان نشانه گذاری گسترش پذیر بر اساس دو کارکرد مدیریتی و توصیفی با استاندارد متس ترسیم شده است.

به عنوان فراداده توصیفی

عنصر ریشه (record) پیشینه مبتنی بر مارک 21 برای پشتیبانی از کارکرد توصیفی در بخش فراداده های توصیفی متس با برچسب بر پایه روش درونی در برچسب، و پیوند به پیشینه مارک بر پایه روش برونی در عنصر جاسازی می شود چنان چه روش درونی مد نظر باشد، داده های کدگذاری شده در قالب ایکس ام ال در برچسب، و داده های در قالب دودویی یا متن خام در برچسب binData قرار میگیرند دیگر استانداردهای فراداده ای با کارکرد توصیفی نیز به همین صورت با استاندارد متس تعامل برقرار میکنند.

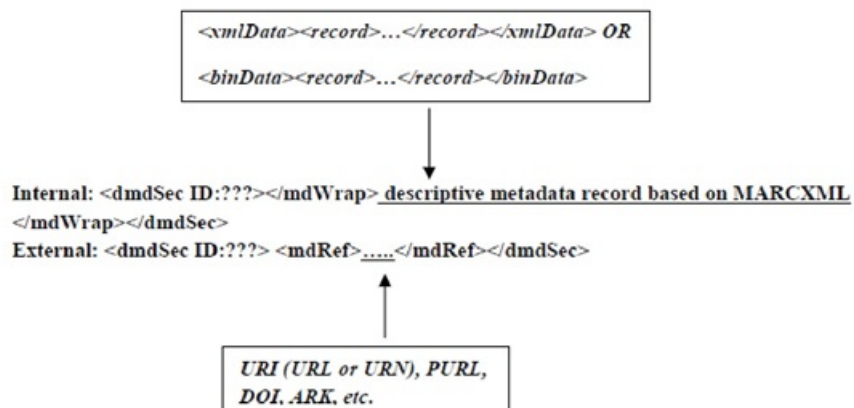
ص: 113

الگوی 1. شیوه درج پیشینه مارک 21 با کارکرد توصیفی در بخش فراداده‌های توصیفی پیشینه متس

الگوی 1 نشان می‌دهد عنصر ریشه یک پیشینه مارک 21 برای پشتیبانی از کارکرد توصیفی قادر است به ترتیب در عناصر ???
<mdWrap> یا یا

به عنوان فراداده مدیریتی

عکس



الگوی ۱. شیوه درج پیشینه مارک ۲۱ با کارکرد توصیفی در بخش فراداده‌های توصیفی پیشینه متس

الگوی ۱ نشان می‌دهد، عنصر ریشه یک پیشینه مارک ۲۱ برای پشتیبانی از کارکرد توصیفی قادر است به ترتیب در عناصر `<dmdSec ID:???'</mdWrap>`، `<mdRdf>` یا `<xmlData>` یا `<binData>`، و `<record>` جاسازی شود.

به عنوان فراداده مدیریتی

عنصر ریشه پیشینه مارک ۲۱ با کارکرد مدیریتی در بخش فراداده‌های مدیریتی استاندارد متس (با برچسب `<amdSec>`)، بر پایه روش درونی در برچسب `<mdWrap>`، و پیوند به پیشینه مارک بر پایه روش برونی در عنصر `<mdRef>` درج می‌گردد. پیشینه‌های فراداده‌ای مبتنی بر استانداردهای با کارکرد مدیریتی در عناصر `<techMD>` (برای استانداردهای فراداده‌ای با کارکرد مدیریت فنی)، `<rightsMD>` (برای استانداردهای فراداده‌ای با کارکرد مدیریت حقوق معنوی)، `<sourceMD>` (برای استانداردهای فراداده‌ای با کارکرد مدیریتی و توصیفی مربوط به اشیای آنالوگ)، و `<digiprovMD>` (برای استانداردهای فراداده‌ای با کارکرد مدیریت اشیای با منشاء دیجیتالی) استاندارد متس جاسازی می‌شوند. دیگر استانداردهای فراداده‌ای با کارکرد مدیریتی با توجه به کارکرد فرعی خاص خود همانند پیشینه‌های مارک ۲۱ در ایکس.ام.ال. با استاندارد متس تعامل پیدا می‌کنند.

عنصر ریشه پیشینه مارک 21 با کارکرد مدیریتی در بخش فراداده‌های مدیریتی استاندارد متس (ب) برچسب `amdSec` (ب) بر پایه روش درونی در برچسب `mdWrap` و پیوند به پیشینه مارک بر پایه روش برونی در عنصر درج می‌گردد. پیشینه‌های فراداده‌ای مبتنی بر استانداردهای با کارکرد مدیریتی در عناصر `techMD` برای استانداردهای فراداده‌ای با کارکرد مدیریت فنی `rights` برای استانداردهای فراداده‌ای با کارکرد مدیریت حقوق معنوی `sourceMD` (برای استانداردهای فراداده‌ای با کارکرد مدیریتی و توصیفی مربوط به اشیای آنالوگ و برای) استانداردهای فراداده‌ای با کارکرد مدیریت اشیای با منشاء (دیجیتالی) استاندارد متس جاسازی می‌شوند. دیگر استانداردهای فراداده‌ای با کارکرد مدیریتی با توجه به کارکرد فرعی خاص خود همانند پیشینه‌های مارک 21 در ایکس ام ال با استاندارد

متس تعامل پیدا می کنند.

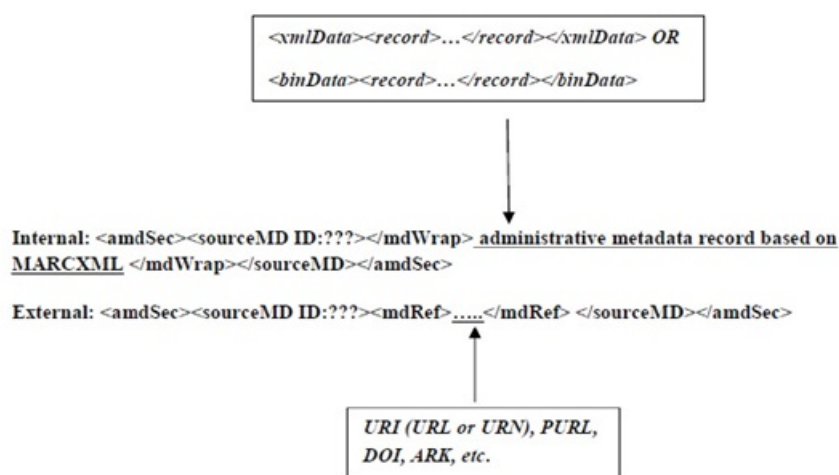
ص: 114

الگوی 2 شیوه درج پیشینه مارک 21 با کارکرد مدیریتی در بخش فراداده های مدیریتی پیشینه متس

چنان چه در الگوی 2 مشاهده میشود پیشینه مبتنی بر مارک 21 برای ایفای کارکرد مدیریتی میتواند به ترتیب درون برچسبهای </mdWrap KamdSec ID:??> یا ، یا

عکس

بررسی تأثیر بستر نحوی ... ۱۱۵



الگوی ۲. شیوه درج پیشینه مارک ۲۱ با کارکرد مدیریتی در بخش فراداده های مدیریتی پیشینه متس

چنان چه در الگوی ۲ مشاهده می شود، پیشینه مبتنی بر مارک ۲۱ برای ایفای کارکرد مدیریتی، می تواند به ترتیب درون برچسب های <amdSec ID:??>، </mdWrap> یا <mdRdf>، <xmlData> یا <binData>، و <record> درج شود. طرح فراداده های توصیف شیء (مودس) این طرح برای مجموعه عناصر کتابشناختی که با اهداف گوناگون، به خصوص کاربردهای کتابخانه-ای استفاده می شوند، در بستر زبان نشانه گذاری گسترش پذیر تهیه شده است. مودس امکان انتقال داده های کتابشناختی گزیده از پیشینه های موجود مارک و ایجاد پیشینه های توصیفی برای اشیای محتوایی جدید را فراهم می آورد. طرح فراداده ای مودس مجموعه ای از عناصر مورد نیاز برای توصیف اشیای دیجیتالی که از فیله های مارک ۲۱ استخراج شده است را در بر می گیرد. نام های برچسب عناصر مودس بر خلاف مارک ۲۱ در قالب ایکس.ام.ال. مبتنی بر واژگان زبان طبیعی هستند (مک کالم^۱، ۲۰۰۴؛ دفتر استاندارد های مارک و توسعه شبکه، ۲۰۱۳). کارکرد اصلی مودس، کارکرد توصیفی است و پیشینه های آن در بخش فراداده های توصیفی متس درج می شود.

این طرح برای مجموعه عناصر کتابشناختی که با اهداف گوناگون به خصوص کاربردهای کتابخانه-ای استفاده میشوند در بستر زبان نشانه گذاری گسترش پذیر تهیه شده. است مودس امکان انتقال داده های کتابشناختی گزیده از پیشینه های موجود مارک و ایجاد پیشینه های توصیفی برای اشیای محتوایی جدید را فراهم می آورد طرح فراداده ای مودس مجموعه ای از عناصر مورد نیاز برای توصیف اشیای دیجیتالی که از فیلهای مارک 21 استخراج شده است را در بر میگیرد نامهای برچسب عناصر مودس بر خلاف مارک 21 در قالب ایکس ام ال مبتنی بر واژگان زبان طبیعی هستند (مک کالم (1)، 2004 دفتر استانداردهای مارک و توسعه شبکه e2013). کارکرد اصلی مودس کارکرد توصیفی است و پیشینه های آن در بخش فراداده های توصیفی متس درج میشود.

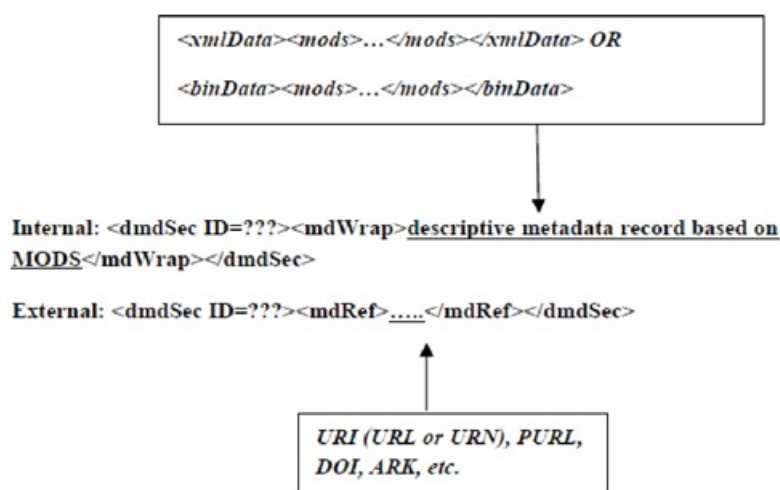
ص: 115

الگوی 3. شیوه درج پیشینه مودس در بخش فراداده‌های توصیفی پیشینه متس

طرح فراداده‌های هسته دوبلین (دی. سی. ام. آی.)

عکس

۱۱۶ مدیریت منابع اطلاعاتی وب



الگوی ۳. شیوه درج پیشینه مودس در بخش فراداده‌های توصیفی پیشینه متس

طرح فراداده‌های هسته دوبلین (دی. سی. ام. آی.)

طرحی بین‌المللی و میان رشته‌ای که مجموعه عناصری ساده و کارآمد برای توصیف طیف گسترده‌ای از اشیای محتوایی شبکه‌ای ارائه می‌دهد. این طرح نخستین تلاش جدی در حوزه طراحی استانداردهای فراداده‌ای پس از تعمیم شبکه جهانی وب محسوب می‌شود. کارکرد اصلی طرح هسته دوبلین نیز همانند طرح مودس، کارکرد توصیفی است. قالب ایکس. ام. ال. یکی از بسترهای نحوی هسته دوبلین است، و امکان پیاده‌سازی پیشینه‌های این طرح در قالب‌هایی دیگر نیز وجود دارد (جانستون^۱ و پاول^۲، ۲۰۰۶: فتاحی و طاهری، ۱۳۸۸).

1. Juhnston
2. Powell

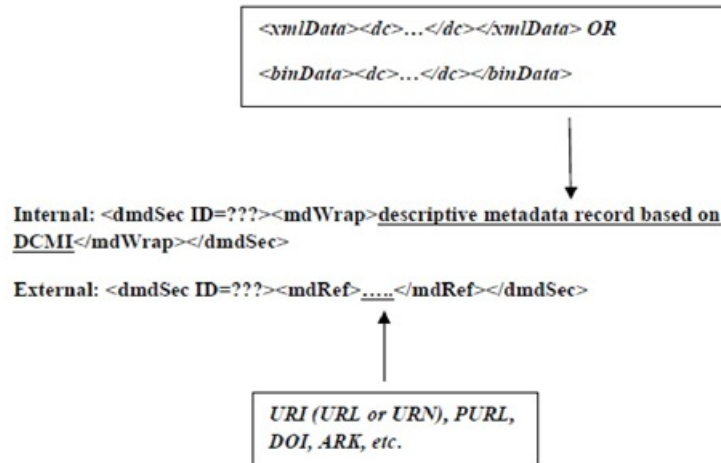
طرحی بین‌المللی و میان رشته‌ای که مجموعه عناصری ساده و کارآمد برای توصیف طیف گسترده‌ای از اشیای محتوایی شبکه‌ای ارائه

میدهد این طرح نخستین تلاش جدی در حوزه طراحی استانداردهای فراداده ای پس از تعمیم شبکه جهانی وب محسوب میشود کارکرد اصلی طرح هسته دویلین نیز همانند طرح مودس کارکرد توصیفی است قالب .ایکس ام ال یکی از بسترهای نحوی هسته دویلین است، و امکان پیاده سازی پیشنهادهای این طرح در قالبهایی دیگر نیز وجود دارد (جانستون(1) و پاول(2)، 2006؛ فتاحی و طاهری، 1388).

ص: 116

Juhnston -1

Powell -2



الگوی ۴. شیوه درج پیشینه هسته دوبلین در بخش فراداده‌های توصیفی پیشینه متس

طرح فراداده‌های توصیف مستند (مدس)

مدس شامل مجموعه عناصری برای توصیف داده‌های مستند مربوط به اشخاص حقیقی و حقوقی، رویدادهای مهم، شناسه‌های موضوعی و جغرافیایی، و نظیر آن است، و به عنوان مکمل طرح فراداده‌ای توصیف شیء طراحی شده است (دفتر استانداردهای مارک و توسعه شبکه، ۲۰۱۳). با این وجود می‌تواند برای مستندسازی ارزش‌های عناصر دیگر طرح‌های فراداده‌ای با کارکرد مدیریتی و توصیفی مانند هسته دوبلین نیز استفاده شود. طرح مدس به صورت مستقیم در پیشینه‌های متس درج نمی‌شود و یا پیوند نمی‌یابد، بلکه به صورت غیر مستقیم و پیوند با طرح‌های فراداده‌ای با کارکرد مدیریتی یا توصیفی با متس تعامل برقرار می‌کند.

```
<?xml version="1.0" encoding="UTF-8"?><mads xmlns=http://www.loc.gov/mads/
xmlns:mods="http://www.loc.gov/mods/v3" xmlns:xlink="http://www.w3.org/1999/xlink"
xmlns:xsi="http://www.w3.org/2001/XMLSchema-instance" xsi:schemaLocation="http://
www.loc.gov/mads/mads.xsd">
  <authority><name><namePart>Smith, John</namePart><namePart
type="date">1995-</namePart></name></authority><variant type
="other"><name><namePart>Smith, J</namePart></name></variant>
  <variant type="other"><name><namePart>Smith, John J</namePart>
```

الگوی ۴. شیوه درج پیشینه هسته دوبلین در بخش فراداده‌های توصیفی پیشینه متس

طرح فراداده‌های توصیف مستند (مدس)

مدس شامل مجموعه عناصری برای توصیف داده‌های مستند مربوط به اشخاص حقیقی حقوقی رویدادهای مهم شناسه‌های موضوعی و

جغرافیایی و نظیر آن است و به عنوان مکمل طرح فراداده ای توصیف شیء طراحی شده است دفتر استانداردهای مارک و توسعه شبکه b2013 با این وجود میتواند برای مستندسازی ارزشهای عناصر دیگر طرحهای فراداده ای با کارکرد مدیریتی و توصیفی مانند هسته دویلین نیز استفاده شود. طرح مدس به صورت مستقیم در پیشنهادها متس درج نمی شود و یا پیوند نمی یابد، بلکه به صورت غیر مستقیم و پیوند با طرحهای فراداده ای با کارکرد مدیریتی یا توصیفی با متس تعامل برقرار میکند.

"xmlns:mods="http://www.loc.gov/mods/v3" xmlns:xlink="http://www.w3.org/1999/xlink

//:xmlns:xsi="http://www.w3.org/2001/XMLSchema-instance xsi:schemaLocation="http

<"www.loc.gov/mads/mads.xsd

,Smith

-type="date " > 1995

John

type

other">Smith, J"=

Smith, John J

ص: 117

```
</name></variant><notetype="history">BiographicalnoteaboutJohnSmith.</note><affiliation><organization>Lawrence Livermore Laboratory</organization><dateValid>1987</dateValid></affiliation></mads>
```

نمونه ۱. نمونه‌ای از یک پیشینه مَدَس مربوط به یک شخص حقیقی

```
Internal: <dmdSec ID="???"><mdWrap><mods><name type="personal"><namePart type="termsOfAddress">Dr.</namePart> <namePart>Smith, John</namePart> </name></mods></mdWrap></dmdSec>
```

پیشینه مبتنی بر استاندارد مَدَس

```
Internal: <dmdSec ID="???"><mdWrap><mods><name type="personal"> <namePart type="termsOfAddress">Dr.</namePart> <namePart>Smith, John</namePart> </name></mods></mdWrap></dmdSec>
```

پیشینه مبتنی بر استاندارد مَدَس

الگوی ۵. شیوه تعامل غیر مستقیم پیشینه‌های مَدَس با پیشینه متس به وسیله پیشینه مودس

پیشینه مبتنی بر استاندارد مَدَس در پایگاه مستند با استفاده از شناسگر پیشینه^۱ با عنصر (فیلد) ارتباطی پیشینه مودس در پایگاه کتابشناختی که فقط ارزش‌های کد شده می‌پذیرد، پیوند می‌یابد، و فرایند کنترل مستندات را پشتیبانی می‌کند. بنابراین میان پیشینه‌های مَدَس و مودس به صورت مستقیم، و میان پیشینه‌های مَدَس و متس ارتباط غیر مستقیم برقرار می‌شود، و بدین گونه کارکرد کنترل مستندات در نظام اطلاعاتی وجود خواهد داشت.

فراداده برای نگهداری اشیای دیجیتالی (پریمیس)

مجموعه عناصر مبتنی بر ایکس.ام.ال. که با هدف ثبت فراداده‌های مربوط به نگهداری شیء دیجیتالی در کتابخانه یا دیگر مجموعه‌های دیجیتالی گسترش یافته است. بنابراین کارکرد اصلی استاندارد پریمیس، نگهداری اشیای دیجیتالی است (هابینگ^۲، ۲۰۰۸). پیشینه‌های مبتنی بر این استاندارد بر اساس نوع

1. RecordID
2. Habing

.BiographicalnoteaboutJohn Smith

iliation> 1987

<dateValid

نمونه 1. نمونه ای از یک پیشینه مدس مربوط به یک شخص حقیقی

:Internal

<??=ID

Dr: Smith, John

<namePart

پیشینه مبتنی بر استاندارد مدس

:Internal

type="termsOfAddress">Dr.Smith, John

پیشینه مبتنی بر استاندارد مدس

الگوی 5 شیوه تعامل غیر مستقیم پیشینه‌های مدس با پیشینه متس به وسیله پیشینه مودس

پیشینه مبتنی بر استاندارد مدس در پایگاه مستند با استفاده از شناسگر پیشینه(1) با عنصر (فیلد) ارتباطی پیشینه مودس در پایگاه کتابشناختی که فقط ارزشهای کد شده می پذیرد، پیوند می یابد، و فرایند کنترل مستندات را پشتیبانی می کند بنابراین میان پیشینه‌های مدس و مودس به صورت مستقیم و میان پیشینه های مدس و متس ارتباط غیر مستقیم بر قرار میشود و بدین گونه کارکرد کنترل مستندات در نظام اطلاعاتی وجود خواهد داشت.

فراداده برای نگهداری اشیای دیجیتالی (پریمیس)

مجموعه عناصر مبتنی بر ایکس ام ال که با هدف ثبت فراداده‌های مربوط به نگهداری شیء دیجیتالی در کتابخانه یا دیگر مجموعه های دیجیتالی گسترش یافته است بنابراین کارکرد اصلی استاندارد پریمیس نگهداری اشیای دیجیتالی است (هابینگ(2)، 2008) پیشینه‌های مبتنی بر این استاندارد بر اساس نوع

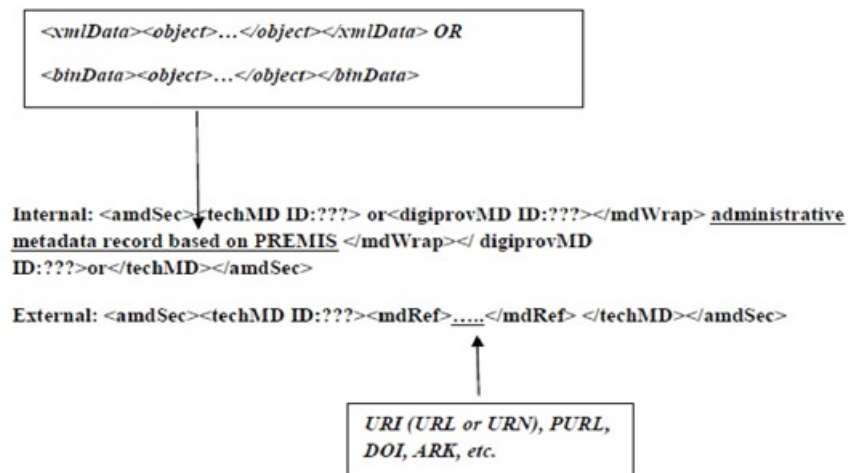
ص: 118

RecordID-1

Habing-2

بررسی تأثیر بستر نحوی ... ۱۱۹

موجودیتی که در بر می‌گیرند، باید در بخش فراداده‌های مدیریتی متس و در عناصر <techMD> و <digiprovMD> درج گردند.



الگوی ۶. شیوه درج پیشینه پرمیس با کارکرد مدیریتی در بخش فراداده‌های مدیریتی پیشینه متس

```

<mets:techMD ID="file-2"><mets:mdWrap MDTYPE="PREMIS">
  <mets:xmlData>
    <premis:object>
      <premis:objectIdentifier>
        <premis:objectIdentifierType>uri</premis:objectIdentifierType>
        <premis:objectIdentifierValue>info:lna/pic/vn3579101-c</premis:objectIdentifierValue>
      </premis:objectIdentifier>
      <premis:preservationLevel>unknown</premis:preservationLevel>
      <premis:objectCategory>file</premis:objectCategory>
      <premis:objectCharacteristics>
        <premis:format>
          <premis:formatDesignation>
            <premis:formatName>image/tiff</premis:formatName>
            <premis:formatVersion>6.0</premis:formatVersion>
          </premis:formatDesignation>
        </premis:format>
      </premis:objectCharacteristics>
      ...
    </premis:object>
  </mets:xmlData>
</mets:mdWrap>
</mets:techMD>
<mets:digiprovMD ID="event-1"><mets:mdWrap MDTYPE="PREMIS">
  <mets:xmlData>
    <premis:event>
      <premis:eventIdentifier>
        <premis:eventIdentifierType>internal</premis:eventIdentifierType>
        <premis:eventIdentifierValue>20903-1</premis:eventIdentifierValue>
      </premis:eventIdentifier>
      <premis:eventType>creation</premis:eventType>
      <premis:eventDateTime>2005-11-03T12:15:59</premis:eventDateTime>
    </premis:event>
  </mets:xmlData>
</mets:mdWrap></mets:digiprovMD>
  
```

شکل ۱. نمونه‌ای از پیشینه‌های مبتنی بر پرمیس جاسازی شده در پیشینه متس

موجودیتی که در بر میگیرند باید در بخش فراداده‌های مدیریتی متس و در عناصر درج گردند.

الگوی ۶. شیوه درج پیشینه پرمیس با کارکرد مدیریتی در بخش فراداده‌های مدیریتی پیشینه متس

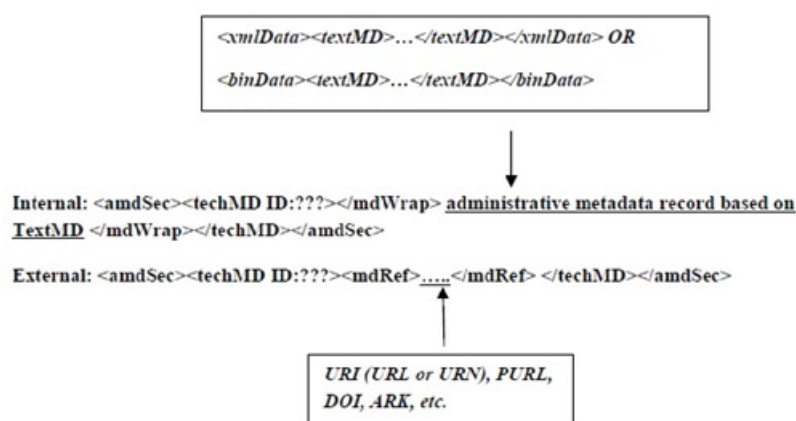
<Spremis:object

```
<Spremis:objectIdentifierType>uris/premis:objectIdentifierType
<Spremis:objectIdentifierValue>info:nla/nla.pic-vn 3579101-c premis:objectIdentifierValue
    <premis:objectIdentifier
        unknown premis:preservationLevel> file
    <Spremis:objectCharacteristics
        ?Spremis:format
        ?Spremis:formatDesignation
        image/tiff
        6.0
    <premis:format
    <premis:objectCharacteristics
    <premis:object
    <Spremis:event
    <Spremis:eventIdentifier
        internal
    <premis:eventIdentifierValue 28903-1
    <Spremis:eventIdentifier
        creation
    <Spremis:eventDateTime>2005-11-03T12:15:59
    <premis:event
```

شکل 1. نمونه ای از پیشینه‌های مبتنی بر پرمیس جاسازی شده در پیشینه متس

فرا داده فنی برای اشیای دیجیتالی متنی (تکست ام. دی.)

استانداردی فراداده‌ای و مبتنی بر زبان نشانه‌گذاری گسترش‌پذیر که عناصری برای توصیف جنبه‌های فنی اشیای دیجیتالی متنی ارائه می‌دهد. پیشنهادها مبتنی بر این استاندارد می‌توانند مستقیماً در بخش فراداده‌های مدیریتی متس، و یا به طور غیر مستقیم در عنصر `<additionalTechnicalCharacteristics>` مربوط به موجودیت شیء^۱ استاندارد پریمیس درج گردند (دفتر استانداردهای مارک و توسعه شبکه، ۲۰۱۳).



الگوی ۷. شیوه درج پیشنهاد تکست ام. دی. در بخش فراداده‌های مدیریتی پیشنهاد متس

فرا داده برای مدیریت تصاویر در بستر ایکس. ام. ال. (میکس)

میکس محصول دفتر استانداردهای مارک و توسعه شبکه، با همکاری کمیته استانداردهای فراداده‌ای فنی برای تصاویر ثابت وابسته به «سازمان استانداردهای اطلاعات ملی (NISO)» است که به منظور مدیریت تصاویر دیجیتالی توسعه یافته است. کارکرد اصلی این طرح مدیریت فنی تصاویر ثابت دیجیتالی است (دفتر استانداردهای مارک و توسعه شبکه، ۲۰۱۳)، و پیشنهادها آن در بخش فراداده‌های مدیریتی متس در عنصر `<techMD>` قرار می‌گیرند.

فرا داده فنی برای اشیای دیجیتالی متنی (تکست ام. دی.)

استانداردی فراداده‌ای و مبتنی بر زبان نشانه‌گذاری گسترش‌پذیر که عناصری برای توصیف جنبه‌های فنی اشیای دیجیتالی متنی ارائه می‌دهد پیشنهادها مبتنی بر این استاندارد می‌توانند مستقیماً در بخش فراداده‌های مدیریتی متس، و یا به طور غیر مستقیم در عنصر `additionalTechnicalCharacteristics` مربوط به موجودیت شیء⁽¹⁾ استاندارد پریمیس درج گردند (دفتر استانداردهای مارک و

الگوی 7 شیوه درج پیشنهاد تکست. ام دی در بخش فراداده‌های مدیریتی پیشنهاد متس

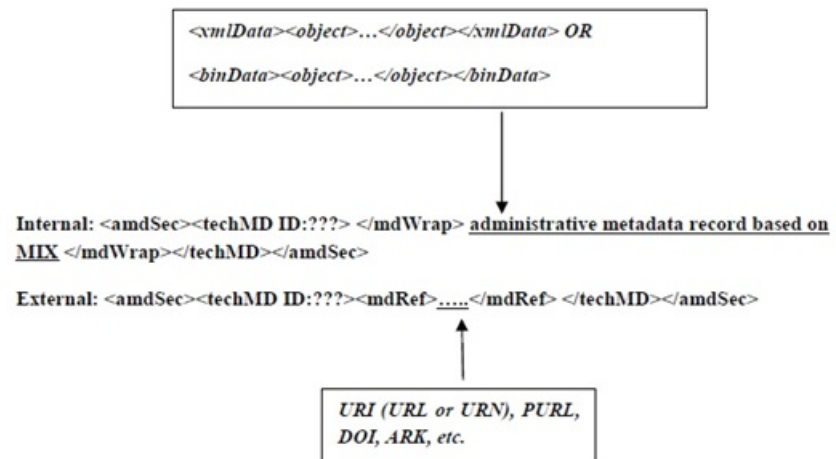
فراداده برای مدیریت تصاویر در بستر ایکس ام ال (میکس)

. میکس محصول دفتر استانداردهای مارک و توسعه شبکه با همکاری کمیته استانداردهای فراداده ای فنی برای تصاویر ثابت وابسته به سازمان استانداردهای اطلاعات ملی (NISO) است که به منظور مدیریت تصاویر دیجیتالی توسعه یافته است. کارکرد اصلی این طرح مدیریت فنی تصاویر ثابت دیجیتالی است (دفتر استانداردهای مارک و توسعه شبکه d2013) و پیشنهادهای آن در بخش فراداده های مدیریتی متس در عنصر techMD قرار میگیرند.

ص: 120

Object -1

بررسی تأثیر بستر نحوی ... ۱۲۱



الگوی ۸ شیوه درج پیشینه میکس در بخش فراداده‌های مدیریتی پیشینه متس

```
<metsHdr CREATEDATE="2013-01-05T14:00:00" RECORDSTATUS="Complete">
<agent ROLE="CREATOR" TYPE="INDIVIDUAL"><name>Sayyed Mahdi Taheri</name></agent></metsHdr>
<dmdSec ID:??></mdWrap> descriptive metadata record based on MARCXML
</mdWrap></dmdSec><dmdSec ID:??></mdWrap> descriptive metadata record
based on MODSExamples/mods99042030.xml</mdWrap></dmdSec><dmdSec
ID:??></mdWrap> descriptive metadata record based on DCMIExamples/
mods99042030.xml</mdWrap></dmdSec><amdSec><sourceMD ID:??></mdWrap>
administrative metadata record based on MARCXML </mdWrap></sourceMD></amdSec>
<amdSec><techMD ID:??> or<digiprovMD ID:??></mdWrap>
administrative metadata record based on PREMIS </mdWrap></digiprovMD ID:??>or</
techMD></amdSec><amdSec><techMD ID:??></mdWrap> administrative metadata
record based on TextMD </mdWrap></techMD></amdSec><amdSec><techMD
ID:??> </mdWrap> administrative metadata record based on MIX </mdWrap></
techMD></amdSec><fileSec><fileGrp ID="VERSI"><file ID="FILE001"
MIMETYPE="application/xml" SIZE="257537" CREATED="2013-01-05"><FLocat
```

الگوی 8 شیوه درج پیشینه میکس در بخش فراداده‌های مدیریتی پیشینه متس

Sayyed Mahdi Taheri

<name

descriptive metadata record based on MARCXML_

ID=???>descriptive

record

based on MODSExamples/mods99042030.xml

ID=???>descriptive

metadata record based on

/DCMIExamples

mods99042030.xml

administrative metadata record based on MARCXML

amdSec> or

administrative metadata record based on PREMIS or

techMD> administrative metadata

record based on TextMD

ID:???> administrative metadata record based on MIX

<techMD

"ID="FILE001

<"MIMETYPE="application/xml" SIZE="257537" CREATED="2013-01-05

LOCTYPE="URL">http://dlib.nyu.edu/tamwag/beame.xml

<fileGrp

<"Introduction" ORDER="1

"BEGIN="INTVWBG

<div

to

<"2

"END="INTVWND

"FILEID="FILE001

</BETYPE="IDREF

<"1

"LABEL="Page TYPE="page" ID="P1

</"FILEID="HTMLF1

"ID="DISS1.1

"STRUCTID="S1.1" BTYPE="uva-bdef:stdImage" CREATED="2002-05-25T08:32:00

LABEL="UVA Std Image Disseminator" GROUPID="DISS1" ADMID="AUDR

<"EC1

-NEW AND IMPROVED Image Mechanism" LOCTYPE="URN" xlink:href="uva

</"bmech:BETTER-imageMech

الگوی 9. برآیند الگوهای ارائه شده یک پیشینه کامل متس که تمامی استانداردهای مورد مطالعه را در بر میگیرد

همان طور که در الگوهای طراحی شده مشاهده میشود بستر نحوی استانداردهای فراداده ای مورد امکان میانکنش پذیری میان آنها را

فراهم آورده است. هر یک از این استانداردها میتوانند با

بررسی یکدیگر و با استاندارد هسته متس ارتباط برقرار کنند و با پشتیبانی از کارکردهای گوناگون، یکپارچگی نظام های اطلاعاتی را میسر سازند به علاوه میتوان بیش از یک استاندارد فراداده ای با کارکرد یکسان درون پیشینه های متس جاسازی نمود هنگامی که بیش از یک استاندارد و یا گزیده ای از عناصر هر استاندارد مورد نیاز است و نیز میتوان پروفایل کاربردی یک مرکز یا محیط اطلاعاتی خاص را درون پیشینه های متس بسته بندی نمود. استاندارد های مورد بررسی تنها بخشی مهمترین و پرکاربردترین از استانداردهای فراداده ای بودند. بدیهی است با استفاده از بستر نحوی، مناسب میانکنش پذیری دیگر استانداردها نیز میسر خواهد بود این تاثیر بستر نحوی بر فرایند میانکنش پذیری فراداده ای را نشان می دهد.

نتیجه گیری

ضرورت استفاده از چند استاندارد فراداده ای در یک نظام اطلاعاتی به منظور پشتیبانی از کارکردهای گوناگون مورد نیاز آن نظام و اقبال ویژه نظامهای اطلاعاتی به مقوله یکپارچگی و ارزشهای افزوده مرتبط با آن بیانگر اهمیت فرایند میانکنش پذیری فراداده ای است طراحی پروفایل های کاربردی متناسب

ص: 122

با نیازهای مراکز یا محیطهای اطلاعاتی خاص نیز این اهمیت را دو چندان نموده است. بستر نحوی یکی از ارکان اصلی فرایند میانکنش پذیری میان استانداردها و نظامهای فراداده ای است استانداردهای فراداده ای با انتخاب بستر نحوی مناسب سطح تعامل پذیری خود با دیگر استانداردهای فراداده ای را افزایش می دهند، و بدین گونه علاقه نظامهای اطلاعاتی به انتخاب آنها را بر می انگیزانند گرایش استانداردهای فراداده ای به انتخاب زبان نشانه گذاری گسترش پذیر به دلیل قابلیت های منحصر به فرد آن از جمله خود توصیف بودن که فرایند میانکنش پذیری را در هر دو سطح نحوی و معنایی تسهیل می نماید، به عنوان قالب اصلی و یا یکی از قالبهای پیاده سازی پیشینهها در همین راستا بوده است (حریری و دیگران، 1391) بستر نحوی زمینه را برای پشتیبانی کارکردهای مورد نظر استانداردهای فراداده ای و ارتباط پیشینه های فراداده ای با یکدیگر را فراهم نموده افزودن بر بهبود یکپارچگی درونی نظامهای اطلاعاتی، میانکنش پذیری آنها با دیگر نظامهای اطلاعاتی از جمله موتورهای کاوش وب به عنوان پرکاربردترین ابزار جستجو و بازیابی اطلاعات در شبکه وب (کین، 2008؛ طاهری و حریری، 2012) یا به عبارت دیگر یکپارچگی برونی یا آنها را نیز افزایش میدهد امکان دسترسی یکپارچه به اشیای محتوایی منتشر شده در نظامهای اطلاعاتی مختلف از طریق ابزارهایی چون دروازه های اطلاعاتی و درگاهها از پیامدهای نیک قابلیت های بستر نحوی میباشد یکی دیگر از ارزشهای افزوده ای که بستر نحوی مناسب تولید خواهد نمود تولید دانش بر اساس ارتباط میان پیشینه های فراداده ای است که جلوه ای دیگر از یکپارچگی درونی و برونی نظامهای اطلاعاتی به شمار می آید.

منابع

حریری نجلا سید مهدی طاهری سید رحمت الله فتاحی 1391. میانکنش پذیری نظامهای فراداده ای و موتورهای کاوش: وب تحولات و رویکردهای جاری پژوهشنامه کتابداری و اطلاع رسانی، 2 (2) طاهری، سید مهدی 1387 الف. طراحی یک کتابخانه دیجیتالی استاندارد. در مجموعه مقالات نخستین همایش کتابخانه های دیجیتالی به کوشش شرکت پارس آذرخش تهران سبزان

طاهری سید مهدی 1387 ب. مقایسه کارایی طرح فراداده‌های هسته دویلین و قالب فراداده مارک 21 در سازماندهی منابع اطلاعاتی شبکه جهانی وب فصلنامه کتابداری و اطلاع رسانی، 43 (پاییز 1387)

طاهری سید مهدی 1391 بررسی امکان نمایه سازی و پیدانمایی نامهای برچسب عناصر فراداده ای هسته، دویلین مارک، 21 و طرح فراداده ای توصیف شیء توسط موتورهای کاوش عمومی و ارائه الگوی مناسب رساله دکترا گروه کتابداری و اطلاع رسانی دانشگاه آزاد اسلامی، واحد علوم و تحقیقات تهران

فتاحی، سید رحمت الله طاهری سید مهدی (1388) فهرستتویسی رایانهای مفاهیم، شیوهها و کاربرد نرم افزارهای رایانهای در سازماندهی اطلاعات با همکاری فرشته ناقدی احمدی تهران کتابدار

.Bray, Tim, Jean Paoli, C. M. Sperberg-McQueen, Eve Maler, François Yergeau. 2008

Extensible Markup Language (XML) 1.0 (Fifth Edition): W3C Recommendation 26

ص: 123

/November 2008. Retrieved 7 Jun. 2013 from: <http://www.w3.org/TR/xml>

//:Dublin Core Metadata Initiative 2013. Glossary. Retrieved 14 Jun. 2013 from: <http://dublincore.org/documents/2003/08/26/usageguide/glossary.shtml>

//:Gigee, Grant, Kelly 2006. MARC and MARCXML. Retrieved 5 Nov. 2011 from: <http://threegee.files.wordpress.com/2006/05/marcxml.pdf>

Gill, Toney 2008. Metadata and the web: Introduction to metadata. Retrieved 5 Jun. 2013
/from: http://www.getty.edu/research/publications/electronic_publications/intrometadata
metadata.pdf

Habing, Tom 2007. METS, MODS and PREMIS, Oh My!: Integrating Digital Library
.Standards for Interoperability and Preservation. Retrieved 5 Jun. 2013 from: <http://www.loc.gov/standards/mods/presentations/habing-ala07>

:(Hirwade, Mangala Anil .2011. A study of metadata standards. Library Hi Tech News, 28(7
.18-25

Juhnston, pete, Andy Powell.2006. Expressing Dublin Core Metadata Using XML. Retrieved
.Jun. 2013 from: <http://dublincore.org/documents/dc-xml> 5

Maarof, M.H.B.S., Y. Yahya 2009. Digital libraries interoperability issues. Electrical
.Engineering and Informatics. ICEEI '09. International Conference on. Retrieved 5 Jun
from: http://ieeexplore.ieee.org/xpl/freeabs_all.jsp?tp=arnumber=5254718&isnu 2013
.mber=5254684

McCallum, H. Sally .2004. An Introduction to the Metadata Object Description Schema
-MODS). Retrieved 5 Jun. 2013 from: <http://dlcd.lib.uchicago.edu/talks/2004/lita2004>)

Network Development and MARC Standards Office (NDMSO) 2013a. MARC 21 XML

.Schema. Retrieved 5 Jun. 2013 from: <http://www.loc.gov/standards/marcxml>

Network Development and MARC Standards Office (NDMSO) .2013b. Metadata Authority

/Description Schema (MADS). Retrieved 5 Jun. 2013 from: <http://www.loc.gov>

[standards/mads](http://www.loc.gov/standards/mads)

Network Development and MARC Standards Office (NDMSO) .2013c. Metadata Encoding

/Transmission Standard (METS). Retrieved 5 Jun. 2013 from: <http://www.loc.gov>

[.standards/mets](http://www.loc.gov/standards/mets)

Network Development and MARC Standards Office (NDMSO) .2013d. Metadata for Images

.in XML Standard (MIX). Retrieved 5 Jun. 2013 from: <http://www.loc.gov/standards/mix>

Network Development and MARC Standards Office (NDMSO) 2013e. Metadata Object
/Description Schema (MODS). Retrieved 5 Jun. 2013 from: <http://www.loc.gov/standards/mods>

.Network Development and MARC Standards Office 2013f. Technical Metadata for Text
Retrieved 5 Jun. 2013 from: <http://www.loc.gov/standards/textMD>

.Network Development and MARC Standards Office 2013g. Understanding PREMIS
.Retrieved 5 Jun. 2013 from: <http://www.loc.gov/standards/premis/understanding-premis.pdf>

Qin, Jian. 2000. Representation and organization of information in the web space: From
.MARC to XML. Retrieved 5 Jun. 2013 from: <http://inform.nu/Articles/Vol3/v3n2p83-88.pdf>

Taheri, S. M., Nadjla Hariri.2012. A comparative study on the indexing and ranking of the
content objects including the MARCXML and Dublin Core's Metadata elements by
(general search engines. Electronic Library. 30(4

//:Word Wild Web Consortium.2013. HTML CSS. Retrieved 05 Jun. 2013 from: <http://www.w3.org/standards/webdesign/htmlcss>

.Wikipedia 2013. Portable document format. Retrieved 05 Jun. 2013 from: http://en.wikipedia.org/wiki/Portable_Document_Format

آرشیو سازی وب به طور خودکار توسط خزشگرهای وب انجام میگیرد. این خزشگرهای صفحات را به صورت ادواری بازبینی و آرشیوها را با عکسهای جدید و تازه، روزآمد میکنند. مقاله حاضر به موضوع آرشیوسازی صفحات وب به طور کارآمد و بهسازی کیفیت آن اشاره دارد و یککرد پیشنهادی در این مقاله سه مفهوم را با هم تلفیق میکند بخش بندی صفحه، دیداری شناسایی تغییر و اهمیت بلاکهای صفحه های وب برای تشخیص بهتر تغییرات مهم میان نسخهها چالش اصلی در این مقاله بهبود کیفیت آرشیو. است هدف ما این است که راهبردهای خزش را طوری تنظیم کنیم که خزش وبگاههای آینده تا آنجا که ممکن است به طور جامع و منسجم انجام شود بنابراین در مقاله حاضر نوعی راهبرد خزش مبتنی بر الگو و دو سنجه کیفیت جامعیت و انسجام برای ارزیابی راهبرد خزش پیشنهاد شده است.

نوشته: میریام بن سعد (1)

ترجمه: مهشید برجیان (2) - ساناز باغستانی (3)

مقدمه

با توجه به اهمیت روزافزون شبکه جهانی وب آرشیوسازی اطلاعات آن به منظور حفظ منابع اطلاعاتی مفید بسیار حائز اهمیت است. به همین دلیل حفظ وبگاههای مفید در نظر بسیاری از مؤسسه ها و سازمانهای آرشیوی ملی در سراسر دنیا به مسئله مهمی تبدیل شده است. اغلب اوقات، آرشیو سازی وب، به طور خودکار توسط خزشگرهای وب انجام میگیرد این خزشگرها، صفحه ها را به صورت ادواری بازبینی و آرشیوها را با عکسهای جدید و تازه روزآمد می کنند دریافت تمام نسخه ها از کل وبگاه ها و حفظ کیفیت آرشیوها کار بی اهمیت و پیش پا افتاده ای نیست. در حقیقت، حفظ یک آرشیو کامل از کل وبگاه شامل تمامی نسخه های کل صفحه ها غیر ممکن است؛ زیرا وبگاهها دائماً در حال تکامل و گسترش هستند و امکانات و منابع (4) تخصیص یافته نیز معمولاً محدودند (مثل پهنای نوار، فضای ذخیره سازی و قوانین اخلاقی سایت). همچنین خزش یک وبگاه بزرگ ممکن است ساعتها و حتی روزها طول بکشد. این امر، باعث افزایش احتمال تغییرات صفحه در حین خزش میشود و در نهایت

ص: 127

Myriam Ben Saad -1

2- کارشناس ارشد کتابداری و اطلاع رسانی سازمان اسناد و کتابخانه ملی ایران

3- کارشناس ارشد کتابداری و اطلاع رسانی

resources -4

عدم انسجام میان صفحه های آرشیو شده را به دنبال دارد به عنوان مثال فرض میکنیم که خزشگر در حال بارگذاری دو صفحه از وبگاه France است. اولین صفحه حاوی پیوند فرامتنی موسوم به obama است که ما را به خزش صفحه دومی ارجاع میدهد که درباره رئیس جمهور آمریکا صحبت می کند. در حین خزش صفحه، اول محتوای صفحه دوم به روز میشود و اکنون در مورد رئیس جمهور فرانسه یعنی سارکوزی صحبت می کند، بنابراین نسخه ذخیره شده صفحه اول به obama اشاره دارد با که نسخه صفحه ای که عمل خزش در آن انجام شده است و در مورد سارکوزی صحبت می کند، متناقض (1) میشود. مسئله تناقض بسیار رایج است و معمولاً حین آرشیوسازی وب شناسایی میشود. بنابراین، سیستمهای آرشیو وب باید دو مسئله را در نظر بگیرند برای داشتن آرشیو جامع، عمل خزش را چگونه بهبود بخشند و چگونه کیفیت آن را با به حداقل رساندن عدم انسجام میان صفحات جمع آوری شده، حفظ کنند؟

یکی از روشهای ایده آل برای آرشیوسازی وب این است که تمام صفحه های سایت را در یک زمان و با هر تغییر خزش کنیم و یا از تغییر محتوای صفحات در طول عمل خزش جلوگیری نماییم، البته، این مورد به دلیل محدودیتهای امکانات و منابع عملاً غیر ممکن است حصول اطمینان 100 درصدی از جامعیت و نیز انسجام در آرشیو عملاً غیر ممکن است. با وجود این هدف ما این است که راهبرد خزش خود را طوری تنظیم کنیم که خزشهای وبگاه در آینده بتوانند تا آن جایی که امکان دارد آرشیوها را جامع سازد و تا حد امکان به آن هماهنگی و انسجام ببخشند. در این مقاله این دو چالش را برای بهبود کیفیت آرشیو بیان میکنیم.

یکی از ایده هایی که برای بهبود جامعیت آرشیو مطرح شده این است که مهمترین نسخه ها از صفحه ها را طوری بارگذاری کنیم که اطلاعات مفید از دست رفته به حداقل ممکن برسد. نسخه مهم، نسخه ای است که بعد از آرشیو آخرین نسخه واجد تغییر مهمی باشد. تغییرات بی اهمیت صفحه مثل تبلیغات را میتوان نادیده گرفت.

پژوهشهای بسیاری هستند که بر تخمین فراوانی تغییرات صفحه های وب برای بهبود خزشگرها

تأکید می کنند. اما هیچ روش مؤثری وجود ندارد که بتواند مشخص کند تغییرات مهم میان نسخه ها دقیقاً چه زمانی و هر چند وقت یک بار رخ میدهد تا امروز رویکردهایی که فراوانی تغییرات را برآورد کرده اند فقط میزان تغییرات شناسایی شده را مورد توجه قرار داده اند ولی اهمیت تغییرات رخ داده را در نظر نگرفته اند. اگر بتوانیم فراوانی تغییرات مهم را به طور دقیق تری پیش بینی کنیم، آنگاه می توانیم عمل خزش وب را بهبود بخشیم.

برای تخمین تعداد روز آمدسازیهها باید تغییرات میان نسخه های بازیابی شده اسناد شناسایی شوند. بسیاری از الگوریتمهای موجود [21]، [14]، [28] به منظور شناسایی تغییرات اسناد نیمه ساختار یافته 14 (فرمت XML و فرمت HTML) طراحی شده اند. با وجود این هیچ روشی وجود ندارد که بتواند تغییرات مرتبط نامرتبط را از اطلاعات مفید زائد اطلاعاتی که موجب اختلال میشوند شناسایی کند و آنها را

ص: 128

هدف ما ارائه رویکردی است برای: 1) شناسایی تغییرات مهم میان نسخه ها و به کارگیری آنها، 2) بهینه سازی خزش، وب 3) بهبود کیفیت نسخه‌های صفحه های آرشیو شده 4) نمایه سازی /ذخیره سازی مؤثر نسخه‌های صفحات رویکرد ما در یکی از مجموعه های (1) سازمان INA (مؤسسه شنیداری - دیداری ملی (فرانسه) به کار برده شده است. یکی از مأموریت‌های INA ایجاد مجموعه ای (2) قانونی است که صفحه‌های وب رادیویی و تلویزیونی و نیز صفحه های مربوط به آنها را نگهداری کند یکی از الزامات مهم این، پروژه حفظ جنبه های دیداری صفحه هاست. بنابراین ما میخواهیم برای تعیین اهمیت بخشهای صفحه های وب با توجه به جایگاه نسبی آنها در صفحه، از روش آنالیز صفحه تصویری استفاده کنیم.

پژوهشهای پیشین [7،9] نشان میدهند که یک صفحه را میتوان به چند بخش یا چند بلاک تقسیم کرد این بلاکها غالباً در صفحه از اهمیت متفاوتی برخوردارند. در حقیقت بخشهای مختلف موجود در یک صفحه وب با توجه به مکان اندازه، بخش و محتوا از وزنهاى اهمیت متفاوتی برخوردار هستند معمولاً مهم ترین اطلاعات در مرکز صفحه تبلیغات در بالای صفحه (3) یا در سمت چپ، و بخش حق مؤلف در قسمت پایین صفحه (4) قرار دارند. هنگامی که یک صفحه بخش بندی می شود، برای هر بلاک باید یک اهمیت نسبی تعیین شود با استفاده از یک الگوریتم [26] و روش یادگیری ماشینی تحت نظارت میتوان به طور خودکار این فرآیند را انجام داد، سپس میتوانیم میزان اهمیت تغییرات میان دو نسخه یک صفحه را براساس موارد زیر محاسبه کنیم: 1) اهمیت نسبی بلاکها و 2) اهمیت نسبی عملیات (درج، حذف، روزآمدسازی و مانند آن) رخ داده در بلاکها که با مقایسه این دو نسخه شناسایی شده اند. از این رو در پژوهش حاضر در صددیم که این مفاهیم را با هم ترکیب کنیم تا مسائل مربوط به شناسایی آنالیز تغییرات مهم صفحات وب را بیان کنیم سپس نتایج حاصل از این تجزیه و تحلیل را می توانیم در آرشیو سازی مؤثر صفحات وب به کار ببریم و موجب بهبود کیفیت وبگاههای آرشیو شده شویم. سؤالهای اصلی و درخور توجه این پژوهش که در آینده را بیان خواهیم کرد نیز مورد بحث قرار گرفته اند.

مسئله پژوهش

طراحی سیستمهای آرشیو وب چالشهای بسیاری را به همراه دارد؛ 1) بهینه سازی خزش وب، به منظور بهبود کیفیت آرشیو 2) نمایه سازی ذخیره سازی مناسب و 3) پرسش مؤثر در لحظه جست و جویا در نظر گرفتن مسئله بارگذاری نسخه های صفحه ای بی اهمیت بی فایده به طور قابل توجهی میتوان اثر بخشی این سه نکته را بهبود بخشید. بنابراین خزشگرهای وب باید به دقت مشخص کنند که کدام

ص: 129

repository -1

header -2

repositite -3

footer -4

صفحه و با چه الیوتی بازسازی (1) / آرشیو شود. آنها، همچنین، باید مشخص کنند که برای به روز نگه داشتن آرشیو وب باید هر چند وقت یکبار صفحه ها را مورد بازبینی قرار داد. در حقیقت، حفظ یک آرشیو کامل از کل وب و یا حتی بخشی از آن که شامل تمام نسخه ها از کل صفحه ها می شود، کاری غیر ممکن است؛ زیرا وب دائماً در حال تکامل و گسترش است و امکانات و ابزارهای تخصیص یافته نیز معمولاً محدود هستند. بنابراین این مسئله را میتوان به این صورت مطرح کرد: چگونه عمل خزش وب را بهبود بخشیم تا مهمترین نسخهها را بارگذاری کنیم طوری که اطلاعات مفید از دست رفته باشد البته این مسئله باید بدون هیچ کمکی از سوی مدیران وبگاه ها حل شود. مهمی بعد از آخرین نسخه آرشیو شده در آن صورت گرفته باشد البته این مسئله باید بدون هیچ کمکی از سوی مدیران وبگاه ها حل شود. بنابراین، به روشی مؤثر و مفید نیاز داریم تا به واسطه آن بدانیم که تغییرات مهم در چه زمانی و هر چند وقت یکبار بین نسخه ها در وبگاهها رخ میدهند تاکنون، رویکردهای موجودی که فراوانی خزنده ها را برآورد میکنند اهمیت تغییرات میان نسخه ها را در نظر نگرفته اند در حقیقت اغلب اتفاق می افتد که خزشگرها صفحه هایی را بارگذاری میکنند که دارای اطلاعات بی اهمیت هستند (به عنوان مثال تبلیغاتی که به روز میشوند) برای برآورد تعداد مناسب خزنده ها، تغییرات میان نسخه های بازایی شده باید شناسایی شوند و مورد تجزیه و تحلیل قرار گیرند با وجود اینکه برای شناسایی تغییرات میان اسناد، الگوریتمهای مختلفی طراحی شده اند هیچ روشی وجود ندارد که بتواند تغییرات مهم بی اهمیت را از اطلاعات مفید بی فایده شناسایی کند و آنها را از هم متمایز سازد.

در این مقاله، ما برخی چالشهای مهم را بیان میکنیم: (1) سیستمهای آرشیوسازی چگونه میتوانند تغییرات مفید میان نسخههای آرشیو شده را شناسایی کنند و چگونه میتوانند اهمیت آنها را تعیین کنند؟ (2) با توجه به امکانات و ابزار محدود و تعداد زیاد اسنادی که باید آرشیو شوند خزشگرها چگونه می توانند ضروری ترین / مهمترین نسخه بازسازی شده را انتخاب کنند؟ (3) برای بهبود کیفیت (جامعیت و انسجام) وبگاههای آرشیو شده چگونه میتوان از نتایج آنالیز اهمیت تغییرات بهره برد؟

پژوهشهای مرتبط

در ادامه به مطالعات مرتبط با این تحقیق میپردازیم. موضوع های اصلی عبارت اند از: آرشیوسازی وب تجزیه و تحلیل دیداری صفحه ،وب شناسایی تغییرات و خزش وب.

آرشیوسازی وب مؤسسههای آرشیوی متعددی (کتابخانه های ، ملی آرشیوهای داده های تاریخی، و مانند آن) در سراسر دنیا، برای حفظ میراث وب کشور خود چندین پروژه را راه اندازی کرده اند. برخی مطالعات بر تعیین محدوده وب برای انتخاب صفحات برای آرشیو شدن تأکید دارند. پژوهشهای دیگری نیز بر روی مدل سازی و ارزیابی فراوانی تغییرات وب کار کرده اند. آنها برای بهبود بازسازی آرشیو، تخمین زندهای فراوانی تغییرات و سیاستهای بازسازی گوناگونی را ارائه میدهند. برخی محققان نیز مسائل مربوط به فرمت اطلاعات ذخیره سازی نمایه سازی شده را با ارائه سیستم ذخیره سازی خودشان

ص: 130

بیان میکنند مطالعات دیگر بر کنترل و نمایش تغییرات تأکید میکنند این مطالعات، برای پرسش (1) و ذخیره سازی مؤثر آرشیو، وب الگوریتم شناسایی تغییر و یا فرمت دلتا ارائه میدهند. کارهای اخیر نیز مسئله انسجام و کیفیت آرشیو را توسط راهبردی خزش بیان میکنند.

جالب اینجاست که این روشها و رویکردها اهمیت تغییرات صفحه ها برای آرشیوسازی وب به طور مؤثر را در نظر نمی گیرند؛ در حالی که محور اصلی روش ما همین اهمیت است.

تجزیه و تحلیل دیداری صفحه وب برای تجزیه و تحلیل نمایش دیداری صفحات وب از چندین روش استفاده شده است. بیشتر این روشها، ساختار منطقی صفحه را با آنالیز اسناد ارائه شده یا آنالیز کد اسناد، کشف می کنند. گو (2) و همکارانش [20]، نوعی الگوریتم بالا به پایین ارائه دادند که الگوریتم ساختار محتوای وب را مبتنی بر اطلاعات صفحه آرایشی شناسایی میکند. کوواشونیک (3) و همکارانش [17]، برای شناسایی بخشهای رایج یک (صفحه بالای، صفحه پایین صفحه مرکز صفحه) فرآیندهای مکاشفه ای مبتنی بر اطلاعات دیداری را تعیین کردند. کای (4) و همکارانش [9] الگوریتم VIPS را مطرح کردند. این الگوریتم براساس اطلاعات دیداری بازیابی شده توسط مرورگر صفحه وب را به چندین بلاک معنایی تقسیم میکند کوسولشی (5) و همکارانش [15] رویکردی را مطرح کردند که میزان تشابه بلاکها را در صفحه های وب با استفاده از اطلاعات موقعیتی عناصر درخت DOM محاسبه میکند. به نظر میرسد که روش VIPS در مقایسه با روشهای موجود مناسبترین روش برای رویکرد ماست زیرا دانگی (6) مناسبی برای بخش بندی صفحه ایجاد میکند منظور) از دانگی، اندازه قطعه های حافظه در سیستم مجازی است این، روش سلسله مراتبی از بلاکهای معنایی صفحه را ایجاد می کند. این سلسله مراتب چگونگی درک کاربر از ساختار صفحه آرایشی، وب مبتنی بر درک دیداری وی را بهتر شبیه سازی می کند. از این رو برای ایجاد ساختار دیداری اسناد از VIPS استفاده شد.

شناسایی تغییرات برای شناسایی تغییرات میان دو نسخه یک سند نیمه ساختاری (XML و HTML)، چندین الگوریتم طراحی شده است. این الگوریتمها حداقل مجموعه ای از عملیات تغییر (درج، حذف، و) را پیدا می کنند که یک درخت دادهها را به درخت دادههای دیگر تبدیل میکند. این عملیات تغییر غالباً در یک متن دلتا و یا یک فایل دلتا گردآوری میشوند. طراحی الگوریتمهای مختلف به اهداف و الزامات (7) آنها پیچیدگی، زمانی عملیاتی که قرار است به کار برده شوند کیفیت دلتا و مانند آن بستگی دارد. کوبنا (8) و همکارانش [14] برای بهبود مدیریت حافظه و مدیریت زمان الگوریتم XyDiff را مطرح کردند. الگوریتم XyDiff عملیات انتقال (9) را پشتیبانی میکند و پیچیدگی زمانی $O(n \log(n))$ را به دست می آورد.

ص: 131

query -1

Giu -2

Kovacevic -3

Cai -4

Cosulshi -5

granularity -6

requirement -7

Cobena -8

move -9

این الگوریتم، با وجود عملکرد بالایی که دارد همیشه نمیتواند نتیجه بهینه و مطلوبی را تضمین کند (منظور از نتیجه بهینه و مطلوب حداقل ویرایش برای متن است). وانگ (1) و همکارانش [28]، الگوریتم XyDiff را مطرح کردند این الگوریتم میتواند تفاوت‌های مطلوب میان دو درخت سامان نیافته XML در معادله زمانی درجه دوم $O(n^2)$ را شناسایی کند؛ ولی هیچ انتقالی را پشتیبانی نمیکند الگوریتم Delta [21] این اسناد XML را برای درختان سامان یافته و سامان نیافته با حمایت از عملیات اصلی میتواند مقایسه ادغام و هماهنگ میکند؛ اما انتقال را شناسایی نمیکند الگوریتم‌های دیگری همانند الگوریتم [Diff – DTD]22 و غیره نیز وجود دارند پس از مطالعه و بررسی این الگوریتم‌ها تصمیم گرفتیم که برای رویکرد آرشیو سازی وب خود از روش‌های موجود استفاده نکنیم زیرا هدف آنها کلی است. از آنجا که الزامات خاص متفاوتی در ارتباط با ساختار صفحه آرایی دیداری اسناد وجود دارد، ترجیح میدهم که از الگوریتم ویژه و متناسب با کار خود استفاده کنیم (الگوریتم Vi DIFF). این الگوریتم امکان ارزیابی بهتر پیچیدگی و جامعیت مجموعه عملیات شناسایی شده را فراهم می‌کند.

خزش وب

تعدادی از پژوهش‌های موجود مسئله بهینه سازی خزش وب را از طریق ایجاد راهبردیهای زمان بندی [23 و 10 و یا از طریق برآورد فراوانی، تغییرات بیان میکنند [27 و 16] مطالعات اخیر [13 و 12]، مسئله کیفیت و انسجام آرشیو را با ارائه راهبردی خزش وبگاه، بیان می‌کنند. اما، از آنجا که راهبردهای آنها مبتنی بر فرآیند پواسون است برای صفحاتی که زود به زود تغییر میکنند مفید نیستند (مانند صفحات رادیویی و تلویزیونی همچنین تا آنجا که ما میدانیم؛ مطالعات خزش، موجود اهمیت تغییرات رخ داده در نسخه‌های تجزیه و تحلیل شده را در نظر نمی‌گیرند. اگر بتوانیم فراوانی تغییرات مهم را به طور دقیق تری پیش بینی کنیم شاید بتوانیم از نمایه سازی و ذخیره اطلاعات بی اهمیت و غیر ضروری جلوگیری کنیم و اثر بخشی و کیفیت آرشیو وب را بهبود بخشیم.

2. رویکرد آرشیوسازی وب

رویکرد آرشیوسازی وب ما، آنالیز آرشیوسازی ساختار دیداری اسناد و تعیین ارزشهای (3) اهمیت برای بلاکهای صفحه های وب با توجه به جایگاه نسبی آنها در صفحه است. به عبارت دیگر، نسخه های یک صفحه مطابق با نمایش دیداری شان بازسازی میشوند شناسایی تغییرات در چنین نسخه های صفحه ای بازسازی شده اطلاعات مناسبی را برای درک دینامیک وبگاهها ارائه میدهد همچنین امکان تشخیص تغییرات مرتبط نامرتبط را از اطلاعات مفید بی فایده فراهم می‌آورد، بنابراین، روش مطرح شده، سه مفهوم زیر را با هم ترکیب میکند آنالیز دیداری صفحه (بخش بندی)، شناسایی تغییرات، و اهمیت بلاکهای صفحه های وب به منظور بهینه سازی خزش وب این مفاهیم جدید نیستند؛ اما تا آنجا که ما

ص: 132

Wang – 1

importance values – 2

Poisson Process – 3

میدانیم هرگز این مفاهیم را برای آرشیوسازی وب با هم ترکیب نکرده اند. معماری آرشیو وب به طور مفصل تری در [6] شرح داده شده است.

بخش بندی صفحات دیداری

ما مدل بخش بندی دیداری موجود [9] VIPS را برای ایجاد ساختار دیداری صفحه های وب گسترش دادیم از مدل VIPS برای بخش بندی صفحه وب به بلاکهای معنایی تو در تو مبتنی بر گره های مناسب در درخت HTML DOM، صفحه استفاده شد. این مدل جدا کننده های افقی و عمودی را در صفحه وب شناسایی می کند همچنین این مدل براساس جداکننده ها درخت معنایی صفحه وبی را ایجاد میکند که به چندین بلاک تقسیم بندی شده است. اساس کار کل صفحه است. هر بلاک به عنوان یک گره در درخت نشان داده می شود. با استفاده از استخراج پیوندها، تصاویر، و متن برای هر بلوک، الگوریتم VIPS را برای تکمیل درخت معنایی کل صفحه گسترش دادیم الگوریتم VIPS توسعه یافته ما نوعی سند Vi-XML به عنوان خروجی تولید می کند این سند ساختار سلسله مراتبی کامل صفحه وب را توصیف می نماید.

شناسایی تغییرات

برای شناسایی تغییرات میان دو صفحه وب براساس بعد دیداری الگوریتم Vi Diff را مطرح کردیم. این الگوریتم دو نوع تغییر را شناسایی میکند تغییرات ساختاری و تغییرات محتوایی. تغییرات ساختاری (درج، حذف و جابه جایی)، معمولاً ساختار سند XML را در سطح بلاکها تغییر میدهند؛ در حالی که تغییرات محتوایی، درج، حذف روزآمدسازی و انتقال محتوای متنی را در سطح پیوندها، تصاویر، متنها تغییر میدهند تمامی عملیات تغییر شناسایی شده در یک فایل Vi-Delta توصیف میشوند. اگر فرض کنیم که هیچ تغییری در ساختار وجود ندارد میزان پیچیدگی Vi-Delta لگاریتم خطی $O(\log(n))$ است که در آن همان تعداد کلی گره هاست. اگر تغییرات ساختاری وجود داشته باشند، در بدترین حالت حالتی که تمام ساختار تغییر (کند که پیچیدگی به صورت معادله درجه دوم $O(n^2)$ است؛ ولی ارزش دارد که بینیم n همیشه اندازه اش کوچک باقی میماند.

اهمیت تغییرات

با توجه به Vi-Delta ایجاد شده توسط Vi Diff تابعی [4] را ارائه میدهم که اهمیت تغییرات شناسایی شده را ارزیابی میکند. این تابع به سه پارامتر اصلی بستگی دارد:

اهمیت بلوک روزآمد شده معمولاً مهمترین اطلاعات در مرکز صفحه وب و تبلیغات در قسمت بالای صفحه و مانند آن قرار میگیرند بنابراین میتوان اهمیت بلاکها را با توجه به جایگاه نسبی آنها در صفحه تعیین کرد به عنوان مثال میتوان از روش سانگ (1) و همکارانش [26] برای به دست آوردن آن

استفاده کرد. آنها از الگوریتمهای یادگیری ماشینی تحت نظارت مبتنی بر ویژگیهای محتوایی و فضایی استخراج شده بلاکها استفاده میکنند تا میزان اهمیت هر بلاک به طور خودکار تعیین شود. همچنین، میتوانیم پارامترهای دیگری را برای ارزشیابی اهمیت هر بلاک با توجه به تاریخچه تغییرات آن در نظر بگیریم. به عنوان مثال، می توانیم فرض کنیم که هر چه یک بلاک بیشتر تغییر کند، میزان اهمیت آن کمتر است. در حال حاضر به دنبال بهترین تکنیک برای تخمین میزان اهمیت بلاکها هستیم.

اهمیت عملکرد اهمیت عملکردها به نوع عملیات (انتقال درج و مانند آن) و عنصر تغییر یافته پیوند، تصویر، و مانند آن بستگی دارد؛ مانند عملیات درج و حذف که میتوانند مهمتر از عملیات انتقال محسوب شوند. همچنین درج یک تصویر میتواند مهمتر از درج یک پیوند و یا متن باشد. ما می خواهیم برای انتخاب بهترین پارامترها برای هر کدام از انواع عملیات به مطالعه روشهای یادگیری ماشینی پردازیم.

میزان تغییر هر بلوک میزان تغییر عملیاتی (حذف، درج و مانند آن) که برای هر عنصر (پیوند، تصویر و متن) در هر بلاک ایجاد میشود از VI-Delta ایجاد شده استنتاج میشود. این میزان درصد تغییر عملیات مشخص شده در هر بلاک (این بلاکها خود به تعدادی عنصر تقسیم شده اند) را نشان میدهد. عملیات پیشنهادی با توضیحات دقیق تر در [4] توصیف شده است.

آزمایشها

با استفاده از الگوریتم VIPS گسترش داده شده، آزمایشهای بخش بندی دیداری بر روی صفحه های وب HTML انجام شدند ما کارایی های بخش بندی دیداری در طول زمان و اندازه خروجی را اندازه گیری کردیم همچنین آزمایشهایی برای آنالیز کارایی الگوریتم Vi-DIFF پیشنهادی در طول زمان و اندازه برون داد، انجام شدند. آزمونها نشان میدهند که مدت زمان امیدوار کننده است. مدت زمان کلی رضایت بخش است زیرا این زمان امکان پردازش بیش از 100 صفحه (اندازه گیری فعلی در هر لحظه (ثانیه) و در هر پردازشگر شرایط پروژه کارتک CARTEC را فراهم می. کند به هر حال، زمان بخش بندی بسیار بیشتر از زمان مقایسه است برای بهینه سازی بیشتر، سیستم باید بر کاهش زمان بخش بندی و یا جلوگیری از بخش بندی تمامی نسخه های صفحه تمرکز کنیم.

کیفیت آرشیو وب

یکی از اهداف ما بهره وری از اهمیت تغییرات شناسایی شده برای بهسازی کیفیت آرشیو است رویکرد ایده آل آرشیوسازی، وب خزش تمامی صفحه های وبگاه به طور همزمان در هر تغییر و یا جلوگیری از تغییر محتویات صفحه در طول خزش است البته با توجه به تعداد زیاد صفحه های هر سایت و محدودیتهای امکانات و منابع این کار به طور عملی غیر ممکن است بنابراین، نمی توانیم آرشیو جامعی داشته باشیم آرشیوی که تمامی نسخه های تمامی صفحه های سایت را در بردارد. همچنین اطمینان از انسجام کل آرشیو صفحه های گردآوری شده وضعیت واقعی یک سایت را در یک لحظه از زمان

منعکس می‌کنند) غیرممکن است. با وجود این، قصد داریم که راهبردهای خزش را طوری تنظیم کنیم که خزش وبگاه در آینده تا آنجا که ممکن است جامع و منسجم باشد. تعدادی از راهبردهای خزش بر اساس مدل فرآیند پواسون هستند. این مدل، زمانی مناسب است که تقسیم‌بندی زمانی بیش از یک ماه باشد. به هر حال، کار ما بر روی صفحه‌های وبی است که بارها تغییر می‌کنند (چندین بار در روز) مانند صفحه‌های وب تلویزیون، رادیو، و صفحه‌های خبری. تقسیم‌بندی زمانی برای این صفحات کمتر از یک ماه است. بنابراین، مدرک مهمی وجود دارد که مدل Poisson برای این مورد اعتبار ندارد. [۷۹، ۲۵] با بازنگری تغییرات صفحه مشاهده کردیم که روزآمد سازی صفحه به ساعت روز و روز هفته بستگی دارد. همچنین، مشاهده کردیم که تغییرات در طول روز به میزان قابل توجهی بیشتر از شب و در روزهای کاری بیشتر از آخر هفته هستند.

بنابراین، با بازنگری تغییرات صفحه از یک خزش ناپیوسته^۱ الگوهای را کشف می‌کنیم. یک الگو، رفتار تغییرات صفحه و اهمیت آنها را در طول زمان به‌عنوان مثال، در طول یک روز طراحی می‌نماید. این مسئله، به روز هفته بستگی دارد. الگو باید پی‌درپی روزآمد شود تا همیشه بتواند تغییرات جاری صفحات وب را منعکس کند. صفحات یک سایت با الگوی مشابه می‌توانند برای داشتن یک الگوی مشترک گروه‌بندی شوند. بر اساس این الگوها، راهبرد این خزشگرهای وبی تنظیم می‌شوند و به‌طور مؤثری در صفحه‌های وب می‌خزند و کیفیت آرشیو را بالا می‌برند.

تعریف: یک الگو از صفحه P_i با طول فاصله l دنباله

$$Patt(P_i) = \{(W_1, T_1); (W_2, T_2); \dots; (W_k, T_k)\}$$

است به طوری که W_k میانگین اهمیت تغییرات در زمان T_k است. مجموع زمانها $\sum_{j=1}^k T_j$ برابر است با l .

ما l را برابر یک روز، به‌عنوان طول الگو، در مدل مورد نظمان انتخاب کردیم.

Page Changes Pattern

| Periods
T | Workdays
ω_k | Saturday
ω_k | Sunday
ω_k |
|----------------|------------------------|------------------------|----------------------|
| [0:00-6:00] | 0,2 | 0,1 | 0,2 |
| [6:00-12:00] | 0,4 | | |
| [12:00-18:00] | 0,6 | 0,4 | 0,35 |
| [18:00-24:00] | 0,1 | 0,2 | 0,13 |

شکل ۱. نمونه الگو

مثال. همانطور که در شکل ۱ نشان داده شده است، دنباله زیر الگوی دوره‌ای صفحه P_i برای روزهای هفته است.

$$Patt(P_i) = \{(0,2, [0h-6h]); (0,4, [6h-12h]); (0,6, [12h-18h]); (0,1, [18h-24h])\}$$

1. off-line

منعکس می‌کنند غیر ممکن است با وجود این قصد داریم که راهبردهای خزش را طوری تنظیم کنیم که خزش وبگاه در آینده تا آنجا که ممکن است جامع و منسجم باشد تعدادی از راهبردهای خزش بر اساس مدل فرآیند پواسون هستند این مدل زمانی مناسب است که تقسیم بندی زمانی بیش از یک ماه باشد به هر حال کار ما بر روی صفحه های وبی است که بارها تغییر میکنند چندین بار در روز مانند صفحه های وب ،تلویزیون رادیو و صفحه های خبری تقسیم بندی زمانی برای این صفحات کمتر از یک ماه است بنابراین مدرک مهمی وجود دارد که مدل Poisson برای این مورد اعتبار ندارد. [79، 25] با بازنگری تغییرات صفحه مشاهده کردیم که روزآمد سازی صفحه

به ساعت روز و روز هفته بستگی دارد. همچنین، مشاهده کردیم که تغییرات در طول روز به میزان قابل توجهی بیشتر از شب و در روزهای کاری بیشتر از آخر هفته هستند.

بنابراین با بازنگری تغییرات صفحه از یک خزش ناپیوسته (1) الگوهای را کشف میکنیم یک الگو را رفتار تغییرات صفحه و اهمیت آنها را در طول زمان به عنوان مثال در طول یک روز طراحی مینمایند. این مسئله به روز هفته بستگی دارد. الگو باید پی در پی روزآمد شود تا همیشه بتواند تغییرات جاری صفحات وب را منعکس کند صفحات یک سایت با الگوی مشابه میتوانند برای داشتن یک الگوی مشترک گروه بندی شوند. بر اساس این الگوها راهبرد این خزشگرهای و بی تنظیم میشوند و به طور مؤثری در صفحههای وب می خزند و کیفیت آرشیو را بالا میبرند.

تعریف یک الگو از صفحه . با طول فاصله 1 دنباله

$$\{ (Patt (P=\{(W,,T,);(W,,T,); \dots; (W,T)$$

است به طوری که میانگین اهمیت تغییرات در زمان T است. مجموع زمانها TW:

است با 1.

ما 1 را برابر یک، روز به عنوان طول الگو در مدل مورد نظرمان انتخاب کردیم.

شکل 1. نمونه الگو

مثال همانطور که در شکل 1 نشان داده شده است دنباله زیر الگوی دوره ای صفحه برای روزهای هفته است.

$$\{ ([Patt(P)-([h-th]); (., [h-h]); (., [h-Ah]); (., [Ah-th$$

ص: 135

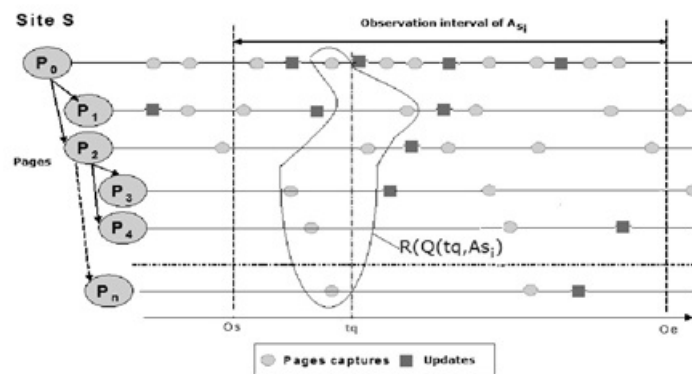
مقادیرهای عددی ۰/۲، ۰/۴، ۰/۴، ۰/۱ به ترتیب میانگین اهمیت تغییرات برای هر بازه زمانی

$$T_1=[0h-6h], \dots, T_i=[(i-1)h-ih]$$

هستند. همچنین، می‌توان همانطوری که در شکل یک نشان داده شد، می‌توان به‌طور جداگانه برای روزهای شنبه و یکشنبه (آخر هفته) تعریف کرد.

مدل آرشیو وب

در مدل آرشیو وب مورد نظر، تمامی صفحه‌های وبگاه، بارها به‌صورت جداگانه ذخیره شده‌اند. خزشگر به‌طور معمول، بر روی تقسیم‌بندی صفحه‌ها کار می‌کند، اما بر روی تقسیم‌بندی یک وبگاه کاری انجام نمی‌دهد. خزشگر، مهم‌ترین صفحه‌های بازسازی شده را انتخاب می‌کند. به‌عنوان مثال، یک صفحه در هر ثانیه مورد خزش قرار می‌گیرد. صفحه‌هایی که بارها تغییرات قابل ملاحظه‌ای را در بردارند، بیشتر بازیابی می‌شوند. آرشیو وبی ما، در AS_i مجموعه‌ای از نسخه‌های صفحه جداگانه‌ای تعریف می‌شود، که از یک سایت S_i بازگذاری شده‌اند. همانطوری که در شکل ۲ نشان داده شده است، یک بازه مشاهده برای تعیین شروع (O_i) و پایان (O_i) مشاهده آرشیو تعریف می‌شود.



شکل ۲. مدل آرشیو وب

تصویربرداری از صفحات i

به روزرسانی i

آرشیو، توسط پرسش استفاده‌کننده‌ای $Q(t_i, A_i)$ که آخرین نسخه‌های در دسترس صفحه‌های وبگاه سایت S را در زمان ارسال پرسش t_i جست‌وجو می‌کند، قابل دسترس است. توجه: فرض می‌کنیم که دنباله $\{S_1, S_2, \dots, S_n\}$ فهرست سایت‌هایی است که باید مورد خزش قرار گیرند. هر سایت از n صفحه $\{P_1, \dots, P_n\}$ تشکیل شده است. هر صفحه P_i الگوی

مقادیرهای عددی ۰/۲، ۰/۴، ۰/۴، ۰/۱ به ترتیب میانگین اهمیت تغییرات برای هر بازه زمانی

هستند. همچنین می‌توان همانطوری که در شکل یک نشان داده شد می‌توان به‌طور جداگانه برای

روزهای شنبه و یکشنبه (آخر هفته) تعریف کرد.

در مدل آرشیو وب مورد نظر تمامی صفحه های وبگاه بارها به صورت جداگانه ذخیره شده. اند خزشگر به طور معمول بر روی تقسیم بندی صفحه ها کار میکند اما بر روی تقسیم بندی یک وبگاه کاری انجام نمی دهد خزشگر مهمترین صفحه های بازسازی شده را انتخاب میکند به عنوان مثال، یک صفحه در هر ثانیه مورد خزش قرار میگیرد صفحه هایی که بارها تغییرات قابل ملاحظه ای را در بردارند بیشتر بازبینی می شوند آرشیو وبی ما در AS مجموعه ای از نسخه های صفحه جداگانه ای تعریف می شود، که از یک سایت S: بارگذاری شده اند همانطور که در شکل 2 نشان داده شده است یک بازه مشاهده برای تعیین شروع (o) و پایان (o) مشاهده آرشیو تعریف می شود.

شکل 2 مدل آرشیو وب

تصویر برداری از صفحات :

به روزرسانی 1

آرشیو، توسط پرسش استفاده کننده ای که آخرین نسخه های در دسترس صفحه های وبگاه

سایت S را در زمان ارسال پرسش t جست و جو میکند قابل دسترس است.

توجه فرض میکنیم که دنباله $\{s_1, s_2, s_3, \dots\}$ فهرست سایتهایی است که باید مورد خزش قرار گیرند هر سایت از 1 صفحه $\{p_1, \dots, p_n\}$ تشکیل شده است. هر صفحه : الگوی P

ص: 136

$\{ (WT); (W,T,); \dots (WT) \}$ را دارد. W ، میانگین اهمیت تغییرات در زمان است. ما به نسخه صفحه P که در زمان t توسط v ذخیره شده است توجه میکنیم فرض میکنیم (P) کپی واقعی صفحه است. که بدون مکث تغییری را در زمان t دنبال میکنند. تعریف آرشیو AS توسط مجموعه ای از نسخه های صفحه ای (P) تعریف می شود. که از سایت S در طول بازه مشاهده $[0]$ ذخیره شده اند؛ به طوری که $in \geq 1$ و $jk \geq 1$.

تعریف $RQt(AS)$ ، نتیجه پرسش $Q(t, AS)$ استفاده کننده تعریف می شود. همانطور که در شکل 2 نشان داده شد $RQt(AS)$ مجموعه ای از n نسخه صفحه های ذخیره شده های $V(P)$ را نشان میدهد که نزدیکترین [صفحه] به زمان داده شده t هستند.

$$R\langle Q(t, A_j) \rangle = \{ V(P) e A - VP \} e A t - t q$$

ما قصد داریم که راهبرد خزش را طوری تنظیم نمایم که کیفیت نتایج بازیابی شده $Q(t, AS)$ را در هر زمان پرسش t افزایش دهد.

کیفیت سنجها

در اینجا، به تعریف دو سنج جامعیت (1) و انسجام میپردازیم که برای ارزشیابی کیفیت آرشیو به کار می روند. جامعیت جامعیت توانایی آرشیو را برای در بر گرفتن تمامی نسخه های کل سایتها را در مقایسه با تعداد کل نسخههایی که میتوان در یک خزش ایده آل ذخیره کرد، اندازه می گیرد.

توجه آرشیو زمانی دارای جامعیت است که تمامی کپیهای واقعی $V(P)$ تمامی صفحه هایی را که میتوان در یک خزش ایده آل ذخیره کرد، در برگرد. بنابراین، هیچ نسخه ای از دست نرفته است.

انسجام انسجام سنجهای است که درجه مجموعه ای از نسخه های صفحه ها را برای انعکاس وضعیت واقعی وبگاه بدون حضور اطلاعات متناقض می سنجد.

توجه: مجموعه ای از صفحات آرشیو شده زمانی دارای انسجام هستند که وضعیت واقعی را حداقل در یک نقطه از زمان منعکس کنند.

دو سنج وزن دار و بدون وزن برای هر دو سنج جامعیت و انسجام تعریف کرده ایم. سنجهای وزن دار اهمیت صفحه و تغییرات مرتبط را مورد توجه قرار میدهند؛ در حالی که سنجهای بدون وزن به نسبت تغییرات توجه دارند این سنجها به دلیل محدودیت جا در این مقاله به طور کامل توضیح داده نشده اند.

ص: 137

با توجه به محدودیت امکانات و ابزارهای قابل دسترس برای ذخیره صفحه ها هدف راهبردی خزش ما تعیین صفحه ها به گونه ای است که کیفیت (جامعیت و انسجام) آرشیو به حداکثر برسد. هر خزشگر میتواند کل M صفحه را در هر بازه زمانی T بارگذاری کند بر اساس الگوها، صفحه ها، براساس ارجحیت و یا ضرورت تعیین میشوند. به هر صفحه، یک ارزش ضرورت (UP) متناسب با اهمیت تغییرات مورد نظر اختصاص داده میشود که توسط الگو در بازه زمانی T مشخص میگردد در هر بازه زمانی MT صفحه با بالاترین ارجحیت ذخیره میشوند ضرورت، صفحه با گذشت زمان تغییر می کند. این ضرورت به زمان آخرین بازسازی و میانگین اهمیت تغییراتی بستگی دارد که توسط الگو مشخص میشود:

P الگوی زیر را دارد:

t زمان جاری است. ($te T$)

W میانگین اهمیت تغییراتی است که توسط الگوی صفحه در بازه زمانی T است.

t آخرین بازسازی، آخرین زمان بازسازی صفحه است.

a ضریب نرمال سازی است. از آنجا که واحدهای اندازه گیری W و زمان t متفاوت هستند، برای کاهش میزان تأثیر ضریب a معرفی شده است.

M صفحه منتخب در یک ترتیب نزولی و براساس میانگین اهمیت تغییر W مرتب و بارگذاری شدند. ریسک روز آمدسازی در طول خزش انسجام) ممکن است توسط ذخیره با اهمیت ترین صفحه ها در حال تغییر در آغاز هر بازه زمانی کاهش داده شود در حقیقت احتمال اینکه یک صفحه نامنسجم باشد، به جایگاه نسبی آن هنگام خزش سایت نیز بستگی دارد. بعد از آن هر صفحه ذخیره شده در زمان (t) با وضعیت قبلی آن صفحه در زمان (1) برای شناسایی تغییرات مقایسه شد. شاید اهمیت تغییراتی که بین دو نسخه شناسایی شدند، به بهره برداری از الگوی روز آمدسازی مربوط شود. به عنوان مثال اگر هیچ تغییری در بازه زمانی [] شناسایی نشود، میانگین اهمیت تغییر توسط الگو در زمان که به بازه [1] متعلق است، می تواند دوباره حساب شود در مقابل اگر تغییری بین دو نسخه ایجاد شود الگوی صفحه نمی توانسته مستقیماً روزآمد شود؛ زیرا ما دقیقاً نمی دانیم که در کدام زمان T تغییر حاصل شده است. به همین دلیل ممکن است جالب باشد که گاهی ابزاری را به نظارت بر تغییرات صفحه و روزآمدسازی الگو اختصاص دهیم.

آزمایش

ما به طور تجربی کارآیی رویکرد خزش پیشنهادی خود را از طریق مقایسه آن با راهبردهای مربوط ارزیابی کردیم آزمایشها بر روی دادههای ترکیبی برای ارزیابی عملکرد راهبردهای مختلف در شرایط آزمایشگاهی کنترل شده انجام شدند میزان امکانات و ابزارهای اختصاص داده شده، نسبت و اهمیت

تغییرات و مانند آن) به ویژه با استفاده از شیبه سازی میزان جامعیت و انسجام به دست آمده (که در بخش 2,3 توصیف شدند) از راهبردهای زیر را با هم مقایسه کردیم:

فراوانی در این راهبرد خزشگر، با توجه به فراوانی تغییرات، صفحات آنها را برای آرشیو کردن انتخاب میکند فراوانی تغییرات هر صفحه بر اساس برآورد کننده [12] Cho ارزیابی شدند. سپس، صفحات منتخب بر اساس زمان آخرین بازسازی در یک ترتیب نزولی مرتب شدند و فقط M صفحه اول بارگذاری شدند. در مواردی که تعداد صفحات انتخابی از M کمتر بود، خزش در بقیه صفحات غیر منتخب براساس آخرین زمان بازسازی انجام شد.

انسجام پیشرفته، این راهبرد با رویکرد پیشنهادی منبع [27] برای گسترش انسجام آرشیو مطابقت دارد. در این راهبرد خزشگر بر تقسیم بندی سایت کار میکند خزشگر مکرراً کل سایت را بارگذاری میکند در حالی که محدودیت صفحه ای را که در طول زمان ذخیره میشود در نظر دارد صفحه ها بر اساس احتمال تغییراتشان به ترتیب نزولی مرتب میشوند برای هر صفحه احتمال عدم انسجام در طول مدت خزش محاسبه میشود، ابتدا صفحه ها با ریسک کمتر بارگذاری می شوند. سپس خزشگر کار ذخیره را در صفحه های جا افتاده ادامه میدهد.

شارک (1) این راهبرد با رویکرد منبع [16] پیشنهاد شده است که با گسترش میزان هشیاری آرشیو مطابقت دارد. همانند رویکرد انسجام، پیشرفته خزشگر بر روی تقسیم بندی سایت کار میکند و مکرراً کل سایت را بارگذاری می کند. در این راهبرد صفحهها با توجه به میانگین نسبت تغییری (8) که توسط مدل Poission محاسبه میشود با ترتیب صعودی مرتب میشوند سپس صفحه هایی ذخیره می شوند که بیشترین تغییر را دارند تا حد امکان به وسط بازه مشاهده نزدیک هستند.

الگوها این راهبرد مبتنی بر الگوهایی است که رفتار تغییر را برای هر صفحه توصیف می کنند راهبرد حاضر، نخستین راهبرد پیشنهادی ماست که در آن الگوی هر صفحه به جای میانگین اهمیت تغییرات به نسبت تغییر - همانند Imp - Pattern - بستگی دارد، بنابراین این راهبرد اهمیت تغییرات بین صفحه ها را در نظر نمی گیرد.

الگوهای ایمپ (2). این دومین راهبرد الگو محور پیشنهادی است که در بخش 3/3 توضیح داده شد. صفحه ها براساس رفتار اهمیت تغییرشان که توسط الگوها تعریف میشوند، بارگذاری میشوند.

تمامی این راهبردها به محدودیتهای امکانات و ابزارها نیز بستگی دارند حداکثر M صفحه در هر زمان T ذخیره میشود تمامی آزمایشهای انجام شده در شرایط یکسان انجام گرفته اند. پیشرفت راهبرد الگوهای ایمپ را با در نظر گرفتن وزن انسجام و جامعیت آرشیو ارزیابی کردیم نتایج نشان دادند که راهبرد ما میزان جامعیت را در مقایسه با الگو 4 درصد، در مقایسه با شارک و انسجام پیشرفته 6 درصد و در مقایسه با فراوانی 25 درصد گسترش میدهد همچنین این راهبردها میزان انسجام را حدود 3 درصد در مقایسه با الگو شارک و انسجام پیشرفته و

مقاله حاضر به موضوع آرشیوسازی صفحات وب به طور کارآمد و به سازی کیفیت آن اشاره دارد. رویکرد پیشنهادی ما سه مفهوم را با هم تلفیق میکند بخش بندی صفحه دیداری شناسایی تغییر، و اهمیت بلاکهای صفحه های وب برای تشخیص بهتر تغییرات مهم میان نسخه ها رویکردهای دیگر آرشیو وب، فقط براساس فراوانی تغییرات هستند. نخستین قدم رویکرد ما ایجاد ساختار دیداری سند بر اساس بلاکهای معنایی خاصی است که الگوریتم VIPS را بسط میدهند [9] قدم دوم، شناسایی تغییرات بین آخرین نسخه صفحه آرشیو شده و نسخه ماقبل آن است الگوریتم Vi-Diff که برای این بخش از کار طراحی شده از روشهای کلی موجود برای ساختار دیداری اسناد مناسبتر است. سپس به موضوع ارزیابی اهمیت تغییرات پرداختیم برای اینکار راهبرد زمان بندی خزشگر را طراحی و با استفاده از آنالیز اهمیت تغییرات به بهینه سازی خزش وب پرداخته شد. آزمایشهای اولیه بر روی بخش بندیها و فازهای مختلف نشان دادند که زمان اجرا امیدوار کننده است. به هر حال زمان بخش بندی بسیار بالاتر از زمان مقایسه است برای بهینه سازی بهتر، سیستم باید بر روی کاهش مدت زمان بخش بندی توجه میکردیم. چالش اصلی در این مقاله بهبود کیفیت آرشیو است. هدف ما این است که راهبردهای خزش را طوری تنظیم کنیم که خزش و نگاههای آینده تا آنجا که ممکن است به طور جامع و منسجم انجام شود. بدین منظور، نوعی راهبرد خزش مبتنی بر الگو را پیشنهاد کردیم. یک الگور رفتار تغییرات صفحه و اهمیت آنها را با گذشت زمان در طول روز طراحی میکند دو سنجه کیفیت جامعیت و انسجام نیز برای ارزیابی راهبرد پیشنهادی تعریف شد سودمندی راهبرد الگو محور خود را با مقایسه آن با راهبردهای مرتبط نشان دادیم در شرایط شبیه سازی شده کامل جامعیت و انسجام کلی هر راهبرد در آرشیوسازی مقایسه شد نتایج نشان دادند که راهبرد ما بهتر از راهبردهای موجود، جامعیت و انسجام را بهبود میبخشد. با وجود، این علاقه مندیم که جامعیت و انسجام به دست آمده را بیشتر گسترش دهیم. کار دیگری که در حال انجام است این است که ما راهبرد خود را بر روی صفحه های وب واقعی بررسی میکنیم. ما میخواهیم که الگوهایی را از صفحه های وب رادیو و تلویزیون کشف کنیم و بعد از این الگوها در تنظیم خزشگرهای وب استفاده کنیم و کیفیت آرشیو را بالا ببریم.

- .The Web archive bibliography, <http://www.ifs.tuwien.ac.at/aola/links/webarchiving.html> [1]
- S. Abiteboul, G. Cobena, J. Masanes, and G. Sedrati. A First Experience in Archiving the [2]
French Web. In ECDL '02: Proceedings of the 6th European Conference on Research
.and Advanced Technology for Digital Libraries, 2002
- H. Artail and K. Fawaz. A fast HTML web page change detection approach based [3]
,on hashing and reducing the number of similarity computations. Data Knowl. Eng
.2008 ,326–337:(2)66
- .M. Ben Saad and S. Gançarski. Using visual pages analysis for optimizing web archiving [4]
.In EDBT/ICDT PhD Workshops, Lausanne, Switzerland, 2010
- M. Ben Saad, S. Gançarski, and Z. Pehlivan. Archiving web pages based on visual [5]
(analysis and diff. In 25^e me Journé es des Bases de Donné es Avancé es (BDA
.Demonstration, Poster), Namur, Belgium, 2009)
- M. Ben Saad, S. Gançarski, and Z. Pehlivan. A novel web archiving approach based on [6]
,visual pages analysis. In the 9th International Web Archiving Workshop (IWA), Corfu
.Greece, 2009
- B. Brewington and G. Cybenko. How dynamic is the web? In In World Wide Web [7]
.conference (WWW'2000), pages 257–276, 2000
- D. J. C. Lampos, M. Eirinaki and M. Vazirgiannis. Archiving the greek web. In 4th [8]
.International Web Archiving Workshop (IWA04), Bath, UK, 2004
- D. Cai, S. Yu, J.–R. Wen, and W.–Y. Ma. VIPS: a Vision–based Page Segmentation [9]

.Algorithm. Technical report, Microsoft Research, 2003

.C. Castillo and B. Sp. Scheduling algorithms for web crawling, 2004 [10]

W. Cathro. Development of a digital services architecture at the national library of [11]

.Australia. EduCause, 2003

J. Cho and H. Garcia-Molina. The Evolution of the Web and Implications for an [12]

Incremental Crawler. In VLDB '00: Proceedings of the 26th International Conference on

.Very Large Data Bases, 2000

J. Cho and H. Garcia-Molina. Estimating frequency of change. ACM Trans. Interet [13]

.Technol., 3(3), 2003

ص: 141

- G. Cobena, S. Abiteboul, and A. Marian. Detecting changes in XML documents. In [14]
.ICDE '02: Proceedings of 18th International Conference on Data Engineering, 2002
- C. N. Cosulschi M. and G. M. Classification and comparison of information structures [15]
.from a web page. In The Annals of the University of Craiova, 2004
- D. Denev, A. Mazeika, M. Spaniol, and G. Weikum. Share: framework for ualityconscious [16]
.web archiving. Proc. VLDB Endow., 2(1):586-597, 2009
- M. K. Evi, M. Diligenti, M. Gori, M. Maggini, and V. Milutinovi. Recognition of [17]
Common Areas in a Web Page Using Visual Information: a possible application in a page
classification. In the proceedings of 2002 IEEE International Conference on Data Mining
.ICDM'02, 2002
- D. Gomes, A. L. Santos, and M. J. Silva. Managing duplicates in a web archive. In SAC [18]
.Proceedings of the 2006 ACM symposium on Applied computing, 2006 :06'
- D. Gruhl, R. Guha, D. Liben-nowell, and A. Tomkins. Information diffusion through [19]
.blogspace. In In WWW '04, pages 491-501. ACM Press, 2004
- X.-D. Gu, J. Chen, W.-Y. Ma, and G.-L. Chen. Visual Based Content Understanding [20]
towards Web Adaptation. In Second International Conference on Adaptive Hypermedia
.and Adaptive Web-based Systems (AH2002), 2002
- R. La-Fontaine. A Delta Format for XML: Identifying Changes in XML Files and [21]
.Representing the Changes in XML. In XML Europe, 2001
- E. Leonardi, T. T. Hoai, S. S. Bhowmick, and S. Madria. DTD-Diff: A change detection [22]
.algorithm for DTDs. Data Knowl. Eng., 61(2), 2007

- C. Olston and S. Pandey. Recrawl scheduling based on information longevity. In WWW [23]
Proceeding of the 17th international conference on World Wide Web, pages 437 :08'
New York, NY, USA, 2008. ACM ,446
- Z. Pehlivan, M. Ben Saad, and S. Gançarski. Vi-diff: Understanding web pages changes [24]
In 21st International Conference on Database and Expert Systems Ap- plications
.DEXA'10), Bilbao, Spain, 2010)
- K. C. Sia, J. Cho, and H.-K. Cho. Efficient monitoring algorithm for fast news alerts [25]
.IEEE Transactions on Knowledge and Data Engineering, 19:950-961, 2007
- R. Song, H. Liu, J.-R. Wen, and W.-Y. Ma. Learning block importance models for web [26]
pages. In WWW '04: Proceedings of the 13th international conference on World Wide
.Web, 2004

M. Spaniol, D. Denev, A. Mazeika, G. Weikum, and P. Senellart. Data quality in web [27]
archiving. In WICOW '09: Proceedings of the 3rd workshop on Information credibility
.on the web, pages 19–26, New York, NY, USA, 2009. ACM

Y. Wang, D. DeWitt, and J.-Y. Cai. X-Diff: an effective change detection algorithm for [28]
XML documents. In ICDE '03: Proceedings of 19th International Conference on Data
Engineering, March 2003

ص: 143

رشد فزاینده وب جهانگستر چالشهایی را در حوزه حفاظت مؤثر داده های وب مطرح ساخته است. ابزارهایی که امروزه آرشیو سازان وب به کار میبرند، در وب به صورت کورکورانه خزش و بدین طریق بدون توجه به نوع صفحه هایی که در دسترس قرار دارند (که به راهبردهای غیر بهینه (1) خزش منتهی میشود) و بدون توجه به ساختاری که محتوای صفحه هایی در خود دارند که منجر به گردآوری منابعی در سطح صفحه می شود که محتوای آن را به دشواری می توان مورد بهره برداری قرار داد اقدام به گردآوری میکنند این نوشته، معطوف به خزش و آرشیوسازی برنامه های کاربردی وبی به ویژه وب اجتماعی است که برای عموم قابل دستیابی اند. برنامه کاربردی وب به هر نوع برنامه ای گفته میشود که از استانداردهای، وب نظیر HTML و HTTP برای انتشار اطلاعات در وب استفاده کنند و توسط مرورگر (2) های وبی قابل دسترسی اند. تالارهای گفت و گوی وبی شبکه های اجتماعی، خدمات مکانهای جغرافیایی (3) و مانند آن از آن جمله اند. ادعای ما این است که بهترین راهبرد برای خزش برنامه های کاربردی آن است که خزشگر وب را طوری طراحی کنیم که نوع برنامه کاربردی وبی را که در حال خزش آن است بشناسد و بتواند فهرست یوآرال هایی را که در صف خزش قرار دارند پالایش نموده اطلاعاتی درباره ساختار محتوای خزش شده به آرشیو ارائه دهد. برای این، کار ویژگیهایی به یک خزشگر مختص آرشیو وب میافزاییم که عبارت است از: توان تشخیص این که یک صفحه متعلق به چه نوع برنامه کاربردی وبی است و کاربرد روش شناسی خزش و برداشت متناسب با آن.

مقوله ها و توصیفگرهای موضوعی

کلیدواژه ها برنامه کاربردی، وب، آرشیووب، خزش برداشت داده XPath

ص: 144

Suboptimal -1

Browser -2

Geolocation services -3

*خزش هوشمند در برنامه های کاربردی وب(1)

محمد فهیم زیر نظر پیر سنلار(2) | ترجمه فرزانه شادان پور(3)

1. مسئله

از زمان معرفی وب 2 وب اجتماعی(4) منبع مهمی برای برداشت(5) محتوا شده است. میلیونها کاربر از وب اجتماعی به عنوان وسیله ای برای انتشار اطلاعات بحث درباره موضوعات، سیاسی به اشتراک گذاشتن محتوای ویدئویی ارسال نظرات مدیریت وب نوشت و بیان نظرات(6) شخصی خود درباره مباحث روز استفاده می کنند فقط کاربران عادی از وب اجتماعی استفاده نمی کنند، بلکه هر روز توجه رهبران سیاسی بیش از پیش به این پدیده جلب میشود. امروزه در امریکا و انگلستان پاسخ به پرسشهای پارلمانی با استفاده از توئیتر امری عادی شده است. به تازگی در 6 جولای، 2011 باراک اوباما نام خود را به عنوان اولین رئیس جمهوری که از توئیتر به عنوان ابزاری برای ارتباط جمعی استفاده کرده است، ثبت نمود [13]. بنابراین وب اجتماعی به صورت بخشی از مبارزات سیاسی و هدایتگر برنامه آینده سیاسی درآمده است.

ص: 145

Intelligent crawling of Web applications for Web archiving –1

Muhammad Faheem, supervised by Pierre Sellenart –2

3- مربی عضو هیئت علمی سازمان اسناد و کتابخانه ملی ایران

4- Social Web

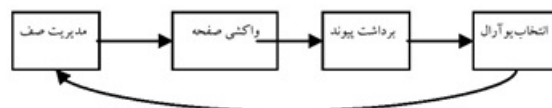
5- Extraction

6- Comments

این مسئله، لزوم حفظ داده‌های اجتماعی را تشدید می‌کند.

ولی آرشيو کردن داده‌ها از وب اجتماعی به روشی هوشمند، هنوز چالشی برای اهل فن است. هدف ما پرداختن به این چالش با معرفی رویکردی انطباقی است که عمدتاً بر شناخت انواع گوناگون برنامه‌های کاربردی مورد استفاده در وب اجتماعی مبتنی است. برنامه کاربردی وب، به هر برنامه مبتنی بر HTTP گفته می‌شود که از وب جهانگستر برای انتشار اطلاعات استفاده می‌کند، در این نوشتار به ویژه بر وجوه اجتماعی وب تأکید خواهیم کرد زیرا، همچنان که به‌عنوان نمونه در تالارهای گفت‌وگوی وب^۱، وب‌نوشت یا تویتر مشاهده می‌شود؛ تکیه زیادی بر محتوای تولیدشده توسط کاربران، تعاملات اجتماعی، و شبکه می‌شود.

برای آگاهی از اینکه چگونه توانمندی ابزارهای فعلی آرشيو سازی وب در حد کارکرد آرشيو کردن وب اجتماعی نیست، معماری ساده شده یک خزشگر وب سنتی (مثل هریتریکس^۲ [۲۳]) را که در تصویر ۱ نشان داده شده است، در نظر بگیرید. یک خزشگر وب (که آن را عنکبوت یا روبات نیز می‌نامند) نوعی برنامه رایانه است که وب را به گونه‌ای روشمند و ارسی می‌کند و مدارک مورد نظر را می‌یابد. خزشگرهای سنتی از حیث مفهومی به شیوه‌ای بسیار ساده در وب خزش می‌کنند. آنها کار خود را از یک فهرست یوآرال هسته^۳ آغاز می‌کنند که در صفی^۴ نگه داری می‌شوند، (یوآرال ممکن است صفحه نخست یک وبگاه باشد). سپس، صفحه‌های وب یکی پس از دیگری از صف انتخاب و واکنشی^۵ می‌شوند. پیوندها [یا همان یوآرال‌ها] موجود در محتوای این صفحه‌ها واکنشی شده استخراج و برداشت می‌شوند. اگر این پیوندها در دامنه کار آرشيو قرار بگیرند، یو آر ال‌های تازه برداشت شده به صف اضافه می‌شوند. این فرآیند پس از زمان مشخصی، یا هنگامی که دیگر یوآرال مناسب دیگری در صف نباشد متوقف می‌شود.



تصویر ۱- زنجیره فرایند سنتی خزش در یک خزشگر وب

در رویکرد بالا، با چالش‌های خزش برنامه‌های کاربردی وب مواجه نخواهیم بود. ماهیت برنامه کاربردی وب که مورد خزش واقع شده است، یا محتوای مورد نظر، در تصمیم‌گیری‌های راهبردی

1. Web forums
2. Heritrix
3. Seed URLs (Uniform Resource Locator)
4. Queue
5. Fetch

این مسئله لزوم حفظ داده‌های اجتماعی را تشدید میکند

ولی آرشيو کردن داده‌ها از وب اجتماعی به روشی هوشمند، هنوز چالشی برای اهل فن است. هدف ما پرداختن به این چالش با معرفی رویکردی انطباقی است که عمدتاً بر شناخت انواع گوناگون برنامه‌های کاربردی مورد استفاده در وب اجتماعی مبتنی است برنامه کاربردی وب، به هر برنامه مبتنی HTTP گفته میشود که از وب جهانگستر برای انتشار اطلاعات استفاده میکند، در این نوشتار به ویژه بر وجوه

اجتماعی وب تأکید خواهیم کرد زیرا همچنان که به عنوان نمونه در تالارهای گفت وگویی وب(1)، وب نوشت یا توییتر مشاهده میشود تکیه زیادی بر محتوای تولید شده توسط کاربران تعاملات، اجتماعی و شبکه می شود.

برای آگاهی از اینکه چگونه توانمندی ابزارهای فعلی آرشیوسازی وب در حد کارکرد آرشیو کردن وب اجتماعی نیست معماری ساده شده یک خزشگر وب سنتی (مثل هریتریکس(2) [23]) را که در تصویر 1 نشان داده شده است، در نظر بگیرید یک خزشگر وب (که آن را عنکبوت یا روبات نیز می نامند) نوعی برنامه رایانه است که وب را به گونه ای روشمند واریسی میکند و مدارک مورد نظر را می یابد. خزشگرهای سنتی از حیث مفهومی به شیوه ای بسیار ساده در وب خزش میکنند آنها کار خود را از یک فهرست یوآرال هسته(3) آغاز میکنند که در صفی(4) نگه داری میشوند (یو آرال ممکن است صفحه نخست یک وبگاه باشد)، سپس صفحه های وب یکی پس از دیگری از صف انتخاب و واکشی(5) میشوند پیوندها یا همان یوآرالهای موجود در محتوای این صفحه ها واکشی شده استخراج و برداشت میشوند. اگر این پیوندها در دامنه کار آرشیو قرار بگیرند یو آرالهای تازه برداشت شده به صف اضافه میشوند. این فرآیند از زمان مشخصی یا هنگامی که دیگر یوآرال مناسب دیگری در صف نباشد متوقف می شود.

تصویر 1- زنجیره فرایند سنتی خزش در یک خزشگر وب

در رویکرد بالا- با چالشهای خزش برنامه های کاربردی وب مواجه نخواهیم بود. ماهیت برنامه کاربردی وب که مورد خزش واقع شده است یا محتوای مورد نظر در تصمیم گیریهای راهبردی

ص: 146

Web forums -1

Heritrix -2

(Seed URLs (Uniform Resource Locator -3

Queue -4

Fetch -5

خزش، مدنظر گرفته نمی شود. برنامه های کاربردی وب با محتوای پویا(1) (مثل تالارهای گفت وگویی وب، وب نوشت ها و مانند آن) ممکن است به گونه ای ناکارآمد خزش شوند؛ اگر محتوا به گونه ای باشد که فقط با عملیات پیچیده خزش (مثل پرس و جوی AJAX ارسال فرم) قابل دسترس باشد، بعضی از بخشهای آن ممکن است از دست برود.

به عنوان نمونه تالارهای گفت وگویی وب دارای ویژگی پویایی هستند به این معنا که برای برداشت محتوای معنایی یا بهبود عملکرد خزش در مورد آنها ماهیت آنها باید شناخته شود. محتوای تالارهای گفت وگویی وب اغلب در یک پایگاه داده نگهداری می شود. هنگامی که یک کاربر، پرس و جویی ارسال می کند صفحه پاسخ به صورت خودکار و با استفاده از یک قالب(2) از پیش تعریف شده تولید می شود. وقتی دو پرس و جو بخش واحدی از محتوای چنین صفحه ای را تقاضا میکنند سرور دو صفحه پویا با محتوای یکسان یا مشابه، ولی با یوآرال مختلف تحویل کاربر می دهد. ولی این صفحه های پویا موجب پدیده افزونگی(3) میشوند که ممکن است زیان آور باشد؛ از این جهت که موجب میشوند منابع بیشتری برای خزش وجود داشته باشند و در نتیجه آرشیو نهایی از کیفیت خوبی برخوردار نخواهد بود. خدمات وب نوشت نیز اطلاعات مکرر در خود دارند به عنوان مثال به صورت ماهانه و سالانه آرشیو می شوند و دربردارنده محتواهایی هستند که دارای تغییرات جزئی است و مکرر محسوب میشود در صورت خزش تالارهای گفت وگو و وب نوشت ها با رویکرد سنتی، خزش با موارد بسیار زیادی از منابع مکرر مواجه میشویم در نهایت خزشگر در تله خزش(4) گیر میافتد چرا که باید بینهایت پیوند را خزش کند. همچنین پیوندهای دارای اختلال(5) نظیر صفحه های مناسب، چاپ یا، تبلیغات و مانند آنها وجود دارند که بهتر است از گردآوری و ذخیره آنها هنگام ایجا آرشیو خودداری کرد. با رویکردهای سنتی خزش، از برنامه های کاربردی وب نیز که در حد بالایی پردازش نویسی(6) شده اند یا از وب عمیق(7) (داده های قابل دسترس از ویرای فرمها) داده نمیتوان برداشت کرد.

دست آخر اینکه خزشگرهای آرشیو وب که در رویکرد سنتی طراحی شده اند، برای برداشت داده به طرق گوناگون تلاشی نمیکنند؛ حال آنکه آرشیو داران و کاربران آرشیوهای وب دوست دارند به ابعاد معنایی بیشتری درباره محتوای آرشیو نظیر برچسبهای زمان(8) پیامهای وب نوشتهها، و توصیفگرهای نویسندگان دست یابند حتی اگر این کار با استفاده از ناوبری های(9) پیچیده در صفحه ها انجام شود که اطلاعات مربوط از صفحه ها مختلف را سازمان می دهند. به عنوان مثال، یک برنامه کاربردی وب که محتوای خود را به گونه ای سامان میدهد که برای برداشت اطلاعات مورد نظر مرور تقویم ضروری باشد،

ص: 147

Dynamic -1

Template -2

Redundancy -3

Spider trap -4

Noisy -5

Scripting -6

Deep Web -7

Timestamps -8

Navigation -9

یا در مورد یک تالار گفت و گوی وب جایی که یک روند(1) میتواند ورودیهای مختلفی متشکل از صفحه ها مختلف داشته باشد و برای برداشت این اطلاعات ناوبری مؤثر صفحه ضروری باشد رویکرد سنتی خزش برای انجام کارآمد این فرآیندها با محدودیتهایی مثل صرف زمان بیشتر در خزش برای تعداد اندکی صفحه مناسب بدون در برداشتن اطلاعات معنایی در محتوا، مواجه خواهد بود.

در ادامه به طور خلاصه وضعیت فعلی فناوری خزش در وب را بررسی میکنیم سپس در بخش 3 رویکرد پیشنهادی ارائه میشود که عبارت است از وارد کردن یک راهنمای آگاه از برنامه کاربردی(2) در فرآیند خزش که به خزشگر در تمام فرآیند خزش کمک و خزش مؤثر داده ها را تضمین میکند. در بخش 4 با توصیف الگوهای اکتشاف(3) برنامه کاربردی وب پیشنهادی مان که ساختاری برای پایگاه دانش(4) یک برنامه کاربردی وب است میکوشیم تا روش شناسی خود را به تفصیل شرح دهیم. همچنین توضیح خواهیم داد که چه نوع از عملکرد خزش مورد نظر ماست. در بخش 5 نتایج اولیه را بر خواهیم شمرد و مقاله را با بحثی درباره پژوهشهای آینده به پایان میبریم.

2. وضعیت فعلی فناوری خزش

خزش در وب مسئله ای است که مطالعات خوبی درباره آن صورت گرفته، ولی همچنان چالش برانگیز است. ژولین ماسانه(5) در [18] مروری بر حوزه آرشیو وب و خزش برای ایجاد آرشیو وب انجام داده است. او به ویژه درباره خزش وب عمیق بحث کرده و بر نیاز به آرشیو داده از سطح و عمق وب برای حفاظت وب تأکید میکند.

یک خزشگر کانونی(6) بر اساس مجموعه ای از موضوعات از پیش تعیین شده خزش می کند [7]. در این روش رفتار خزشگر نه بر مبنای ساختار برنامه کاربردی- وب که هدف ماست - بلکه بر مبنای محتوای صفحه های وب تنظیم میشود در رویکرد ما هدف خزش کانونی نیست، بلکه جای آن را رویکرد بهتری برای برنامه های کاربردی وب میگیرد. هر دو راهبرد، روشهای تکمیل کننده ای برای ارتقای عملکرد خزشگر سنتی به حساب می آیند.

در یک برنامه کاربردی وب یا یک سامانه مدیریت محتوا محتوا با توجه به یک قالب(7) (اجزای قالب به عنوان مثال شامل نوار سمت چپ یا راست صفحه، وب نوار ناوبری، صفحه سرصفحه و پاصفحه منوی اصلی و... است).

از میان آثار متعدد درباره برداشت قالب(8)، گیسون پونرا و تامکینز [10] زمینه محتوای مبتنی بر قالب را در وب مورد بررسی قرار داده اند آنها دریافته اند که 40 تا 50 درصد محتوای وب در (2005) مبتنی

ص: 148

Thread -1

Application-aware helper -2

Detection patterns -3

Knowledge base -4

Julien Masanes -5

Focused or goal-directed crawler -6

Template -7

بر قالب است یعنی بخشی از یک برنامه کاربردی وب است. یافته های آنها همچنین نشان میدهد که صفحه های وب با نرخ معادل 6 تا 8 درصد در سال در حال رشد هستند. این پژوهش با قوت به منافع خزشگری که بتواند به شیوه ای خاص برنامه های کاربردی وب را خزش کند اشاره دارد.

گرچه جایی که ما میدانیم هنوز خزش آگاه از برنامه های کاربردی به صورت عام مورد توجه قرار نگرفته است، تلاشهایی برای برداشت محتوا از تالارهای گفت وگویی وب صورت پذیرفته است [6، 11]. اولین مورد که خزش تالار گفت و گو (1) نامیده میشود، [11] ساختار سازمان یافته تالارهای گفت و گوی وب را هدف خزش قرار میدهد و رفتار کاربر را در فرآیند برداشت اطلاعات شبیه سازی میکند. خزش تالار گفت و گو با مسئله به طور بسیار مؤثری سروکار پیدا میکند ولی هنوز دچار محدودیتهایی است که ناشی از داشتن قواعدی ساده است و تنها میتواند در مورد تالارهایی با ساختار ساماندهی مشخص به کار رود رویکرد دوم [6] وابسته به ساختار تالار گفت و گوی وب نیست در این روش سامانه iRobot فرآیند برداشت را با فراهم کردن نقشه سایت (2) برنامه کاربردی وب که باید خزش شود، همراهی می کند. نقشه سایت از طریق خزش تصادفی چند صفحه از برنامه کاربردی ساخته میشود. این فرآیند به شناسایی مناطقی که به شدت تکرار پذیرند کمک و بعد آنها را بر مبنای آرایششان خوشه بندی می کند [25]. بعد از تولید نقشه سایت iRobot ساختار تالار گفت و گوی وب را در شکل یک گراف جهت دار (3) متشکل از رئوس (4) (صفحه ها وب) و یالهای (5) جهت دار (پیوندهای میان صفحه ها مختلف وب) می سازد. علاوه بر این مسیر نیز به منظور فراهم کردن بهترین مسیر عبور که فرآیند برداشت را هدایت کند، و از واکنشی صفحه های مکرر و بی ارزش نیز جلوگیری نماید مورد تجزیه و تحلیل قرار می گیرد. هدف ما توسعه فناوریهای مشابه برای برنامه های کاربردی اختیاری (6) بوده است و نه فقط برای تالارهای گفت و گوی وب در اینجا توضیح این نکته لازم است که رویکرد ما برای کشف نوع برنامه کاربردی وب که خزش شده است، بر مکانیسمی عام (7) استوار است آثاری درباره وب نوشت ها و تالارهای گفت و گوی خاص نگارش یافته اند. در [14] به ویژه از مدل (8) SVM برای کشف صفحه ای که متعلق به یک وب نوشت است استفاده شده است. [19 5 SVM] برای دسته بندی متون بسیار مورد استفاده قرار می گیرد. در [14] SVM ها با استفاده از بردارهای مختلف سنتی که از دسته های (9) واژگان یا n-grams موجود در در محتوا تشکیل شده اند آموزش داده میشوند. پاره ای از ویژگیهای (10) جدید نیز برای کشف وب نوشت، نظیر دسته های یو آر ال پیوند داده شده و دسته های گرانوازه (11) معرفی می شوند. برای انتخاب ویژگیها

ص: 149

Board Forum Crawling -1

Sitemap -2

Directed graph -3

Vertices -4

Arcs -5

Arbitrary -6

Generic mechanism -7

Support Vector Machine -8

Bags -9

Features -10

از آنتروپی (1) نسبی استفاده می شود. براساس نتایجی که در این نوشته آمده است، شناسایی وب نوشت با استفاده از ویژگی‌هایی مرکب از یو آر ال‌ها گرانوازه ها و ابر برچسبها بهتر انجام شده است. با این حال در موضوع کشف برنامه های کاربردی وب روشهایی برای کشف منابع وب پنهان وجود دارد [3 و 4] از جمله تجربیات در این زمینه خزشگر فرم کانون (2) - [3] برای کشف پایگاههای اطلاعاتی برخط بوده است. در این روش، از یک خزشگر کانونی برای بارگذاری منابع در موضوعات مورد نظر به کمک یک برنامه دسته بندی پیوند (3)، استفاده میشود برنامه دسته بندی، پیوند فرمهای قابل جست و جورا کشف و پیوندهای موجود در آنها را الویت بندی میکند در این روش محدودیتهایی نیز وجود دارد که از آن جمله اند: تنظیم دستی وابستگی به برنامه دسته بندی پیوندها و نتایج ناهمگن. در روشی که اخیراً به کار رفته [14] و سازگار یافته تر است خزشگر سازگار یافته برای منابع وب پنهان در مورد این محدودیتهای ملاحظاتی صورت گرفته است در این روش نقاط ورود به وب پنهان به گونه ای کارآمد و با الگویی ناشناخته اکتشاف میشود و این تجربه به طور خودکار به فرآیند یادگیری افزوده می شود.

آثاری [1، 15، 16] نیز به طور عام معطوف به شناسایی دسته بندیهای عام (مثل وب نوشت، وبگاههای دانشگاهی شخصی و مانند آن) برای یک وبگاه با استفاده از برنامه دسته بندی مبتنی بر ویژگیهای ساختاری صفحه وب است که میکوشد تا کارکرد پذیریهایی (4) این صفحه های وب را کشف کند این امر مستقیماً در تنظیمات کار ما قابل به کارگیری نیست؛ نخست به این علت که این متون در سطح وبگاه قابل به کارگیری اند (مثلاً) مبتنی بر صفحه نخست (اند و نه در سطح صفحه های منفرد کاملاً روشن است که همه این فنون شناسایی برنامه های کاربردی بر برنامههای دسته بندی آموزش دیده متکی هستند. این مسئله میتواند یکی از جهات ممکن برای مکانیسم کشف برنامه کاربردی وب در این پژوهش باشد.

3. رویکرد پیشنهادی

مدعای اصلی ما این است که برای انواع برنامه های کاربردی وب فنون خزش متناسب و متفاوتی به کار برده شود؛ یعنی داشتن راهبردهای خزش متفاوت برای انواع وبگاههای اجتماعی (وبنوشتها، ویکیها، شبکه های اجتماعی نشانکهای اجتماعی (5)، ریز بلاگها (6)، شبکه های موسیقی (7)، تالارهای گفت و گوی، وب، شبکههای عکس (8)، شبکه های ویدئویی (9)، و مانند آن) برای سیستمهای ویژه مدیریت محتوا (مانند Wordpress, php BB) و برای سایتهای خاص مانند توئیتر و فیس بوک در رویکردی که ما در این پژوهش در پیش گرفته، ایم نوع برنامه کاربردی وب (نوع به طور کلی سیستم مدیریت محتوا، یا

ص: 150

Relative entropy -1

Form-focused crawler -2

Link classifier -3

Functionalities -4

Social bookmarks -5

Microblogs -6

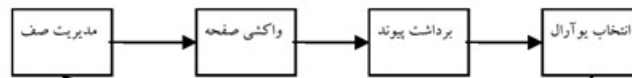
Music networks -7

Photo networks -8

خزش هوشمند در برنامه های کاربردی وب ۱۵۱

سایت) حین عملکرد خزش، و نوع صفحه وب درون برنامه کاربردی (مانند پروفایل کاربر در یک شبکه اجتماعی) کشف و در مورد عملیات مورد نیاز برای خزش بعدی (تعقیب پیوند، استفاده از یک API، ارسال فرم، برداشت محتوای ساختارمند) تصمیم گیری می شود. برای سازگار کردن یک خزشگر سنتی با الزامات و نیازمندی های این پژوهش، ما معماری سنتی یک خزشگر وب را آن گونه که در شکل ۲ نشان داده شده است بسط دادیم. ماژول واکنشی^۱ صفحه (نگاه کنید به شکل ۱) به ماژول واکنشی دقیق تری بسط داده شده است که قادر است منابعی را بازیابی کند که با پرس و جوی HTTP GET قابل دسترسی نیستند (بلکه با یک رشته پرس و جو، یا با استفاده از یک API)، یا اشیای وبی منفردی در یک صفحه وب هستند (مانند یک پست و بنوشت، یک دیدگاه^۲، یا نام یک پست گذار^۳). در واقع، به جای کارکرد معمول برداشت پیوند، ماژول راهنمای آگاه از برنامه کاربردی به کار برده می شود تا در اثنای خزش بتواند برنامه کاربردی را شناسایی کند و ضمن تقسیم بندی عملکردهای خزشی، که می توان برای هر برنامه انجام داد، در این مورد تصمیم گیری کند.

این تغییرات در دو خزشگر اعمال می شوند: خزشگر تحت مالکیت بنیاد حافظه اینترنت^۴، که ما با آنها همکاری نزدیکی داریم؛ و نسخه سفارشی خزشگر هریتریکس^۵ [۲۳]، که توسط آزمایشگاه پژوهشی ATHENA در چارچوب طرح ARCOMEN [۲] تهیه شده است. در بخش بعدی این نوشتار ماژول راهنمای آگاه از برنامه های کاربردی با جزئیات بیشتری تشریح می شود.



شکل ۲- معماری بسط داده شده خزشگر

۴. روش شناسی

در این بخش، ماژول راهنمای آگاه از برنامه های کاربردی معرفی می شود. این ماژول خزشگر آرشیو وب را در فراهم آوری محتوا از وب اجتماعی به شیوه ای هوشمند و سازگار یافته یاری می رساند. این ماژول قابلیت های کارکردی خزشگر را افزایش می دهد و فرآیند خزش را کارآمدتر می کند.

پایگاه دانش در برنامه کاربردی وب، خزشگر با پایگاهی از دانش در مورد برنامه های کاربردی وب یاری می شود که در آن نحوه خزش هوشمند و بگاهها توصیف می شود. این پایگاه دانش چگونگی اکتشاف برنامه های ویژه کاربردی وب و عملکرد خزشی را که برای آن لازم است به اجرا گذاشته شود،

1. Fetching
2. Comment
3. Poster
4. Internet Memory Foundation
5. Heritrix

سایت) حین عملکرد خزش و نوع صفحه وب درون برنامه کاربردی (مانند پروفایل کاربر در یک شبکه اجتماعی) کشف و در مورد عملیات مورد نیاز برای خزش بعدی تعقیب پیوند استفاده از یک API ارسال فرم برداشت محتوای ساختارمند تصمیم گیری میشود برای سازگار کردن یک خزشگر سنتی با الزامات و نیازمندی های این پژوهش ما معماری سنتی یک خزشگر وب را آن گونه که در شکل ۲ نشان داده شده است بسط دادیم. ماژول واکنشی^(۱) صفحه (نگاه کنید به شکل ۱) به ماژول واکنشی دقیق تری بسط داده شده است که قادر است منابعی را بازیابی کند که با پرس و جوی HTTP GET قابل دسترسی نیستند بلکه با یک رشته پرس و جو یا با استفاده از یک (API) یا

اشیای و بی منفردی در یک صفحه وب هستند (مانند یک پست وب نوشت، یک دیدگاه(2)، یا نام یک پست گذار(3)). در واقع، به جای کارکرد معمول برداشت پیوند ماژول راهنمای آگاه از برنامه کاربردی به کار برده میشود تا در اثنای خزش بتواند برنامه کاربردی را شناسایی کند و ضمن تقسیم بندی عملکردهای خزشی که میتوان برای هر برنامه انجام داد، در این مورد تصمیم گیری کند.

این تغییرات در دو خزشگر اعمال میشوند خزشگر تحت مالکیت بنیاد حافظه اینترنت(4)، که ما با آنها همکاری نزدیکی داریم؛ و نسخه سفارشی خزشگر هریتریکس(5) [23] که توسط آزمایشگاه پژوهشی ATHENA در چارچوب طرح [ARCOMEN 2] تهیه شده است. در بخش بعدی این نوشتار ماژول راهنمای آگاه از برنامههای کاربردی با جزئیات بیشتری تشریح میشود.

شکل-2 معماری بسط داده شده خزشگر

4. روش شناسی

در این بخش ماژول راهنمای آگاه از برنامه های کاربردی معرفی میشود این ماژول خزشگر آرشیو وب را در فراهم آوری محتوا از وب اجتماعی به شیوهای هوشمند و سازگار یافته یاری میرساند. این ماژول قابلیت های کارکردی خزشگر را افزایش میدهد و فرآیند خزش را کارآمدتر میکند.

پایگاه دانش در برنامه کاربردی وب خزشگر با پایگاهی از دانش در مورد برنامه های کاربردی وب یاری میشود که در آن نحوه خزش هوشمند و بگاهاها توصیف می.شود این پایگاه دانش چگونگی اکتشاف برنامه های ویژه کاربردی وب و عملکرد خزشی را که برای آن لازم است به اجرا گذاشته شود،

ص: 151

Fetchng -1

Comment -2

Poster -3

Internet Memory Foundation -4

Heritrix -5

مشخص می. کند پایگاه دانش از تقسیم بندیهای عمومی گرفته تا موارد خاص (وبگاهها) یک برنامه کاربردی به شیوه سلسله مراتبی تنظیم میشود برای مثال وبگاههای شبکه های اجتماعی را میتوان به وب نوشتهها تالارهای گفت و گوی، وب. ریز بلاگها شبکه های، ویدئویی و مانند آن تقسیم بندی کرد غیر از این میتوان برنامه های کاربردی را براساس سامانه های مدیریت محتوا به گونه دیگری تقسیم بندی کرد. مثلاً Wordpress و Movable Type نمونه هایی از سامانه های مدیریت وب نوشت هستند. phpBB و vBulletin نیز نمونه هایی از سامانه های مدیریت محتوای تالارهای گفت و گوی وب.

علاوه بر این هر برنامه کاربردی وب معمولاً از انواع مختلف صفحه وب تشکیل می شود. در یک تالار گفت و گوی، وب صفحه ها وجود دارند که فهرستی از تالارها را ارائه میکنند صفحه هایی که فهرستی از پستهای مربوط به تالارهای گفت و گوی خاص را نشان میدهند، صفحه هایی که به پستهای منفرد، همراه با دیدگاههای آنها ارجاع می دهند. بنابراین، پایگاه دانش، انواع مختلف صفحه های وب تحت یک برنامه کاربردی را توصیف میکند و بر این مبنا میتوانیم عملکردهای مختلف خزشی را که باید در قبال هر نوع از صفحه به اجرا بگذاریم - تعیین کنیم.

پایگاه دانش باید در زبانی روان و گزاره ای بیان شود تا به آسانی قابل روز آمدسازی و اشتراک باشد ناآشنایان به برنامه نویسی هم بتوانند با آن کار کنند و حتی اگر بشود از مثالها به طور خودکار یادبگیرند کنسرسیوم وب جهانی(1)، نوعی زبان توصیف برنامه کاربردی به نام [12] (2) WADL را تهیه و برای بهره برداری ارائه کرده است که توصیف منابع سازگار با انتقال در پروتکل HTTP را در قالبی با قابلیت پردازش، ماشینی ممکن می سازد WADL برای توصیف مجموعه های منابع روابطشان با یکدیگر روشی که برای هر منبع باید به کار رود و قالبهای نمایش منبع به کار می رود WADL را می توان به عنوان جزء قابل کاربرد و قالب ارسال برای پایگاه دانش به کار برد، ولی همه نیازهای ما را مرتفع نمی سازد؛ از جمله توصیف الگوهای بازشناسی برنامه کاربردی وب و تعاملات برنامه کاربردی که ورای یک تقاضای GET و POST ساده جریان دارند در نتیجه پایگاه دانش ما باید در قالب XML توصیف شود و به خوبی با ساختار درختی سلسله مراتب برنامه کاربردی وب و سطوح مختلف صفحه ها انطباق یابد.

ماژول کشف برنامه کاربردی وب تشخیص برنامه کاربردی وب و پس از آن در پیش گرفتن بهترین راهبرد خزش متناسب با آن مهمترین چالش در خزش و برداشت محتواست. در مورد شناسایی برنامه های کاربردی وب کار زیادی انجام نشده ولی تلاشهایی برای دسته بندی صفحه ها وب تحت برنامه های کاربردی وب مختلف صورت گرفته است، [11، 15، 16] برای کشف یک برنامه کاربردی وب خاص پایگاه دانش با توصیف قواعد مختلف بر مبنای الگوهای یوآرال فراداده های HTTP، محتوای متنی، الگوهای XPath(3) ارجاعات به سامانه دسته بندی و در صورت امکان ویژگیهای مبتنی بر گراف

ص: 152

World Wide Web Consortium -1

Web Application Description Language -2

3- زبان مسیر، xml زبان پریشی برای انتخاب گره ها از سند xml است که توسط WWC تعریف شده است.

وب(1) را میسر می سازد. شناسایی سطح صفحه درون هر برنامه کاربردی وب ممکن است با طبقه بندی صفحه متناسب با ویژگیهای ساختاری انجام شود.

بد نیست در اینجا سامانه مدیریت محتوای تالار گفت وگویی وب موسوم به vBulletin را ذکر کنیم که به عنوان نمونه از طریق جست و جوی یک ارجاع به پردازش نوشته(2) جاوای vbulletin_global.js با استفاده از عبارت XPath ساده //script/@src قابل شناسایی است.

صفحه های سطح «فهرست تالار گفت و گو هنگامی که با عبارت XPath //a[@class="forum"]@href در بیانند قابل شناسایی خواهند بود.(3)

خزش و برداشت. بعد از کشف برنامه کاربردی که صفحه وب متعلق به آن است، مرحله بعدی این است که عملکردهای مناسب خزش را تعیین کنیم دامنه عمل خزش فراتر از افزودن فهرست یوآرال به صف یوآرلهای در دست خزش است. این امر شامل هر عملکردی است که در فرآیند خزش به نحوی دخیل است استفاده از API برای برداشت دادههای مرتبط از سایتهای شبکه های اجتماعی نظیر توییتر یا انجام تعاملات پیچیده با برنامه های کاربردی بر مبنای AJAX یا شناسایی موجودیتهای وی، به ویژه برنامه های کاربردی وب عملکردهای خزش به معنای خاص آن به دو نوع اند:

عملکردهای ناوبری ناوبری به یک صفحه وب یا منابع وی.

عملکردهای برداشت برداشت موجودیتهای معنایی منفرد از صفحه های وب (مانند برچسبهای

زمان(4)، پستهای وب نوشت، دیدگاهها).

مشابه همین ما یک زبان اعلانی(5) برای توصیف همه عملکردهای خزش (مطلوب است یک پایگاه دانش که به سادگی قابل نگهداری باشد از جمله نگهداری ماشینی داشته باشیم). میخواهیم بنابراین به یک زبان برنامه نویسی برای ناوبری و برداشت اطلاعات نیاز داریم که قادر باشد به داده های وب عمیق نیز دست یابد همان گونه که به یوآرلهای معمولی دسترسی پیدا میکند.

ما از [oxPath9] استفاده میکنیم که بسط یافته XPath است و دارای امکاناتی برای تعامل با برنامه های کاربردی وب و برداشت دادههای مرتبط است. این زبان شبیه سازی عملکردهای کاربر را در تعامل با رابطهای کاربری چند صفحه ای پردازش نویسی شده(6) برنامه های کاربردی وب (ارزیاب(7) یا با مرورگر موزیلا- و یا با مرورگر Webkit- کار میکند) را ممکن می سازد، oxPath خصوصیتی مشابه XPath دارد و استفاده از گزینشگرهای مبتنی بر الگوهای آبخاری(8) با آن میسر است و میتوان با آن در چندین صفحه مختلف با کلیک ناوبری کرد و حتی از صفحه ها قبلی اطلاعات برداشت پیاده سازی شده

ص: 153

Web-graph-based-features -1

Script -2

3- مثال برای ارائه ساده شده است. ولی در واقع ما با چندین چینش (layout) سروکار داریم که vBulletin می تواند تعریف کند.

Timestamp -4

Declarative language -5

Scripted multipage interface -6

Evaluator -7

CSS: Cascading style sheet -8

این برنامه با کد منبع باز در دسترس است که در سامانه ما نیز مورد استفاده قرار گرفته است. نمونه ای از یک عملکرد ساده XPath که میتوان آن را بر vBulletin اجرا کرد عبارت است از: `a.forum/@href// click/` که به هر پیوند به تالار گفت وگویی وب کلیک میکند برای نمونههای بهتر از متغیرهای XPath به [9] نگاه کنید.

همچنین چندین جایگزین برای XPath را بررسی خواهیم کرد؛ به ویژه مواردی که در شماره های 17، 20، 21، 22 و 23 منابع این مقاله ذکر شده اند. این روشها به جز مورد مطرح در شماره [24] با تعامل با وب، پنهان نظیر ارسال فرم و ناوبری در صفحه ها کاری ندارند. در [22] از Datalog به عنوان زبان برنامه نویسی برای برداشت داده از صفحه ها وب استفاده شده است. در این روش زبان Xlog به عنوان برنامه کاربردی Datalog که شامل محمولات از پیش تعریف شده ای در برداشت داده است معرفی می شود. این روش پنجره ای برای محققان در استفاده از Datalog به عنوان پایه ای برای فرآیند برداشت می، گشاید ولی هنوز تلاشهای زیادی تا وب پنهان در آن لازم است [21] نیز با همین چالش روبه روست در این روش محتوا از یک صفحه ساده یا از صفحه های کتابشناختی برداشت می شود نیز در این روش هیچ عملی برای پر کردن فرم یا ناوبری در صفحه ها شبیه سازی نمی شود و در آن از یک زبان با نام Wraplet استفاده شده است که دادههای ساختارمند را از صفحه ها وب در قالب HTML برداشت می کند این زبان با عبارتهای پردازنده نویسی Wraplet عبارتهای برداشت (داده نوشته شده است و یک سند HTML را به عنوان ورودی گرفته خروجی را در XML میسازد. متأسفانه این روش فقط برای صفحه ها تک قابل استفاده است که دارای ویژگی پویایی وب نیستند روشی که از حیث امکانات و کارکرد، نظیر ارسال فرم و ناوبری بسیار به XPath شبیه است در [24] معرفی شده است. در نوشته اخیر نمونه های مختلفی برای روشن شدن مفاهیم ارائه شده اند. نویسندگان در آن سامانه ای را معرفی کرده اند و آن را سامانه برداشت داده مرور-گرا (1) نامیده اند که اطلاعات را از صفحه های وب برداشت و از پیوندها برای ناوبری به صفحه بعدی برای برداشت اطلاعات استفاده می کند سامانه مذکور برنامه های کاربرد را حتی هنگامی که در حال انجام کارکردهای پردازنده نویسی مانند جاوا یا AJAX هستند گردآوری میکند تا به محتوای صفحه بعدی دست پیدا کند این سامانه همچنین عملکردهای کاربران را برای تعامل با وب پنهان شبیه سازی می کند علیرغم همه این مزایا این سامانه محدودیتهای ناشی از مدیریت حافظه را مد نظر نمی گیرد و به همین علت برای سامانه مورد نیاز ما که لازمه اش خزش و آرشیو کردن مداوم مقادیر عظیم داده است مناسب نیست در این سامانه برای ناوبری متعدد صفحه ها کار مرورگر تکرار میشود تا عملکرد بهینه شود. در سامانه ای که ما تدارک دیده ایم با برداشت داده در مقیاس وسیع سروکار داریم. همین علت به سامانه ای نیاز داریم که مراقب عملکرد و مدیریت حافظه باشد XPath از حیث حافظه زمان به خوبی از عهده این حجم کار بر می آید.

ص: 154

از آنجا که این رساله دکتری تنها چند ماه پیش آغاز شده است، ما بر توسعه یک معماری متمرکز بودیم که روش شناسی آن در قسمتهای قبل ذکر شد. با این پیش فرض که این کار شدنی است و با بازشناسی آن در دو برنامه کاربردی وب، یک پیش نمونه اولیه از راهنمای آگاه از برنامه کاربردی وب پیاده سازی شد. بدین ترتیب که مثلاً پیش نمونه در مواجهه با برنامه کاربردی vBulletin و وبگاههای متعددی که از این سامانه مدیریت محتوا استفاده میکردند مورد ارزیابی قرار گرفت. در حال حاضر سامانه قادر است نوع برنامه کاربردی وب و سطح [صفحه] در برنامه کاربردی را تشخیص دهد و عملکردهای مناسب خزش را به اجرا بگذارد. به تازگی سامانه 8 [Yfilter] (سامانه فیلترینگ مبنی بر INFA1) را برای نمایه سازی کارآمد الگوهای اکتشاف و به منظور یافتن برنامه های کاربردی وب، مرتبط به سامانه ملحق کرده ایم هم اکنون، سامانه قادر است صفحه ها مرتبط بیشتری را با صرف حداقل، توان در مقایسه با روشهای قبلی خزش، برداشت کند هنوز لازم است بتوانیم عملکردهای خزش را با استفاده از یک ارزیاب XPath به منظور تبدیل آنها یوآرال یا عبارتهای XPath، در صورت امکان به اجرا درآوریم و راهنمای سامانه را رابط خزشگر قرار دهیم ولی نتایج اولیه برای پژوهشهای آینده راضی کننده و امیدبخش بوده اند.

5. پژوهشهای آینده

چالشهای جالبی در این حوزه وجود دارد که پژوهشهای بیشتری را طلب میکند که در ادامه این رساله دکتری در دستور کار ما خواهد بود:

1. استفاده از عبارتهای XPath 100 برای کشف الگوها با برخی محدودیتها در تبیین (2) مواجهه است:

مثلاً در بعضی موارد ممکن است عبارتهای معمول برای شناسایی یک برنامه کاربردی وب مورد نیاز باشد میتوانیم به عبارت های 20 XPath برگردانیم و از آنها استفاده کنیم یا کارکردهای بسط به آنها برای این منظور بیفزاییم ولی باید با اهداف بهینه سازی باید بکوشیم زبان برنامه نویسی را حتی الامکان اعلانی نگه داریم.

2. از چالشهای مهم تحقیق در امکان خودکار کردن بدون نظارت آموزش برنامه های کاربردی وب جدید ب) واسط الگوهای همگانی و انطباق با تغییرات اندک در قالبهایی است که پوششها را غیر قابل استفاده میکند.

3. همچنین به طور قطع باید در سراسر این کار تلفیق دقیق خزشگرها را در سامانه با تهیه و توسعه مکانیزمی برای تعامل با سایر اجزا صورت گیرد از جمله چالشهای این حوزه اینکه که به علت لزوم

ص: 155

1- Nondeterministic Finite در تئوری محاسبات اتوماتون تعیین ناپذیر متناهی یا اتوماتون غیر قابل تعیین متناهی یا اتوماتون پشته ای یا آن اف ای به او توماتونهایی گفته میشود که در مورد آنها برای برخی از دوتاییهای حالت و سمبل ورودی امکان عبور به بیشتر از یک حالت جدید اجازه داده شده باشد. (NFA Automaton)

2- Expressiveness

پایبندی به مسئله «اخلاقیات» (1) هنوز خزشگر همه تعاملات وبی را برعهده دارد، آن هم در جایی بعضی عملکردهای خزش ممکن است مستلزم ورود به یک برنامه خارجی (2) (یک خزشگر که مثلاً API یا یک ارزیاب oxPath) باشد.

بدیهی است که چالش دیگر ذکر مقیاسهایی است که با آنها عملکرده سامانه ما، هم از حیث کارآمدی و هم از حیث اثربخشی، با توجه به رویکردهای کلاسیک خزش مورد ارزیابی قاعده مند قرار می گیرند.

قدردانی

این پژوهش با پشتیبانی مالی هفتمین برنامه اتحادیه اروپایی (3) تحت تفاهم نامه مالی شماره 270239 (ARCOMEM) انجام شده است. نگارندگان همچنین مراتب قدردانی خود را به ژولین ماسانه (از مؤسسه حافظه اینترنت) برای ارائه مباحثی در موضوع این پایان نامه ابلاغ میکنند.

منابع

1. E. Amitay, D. Carmel, A. Darlow, R. Lempel, and A. Soffer. The connectivity sonar.

.Detecting site functionality by structural patterns. In HT, 2003

2. ARCOMEM Project. <http://www.arcomem.eu/>, 2011–2014

3. L. Barbosa and J. Freire. Searching for hidden-Web databases. In WebDB, 2005

4. L. Barbosa and J. Freire. An adaptive crawler for locating hidden-web entry points. In

.WWW, 2007

5. B. E. Boser, I. Guyon, and V. Vapnik. A training algorithm for optimal margin classifiers

.In COLT, 1992

6. R. Cai, J.-M. Yang, W. Lai, Y. Wang, and L. Zhang. iRobot: An intelligent crawler for

.Web forums. In WWW, 2008

7. S. Chakrabarti, M. van den Berg, and B. Dom. Focused crawling: A new approach to

.topic-specific Web resource discovery. Computer Networks, 31(11-16), 1999

8. Y. Diao, M. ALTINEL, M. J. Franklin, H. Zhang, and P. Fischer. Path sharing and

.predicate evaluation for high-performance XML filtering. ACM TODS, 2003

T. Furche, G. Gottlob, G. Grasso, C. Schallhart, and A. J. Sellers. OXPath: A language for .9

.scalable, memory-efficient data extraction from web applications. PVLDB, 4(11), 2011

.D. Gibson, K. Punera, and A. Tomkins. The volume and evolution of web page templates .10

ص: 156

Politeness -1

External programm -2

.(European Union Seventh Framework Programm (FP7/2007-2013 -3

- In WWW, 2005. [11] Y. Guo, K. Li, kai Zhang, and G. Zhang. Board forum crawling: A .Web crawling method for Web forums. In WIC, 2006
- .M. Hadley. Web application description language. <http://www.w3.org/Submission/wadl> .11
- /International Business Times. <http://www.ibtimes.com/articles/175488> .12
- .obama–twitter–townhall.htm, 2011/20110706 .13
- P. Kolari, T. Finin, and A. Joshi. Svms for the Blogosphere: Blog Identification and Splog .14
- .Detection. In AAI, 2006
- C. Lindemann and L. Littig. Coarse–grained classification of Web sites by their structural .15
- .properties. In CIKM, 2006
- .C. Lindemann and L. Littig. Classifying Web sites. In WWW, 2007 .16
- .M. Liu and T. W. Ling. A rule–based query language for HTML. In DASFAA, 2001 .17
- .J. Masanè s. Web archiving. Springer, 2006 .18
- E. Osuna, R. Freund, and F. Girosi. An improved training algorithm for support vector .19
- .machines. In Workshop on Neural Networks for Signal Processing, 1997
- A. Sahuguet and F. Azavant. Building light–weight wrappers for legacy Web data–sources .20
- .using W4F. In VLDB, 1999
- N. Sawa, A. Morishima, S. Sugimoto, and H. Kitagawa. Wraplet: Wrapping your Web .21
- .contents with a lightweight language. In SITIS, 2007
- W. Shen, A. Doan, J. F. Naughton, and R. Ramakrishnan. Declarative information .22
- .extraction using Datalog with embedded extraction predicates. In VLDB, 2007
- .K. Sigurðsson. Incremental crawling with Heritrix. In IAWW, 2005 .23

J.-Y. Su, D.-J. Sun, I.-C. Wu, and L.-P. Chen. On design of browser-oriented data .24

.extraction system and plug-ins. In JMST, 2010

S. Zheng, R. Song, J.-R. Wen, and D. Wu. Joint optimization of wrapper generation and .25

.template detection. In SIGKDD, 2007

ص: 157

امروزه دسته بندی صفحات وب برای بازیابی و مدیریت اطلاعات وب، ساختار، حفاظت یا گسترش فایل‌های وب (توالی وب)، بهبود کیفیت نتایج تحقیق، و کاهش زمان جست و جو استفاده در موتورهای جست و جوی عمودی و تمرکز خاص ناحیه ای در این موتورها و فیلتر کردن محتوای وب با استفاده از مرور کردن وب کاربرد و ضرورت فراوانی دارد ما برای بهبود دسته بندی مفهومی صفحات وب از رویکرد ارزش واژه در روشی مبتنی بر اتوماتای یادگیر توزیع شده استفاده می.کنیم در روش پیشنهادی به هر صفحه یک اتوماتای یادگیر تخصیص داده می شود که وظیفه آن استفاده از وزن کلمات کلیدی صفحات به منظور یادگیری میزان ارتباط آن صفحه با سایر صفحات وب دیگر است که این امر موجب دسته بندی سلسله مراتبی صفحات میشود برای ارزیابی روش پیشنهادی آن را بر روی داده های مختلفی آزمایش کردیم که نتایج خوبی حاصل شد و همچنین الگوریتم پیشنهادی از سرعت خوبی برخوردار بود.

کلیدواژه: اتوماتای یادگیر اتوماتای یادگیر توزیع شده دسته بندی صفحات وب وزن واژه.

روح اله گودرزی (1)، مصطفی پیرهادی (2)

1 - مقدمه

با افزایش روزافزون تعداد صفحات وب، دستیابی به صفحه‌های مورد نیاز و همچنین تفسیر آنها به عنوان یک چالش فراروی بازیابی اطلاعات و داده کاوی مورد توجه قرار گرفته است. بنابراین دسته بندی کردن صفحات وب میتواند نقش مهمی در افزایش جست و جو تفکیک خلاصه سازی و تفسیر وب داشته باشد. دسته بندی صفحات وب نوع نظارت شده ای از مسئله آموزشی است که به منظور دسته بندی این صفحات به مجموعه ای از دسته‌های از پیش تعریف شده به کار میرود که بر اساس داده های آموزشی برچسب دار می باشند. دسته بندی وظایف شامل اختصاص یافتن اسناد بر اصول موضوع، عملکرد، نظر، نوع و غیره میباشد بر خلاف بسیاری از دسته بندیهای متنی کلی، روش های دسته بندی صفحات وب دارای مزیت محتوایی یکسان ساختاری و ارتباطی با دیگر صفحات در وب است. در سالهای اخیر کارهایی در زمینه دسته بندی صفحات وب گزارش شده است که برخی از آنها عبارتند از: استفاده از نزدیکترین درخت مجاورت (K-NN) (2003.Kwon et al) استفاده از شبکه های عصبی

ص: 159

1- کارشناسی ارشد مهندسی کامپیوتر مدرس گروه کامپیوتر آموزشکده سما واحد بروجرد ruhollahgudarzi@yahoo.com

2- کارشناسی ارشد مهندسی کامپیوتر، کارشناس سازمان فناوری اطلاعات شهرداری بروجرد mpirhadi89@yahoo.com

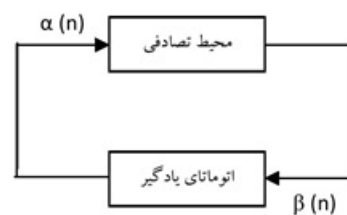
(Anagnos et al. 2004) (Selamat et al. 2004) ، استفاده از تنوری‌های ناهنجار و نامعلوم و سیستم های فازی برای کاهش داده های زائد (Jensen et al. 2004) ، استفاده از مدل SVM دستگاه حفاظتی بردار (Wakaki et al., 2006) (Chen et al. 2006) ، استفاده از وزن کلمات با الگوریتم TFIDF در صفحات وابسته به یک دیکشنری (Liang et al. 2006) (Ulmer et al. 2010) ، استفاده از درخت های تعیین کننده مبتنی بر فاصله (Estruh et al. 2006) و استفاده از الگوریتم ژنتیک (Ozel et al. 2010).

در این مقاله روشی مبتنی بر اتوماتای یادگیر توزیع شده که برای دسته بندی صفحات وب از وزن کلمات کلیدی که بر مبنای الگوریتم TFIDF محاسبه می شوند، پیشنهاد می گردد. در این روش به هر صفحه وب یک اتوماتای یادگیر اختصاص داده می شود که وظیفه اش یادگیری ارتباط آن صفحه با دیگر صفحات و در نتیجه تعیین دسته مطلوب آن صفحه می باشد که اتوماتا برای این کار باید از وزن کلمات کلیدی صفحه که با استفاده از الگوریتم TFIDF به دست می آید استفاده نماید.

روش پیشنهادی ضمن داشتن کارایی مناسب و صحت مدل ، پارامترهای یادگیری در اتوماتای یادگیر توزیع شده را با توجه به تعداد اسناد وب و با توجه به تعداد کلمات کلیدی به صورت پویا تنظیم می کند و اجرای آن بر روی صفحات وب آزمایشی نتایج خوبی را نشان می دهد. ادامه مقاله به دین صورت سازماندهی شده است : در بخش ۲ اتوماتای یادگیر و اتوماتای یادگیر توزیع شده به اختصار بیان می شود. در بخش ۳ روش TFIDF برای محاسبه شباهت سند و وزن یابی واژه مورد بحث قرار می گیرد. در بخش ۴ روش پیشنهادی ارائه خواهد شد و در بخش ۵ نتایج حاصل از ارزیابی روش پیشنهادی بیان می شود و بخش نهایی هم نتیجه گیری است.

۲- اتوماتای یادگیر (Narendra et al. 1989 , Lakshmirarahan et al. 1981)

اتوماتای یادگیر یک مدل انتزاعی است که تعداد محدودی عمل را می تواند انجام دهد. هر عمل انتخاب شده توسط محیطی احتمالی ارزیابی شده و پاسخی به اتوماتای یادگیر داده می شود. اتوماتای یادگیر از این پاسخ استفاده نموده و عمل خود را برای مرحله بعد انتخاب می کند. شکل ۱ ارتباط بین اتوماتای یادگیر و محیط را نشان می دهد.



شکل ۱: ارتباط بین اتوماتای یادگیر و محیط

(Anagnos et al. 2004) (Selamat et al. 2000) ، استفاده از تنوریهای ناهنجار و نامعلوم و سیستم های فازی برای کاهش دادههای زائد (Jensen et al. (Wakaki et al., 2006) (Chen et al. 2006) ، استفاده از مدل SVM دستگاه حفاظتی بردار (Wakaki et al., 2006) (Chen et al. 2006) ، استفاده از وزن کلمات با الگوریتم TFIDF در صفحات وابسته به یک دیکشنری (Liang et al. (Ulmer et al. 2010) ، استفاده از درختهای تعیین کننده مبتنی بر فاصله (Estruh et al. 2006) و استفاده از الگوریتم ژنتیک (Ozel et al. 2010).

در این مقاله روشی مبتنی بر اتوماتای یادگیر توزیع شده که برای دسته بندی صفحات وب از وزن کلمات کلیدی که بر مبنای الگوریتم TFIDF محاسبه میشوند پیشنهاد میگردد. در این روش به هر صفحه وب یک اتوماتای یادگیر اختصاص داده میشود که وظیفه اش یادگیری ارتباط آن صفحه با دیگر صفحات و در نتیجه تعیین دسته مطلوب آن صفحه میباشد که اتوماتا برای این کار باید از وزن کلمات کلیدی صفحه که با استفاده از الگوریتم TFIDF به دست می آید استفاده نماید.

روش پیشنهادی ضمن داشتن کارایی مناسب و صحت مدل پارامترهای یادگیری در اتوماتای یادگیر توزیع شده را با توجه به تعداد اسناد وب و با توجه به تعداد کلمات کلیدی به صورت پویا تنظیم میکند و اجرای آن بر روی صفحات وب آزمایشی نتایج خوبی را نشان میدهد. ادامه مقاله به دین صورت سازماندهی شده است: در بخش 2 اتوماتای یادگیر و اتوماتای یادگیر توزیع شده به اختصار بیان میشود. در بخش 3 روش TFIDF برای محاسبه شباهت سند و وزن یابی واژه مورد بحث قرار میگیرد در بخش 4 روش پیشنهادی ارائه خواهد شد و در بخش 5 نتایج حاصل از ارزیابی روش پیشنهادی بیان میشود و بخش نهایی هم نتیجه گیری است.

2- اتوماتای یادگیر (Narendra et al. 1989, Lakshmiarahan et al. 1981)

اتوماتای یادگیر یک مدل انتزاعی است که تعداد محدودی عمل را میتواند انجام دهد. هر عمل انتخاب شده توسط محیطی احتمالی ارزیابی شده و پاسخی به اتوماتای یادگیر داده میشود. اتوماتای یادگیر از این پاسخ استفاده نموده و عمل خود را برای مرحله بعد انتخاب میکند. شکل 1 ارتباط بین اتوماتای یادگیر و محیط را نشان می دهد.

شکل 1: ارتباط بین اتوماتای یادگیر و محیط

اتوماتای یادگیر با ساختار متغیر توسط 4 تایی app نشان داده میشود که در آن a, a, a, \dots مجموعه عملهای اتوماتا. $\{ \dots, \dots \}$ $B =$ مجموعه ورودیهای اتوماتا $p = \{ \dots, \dots, p \}$ بردار احتمال انتخاب هر یک از عملها و $(Tan(n(1))$ الگوریتم یادگیری میباشد. در این نوع از اتوماتاها اگر عمل در مرحله n ام انتخاب شود و پاسخ مطلوب از محیط دریافت نماید، احتمال (n) افزایش یافته و سایر احتمالها کاهش می یابند و برای پاسخ نامطلوب احتمال (n) کاهش یافته و سایر الگوریتم زیر یک نمونه از الگوریتمهای یادگیری خطی در اتوماتای با ساختار متغیر است.

الف - پاسخ مطلوب

ب - پاسخ نامطلوب

در روابط فوق، a پارامتر پاداش و پارامتر جریمه می باشد.

2-2- اتوماتای یادگیر توزیع شده

یک اتوماتای یادگیر توزیع شده شبکه ای از اتوماتاهای یادگیر است که برای حل یک مسأله خاص با یکدیگر همکاری دارند در این شبکه اتوماتاهای یادگیر همکار در هر زمان تنها یک اتوماتا فعال است تعداد اعمال قابل انجام توسط یک اتوماتا در DIA برابر با تعداد اتوماتاهایی است که به این اتوماتا متصل شدهاند. انتخاب یک عمل توسط اتوماتای یادگیر در این شبکه باعث فعال شدن اتوماتای یادگیر متصل شده به این اتوماتای یادگیر متناظر با این عمل گردد. به عبارت دیگر انتخاب یک عمل توسط یک اتوماتای یادگیر در این شبکه متناظر با فعال شدن یک اتوماتای یادگیر دیگر در این شبکه است. یک DIA توسط یک گراف که هر یک از رئوس آن یک اتوماتای یادگیر است، نشان داده میشود وجود یال $(LALA)$ در این گراف بدین معناست که انتخاب عمل توسط LA باعث فعال شدن LA می گردد. تعداد اعمال قابل انتخاب توسط LA بصورت $pk = \{ p, p, \dots, p \}$ نمایش داده شود. در این مجموعه عدد نشان دهنده احتمال مربوط به عمل $aman$ است انتخاب عمل توسط LA باعث فعال شدن LA میشود. r تعداد اعمال قابل انجام توسط اتوماتای LA را نشان میدهد. برای کسب اطلاعات بیشتر راجع به اتوماتاهای یادگیر توزیع شده و کاربردهای آن میتوان به مراجعه $(Narendra et al. 1989)$ ، $(Lakshmiarahan et al. 1981)$ نمود.

3- شباهت سند در روش TFIDF (Ulmer et al. 2010)

شباهت سند با استفاده از اهمیت اصطلاح، یک روش بازیابی اطلاعات است که به صورت زیر بیان

می شود:

۱. اصطلاحات مهم و معتبر در هر سند مانند گویا ترین اصطلاحات بیشترین اهمیت و اعتبار را دریافت کنند.

۲. هر سند توسط یک بردار اهمیت اصطلاح نشان داده شود.

۳. مقایسه ی اسناد با همدیگر با استفاده از یک مقیاس شباهت در فضای برداری اصطلاح .

یک طرح موثر و شناخته برای هدف، tfidf است. اعتبار tfidf شامل دو قسمت است: idf و tf. تکرار tf ارزیابی می کند که چند وقت یکبار اصطلاح خاص نسبت به همه ی اصطلاحات V در صفحه d رخ می دهد:

$$Tf(t \ll d) = \frac{\text{count}(t, d)}{\sum_{v \in D} \text{count}(v, d)} \quad (3)$$

تکرار سند معکوس (idf)، نسبت اسناد d را که در آنها وجود دارد در مجموعه D ارزیابی می کند.

$$idf(t) = \frac{\log |D|}{|[d \in D : t \in d]|} \quad (4)$$

Tfidf بیشترین اهمیت (وزن) اصطلاح t را که غالباً در سند d وجود دارد تعیین می کند این اعتبار در اسناد دیگر، کمتر وجود دارد.

$$Tfidf(t \ll d) = tf(t \ll d) \cdot idf(t) \quad (5)$$

در مدل فضای بردار، سند d با بردار vd متعلق به tfidf مشخص می شود هر قسمت i از vd، امتیاز نام اصطلاح در مجموعه سند tfidf را می گیرد شباهت بین اسناد d و a با استفاده از کسینوس زاویه θ بین بردارهای va و vb محاسبه می شود.

$$\text{sim}(d, a) = \cos(\theta_{v_d, v_a}) \quad (6)$$

یعنی

$$\text{sim}(d, a) = \frac{v_d \cdot v_a}{\|v_d\| \|v_a\|} \quad (7)$$

یا به طور هم ارز

$$\text{sim}(d, a) = \frac{\sum_{t \in D} tfidf(t, d) \cdot tfidf(t, a)}{\sqrt{\sum_{t \in D} tfidf(t, d)^2} \sqrt{\sum_{t \in D} tfidf(t, a)^2}} \quad (8)$$

۴- الگوریتم پیشنهادی

الگوریتم پیشنهادی برای دسته بندی صفحات وب از روشی مبتنی بر اتوماتای یادگیر توزیع شده و از وزن کلمات کلیدی استفاده می کند. اختلاف این الگوریتم با الگوریتم های پیشین مبتنی بر اتوماتای یادگیر توزیع شده (Anari et al. 2007)، (Baradaran et al. 2007) در این است که در الگوریتم های قبلی در حقیقت خوشه بندی صفحات انجام می گیرد و اتوماتا یاد می گیرد در بین صفحات وبی که در حال

1. اصطلاحات مهم و معتبر در هر سند مانند گویا ترین اصطلاحات بیشترین اهمیت و اعتبار را دریافت کنند.

2. هر سند توسط یک بردار اهمیت اصطلاح نشان داده شود.

3. مقایسه ی اسناد با همدیگر با استفاده از یک مقیاس شباهت در فضای برداری اصطلاح .

یک طرح موثر و شناخته برای ، هدف tfidf . است اعتبار tfidf شامل دو قسمت است tf و idf. تکرار

f ارزیابی میکند که چند وقت یکبار اصطلاح خاص نسبت به همه ی اصطلاحات 7 در صفحه d رخ می دهد:

تکرار سند معکوس (idf) نسبت اسناد d را که در آنها وجود دارد در مجموعه D ارزیابی می

کند.

Tidf بیشترین اهمیت (وزن) اصطلاح t را که غالباً در سند d وجود دارد تعیین می کند این اعتبار در اسناد دیگر کمتر وجود دارد. در مدل فضای بردار سند d با بردار vd متعلق به tfidf مشخص میشود هر قسمت از vd امتیاز نام اصطلاح در مجموعه سند tfidf را می گیرد شباهت بین اسناد d و a با استفاده از کسینوس زاویه 0 بین بردارهای va و gvba در tfidf محاسبه می شود.

یعنی

یا به طور هم ارز

4- الگوریتم پیشنهادی

الگوریتم پیشنهادی برای دسته بندی صفحات وب از روشی مبتنی بر اتوماتای یادگیر توزیع شده و از وزن کلمات کلیدی استفاده میکند اختلاف این الگوریتم با الگوریتمهای پیشین مبتنی بر اتوماتای یادگیر توزیع شده (Anari et al., (Baradaran et al., 2007) در این است که در الگوریتمهای قبلی در حقیقت خوشه بندی صفحات انجام میگردد و اتوماتا یاد میگیرد در بین صفحات وبی که در حال

ص: 162

حاضر موجود است کدام ورودی ها را با هم در یک دسته قرار دهد اما در این الگوریتم عمل دسته بندی بدین معنا انجام میگیرد که برای هر دسته از روی تشابه مفهومی صفحات آن دسته و با توجه به وزن کلمات کلیدی آن صفحات با استفاده از الگوریتم TFIDF ، یک تابع عضویت دسته تعریف می شود و هر صفحه جدیدی که وارد شود در این تابع تست می گردد تا عضویت آن در دسته بررسی شود که البته با توجه به تعداد کلمات انتخاب شده، میزان شباهت صفحات به یکدیگر نیز به صورت سلسله مراتبی قابل تعیین است. روند اجرای الگوریتم پیشنهادی بصورت شکل 2 می باشد.

1 اتوماتای یادگیر توزیع شده متناظر با ساختار اسناد وب ایجاد کن

2. بردار احتمالات اتوماتای یادگیر موجود در اتوماتای یادگیر توزیع شده را مقدار دهی اولیه کن 3. تکرار ft را برای همه اسناد مطابق فرمول 3 محاسبه کن

4. به ازای هر کاربر موجود در لاگ فایل انجام بده

4-1- به ازای هر حرکت به صورت nm در طول مسیر انجام بده

1-1-4- بردار احتمال اتوماتای متناظر با سند m را مطابق روابط زیر بروز کن

2-1-4- تکرار سند معکوس (fdi را طبق فرمول (4) برای مجموعه اسناد محاسبه کند.

ه به ازای صفحات فاقد کاربر در لاگ فایل مطابق فرمول (8) عمل کن و سپس برو به 4-1-1

شکل 2- الگوریتم پیشنهادی

5- شبیه سازی و ارزیابی الگوریتم پیشنهادی

برای شبیه سازی الگوریتم پیشنهادی از مدل معرفی شده در (Liu et al.2004) استفاده می شود که در آن Liu و همکارانش نظم موجود در رفتارهای کاربران در محیط وب را با استفاده از یک مدل مبتنی بر عامل مشخص و اعتبار مدل خود را با استفاده از اطلاعات استفاده از وب چندین سایت وب بزرگ مانند مایکروسافت تایید کرده اند. آنها به جای استفاده از صفحات وب واقعی و داده های واقعی کاربران وب از این مدل استفاده کرده اند این مدل محیطی شامل صفحات وب و کاربران آن را فراهم می کند. مزیت استفاده از این مدل آن است که تشخیص کاربران و بازدیدهای انجام شده از صفحات وب با استفاده از این مدل بسیار دقیق تر میباشد و به عملیات پالایش دادهها نیز احتیاجی نخواهد بود

البته پارامترهای معرفی شده در این مدل بایستی بدقت تنظیم گردند تا نتیجه حاصل از آن مشابه با محیط واقعی گردد. هر صفحه وب در این مدل دارای برداری است که هر مولفه آن میزان ارتباط این صفحه با موضوع متناظر با این مولفه را نشان می دهد (تعداد موضوعات ثابت و قابل تعریف است). میزان ارتباط هر صفحه با یک موضوع به صورت عددی بین صفر و یک بیان می شود به طوری که مجموع آنها برای همه موضوعات برابر با یک است. همچنین هر صفحه دارای پیوندهایی با صفحات دیگر است برای آزمایشها پروفایل علاقه کاربران به صورت توزیع توانی و توزیع محتوای اسناد به صورت توزیع نرمال در نظر گرفته شده است سایر پارامترهای استفاده شده در این مدل برای شبیه سازی انجام شده در این قسمت در جدول (1) نشان داده شده است.

جدول شماره 1

برای ارزیابی الگوریتمهای پیشنهادی تعیین ساختار اطلاعاتی اسناد وب از معیار کورولیشن (1) استفاده می شود. کورولیشن معیاری برای بدست آوردن وابستگی خطی بین دو بردار است و به صورت زیر تعریف میشود:

عکس

البته پارامترهای معرفی شده در این مدل بایستی بدقت تنظیم گردند تا نتیجه حاصل از آن مشابه با محیط واقعی گردد. هر صفحه وب در این مدل، دارای براداری است که هر مولفه آن میزان ارتباط این صفحه با موضوع متناظر با این مولفه را نشان می دهد (تعداد موضوعات ثابت و قابل تعریف است). میزان ارتباط هر صفحه با یک موضوع به صورت عددی بین صفر و یک بیان می شود به طوری که مجموع آنها برای همه موضوعات برابر با یک است. همچنین هر صفحه دارای پیوندهایی با صفحات دیگر است. برای آزمایشها پروفایل علاقه کاربران به صورت توزیع- توانی و توزیع محتوای اسناد به صورت توزیع نرمال در نظر گرفته شده است. سایر پارامترهای استفاده شده در این مدل برای شبیه سازی انجام شده در این قسمت در جدول (۱) نشان داده شده است.

جدول شماره ۱

| | |
|---|------|
| حد آستانه ایجاد اتصال | ۰/۷ |
| تعداد کاربران | ۲۰۰ |
| تعداد اسناد | ۵۰۰ |
| تعداد موضوع ها | ۲۵ |
| T_e مقدار ثابت سند اولیه (صفحه اولیه سایت) در موضوعات مختلف | ۰/۲ |
| a_u پارامتر توزیع قاتون - توانی توزیع احتمال علاقه کاربران | ۱ |
| a ضریب پاداش دریافتی از مشاهده یک سند | ۰/۹ |
| b ضریب جریمه دریافتی از پیمایش یک دور | ۰/۱ |
| λ ضریب جذب اطلاعات از یک سند توسط یک کاربر | ۰/۵ |
| μ_m میانگین توزیع نرمال $M\Delta^2$ | ۵/۹۷ |
| σ_m واریانس توزیع نرمال $M\Delta^2$ | ۰/۲۵ |
| a_p پارامتر توزیع قاتون - توانی توزیع احتمال وزنهای مطالب برای هر سند | ۳ |
| σ_p واریانس توزیع نرمال برای مقدار افزایش یک گره برای یک موضوع خاص | ۰/۲۵ |
| θ ضریب کاهش علاقه کاربر | ۱ |
| حداقل اشتیاق کاربر برای ادامه جستجو | ۰/۲ |

برای ارزیابی الگوریتمهای پیشنهادی تعیین ساختار اطلاعاتی اسناد وب از معیار کورولیشن^۱ استفاده می شود. کورولیشن معیاری برای بدست آوردن وابستگی خطی بین دو بردار است و به صورت زیر تعریف می شود:

1. Correlation

شکل (3) کارایی الگوریتم پیشنهادی را با معیار کورولیشن در مقایسه با الگوریتمهای ارائه شده (Anari et al.2007) نشان میدهد محور افقی تعداد کاربران و محور در (Baradaran et al.2007) عمودی میزان کورولیشن را نشان میدهد. در ارزیابی فوق تعداد صفحات 20, تعداد موضوعات 5 و تعداد کاربران 20000 در نظر گرفته شده است.

شکل (3) - مقایسه روش پیشنهادی

عکس

دسته‌بندی مفهومی صفحات وب... ۱۶۵

$$Corr(p, p') = \frac{\sum p p' - (\sum p \sum p')/n}{\sqrt{(\sum p^2 - (\sum p)^2/n)(\sum p'^2 - (\sum p')^2/n)}}$$

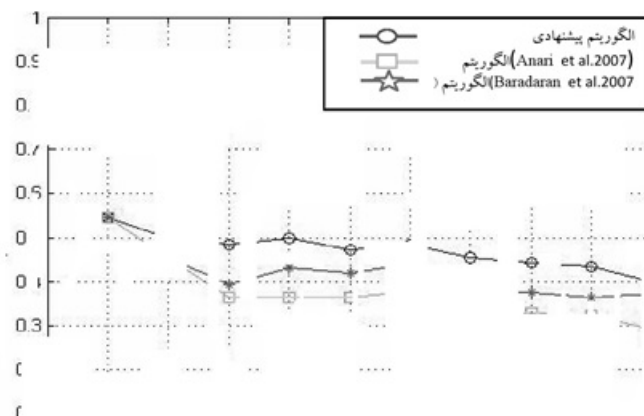
$$p = \{p_{ij} \mid i, j = 1, 2, \dots, n, i \neq j\}$$

$$P_{ij} = \frac{d_{ij}^{-1}}{\sum_{k=1}^n d_{ik}^{-1}} \quad (9)$$

$$d_{ij} = \sqrt{\sum_{k=1}^m (cw_i^k - cw_j^k)} \sqrt{\sum_{k=1}^m (cw_i^k - cw_j^k)}$$

$$p' = \{p'_{ij} \mid i, j = 1, 2, 3, \dots, n, i \neq j\}$$

شکل (3) کارایی الگوریتم پیشنهادی را با معیار کورولیشن در مقایسه با الگوریتمهای ارائه شده در (Baradaran et al.2007) و (Anari et al.2007) نشان می‌دهد. محور افقی تعداد کاربران و محور عمودی میزان کورولیشن را نشان می‌دهد. در ارزیابی فوق تعداد صفحات 20, تعداد موضوعات 5 و تعداد کاربران 20000 در نظر گرفته شده است.



شکل (3) - مقایسه روش پیشنهادی

با توجه به اهمیت دسته بندی صفحات وب در این مقاله روشی مبتنی بر اتوماتای یادگیر توزیع شده و ارزش کلمات پیشنهاد گردید و کارایی آن با توجه به معیار کورولیشن و در نتیجه وابستگی خطی در مقایسه با الگوریتمهای ارائه شده در (Baradaran et al.2007) و (Anari et al.2007) مورد ارزیابی قرار گرفت که الگوریتم پیشنهادی از کورولیشن بالاتری نسبت به این دو روش برخوردار است. ایده اصلی این الگوریتم این بود که صفحاتی که با همدیگر ارتباط مفهومی بیشتری دارند، پاداش بیشتری را دریافت مینمایند. از نتایج بدست آمده از این الگوریتم میتوان به عنوان ابزاری برای یافتن صفحات مشابه و مرتبط با یک موضوع خاص استفاده نمود.

منابع

- Kwon,O.w., J.H. Lee. (2003). Text categorization based on k-nearest neighbor approach for .Web site classification. *Information Processing and Management*, 39:25-44
- Anagnostopoulos, I., C. Anagnostopoulos, V. Loumos, and E. Kayafas. (2004). Classifying .(Web pages employing a probabilistic neural network. *IEEE Proc.Softw*, 151(3
- Selamat, A., S. Omatu.(2004).Web page feature selection and classification using neural .networks. *Information Sciences*, 158:69-88
- Jensen, R., Q. Shen.(2004). Fuzzy-rough attribute reduction with application to web .categorization. *Fuzzy Sets and System*, 141:469-485
- Chen, R.C., C.H. Hsieh. (2006). Web page classification based on a support vector machine .using a weighted vote schema.*Expert System with Applications*, 31:427-435
- Wakaki, T., H. Itkura, M. Tamura, H. Motoda, and T.Washio. (2006).A study on rough set aided feature selection for automatic web-page classification. *Web Intelligence and Agent .System*, 4: 431-441
- Liang, C. Y., L. Guo, Z.J. Xia, F.G. Nie, X.X. Li, L. Su, and Z. Y. Yang.(2006). Dictionary-based ,text categorization of chemical web pages. *Information Processing and Management*

Ulmer, C., M. Gokhale, B. Gallagher, Ph. Top, and T. Eliassi-rad. (2010). Massively parallel

.acceleration of a document -similarity classifier to detect web attacks. J.Parallel Distrib

Comput

Estruh, V., C. Ferri, J. Hernandez-Orallo, M.j. Ramirez-Quintana .(2006). Web categorization

,using distance-based decision trees. Electronic Notes in Theoretical Computer Science

.157:35-40

ص: 166

- Ozel,S.A. (2010). A Web page classification system based on a genetic algorithm using .tagged-terms as features. Expert Systems With Applications
- :Lakshmivarahan,S.(1981). Learning algorithms: Theory and applications. New York .Springer-Verlag
- Narendra,k.S., K.S. Thathachar (1989). Learning automata: An inter oduction. New York .Prentice Hall :
- Anari, B., M.R. Meybodi. (2007). A new method based on distributed learning automata for determining web documents structure.Proceedings of the 12 Annual international CSI Computer Conference, CSICC2007, Tehran, Iran, pp.2281
- Baradaran Hashemi, A., M.R. Meybodi. (2007).Web usage mining using distributed learning autoumata. Preceedings of the 12 Annual International CSI Computer .Conferenc, CSICC2007, Tehran, Iran, pp. 553-560
- Liu, J., S. Zhang, J. Yang. (2004), Characterizing web usage regulaities with information .(foraging agenst. IEEE Transaction in Kniwlodge and Data Engineering, 16(5
- Chao, M., H. Chen (2008). A machine learning approach to web page filtering using content .and structure analysis.Decision Support System, 44:482-494

ایجاد ابزارهای خودکار در سازمانهایی که از فرآیندهایی برخوردارند که به خوبی تعریف شده اند، بهبود فعالیتهای افزایش کارایی، کاهش خطاها، و بهبود سطح خدمات آن سازمان خواهد شد.

در سال 2010 تیم آرشیو وب(1) کتابخانه کنگره با توسعه و ایجاد ابزاری با نام DigiBoard به خود کارسازی جریانهای کاری در آرشیو وب پرداخت. DigiBoard نوعی برنامه کاربردی تحت وب است که به منظور پشتیبانی از جریان کاری آرشیو وب (نامزدی(2)، پردازش مجوزها(3)، و داوری(4)) در کتابخانه کنگره به کار میرود و به افزایش بهره وری تیمها و کاربران کتابخانه کنگره کمک می کند و کاهش میزان خطاهای معمول در روشهای دستی آرشیو وب را میسر می. سازد مقاله حاضر برخی کارکردهای ابزار جدید و سفارشی شده(5) تیم آرشیو وب کتابخانه کنگره را توصیف خواهد کرد.

ص: 168

Web Archiving Team -1

Nomination -2

Permission Processing -3

Reviewing -4

Custom-built -5

DigiBoard*: ابزار افزایش کارایی فعالیتهای پیچیده آرشیو وب در کتابخانه کنگره (1)

آبه گروتک (2) | جینا جونز (3) | ترجمه: سعیده اسلامی (4)

1 - مقدمه

کتابخانه کنگره از 10 سال پیش تا امروز به بایگانی کردن محتوای وب پرداخته است و در این سالها بیشتر از 14/000 وبگاه را در 40 مجموعه آرشیو وب موضوعی (5) و گزینشی (6) گنجانده و روالهای کاری، نامزدی پردازش مجوز تعیین میدان خزش داوری، کیفیت فهرست نویسی و دسترسی را یکپارچه کرده است. البته کتابخانه کنگره به منظور خزش وب حکم و قانون و اسپاری معینی ندارد؛ بنابراین برای بایگانی برخی وبگاهها مانند وب نوشتهها و سایتهای سازمانهای، خبری ملزم به دریافت مجوز است، اما وبگاههای دیگر قابل آرشیو شدن هستند و فقط از طرف کتابخانه کنگره مطلع میشوند که محتوای آنها آرشیو خواهد شد. جهت دسترسی به آرشیو کتابخانه کنگره، کتابخانه ملزم به اخذ مجوز از تمامی وبگاههاست وبگاههای دولتی از این قاعده جدا هستند. اگر چه در سالهای اخیر کتابخانه قابلیت

ص: 169

DigiBoard: A Tool to Streamline Complex Web Archiving Activities at the Library of Congress -1

Abbie Grotke -2

Gina Jones -3

4- کارشناس ارشد مهندسی نرم افزار سازمان اسناد و کتابخانه ملی ایران (s-eslami@nlai.ir)

Thematic -5

Selective -6

خزشگری در جا(1) را پیاده سازی کرده است بیشتر محتوا از طریق یک عامل خزشگر غیر در جا(2) حاصل می شود.

2- ضرورت ایجاد ابزار جدید

در سال 2003 کتابخانه کنگره ابزاری ابتدایی را برای تأمین نیازهای مجوزدهی مستندات ساخت که فقط از جریانهای کاری نامزدی و پردازش مجوزها پشتیبانی میکرد تقاضای ساخت چنین ابزاری از طرف (3) OGC کتابخانه کنگره و اداره حق نشر ایالات متحده (4) ارائه شد این ابزار در بازه زمانی یک ماهه ساخته شد اما با افزایش چندین برابری تعداد کاربران به قدر کافی توسعه پذیر نبود، علاوه بر آن، به موجب تغییر حقوق قانونی اخذ مجوزها مرحله پردازش مجوزها بیش از حد تغییر می کرد. در حال حاضر 5 کارمند تمام وقت، 3 پیمانکار و بیش از 80 کتابدار در واشنگتن دی سی و نقاط مختلف جهان با هدف ساده سازی مراحل نامزدی و پردازش مجوزها از این ابزار استفاده می کنند.

کتابخانه کنگره لازم میدانند از میزان مشارکت کانون وکلا در فعالیتهای آرشیو وب کتابخانه بکاهد؛ بدین ترتیب کاربران مختلف در خدمات کتابخانه (5)، کتابخانه قانون کتابخانه کنگره و خدمات تحقیقاتی کنگره ای (6) به حداقل آموزشها جهت انتخاب کردن وبگاهها نیاز خواهند داشت. علاوه بر این تیم آرشیو وب برای مدیریت بهتر جریانهای کاری اش به ابزاری نیاز دارد جریانهای کاری آرشیو وب کتابخانه و مراحل آن به تدریج توسعه پیدا کرد و در مقایسه با سال 2003 که ابزار اولیه ساخته شد، رسمیت بیشتری پیدا کرده است. بنابراین کتابخانه کنگره نیازمند ابزاری با ویژگیهای زیر است:

(1) داده های فعالیتهای مختلف را مدیریت کند نامزدی، مجوزها، خزش، بررسی کیفیت، گزارشگیری، و مانند آن؛ و

(2) به منظور کاهش زمان پردازش URL برای نامزد کننده و بگاهها و تیم فعالیتهای دستی را خودکار سازی کند.

3- ایجاد و توسعه DigiBoard

DigiBoard، با استفاده از زبان برنامه نویسی PHP و بانک اطلاعاتی MySQL در چند مرحله توسعه یافته است. نخستین ماژول آن برای نامزد کننده وبگاهها در سپتامبر 2009 راه اندازی شد یک پیمانکار مسئولیت توسعه ابزار اولیه را بر اساس نیازمندیهای مشخص شده توسط 5 عضو تیم آرشیو وب عهده دار شد. کتابداران نیازمندیهای دیگری را تدارک دیدند تا بهسازی مراحل کاری خاصشان را فراهم کند.

مدل

ص: 170

On-site -1

Off-site -2

Office of General Counsel -3

United state copyright office -4

Library services -5

Congressional Research Service -6

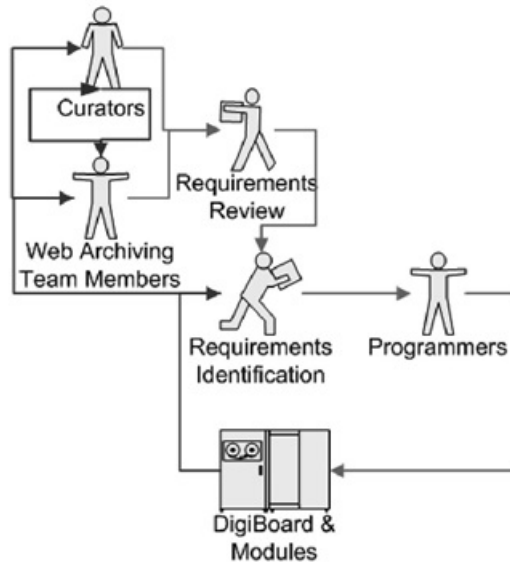
توسعه DigiBoard به صورت چابک(1) و فرآیند توسعه تکراری(2) (شکل 1) میباشد که نیازمندیها و مسائل شناسایی شده توسط کتابداران و اعضای تیم را پوشش میدهد نیازمندیهای هر مرحله از فرآیند تیم آرشیو وب به صورت ترتیبی و افزایشی به سیستم اضافه شده است؛ بنابراین کارآیی مورد اضافه شده توسعه سیستم را تحت تأثیر قرار داده است.

شکل 1. مراحل توسعه دیجی بورد فرآیند چابک و تکراری توسعه

قابلیتهای ناشی از توسعه موفق آمیز DigiBoard در کتابخانه کنگره عبارت اند از لایق بودن تیم توسعه، دارا بودن دانش کافی درباره واسطهای کاربری اشراف بر تواناییهای پایگاه دادههای وب پویا، توانایی انتقال اطلاعات به توسعه دهندگان سیستم، تجربه بالای افراد تیم و علاقه مند به فعالیتهای خودکارسازی فرآیندها.

عکس

توسعه DigiBoard به صورت چابک^۱ و فرآیند توسعه تکراری^۲ (شکل ۱) می باشد که نیازمندی ها و مسائل شناسایی شده توسط کتابداران و اعضای تیم را پوشش می دهد. نیازمندی های هر مرحله از فرآیند تیم آرشیو وب به صورت ترتیبی و افزایشی به سیستم اضافه شده است؛ بنابراین، کارآیی مورد اضافه شده توسعه سیستم را تحت تأثیر قرار داده است.



شکل ۱. مراحل توسعه دیجی بورد: فرآیند چابک و تکراری توسعه

قابلیت های ناشی از توسعه موفق آمیز DigiBoard در کتابخانه کنگره عبارتند: از لایق بودن تیم توسعه، دارا بودن دانش کافی درباره واسط های کاربری، اشراف بر توانایی های پایگاه داده های وب پویا، توانایی انتقال اطلاعات به توسعه دهندگان سیستم، تجربه بالای افراد تیم و علاقه مند به فعالیت های خودکار سازی فرآیندها.

۴- بهبود فرآیند نامزدی و پردازش مجوزها

۴-۱- فرآیند نامزدی

آرشیوهای وب کتابخانه کنگره بر مبنای مفاهیم مجموعه ها ساخته می شود. URL ها برای یک مجموعه و یا چندین مجموعه نامزد می شوند و حد کمینه فراداده برای نامزد کردن یک URL لازم است. قبل از

1. Agile
2. Iterative process

4- بهبود فرآیند نامزدی و پردازش مجوزها

4-1- فرآیند نامزدی

آرشیوهای وب کتابخانه کنگره بر مبنای مفاهیم مجموعه ها ساخته می شود URL ها برای یک مجموعه و یا چندین مجموعه نامزد می شوند و حد کمینه فراداده برای نامزد کردن یک URL لازم است قبل از

Agile -1

Iterative process -2

ارائه این ابزار لازم بود تا نامزد کننده ها دو فرم جداگانه خارج از ابزار را کامل کنند و سپس برای ادامه نامزدی وارد سیستم شوند اما اکنون کلیه مراحل یکپارچه و ساده تر شده است فرمی که توسط کتابداران کتابخانه کنگره استفاده میشود (به این فرم با نام نامزد کننده در این مقاله اشاره شده است) ساده تر بوده و اطلاعات بیشتری مانند موضوع زبان و دیگر اطلاعات درباره URLها را جمع آوری میکند. این فرم اجازه میدهد تا کارهای نیمه تمام کاربران ذخیره شده و در زمان دیگری ادامه داده شوند و امکان کپی رکوردها از یک مجموعه به مجموعه دیگر را فراهم می آورد کلیه دادهها درباره URL نامزد شده در بانک اطلاعاتی ذخیره شده و در هر زمان به راحتی قابل مشاهده و قابل جست و جوست. همچنین، نامزدکننده ها میتوانند پلاگینهایی را به منظور نامزد کردن وبگاهها برای شمول در یک مجموعه بارگذاری و استفاده کنند (شکل 2 را ببینید).

شکل 2. پلاگینهای موجود برای مرحله نامزدی در دیجی بورد

2-4- پردازش مجوزها

از آنجا که به دست آوردن مجوز از مالکان وبگاهها برای آرشیو در کتابخانه کنگره ضروری است فرآیند مجوز درون فرآیند نامزدی ساخته شده است از سال 2003، الزاماً باید تمام وبگاههای آمریکا به غیر از وبگاههای دولتی از قصد کتابخانه کنگره جهت آرشیو محتوای وبگاه شان مطلع شوند. علاوه بر این، OGC کتابخانه کنگره صدور مجوز از طرف مالکان محتوا جهت گردآوری اطلاعات را برای برخی وبگاهها ضروری دانسته است ابزار قدیمی از طریق پست الکترونیکی مجوزها را بر اساس نوع مجوز مجوز گردآوری (مجوز نمایش محتوای غیر درجا) به مالکان محتوا ارسال می کرد در سال 2006، OGC امکان مجوزهای روکشی (1) را فراهم آورد بدین معنی که اگر مالک محتوا مجوز گردآوری یا نمایش را اعطا کند، این مجوز برای مجموعه های آینده هم قابل استفاده مجدد خواهد بود در سالهای اخیر کتابخانه کنگره به ازای هر URL اضافه شده به مجموعه حتی اگر قبلاً مجوز کسب کرده باشد، اطلاعیه ای به مالکان محتوا ارسال میکند.

مدیریت مجوزهای روکشی در ابزار قدیمی به آسانی ممکن نبود و نامزد کننده ها ملزم بودند تا یک

ص: 172

فایل خارجی حاوی فهرست وبگاههایی دارای مجوز کنکاش کنند در برخی موارد، مجوزها مجدداً درخواست میشوند در نتیجه حدود 15 تا 30 دقیقه به ازای هر وبگاه زمان از دست می‌رفت. اگرچه ماژول مجوز DigiBoard امکان ردیابی بهتر و قابلیت استفاده از مجوزهای روکشی را فراهم آورده است به دلیل پیچیدگی وب و نیاز به آموزشهای پیشرفته برای کاربران جهت اعمال مجوزهای روکشی به طور مناسب به نامزدهای جدید این ماژول چالش برانگیزترین ماژول در مرحله پیاده‌سازی بود. به عنوان مثال URLهای مختلف از یک دامنه وبگاه برای مجموعه‌های مختلفی نامزد میشود و لازم است نامزدکننده‌ها آموزش ببینند که گاهی اوقات زیر دامنه دامنه‌ها، دایرکتوریها و فایلها به یک مالک محتوا تعلق دارد و گاهی بنا به ماهیتشان به همان مالک محتوا تعلق ندارد برای مثال، یک وب‌نوشت درباره قانون با میزبانی دانشگاه کلمبیا برای شمول در آرشیو وب دیوان عالی کشور (1) نامزد شده بود و مالک، سایت مجوز آن را برای آرشیو به کتابخانه کنگره اعطا کرده بود وبگاه معماری دانشگاه کلمبیا برای یک آرشیو وب دیگر نامزد شده بود و اطلاعاتی در DigiBoard مبنی بر اینکه مالک وب‌نوشت مجوز اعطا کرده است به کتابدار ارائه شد و آنها مجوز روکشی را به وبگاه کتابخانه اعمال کردند. به دلیل اینکه مالکان سایت به صورت دو موجودیت مختلف در نظر گرفته میشوند آموزشهایی جهت وضوح اینکه بخش‌های مختلف دانشگاه کلمبیا ممکن است چندین مالک داشته باشند لازم است.

3-4- بهبود پاسخ مالک محتوا

به علت نحوه طراحی DigiBoard از زمان پیاده‌سازی آن افزایش قابل توجهی در میزان پاسخهای مالکان محتوا به وجود آمده است. در حالی که ابزار قبلی امکان ارسال پست الکترونیکی و پاسخ دادن از طریق یک فرم وب را برای مالکان محتوا فراهم آورده بود و گاه پست الکترونیکی به اشتباه به صورت اسپم (2) در نظر گرفته میشد علاوه بر این فرم و بی‌پاسخها با سایت آرشیو وب یکپارچه نبود. با ابزار جدید ایمیل‌های تولید شده اجازه میدهند تا مالک محتوا نیز مستقیماً با کلیک بر یک پیوند ساده و یا با رجوع به صفحه وبگاه و با وارد کردن یک شماره یکتا - که امکان پیگیری پاسخها را داشته باشد - پست الکترونیکی خود را ارسال کند (شکل 3).

شکل 3. بازیابی پست الکترونیکی مجوز

ص: 173

supreme court -1

spam -2

میزان پاسخ مجوزها با پیاده سازی دیجی بورد به میزان 50 درصد افزایش پیدا کرده است.

4-4- خصوصیات جدید

خصوصیات جدیدی به ماژولهای نامزدی و ماژول مجوز DigiBoard اضافه شده است که جریانهای کاری و پردازش URLها را بهبود میبخشد. در حال حاضر:

- دیجی بورد روشی برای کپی کردن یک URL در یک مجموعه جدید نگهداری تاریخچه تغییرات و ثبت عملیات انجام شده امکان ثبت تغییرات URL، ها امکان فهرستنویسی یکپارچه و مدیریت مجوزها را فراهم کرده است.

- DigiBoard مکانیزمی برای شناسایی و پیاده سازی آسان خصوصیات مشترک وبگاهها فراهم کرده است که پیچیدگی انتخاب مجوزها از دامنه زیر دامنه دامنه های سطح بالا- و میزبانهای وب را کاهش داده است. به عنوان مثال بیش از 250 وبگاه بلاگ اسپات در ابزار DigiBoard نامزد شده اند که هر کدام متعلق به یک مالک محتوا مختلف است. این ابزار امکاناتی فراهم آورده که به موجب آن، نیاز کمتری به تصمیم گیری نامزد کننده ها در مرتبط و نامرتب بودن وبگاه وجود دارد.

- در DigiBoard نامزد کننده ها به رکوردهای خودشان شامل URLها در مجموعه های فعال، مجموعه های غیر فعال و نشانیهای وبگاههای آرشیو شده دسترسی دارند.

- DigiBoard امکان دسترسی تک - کلیدی را به منظور ارائه اطلاعات مدیریتی درباره مجموعه ها شامل مستندات برنامه ریزی شده مجموعه و اطلاعاتی درباره کنترل مجوزهای هر پروژه برای نامزد کننده ها فراهم می آورد.

- DigiBoard دارای واسط کاربری با قابلیت دسته بندی براساس مجموعه طبقه نامزد کننده بخش و مجوز میباشد.

- DigiBoard امکان ارسال چندین پست الکترونیکی مجوز به یک مالک محتوا را فراهم کرده است.

- DigiBoard امکان انقضای سالیانه و یا شش ماهه URLها را فراهم میکند و طی آن لازم است نامزد کننده ها مجدداً URLها را جهت گردآوری دوباره تأیید کنند.

- DigiBoard امکانی برای داوران فراهم میکند تا نامزدها را به منظور اعمال کار بیشتر و بررسی مجدد مجوزها به نامزد کننده ها برگشت دهند و وبگاههایی را که با راهنماهای انتخاب(1) مجموعه همخوانی ندارند رد کنند.

- DigiBoard امکان انتقال مالکیت یو آر ال از یک نامزد کننده به نامزده دیگر را فراهم می کند که طی آن مسئولیت تصدیق URL و گردآوری نامزدی دوباره نسبت داده شده(2)، با نامزد کننده جدید خواهد بود.

بررسی کیفیت پیش - خزش (2) فعالیتهای مدیریت، هسته و توسعه دستور عملهای تعیین حوزه خزش برای خزشگرهای هر URL نامزد شده را در برمی گیرد. قبل از توسعه این ماژول در ابزار قبلی مرحله مدیریت هسته به صورت خارجی صورت می‌گرفت بدین صورت که به ازای هر مجموعه یک فایل متنی از هسته ها به همراه دستور عملهای تعیین حوزه خزش وجود داشت، بنابراین ویرایش و به روز نگهداشتن فایل کار مشکلی بود و قالب بندی خاص URLها دستی انجام می‌گرفت که متمایل به خطا بود با پیاده سازی ماژول مدیریت هسته DigiBoard چندین فعالیت خودکار شده را فراهم کرده است که به موجب آن به نسبت آن به نسبت 10:1 در مصرف زمان صرفه جویی میشود.

1-5 - واریسی خودکار در نقطه شروع هسته

بیشتر مواقع نامزد کننده ها URL وبگاهی را معرفی میکنند که شامل فایلی با نام index است (برای مثال <http://www.loc.gov/index.html>). در این حالت سیستم در مرحله مرور پیش‌خزش بررسی میکند که آیا تعداد بایتهای <http://www.loc.gov> و <http://www.loc.gov/index.html> برابر هستند یا نه اگر برابر باشند یا نه سیستم <http://www.loc.gov> را به عنوان نقطه شروع بهینه برای خزش معرفی میکند.

2-5 - واریسی تک - کلیک برای پیوندهای سایت و نقشه سایتها

واریسی تک - کلیک برای پیوندها روش آسانی برای مرورگران تیم آرشیو وب فراهم میکند که به موجب آن میتوانند انتخاب و نمایش پیوندها برای تعیین میدانی برای شمول و همچنین برای انتخاب میدان مناسب برای محتوای سایت را به راحتی انجام دهند. برای مثال آرشیو انتخابات را در نظر بگیرید وبگاههای فیس بوک توئیتر، فلیکر مای، اسپیس و یوتیوب متعلق به کاندیدای انتخاب به عنوان وبگاه هایی در نظر گرفته میشوند که با دستور عملهای خزش مرتبط هستند. ماژول مدیریت هسته، شناسایی این URLها و تغییر فرمت مناسب برای خزشگر را آسان می کند.

3-5 - قالب بندی فهرست هسته

قالب بندی URLها در دامنه زیر دامنه مسیر و فایلهایی با فرمت (3) SURT که در حال حاضر در خزشگر هریتریکس (4) پشتیبانی میشود در ابزار DigiBoard به طور خودکار انجام می شود.

ص: 175

Seed management - 1

Pre-crawl - 2

Sort-friendly URI Reordering Transform - 3

Heritrix - 4

DigiBoard با پشتیبانی از امکان مدیریت، هفتگی ماهیانه سه ماهه شش ماهه و سالیانه خزشها، امکان ایجاد فهرست هسته بر اساس تناوب و براساس مجموعه را فراهم میکند این ماژول مراحل جاری قرار داد بستن با عاملهای خزشگر را آسان میکند فهرست هسته صادر شده و به عاملها تحویل داده می شود.

عکس

۱۷۶ مدیریت منابع اطلاعاتی وب

5-4- مدیریت تناوب^۱ هسته

DigiBoard با پشتیبانی از امکان مدیریت هفتگی، ماهیانه، سه ماهه، شش ماهه، و سالیانه خزشها، امکان ایجاد فهرست هسته براساس تناوب و براساس مجموعه را فراهم می کند. این ماژول، مراحل جاری قرارداد بستن با عاملهای خزشگر را آسان می کند، فهرست هسته صادر شده و به عاملها تحویل داده می شود.

6- بررسی پس خزش^۲

DigiBoard از دو نوع بررسی کیفیت پس خزش پشتیبانی می کند: براساس نامزدکننده، براساس تیم آرشیو وب.

6-1- بررسی کیفیت نامزدکننده

امکان انجام بررسی کیفیت برای نامزدکنندهها، باعث می شود تا تیم آرشیو وب از تعداد نامزدکنندههایی که محتوای گردآوری شده را بازرسی می کنند مطلع شوند. واسط کاربری DigiBoard پیوندی به آرشیو هر کدام از ویگاههای نامزد شده برای نامزدکنندهها فراهم می کند و بررسی نامزدها را با ثبت شناسه کاربری، و URL بررسی شده ضبط و ثبت می کنند. با ورود به آرشیو، پیوندی در بنر آرشیو وجود دارد که با کلیک روی آن می توانند با «ضبط موفقیت آمیز است» و «موفقیت آمیز نیست» همراه با توضیحی پاسخ دهند (شکل ۴).



شکل ۴. فرم بررسی کیفیت نامزدکننده

1. Frequency
2. Post-Crawl

DigiBoard از دو نوع بررسی کیفیت پس خزش پشتیبانی میکند بر اساس نامزدکننده بر اساس تیم آرشیو وب.

1-6 بررسی کیفیت نامزد کننده

امکان انجام بررسی کیفیت برای نامزد کننده ها باعث میشود تا تیم آرشیو وب از تعداد نامزدکنندههایی که محتوای گردآوری شده را بازرسی میکنند مطلع شوند واسط کاربری DigiBoard پیوندی به آرشیو هر کدام از وبگاههای نامزد شده برای نامزد کننده ها فراهم میکند و بررسی نامزده کننده ها را با ثبت شناسه کاربری و URL بررسی شده ضبط و ثبت میکنند با ورود به آرشیو پیوندی در بنر آرشیو وجود دارد که با کلیک روی آن میتوانند با ضبط موفقیت آمیز است و موفقیت آمیز نیست همراه با توضیحی پاسخ دهند (شکل 4).

شکل 4. فرم بررسی کیفیت نامزد کننده

ص: 176

frequency -1

post-crawl -2

تأیید عمق و بگاهی که کتابخانه LOC ضبط میکند در طول مرحله بررسی کیفیت پس خزش رخ میدهد به دلیل مرحله پردازش مجوزها، خزشهای LC با دستور عملهای صریح به خزشگر - با توجه به اینکه اجازه داده کجا در وب برود - انجام میشود. بررسی کیفیت پس خزش به سؤالیهای زیر پاسخ میدهد:

آیا قلمرو یا حوزه سایت به درستی مشخص شده است؟

آیا تمام وبگاههای مرتبط دیگر مشخص شده و حوزه آنها تعیین شده است؟

آیا ضبط مشکلی دارد؟

آیا این مشکلات شناخته شده و قابل شناسایی هستند یا به تحقیقات بعدی نیاز است؟

مرورگر حداقل یکبار در زمان گردآوری، سایت به ازای هر URL انتخابی با استفاده از ماژول [DigiQR\(1\)](#) دیجی بورد و پلاگین QR فایر فاکس بررسی کیفیت را انجام میدهد. همان طور که در شکل ه نشان داده شده است مرورگر میتواند مشخص کند که سایت خوب درو (جمع آوری) شده است، یا به منظور لزوم بررسی بیشتر گزینه partial را انتخاب کرده و مسئله ای [\(2\)](#) را عنوان کند.

شکل 5. بررسی کیفیت وبگاه

اگر مسئله آرشیوی باشد مرورگر نوع مسئله و شدت آن را شناسایی می کند. مسئله ها به انواع خاصی از مسئله متعارف طی مرحله بررسی کیفیت طبقه بندی شده. است پلاگین نوعی واسط کاربری را ایجاد کند که مسئله را - در همان صفحه ای که پیدا شده است - گزارش میدهد(شکل 6).

ص: 177

شکل 6. صفحه گزارش کردن مشکلاتی در ماژول DigiQR

در مثال QR (شکل 7) وب سرور مانع ضبط کردن وبگاه توسط خزشگر می شود. متخصص بررسی کیفیت، مسئله را شناسایی میکند و سطح شدت آن را تعیین می کند.

شکل 7. شناسایی مشکلات در ماژول DigiQR

عکس

The Library of Congress QR TOOLS

QR Tools: Issue Report

| | |
|------------|-------------------------------------|
| Issue URL | http://www.randyforcongress.com/ |
| Live | <input checked="" type="checkbox"/> |
| Seed | |
| Collection | U.S. Election 2010 |

No previous entry for this or a similar URL found.

To exit at any time, simply close this window.

- Only work in proxy
- Timeline/Banner issues
- LC template/Banner shows in place of images
- Navigation tabs send users to live site
- Archive issues but looks the same on the live site
- Layout issues
- Can't go beyond resource page
- Image issues
- Video playback issues
- Website not in archive
- Some contents not in archive
- Scoping issues
- Redirect issue
- Website disappeared
- Other

شکل ۶. صفحه گزارش کردن مشکل‌های در مازول DigiQR

در مثال QR (شکل ۷) وب سرور مانع ضبط کردن وبگاه توسط خزشگر می‌شود. متخصص بررسی کیفیت، مسئله را شناسایی می‌کند و سطح شدت آن را تعیین می‌کند.

Other

Recommendation(s) Please be detailed in your description.

Comment

Severity

Status

شکل ۷. شناسایی مشکلات در مازول DigiQR

این مسئله توسط یک عضو تیم آرشیو وب بیشتر بررسی میشود مراحل تحقیقات و نتایج آن ثبت می‌گردد واسط جست و جوی DigiQR امکان جست و جو بر اساس ، مسئله براساس مجموعه، براساس URL و براساس خزش را فراهم می‌سازد.

پیش از توسعه ماژول DigiQR کلیه کارها به صورت مجزا صورت میگرفت و اعضای تیم آرشیو وب از صفحه گسترده‌های خودشان با اصطلاحات فنی متناقض جهت مستندسازی مسئله‌ها استفاده می‌کردند. اکنون DigiQR تیم را قادر می‌سازد تا مسئله‌های محتوای خزش شده را در یک چارچوب مدیریت دانش اشتراکی با کارآیی بیشتر ثبت و پیگیری نمایند. اعضای تیم یاد گرفته اند که مسئله‌های QR پیچیده تری را شناسایی کنند زیرا ابزار امکان دسترسی آسان به اطلاعات مسئله‌ها و نتایج تحقیقات را فراهم می‌آورد این ماژول زمان صرف شده در کپی کردن URL‌ها در صفحه گسترده‌ها را کاهش داده امکان کارآمدتری برای پیگیری مسئله فراهم میکنند این ماژول مکان رخ دادن مشکل را دقیقاً ثبت جست و جوی‌های بعدی را تسهیل می‌بخشد.

7- فعالیتهای اجرایی

ماژول‌های اجرایی از تمام فعالیتهای مدیریتی دیجی‌برد پشتیبانی میکند، اما به طور قابل توجهی امکاناتی برای مدیریت نقشهای کاربران و مدیریت مجموعه‌های آرشیو وب فراهم شده است.

7-1- مدیریت کاربر

با توجه به تعداد تخمینی کاربران DigiBoard توانایی مدیریت فعالیتهای کاربران توسط تیم آرشیو وب ، امری حیاتی است به دلیل اینکه شمول یک URL در فهرست هسته به قیمت زمان کارمندان در مرحله QR، دسترسی و فهرست‌نویسی و صرف اعتبار مالی تمام میشود در سالهای اخیر یک مفهوم دو لایه از نامزد کننده‌ها و داوران به عنوان بخش مهمی از فرآیند جریان کاری به وجود آمد. دیجی‌برد جهت تطبیق با این جریان کاری سه سطح کاربری دارد نامزد کننده داور و مدیر نامزد کننده می‌تواند هر منبعی از وب را برای مجموعه معرفی کند داور قادر است تصدیق کند آیا نامزدها با راهنمای گزینش منبع متناسب است یا نه. البته نامزد کننده مسئولیت فراداده ابتدایی مانند ، موضوع زبان وبگاه و فرآیند مجوزها را نیز دارد. مدیران نیز عضوهای تیم آرشیو وب هستند.

7-2- مدیریت مجموعه

در حال حاضر دیجی‌برد داده‌های 14,151 وبگاه را در 39 مجموعه آرشیو وب مدیریت میکند و گزارشهایی که از بخش مدیریت کتابخانه کنگره درباره پروژه آرشیو وب درخواست میشود از محتوای پایگاههای اطلاعاتی به راحتی قابل تولید است. داده‌های قدیمی از ابزار قدیمی به DigiBoard انتقال یافته است، بنابراین، سوابق گذشته و داده‌های جدید در یک سیستم یکتا قابل نگهداری است.

زمانی که توسعه و ایجاد DigiBoard رشد پیدا کرد ضرورت ارتباط با کاربران به عنوان یک نیاز شناسایی شد. صفحه اصلی به منظور ایجاد فضایی برای اعلانها تغییر کرد و نوعی سیستم پیغام رسانی برای پیگیری نیز فراهم شد.

8- مازولهای ویژه

ماژولهای مختلفی جهت پشتیبانی فعالیتهای فرعی در حال توسعه هستند این ماژولها با هدف مدیریت مجموعه ها و مدیریت محتوای آرشیو وب تدارک دیده شده اند.

8-1- مازول نامزدهای انتخابات ملی آمریکا

ماژول نامزدها در اول جولای 2010 برای نخستین بار شروع به کار کرد این ماژول، آرشیو وب انتخابات را برای علم کتابداری (LS1) فراهم میکند و هر دو سال یکبار سایتهای جدید به آرشیو می پیوندند. در گذشته، داده های مبارزه های انتخاباتی فدرال (2) را به منظور بررسیهای بعدی در یک پایگاه داده اکسس وارد میکردند این پایگاه اطلاعاتی نوعی جریان کاری برای پردازش اطلاعات 2200 کاندیدا، پشتیبانی میکرد اما با ابزار تیم آرشیو وب یکپارچه نبود، بنابراین گرفتن داده از یک پایگاه داده و وارد کردنش به پایگاه داده دیگر مستلزم پردازشهای دستی بود علاوه براین به دلیل اینکه نامزد کننده ها با پایگاه اطلاعاتی اکسس کار میکردند اطلاعات مجوز روکشی برایشان قابل دسترس نبود بنابراین دیگر به تعیین هویت اطلاعات تماس نامزدهایی که قبلاً پاسخ داده بودند نیاز نبود مقصود ماژول نامزدها توسعه ابزاری کارآمد و یکپارچه برای علم کتابداری و تیم آرشیو وب بود.

8-1-1- مجوزهای روکشی نامزدها

مجوزهای روکشی توسط نامزدها مدیریت میشوند؛ بدین ترتیب موجب سهولت مدیریت مجوزها برای نامزدهایی میشوند که در هر انتخابات URL متغیری دارند.

8-1-2- داده های گذشته

دادههای گذشته هر نامزد برای نامزد کنندهای ماژول در دسترس است تمامی داده های خزش شده داده های مبارزه های انتخاباتی فدرال و عملیات انجام شده حفظ خواهند شد و به راحتی در اختیار محققان آینده قرار خواهند گرفت.

ص: 180

این ماژول قابلیت استانداردسازی فراداده هر نامزد را فراهم می آورد داده مبارزه های انتخاباتی فدرال که به صورت دستی در پایگاه داده ها وارد شده است دارای خطاهایی در ورودی است که این ماژول با گذر، زمان رکوردهای پایداری از آنها فراهم خواهد آورد.

-4-1-8- یکپارچگی DigiBoard

اگرچه ماژول نامزدها کاملاً با DigiBoard یکپارچه شده است از دیدگاه بصری متفاوت است و برای به انجام رسانیدن تکلیف ویژه ای طراحی شده است شناسایی نامزدها با وبگاههایی که کتابخانه کنگره تمایل دارد به آرشیو اضافه کند زمانی که کار در ماژول نامزدها تکمیل شد رکورد به عنوان رکوردی جدید به DigiBoard نیز انتقال می یابد تا مراحل آرشیو وب آن تکمیل می شود.

9- ماژولهای آینده

ماژولی با نام مدیریت محتوا در دست توسعه است که امکان ردیابی محتوا در زمان تولید آن توسط عامل خزشگر، انتقال به کتابخانه کنگره، کپی شدن در مخزن ذخیره سازی و دسترسی به سرورها را فراهم میکند این ماژول امکانی برای پاسخ دادن به پرسشهای مدیریتی درباره میزان محتوا آرشیو وب کتابخانه کنگره و مکان ذخیره سازی آن را فراهم می کند این ابزار جهت تکمیل ابزارهای پیچیده ردگیری و ابزارهای سیاهه بندی و قفسه بندی - که توسط کتابخانه کنگره به منظور پشتیبانی از حفاظت رقومی پیاده سازی می شوند - به راحتی قابل اصلاح است.

فعالیتهای توسعه آینده شامل موارد زیر است:

یک ماژول فهرست نویسی که تولید خودکار رکوردهای متس را برای آرشیو وب پشتیبانی میکند؛ بنابراین با در نظر داشتن فراداده ای که توسط ابزار ضبط می، شود از مدت زمانی که توسط فهرست نویسان جهت پردازش این رکوردها مورد نیاز است کاسته میشود. ماژول دیگری نیز مد نظر است که از یک رویکرد انتخابی به منظور شناسایی منابع برای فعالیتهای حفاظت پشتیبانی کند.

10 انجمن آرشیو وب

10-1- ابزار کتابدار وب

ابزار متصدی وب (1)(WCT) برنامه مدیریت جریان کاری متن بازی برای آرشیو وب انتخابی است که در سال 2006 با همکاری مشترک کتابخانه ملی نیوزلند و کتابخانه بریتانیا توسعه یافت و نخستین ابزار کتابدار برای انجمن آرشیو وب شد تا زمانی که WCT برای پتانسیل کاربردی ارزشیابی می شد، کتابخانه کنگره با یک برنامه کاربردی پست الکترونیکی برای اطلاع رسانی مالکان محتوا همکاری می کرد. سرانجام،

ص: 181

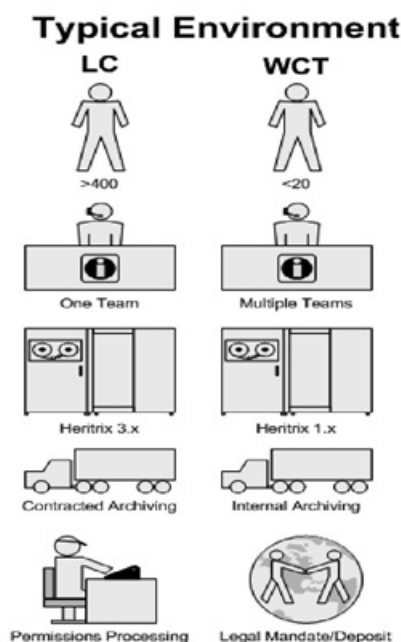
کتابخانه متقاعد شد که WCT به علت برخی تفاوت‌های عملیاتی نیازهای آرشیو وب کتابخانه کنگره را مرتفع نخواهد کرد (شکل 8)

شکل 8. مقایسه محیط کاری LC و WCT

التزامهای پردازش مجوزها در کتابخانه کنگره، تعداد کارکنان درگیر در مرحله نامزدی و طبیعت اینکه چطور آرشیو وب کتابخانه ساخته شده است ضرورت ایجاد DigiBoard و ماژولهای وابسته اش را تعیین کرد.

عکس

کتابخانه متقاعد شد که WCT به علت برخی تفاوت‌های عملیاتی، نیازهای آرشیو وب کتابخانه کنگره را مرتفع نخواهد کرد (شکل ۸).



شکل ۸. مقایسه محیط کاری LC و WCT

التزام‌های پردازش مجوزها در کتابخانه کنگره، تعداد کارکنان درگیر در مرحله نامزدی و طبیعت اینکه چطور آرشیو وب کتابخانه ساخته شده است، ضرورت ایجاد DigiBoard و مازول‌های وابسته اش را تعیین کرد.

۱۰-۲- فعالیت‌های گردآوری درون سازمانی

در حال حاضر، برنامه‌ای جهت توسعه مازولی به منظور مدیریت فعالیت‌های گردآوری درون سازمانی^۱ با استفاده از خزشگر هریتریکس^۲ کتابخانه کنگره وجود ندارد. WCT و نرم‌افزار سوئیت نت آرشیو که توسط کتابخانه ملی و اسپاری و کتابخانه رویال دانمارک توسعه یافته است برای استفاده در مراحل مدیریتی خزش درون سازمانی، ارزیابی خواهد شد.

1. On-site
2. Heritrix Crawler

10-2- فعالیتهای گردآوری درون سازمانی

در حال حاضر برنامه ای جهت توسعه مازولی به منظور مدیریت فعالیتهای گردآوری درون سازمانی(1) با استفاده از خزشگر هریتریکس(2) کتابخانه کنگره وجود ندارد WCT و نرم افزار سوئیت نت آرشیو که توسط کتابخانه ملی و اسپاری و کتابخانه رویال دانمارک توسعه یافته است برای استفاده در مراحل مدیریتی خزش درون سازمانی ارزیابی خواهد شد.

On-site -1

Heritrix Crawler -2

اگر چه DigiBoard به منظور پشتیبانی از جریان کاری کتابخانه کنگره توسعه یافته است، غرض اولیه این ابزار توسعه یک واسط کاربری آسان برای تعداد عظیمی کتابدار است تا فعالیتهای گردآوری را انتخاب و مدیریت کنند نامزدی پردازش، مجوزها مرور کیفیت و مدیریت خزش این جریانهای، کاری با جریانهای، کاری کشورهایی که قانون و اسپاری ندارند. همسوست کتابخانه کنگره تمایل دارد ماژولهای DigiBoard را برای نسخه متن باز آن استاندارد نماید.

11- نتیجه گیری

توسعه ابزار کتابدار که با جریان کاری منحصر به فرد، کتابخانه نیازمندی مجوزها و محیط کاربران متناسب باشد گامی رو به جلو در آرشیو وب کتابخانه کنگره است خودکارسازی فعالیتهای دستی تیم آرشیو وب جهت افزایش قابلیت مقیاس پذیری فعالیتهای آرشیو وب کتابخانه کنگره بسیار حیاتی است. از طرف دیگر همزمان با بهبود DigiBoard و سهولت استفاده آن برای کاربران مختلف، می توان کارکنان کتابخانه کنگره را در پروژههای آرشیو وب بیشتری درگیر کرد.

12- منابع

[1] [/http://archive-access.sourceforge.net/projects/wayback](http://archive-access.sourceforge.net/projects/wayback)

[2] [/http://webcurator.sourceforge.net](http://webcurator.sourceforge.net)

[3] [/http://crawler.archive.org](http://crawler.archive.org)

[4] [/http://netarchive.dk/suite](http://netarchive.dk/suite)

حفاظت و آرشیو کردن وبگاهها چالشی عظیم از حیث فناوری است که باید بی وقفه دنبال شود تحول و تکامل در، فناوری محتوا و رشد مستمر وب ما را به این معنا هدایت می کند که هیچ راه حل قطعی در آرشیوسازی وب نخواهیم یافت. رونوشت برداری از وبگاهها گزینشی در طراحی وب است. سه روش برای رونوشت برداری کامل از هر وبگاه وجود دارد اولین، روش آرشیو از سرور، که از همه روشها دشوارتر است. در این روش باید با مدیران سایت تماس گرفت و آنها را متقاعد کرد که رونوشتی از فایل های سامانه اطلاعاتی داخلی فراموهای پایگاه داده و ویژگیهای سامانه را تهیه کنند روش دوم این است که این کار را در سرور انجام دهیم و تراکنشهای مربوط به عمل آرشیو کردن منابع وب را ضبط کنیم. آخرین روش، گردآوری خودکار اطلاعات به طور مستقیم از وبگاههاست، به همان ترتیبی که یک مرورگر معمولاً آن را انجام می دهد. در مقاله حاضر به بحث رونوشت برداری از وبگاهها از جنبه های مختلف نظیر روشهای رونوشت برداری، تجزیه کننده هسته HTML، تجزیه کننده پردازنده تجزیه کننده های گروه جاوا، واکنشی مدارک، تایید وثوق، اتصال پذیری و مدیریت روزآمد.

*رونوشت برداری از وبگاه ها(1)

نوشته: خاویر روش(2) | ترجمه: فرزانه شادان پور(3)

1 - مقدمه فناوری رونوشت برداری از وبگاه ها

تفاوت اساسی میان رونوشت بردار(4) از داده ساختارها(5) یا سایتهای ftp و وبگاهها در طبیعت ژرف وب جهانی نهفته است. این نه به ویژگی «فهرست راهنما» بودن در پروتکل HTTP مربوط می شود، نه به ویژگی انتقال حجم انبوه و بگاهاها؛ بلکه گزینشی در طراحی وب است که خود مجموعه ای است از منابع ناهمگن که لزوماً با یکدیگر در ارتباط نیستند. به عنوان مثال، مجموعه ای از صفحاتی که در محتوای یک پایگاه داده تولید میشوند قلمرو بی ثباتی از اطلاعات است. پایگاه داده که پاکسازی(6) شود، صفحات مذکور نیز ناپدید میگردند، وب به طور کلی شکلی در حال تغییر است راهنماهای ftp تغییر می کنند ولی شما هر جا که هستید به آسانی میتوانید آنها را همزمان(7) کنید تا روز آمد باشند. این فقط مقداری داده ذخیره شده در یک انباره(8) فایل است. ولی یک صفحه وب بالقوه منحصر به فرد است ممکن است یک

ص: 185

Roche. Xavier (2006). "Copying websites", in Masanes, Julien (ed.), Web Archiving. Berlin Heidelberg – 1
New York: Springer.pp.93-113

Roche. Xavier -2

3- مری، عضو هیئت علمی سازمان اسناد و کتابخانه ملی ایران

Coping -4

File structure -5

Flush -6

Synchronize -7

Repository -8

ساعت شمار(1)، اطلاعات تحویلی بی درنگ بر حسب تقاضا یا نمایی(2) ویژه کاربر یا ویژه تراکنش(3) از یک مجموعه کلیتر از داده باشد میتواند هر چه شما میخواهید باشد منطق درونی آن از ابزارهای کاربری و نیز کاربران این ابزارها پوشیده است و ساده بگویم یک سرور ftp مجموعه ای از فایلهاست که بیشتر به دیسک سختی که عمومی و قابل کنترل از راه دور است شباهت دارد.

یک سرور وب را میتوان مجموعه ای از منابع منطقی(4) دانست که به مشتری ها محتوا تحویل می دهند. این منابع منطقی ممکن است برنامههایی باشند که به پایگاه داده و سایر سامانه ها از طریق تراکنش با ترجیحات و نیازهای کاربر و محیط سرور (پایگاه داده منابع، بیرونی وضعیت فعلی و غیره) متصل اند. مشتری راه دور هرگز این منطق را مشاهده نمی کند و فقط محتوای به دست آمده قابل دسترسی اوست.

بنابراین سه روش برای رونوشت برداری کامل از هر وبگاه وجود دارد نخستین آرشیو از سرور(5) از همه روشها دشوارتر است در این روش باید با مدیران سایت تماس گرفت و آنها را متقاعد کرد که رونوشتی از فایل های سامانه اطلاعاتی، داخلی فرمانهای(6) پایگاه داده و ویژگیهای سامانه را تهیه و همان معماری را از حیث سخت افزار نرم افزار و محیط رایانشی (مانند منابع داده ای که میتوانند مورد استفاده قرار گیرند) در آن ترتیب دهند.

این راه حل را که معمولاً حتی برای خود مدیران سایت نیز بسیار دشوار است نمی توان به عنوان یک راه حل عمومی در نظر گرفت. گزینه دوم این است که این کار را در سرور انجام دهیم و همه تراکنشهای(7) مربوط به عمل آرشیو کردن منابع وب را ضبط کنیم. آخرین روش، گردآوری خودکار اطلاعات به طور مستقیم از وبگاههاست به همان ترتیبی که یک مرورگر معمولاً آن را انجام میدهد (آرشیو از سمت کاربر(8)). این یک راه حل موقت است چرا که رونوشتها هرگز کامل نیستند مثل این است که از یک صفحه متحرک یک عکس بگیریم در این صورت نمیتوانید حرکت آن صفحه را در عکس ایجاد کنید. همچنین نخواهید توانست امتیاز بی درنگ بودن را هنگام مرور در گزارشهای برخط یا جنب و جوش مبادله اطلاعات را داشته باشید ولی این از حیث امکان اجرا و کیفیت در بیشتر موارد ناگزیر پذیرفتنی است. مثل عکس ایستایی(9) از یک وبگاه است که باید در یک آلبوم عکس نگهداری شود؛ عکسی که بارها و بارها بشود آن را، دید حتی بدون نگرانی از این که اصل آن دیگر وجود ندارد. رونوشت برداری از وبگاهها با این روش امری بسیار شهودی است روش دقیقاً همان است که گویی خودتان بخواهید با یک مرورگر از یک وبگاه رونوشت برداری کنید در این صورت از نخستین صفحه آغاز میکردید صفحه و تصاویر آن را ذخیره و سپس بر هر پیوند در صفحه کلیک میکردید تا

ص: 186

Clock Counter -1

View -2

Session -3

Logical -4

Server - side -5

Schemas -6

Transaction -7

Client-side archiving -8

آنها را ببینید و صفحات مرتبط را در یک دیسک ذخیره می‌کردید و آنقدر این کار را ادامه می‌دادید تا از همه صفحاتی که می‌خواهید رونوشت بردارید. بعد تغییراتی در صفحات HTML می‌دادید و همه تگهای مربوط را واری می‌کردید تا با مرورگر سیستم خودتان قابل دیدن شوند. ولی رونوشت برداری از بیش از یکی دو صفحه به صورت دستی خسته کننده است و داشتن ابزاری برای خودکار کردن آن می‌تواند راه حل مؤثری باشد. شماره‌ده خودکار پیوندها را معمولاً تجزیه کننده (1) و نرم‌افزاری را که داده را به طور خودکار و از راه دور بارگذاری میکند خزش گر مینامند. این دو جزء اصلی نقشهای دیگری هم دارند. تجزیه کننده عهده دار بررسی و تضمین این است که پیوندها در یک رونوشت محلی (روی سیستم) نیز کار کنند، و این کار با تغییر نحو (2) یوآرال آن به نحوی صحیح و کاملاً مرتبط انجام پذیر می‌شود و خزش گر مسئول حافظه دم دستی (3) و روزآمد کردنهاست.

دلایل متعددی برای رونوشت برداری از وبگاهها وجود دارد در مدرسه ملی مهندسی کن (4) ما می‌خواستیم وبگاههایی با اندازه کوچک و متوسط را آرشیو کنیم، نه برای آرشیو سازی رایج، بلکه برای گردآوری سایت‌های فنی که اشخاص راه می‌انداختند و سریع هم دستخوش تغییر می‌شدند. همچنین می‌خواستیم وبگاههای بزرگ را که دارای محتوای چند رسانه ای بودند و با استفاده از خط‌های dial-up خانگی قابل دسترسی نبودند، گردآوری و آنها را در رسانه های دائمی مانند CD ذخیره کنیم تا بعد بتوانیم آنها را به صورت غیر برخط ببینیم. در مجموع به ابزاری برای گردآوری اطلاعات بسیار ویژه از وب برای کاربران نهایی نیاز داشتیم.

طرح HTTPTrack برای پاسخ به این نیازها پدید آمد این نرم افزار ابزاری سه‌لایه استفاده است که به کاربران معمولی اجازه می‌دهد از بخشهای کوچک - ولی مهم وب رونوشت برداری کند طراحی این نرم افزار بیشتر به طور تجربی صورت گرفته است. اینترنت و معماری شبکه مربوط به آن عرصه های نسبتاً جدیدی بودند که ما در آنها به اکتشاف می پرداختیم و رونوشت برداری از وبگاه به ویژه و به طور خاص موضوع کاملاً جدیدی برای ما بود تجربیاتی که از تدوین و توسعه این نرم افزار به دست آمد روش - «فناوری (5)» رونوشت برداری از وبگاهها و راه‌های برطرف کردن نقاط ضعف موجود را تشریح می‌کند.

2- تجزیه کننده

2-1- تجزیه کننده هسته (6) HTML

تجزیه کننده HTML یکی از دو جزء هسته در یک ابزار رونوشت برداری (7) از وب است. اگر یک صفحه

ص: 187

Parser -1

Syntax -2

Cache -3

National School of Engineering of Caen -4

Art -5

.Core parser -6

Web copying tool -7

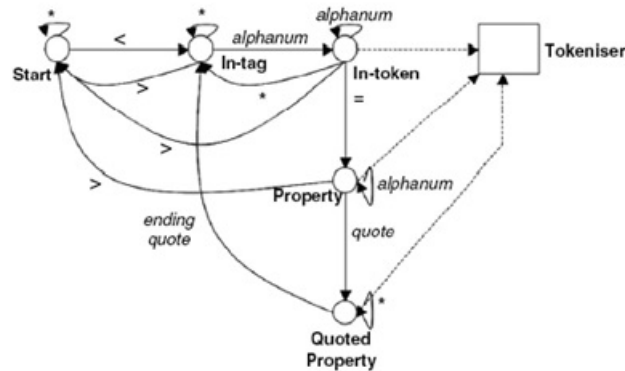
1) HTML را در نظر بگیریم که یک فایل متنی 8 بیتی (2) متشکل از متن ساده (بدون رمز) (3) و تگهای نشانه گذاری (4) است - و اطلاعات همراه آن مانند URL اولیه (5) هدف تجزیه کننده HTML است تا با پویش (6) صفحه پیوندها را گردآوری و تجزیه و تحلیل کند و آنها را به خزشگر بسپارد. ساختار HTML برای گردآوری پیوندها مناسب نیست ما در ابتدا به تعداد محدودی از تگ ها نظیر عناصر «a» یا «img»، علاقه مندیم که بالقوه ممکن است حاوی ابرپیوندهایی به سایر منابع (تصاویر style sheets، صفحه های HTML، و مانند آن) باشند جای قرار گرفتن آنها در صفحه معمولاً مهم نیست، مگر در دامنه های خاصی که هیورستیک پیشرفته (7) میتواند اطلاعات اضافی (مانند موضوعات «حول و حوش») به آن ضمیمه کند. این تگها معین میکنند کدام صفحهها تعقیب شوند و کدام صفحه های نامرتبط تعقیب نشوند. برای یک تجزیه کننده معمولی تنها اطلاعات مهم تگها و متعلقات جاسازی شده آنها هستند.

شکل ساده شده اتوماتون (8) هسته، به خوبی و سادگی قابل درک است پویش خطی (9) دادههای صفحه HTML، شروع از اول کشف تگهای آغازین (<) و بازشناسی (10) عناصر مختلف HTML با نامشان (نگاه کنید به شکل 1).

عکس

HTML^۱ را در نظر بگیریم - که یک فایل متنی ۸ بیتی^۲ متشکل از متن ساده (بدون رمز)^۳ و تگهای نشانه گذاری^۴ است - و اطلاعات همراه آن مانند URL اولیه^۵ هدف تجزیه کننده HTML است تا با پوش^۶ صفحه پیوندها را گردآوری و تجزیه و تحلیل کند و آنها را به خزشگر بسپارد. ساختار HTML برای گردآوری پیوندها مناسب نیست. ما در ابتدا به تعداد محدودی از تگها نظیر عناصر «a» یا «img»، علاقه مندیم که بالقوه ممکن است حاوی ابرپیوندهایی به سایر منابع (تصاویر، style sheets، صفحه‌های HTML، و مانند آن) باشند. جای قرار گرفتن آنها در صفحه معمولاً مهم نیست، مگر در دامنه‌های خاصی که هیورستیک پیشرفته^۷ می‌تواند اطلاعات اضافی (مانند موضوعات «حول و حوش») به آن ضمیمه کند. این تگها، معین می‌کنند کدام صفحه‌ها تعقیب شوند و کدام صفحه‌های نامرتبط تعقیب نشوند. برای یک تجزیه کننده معمولی تنها اطلاعات مهم، تگها و متعلقات جاسازی شده آنها هستند.

شکل ساده شده اتوماتون^۸ هسته، به خوبی و سادگی قابل درک است: پوش خطی^۹ داده‌های صفحه HTML، شروع از اول، کشف تگهای آغازین (<) و بازشناسی^{۱۰} عناصر مختلف HTML با نامشان (نگاه کنید به شکل ۱).



شکل ۱. هسته تجزیه کننده

۱. نگاه کنید به [۱۸۶۶]
۲. به یاد داشته باشید که رمزگذاری کاراکتر صفحه برای نامگذاری (naming) به خصوص در فایل سیستمهای UCS۲ (از جمله ویندوز) اهمیت خواهد داشت.
۳. Plain text
۴. Markup tags
۵. نگاه کنید به [۱۷۳۸]
۶. Scan
۷. Advanced heuristic
- ۸ Automaton : ماشین خودکار، کامپیوتر
۹. Linear Scan
۱۰. Recognizing

شکل ۱. هسته تجزیه کننده

ص: 188

1- نگاه کنید به [1866]

2- به یاد داشته باشید که رمزگذاری کاراکتر صفحه برای نامگذاری (naming) به خصوص در فایل سیستمهای UCS2 از جمله ویندوز اهمیت خواهد داشت.

Plain text -3

Markup tags -4

-5 نگاه کنید به [1738]

Scan -6

Advanced heuristic -7

-8 Automaton : ماشین خودکار کامپیوتر

Linear Scan -9

Recognizing -10

دو دسته از محتویات درون تگ های HTML قابل بازشناسی است نام تگها(1) مانند "img" یا "a" و ویژگیهای تگ مانند href یا "src" این تگها را میتوان به دو گروه اصلی تقسیم کرد تگهایی که امکان جاسازی(2) منابع (نظیر، تصویر style sheetsهایی که در صفحه وجود دارند) و تگهایی که ناوبری به سایر منابع (ابریوند) را میسر می سازند در یک صفحه میتوانید از پیوندهای نامرتبط در گروه دوم صرف نظر کنید (برای مثال پیوندهایی که خارج از دامنه رونوشت آینه ای(3) هستند). در این صورت، این یک محیط برون خطی غیر قابل دسترسی خواهند بود ولی این مسئله باعث تغییر صفحه نمی شود در مورد تگهای گروه نخست باید دقت بیشتری به خرج دهید چرا که در غیر این صورت صفحه در حالت غیر برخط به درستی قابل دیدن نخواهد بود ممکن است تصاویری را از دست بدهید یا چینش(4) صفحات به علت اجزای از دست رفته مانند style sheets یا فایلهای پردازش نویسی(5) جاسازی شده به هم ریخته باشد بنابراین یوآرلهای پیوند تنها اطلاعاتی نیستند که باید به خزشگر داده شوند. بافت تگ(6)، مثلاً اینکه آیا منبع یک منبع جاسازی شده است یا نه؟ هنگام تصمیم گیری در مورد این که این پیوند را بگیر یا خیر اهمیت خواهند داشت.

Tokenizer پیوندها را با تجزیه و تحلیل ویژگیهای شناخته شده شان برداشت میکند. پیوندها با استفاده از URL صفحه اصلی به شکل مطلق(7) تبدیل میشوند که این قسمتهاست پروتکل http، میزبان: "http://www.example.com" و مسیر مربوط index.html. به عنوان مثال پیوندهای نسبی(8) (top.html) در درون صفحه [مثال]

http://www.-example.com/foo/index.html)

تبدیل می شوند به پیوند «http://www.example.com/foo/top.html». سپس جای پیوند بررسی میشود تا از منطبق بودن با دامنه پیش فرض رونوشت آینه ای اطمینان حاصل آید؛ بررسی شامل عبارت(9) معمولی است که مقدار آن به طور پیش فرض پیشوند(10) اصلی URL است. اگر گرفتن رونوشت آینه ای را از «http://www.example.com/foo/index.html» آغاز کرده باشیم دامنه پیش فرض در نحو شبه معمول عبارت خواهد بود:

*http://www.exampel.com/foo

غیر از این پیوندهایی نظیر

http://www.example.com/foo/top.html,

ص: 189

Tag names -1

Embed -2

Mirror -3

Layout -4

Scripting files -5

Tag context -6

4.4 Hierarchical URL and Relative Forms [2396] Sect: نگاه کنید به: Absolute Form -7

8- : نگاه کنید [1808]. Relative links

9- Expression

10- Prefix

به طور پیش فرض در رونوشت آیینهای داخل خواهند شد. البته بسته به سایتی که قرار است رونوشت برداری شود، قواعد بیشتری ممکن است لازم باشند بنابراین عبارت پیش فرض - بسته به نیاز - باید قابل تغییر باشد. سرانجام پیوندهای تکراری نباید به خزشگر منتقل شوند تجزیه کننده باید وضعیت همه URL های شناخته شده را بداند و از چندبار گرفتن پیوندها پرهیزد.

تجزیه کننده همچنین باید بتواند با نحوهای متعددی کار کند که ممکن است مرکب از اشکال نسبی یا مطلق URL، ها گریز (1) HTML (مانند nbsp) گریز یو آر ال (2) (مانند a3 درصد) و به طور کلی هرگونه نحوه ای باشند مرورگرها تا حد زیادی با این نحوها مدارا میکنند حتی وقتی صفحه دچار از هم گسیختگی است (از جمله خطا در نحو نگها) مرورگرها اغلب نهایت تلاش میکنند تا آن را تجزیه و تحلیل کنند و به شکلی که قابل درک باشد بالا بیاورند.

به عنوان مثال پیوند URL مطلق `http://www.Example.com/page2.html` را میتوان با نحوهای متعددی از جمله نحوهای غلط ارجاع داد به هر شکل URL را باید بتوان بازشناخت و آن را به حساب آورد.

در پایان، پیوندها باید بازشناسی شوند تا با ساختار رونوشت آیینهای وبگاه هماهنگ شود. لازم است پیوندهایی که شکل مطلق دارند مانند

" `http://www.example.com/index.html` به شکل نسبی مبدل شوند مانند `index.html` پیوندهایی که خارج از دامنه رونوشت آینه ای هستند پیوندهایی که با دامنه عبارت پیش فرض جور نیستند باید به شکل مطلق تبدیل شوند بنابراین در صفحه های رونوشت آینه ای باید تغییراتی داد تا در یک ساختار محلی قابل استفاده شوند.

2-2- تجزیه کننده پردازنده (3)

چند ماه پس از آغاز تدوین و توسعه HTTPTrack و علیرغم بهبودی که در تجزیه کننده های HTML رخ داده بود ملاحظه میشد که بعضی وبگاهها به درستی رونوشت برداری نشده اند و تعداد زیادی از تصاویر و فایل های صفحات در رونوشتها وجود نداشت و این باعث خطا در ناوبری میشد فقط به این علت که تجزیه کننده پیوندهای مربوط به آن تصاویر و فایلها را ندیده بود.

درون صفحات HTML باید مناطق (4) پردازنده نویسی (5) خاصی مانند جاوا (کد فعال داخل شده در صفحات) در نظر گرفته میشدند که عملکرد تجزیه ویژه ای برایشان لازم بود پرداخت کامل کد پردازنده تقریباً غیر ممکن است کد با تگهای HTML متفاوت است. آنها مشابه تگهای HTML نیستند که بسیار آسانتر میتوان تجزیه و تحلیل شان کرد منطبق نهفته در متغیرها توابع و عبارات ممکن است

ص: 190

[HTML escaping: "Proposed Eutites". 14.sect [7666 -1]

2- نگاه کنید: [URL escaping[1630].

Script arser -3

Zones -4

-5 : نگاه کنید به [ECMA-262 Scripting generalization ECMA Scripting]

بالقوه غیر قابل دریافت باشد نخست اینکه حتی با تفسیر کننده(1) جاوای کامل، عملیاتی که یکباره به علت وضعیت قرار گرفتن موشواره انجام میشوند کلیک بر عناصر و اجزای صفحه، یا محیط (زمان، متغیرهای کاربری و به طور کلی آنتروپی محیط و مانند آن) را نمیتوان به دست آورد. دوم اینکه گرفتن پیوندها با استفاده از تفسیر کننده مسئله دیگر را حل نخواهد کرد که آن تغییر منطق کدها برای انطباق با سایت رونوشت برداری شده است و اگر چه این کار درون تگ HTML ساده است انجام آن درون یک کدپدازه نویسی پیچیده تقریباً غیر ممکن است.

خوشبختانه در بیشتر موارد کدهای جاوای به کار گرفته شده آنقدر ساده هستند که با توان محدود برنامه بتواند جور در بیاید برای بارگیری خودکار تصاویر یا برای گرفتن آنها در تصویر زمینه طراحان نرم افزار معمولاً از ارجاع مستقیم به ویژگی شیء با استفاده از رشته های ایستایی(2) نظیر `foo.src bar grf` استفاده میکنند یا برای گشودن یک پنجره جدید عبارتی مانند `window. Open (foo.html)` را به کار میبرند. حدود 80 درصد از پیوندهای پنهان درون مناطق پردازش را با به کار بردن این نمونه های ساده میتوان کشف کرد و تغییر داد سایر موارد با استفاده از عبارتها یا روشهای ناشناخته همانطور که هستند باقی میمانند نتیجه عالی نخواهد بود و در مورد `HTTrack`، از ابتدا میدانستیم که همه چیز آنطور که میخواهیم پیش نخواهد رفت هدف رسیدن به سطح قابل قبولی از کیفیت برای بیشتر سایت ها بود.

مناطق(3) `CSS` را میتوان با الگوریتمهای مشابهی تجزیه کرد.

اتوماتون ساده شده زیر برداشت رشته های(4) متنی درون ناحیه های پردازش نظیر بخشهای تگ



مرکز تحقیقات رایانگی

اصفهان

گامی

WWW



برای داشتن کتابخانه های تخصصی
دیگر به سایت این مرکز به نشانی

www.Ghaemiyeh.com

www.Ghaemiyeh.net

www.Ghaemiyeh.org

www.Ghaemiyeh.ir

مراجعه و برای سفارش با ما تماس بگیرید.

۰۹۱۳ ۲۰۰۰ ۱۰۹

