

ارائه روشی برای گسترش پرس و جو در بازیابی اطلاعات XML با استفاده از بازخورد کور

فروغ شهابیان^۱، احمد برانی^۲

^۱ دانشجوی کارشناسی ارشد، گروه کامپیوتر دانشکده فنی و مهندسی، دانشگاه اصفهان، اصفهان
shahabian@comp.ui.ac.ir

^۲ استادیار، گروه کامپیوتر دانشکده فنی و مهندسی، دانشگاه اصفهان، اصفهان
ahmadb@eng.ui.ac.ir

چکیده

در بازیابی اطلاعات XML ای که داده‌ها شمای ناهمگون، پیچیده و نامعلومی دارند، اغلب کاربران از پرس و جوهای مبتنی بر محتوا استفاده می‌کنند. بیشتر این پرس و جوها از اطلاعات مهمی که در ساختار XML نهفته است بهره نمی‌برند و نیاز کاربر را نیز به طور کامل بیان نمی‌کنند. از این رو و به منظور بالا بردن کیفیت بازیابی‌ها، می‌توان پرس و جوها را گسترش داد. برای گسترش پرس و جوها می‌توان از بازخورد رابطه کور استفاده کرد که اخیراً در بازیابی اطلاعات XML مورد توجه قرار گرفته است. در روش پیشنهادی این مقاله ابتدا بستر کلمات قابل افزودن به پرس و جو با استفاده از بازخورد کور شناسایی می‌شوند. سپس این بسترها برای یافتن کلمات مناسب مورد کاوش قرار می‌گیرند. بیشتر موتورهای بازخورد موجود، تنها از محتویات بازخورد کور بهره می‌برند در حالیکه روش ارائه شده در این رهاورد از اطلاعات ساختاری این بازخوردها نیز استفاده می‌کند. همچنین برای افزایش دقت بازیابی‌ها و کاهش ابهام پرس و جوها از معنای پرس و جو استفاده شده است. ارزیابی صورت گرفته، افزایش دقت بازیابی‌ها را با استفاده از روش پیشنهادی نشان می‌دهد.

کلمات کلیدی

بازیابی اطلاعات، XML، گسترش پرس و جو، بازخورد رابطه کور، المان، روابط معنایی، شبکه واژگان

کدگذاری شده‌اند. این ویژگی سبب می‌شود تا اطلاعات زمینه‌ی مستند در قالب مسیرهای XML ارائه شود. این امر برای پروسه بازیابی اهمیت بسیاری دارد. پرس و جوها در بازیابی اطلاعات XML به دو بخش عمده تقسیم می‌شوند: پرس و جوهای مبتنی بر محتوا و پرس و - جوهای مبتنی بر محتوا و ساختار. با توجه به پیچیدگی زبان پرس و - جوهای مبتنی بر محتوا و ساختار، ناهمگون بودن داده‌های XML و نیز فقدان دانش از شمای زیرین این داده‌ها برای کاربران غیر خیره، بسیاری از موتورهای جستجو از پرس و جوهای مبتنی بر محتوا استفاده می‌کنند. اما از آنجایی که اغلب این پرس و جوها نمی‌توانند نیاز کاربر را به طور کامل بیان کنند و از اطلاعات ساختاری که موجب بهبود

۱- مقدمه

امروزه XML به عنوان استاندارد جدید برای ذخیره سازی و مبادله اطلاعات در کتابخانه‌های دیجیتال، وب، اینترنت و غیره استفاده می‌شود. از این رو بازیابی اطلاعات روی داده‌های XML به شدت مورد توجه قرار گرفته‌است. از ویژگی‌های شاخص این نوع بازیابی‌ها، بازگرداندن بخشی از مستند به جای کل مستند است [1]. از تفاوت‌های این نوع سیستم‌های بازیابی با سیستم‌های بازیابی متنی رایج، در نوع اطلاعات مورد بازیابی و نوع پرس و جوها است. در بازیابی اطلاعات که روی داده‌های XML صورت می‌گیرد، داده‌ها به صورت XML

درباره محدود مقالاتی که در آنها به نوعی به ساختار XML توجه شده است، [7,8,9] گام مؤثری برای گسترش پرس و جویها برداشته اند. روش ارائه شده در این مقالات به گونه ای است که با استفاده از اطلاعات ساختاری (مسیر المان مرتبط، خصوصیات مستند شامل المان مرتبط) و محتوایی، نوع پرس و جویها از پرس و جوی مبتنی بر محتوا به پرس و جوی مبتنی بر محتوا و ساختار تغییر داده می شود.

روش های ارائه شده در این مقالات باعث افزایش دقت در بازیابی ها می شوند. در [10] پرس و جویها مبتنی بر محتوا با استفاده از بازخورد رابطه کور و اطلاعات ساختاری گسترش می یابند. این اطلاعات ساختاری از همسایگان المان هایی که شامل کلمات پرس و جوی هستند، استخراج می شود. در [11] پیشوند مشترک المان های مرتبط به پرس و جوی افزوده شده است. ضمن اینکه در این مقاله گسترش هر دو نوع پرس و جوی بحث شده است. متاسفانه هیچگونه ارزیابی برای روش پیشنهادی این مقاله صورت نگرفته است.

۳- نگاهی اجمالی به فرایند گسترش پرس و جوی

یک روش برای گسترش پرس و جوی، انتخاب کلماتی است که به مفهوم پرس و جوی نزدیک هستند. برای یافتن کلمات مناسب قابل افزودن به پرس و جوی باید بسترهای مناسبی مورد جستجو قرار گیرند. در این مقاله گسترش پرس و جوی شامل دو فاز اصلی است. در فاز اول بسترهایی که شامل کلمات کاندید برای گسترش پرس و جوی هستند انتخاب می شوند و در فاز دوم کلمات کاندید از این بسترها انتخاب می گردند.

شکل (۱) نگاهی اجمالی به فرایند گسترش پرس و جوی در روش پیشنهادی این مقاله دارد. در ادامه، گام های این فرایند مرور خواهد شد. ۱. در گام اول فرایند، المان ها طبق آنچه کاربر در پرس و جوی آورده است بازیابی می شوند.

۲. در گام دوم، المان هایی که احتمال یافتن کلمات مناسب قابل افزودن به پرس و جوی در آنها بیشتر است شناسایی می گردند. سپس المان ها طبق این معیار رتبه بندی خواهند شد.

۳. در گام سوم، کلمات کاندید وزن دهی و کلماتی با وزن بالاتر برای افزودن به پرس و جوی انتخاب می شوند.

۴. در گام آخر کلمات منتخب به پرس و جوی افزوده می شوند. طبق آنچه در بازخورد رابطه کور مطرح است X المانی که در بالای لیست بازیابی قرار گرفته اند، به پرس و جوی مرتبط تر هستند. بنابراین کافی است به جای آن که کل مستند XML برای یافتن کلمات مورد نظر جستجو شود، این المان ها و المان های شبیه به آنها مورد کاوش قرار گیرد.

بازیابی ها می شود نیز استفاده نمی کنند، نیاز به غنی ساختن این پرس و جویها احساس می شود. یکی از شیوه های مؤثر در غنی ساختن پرس و جویها، گسترش دادن آنها است. برخی از روش های گسترش پرس و جوی مبتنی بر استفاده از بازخورد رابطه کور است. در این نوع بازخوردها فرض می شود تعدادی از المان ها که در مرحله ی اول بازیابی در بالای لیست مرتب نتیجه قرار دارند، مرتبط با پرس و جوی هستند و می توان از این المان ها در غنی سازی پرس و جویها استفاده کرد.

در محدود کارهای انجام گرفته برای گسترش پرس و جوی روی داده های XML با استفاده از بازخورد، اغلب سعی بر آن بوده است که از محتویات این المان ها استفاده شود. اگر چه این کارها در بازیابی ها مؤثر بوده است، اما می توان گفت به دلیل بی توجهی به طبیعت نیم ساخت یافته متون XML ای، کاملاً از بازیابی اطلاعات متنی الهام گرفته شده است.

ایده اصلی این مقاله، گسترش پرس و جوی با استفاده از اطلاعات ساختاری موجود در درخت XML و استفاده از معنای پرس و جویها می باشد. تکنیک های مبتنی بر نزدیکی ساختاری المان ها از جمله تکنیک های مؤثر بازیابی اطلاعات روی داده های ساخت یافته هستند که منجر به بهبود دقت پرس و جوی می شوند. این تکنیک ها روی داده های نیم ساخت یافته کمتر بررسی شده اند.

۲- مروری بر کارهای گذشته

از دیرباز، تکنیک های گوناگون بازخورد رابطه در بازیابی اطلاعات متنی رایج بوده است. این در حالی است که تاکنون تحقیقات چندانی روی بازخوردهای رابطه در بازیابی اطلاعات XML صورت نگرفته است. در ابتدا Rocchio تحقیقات روی گسترش پرس و جوی را با استفاده از بازخوردهای رابطه کاربر در مدل فضا-بردار انجام داد [2]. وی کلماتی را به پرس و جوی افزود که مستندات مرتبط تر را نشان می دهند و از افزودن کلماتی که مستندات غیرمرتبط را نشان می دهند اجتناب کرد. در [3] یک تحقیق کامل از تکنیک های بازخورد رابطه که روی مدل های بازیابی مختلف صورت گرفته، ارائه شده است.

در بازیابی اطلاعات XML ای، اغلب موتورهای بازخورد تنها از محتویات المان هایی که مرتبط یا غیر مرتبط تشخیص داده شده اند استفاده می کنند، بدون آنکه به معنای ضمنی^۱ که در ساختار این المان ها وجود دارد توجه کنند [4,5]. در [6]، چارچوبی برای پالایش پرس و جوی بر اساس دانش آنتالوژی و بازخوردهای کاربر ارائه شده است. در این مقاله ابتدا وزن کلمات پرس و جوی با استفاده از آنتالوژی مفاددهی و سپس با استفاده از بازخوردهای کاربر اصلاح می شوند. هر چند استفاده از بازخورد کاربر تا حدی باعث بهبود کارایی بازیابی شده است اما عدم استفاده از پتانسیلی که در ساختار XML نهفته است از نقاط ضعف این رهیافت محسوب می شود.

۴-۱- مفهوم نزدیکی بر گرفته از ساختار

آنچه که در این رهاورد نزدیکی برگرفته از ساختار المان‌ها نامیده می‌شود، نزدیکی مفهومی است که بین المان‌ها در درخت XML به دلیل چیدمان آن‌ها وجود دارد. به عبارت دیگر آرایش نوده‌ها در درخت XML به صورت اتفاقی نیست و محل قرار گرفتن المان‌ها در XML خود دارای معنا می‌باشد. با بررسی محل قرار گرفتن المان‌ها می‌توان دانش‌های مهمی را استخراج کرد.

از جمله دانش‌های مهمی که از چیدمان المان‌ها در درخت XML استنتاج می‌شود، میزان ارتباط معنایی محتوای المان‌هاست که به فاکتورهای گوناگونی بستگی دارد. یکی از این فاکتورها فاصله بین المان‌هاست که با شمارش تعداد یال‌های بین آن‌ها محاسبه می‌گردد. هر چه فاصله بین این المان‌ها در درخت XML بیشتر باشد گواه این مطلب خواهد بود که در مستند XML نیز محل قرار گرفتن این دو المان از هم فاصله دارد. از طرفی از آن جا که در مستندات، مطالب مرتبط که مضامین مشابه و مرتبطی با هم دارند، پیوسته و نزدیک به هم قرار خواهند گرفت این نتیجه گرفته خواهد شد که مفهوم این دو المان به هم کمتر شبیه خواهند بود. برعکس این مطلب نیز صحیح است یعنی اگر فاصله المان‌ها در درخت XML کم باشد نشان‌دهنده این مطلب خواهد بود که در مستند مذکور نیز این دو المان به هم نزدیک هستند در نتیجه مضمون این دو المان به هم مرتبط خواهد بود.

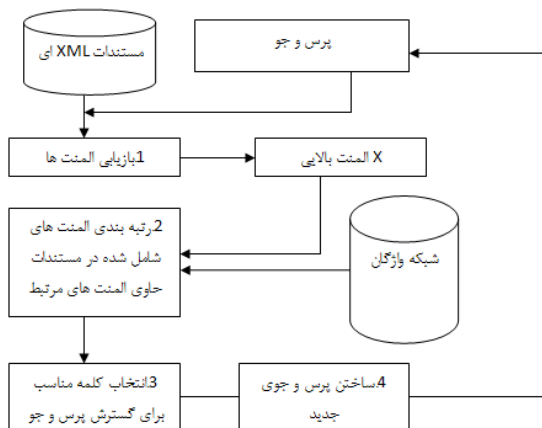
یکی دیگر از مهمترین معیارهای نزدیکی المان‌ها، کوچکترین جد مشترکشان^۲ است [11]. تحقیقات نشان داده است هر چه اختلاف ارتفاع کوچکترین جد مشترک دو المان با ارتفاع هر دو المان کمتر باشد، ارتباط معنایی آن دو بیشتر خواهد بود. عکس این مطلب نیز صحیح است یعنی هر چه اختلاف ارتفاع کوچکترین جد مشترک دو المان با ارتفاع المان‌های مورد نظر بیشتر باشد ارتباط معنایی آن دو المان کمتر است.

بر مبنای آنالیز صورت گرفته و در نظر گرفتن عوامل موثر بالا، رابطه‌ی زیر برای سنجیدن میزان نزدیکی المان‌ها از روی ساختار درخت XML به صورت زیر ارائه شده است:

(۲)

$$NS(e_i, e_j) = \alpha \times \frac{\text{depth}(lca(e_i, e_j))}{\text{depth}(e_i) + \text{depth}(e_j)}$$

ارتفاع المان e_i با $\text{depth}(e_i)$ ، ارتفاع المان e_j با $\text{depth}(e_j)$ ، کوچکترین جد مشترک با $lca(e_i, e_j)$ و فاصله بین دو المان با $\text{distance}(e_i, e_j)$ نشان داده می‌شود. α نیز ضریبی است که توسط کاربر تعیین می‌شود.



شکل ۱. فرایند گسترش پرس و جو

۴-۲ مفهوم شباهت در XML

معمولاً در داده‌های XML ای شباهت المان‌ها با استفاده از اطلاعات ساختاری برآورد می‌شود. منظور از اطلاعات ساختاری، اطلاعاتی است که از روی ساختار درخت XML بدست می‌آید. به عنوان مثال تعداد زیرشاخه‌های یک المان خاص، مسیر رسیدن به یک المان از ریشه، تعداد فرزندان و اجداد یک المان و اطلاعاتی از این دست اطلاعات ساختاری هستند که می‌توان تا حدودی از آن‌ها برای تشخیص المان‌هایی که مضامین مرتبط و مشابهی دارند استفاده کرد.

اما وقتی سخن از بازیابی اطلاعات به میان می‌آید، انتظار می‌رود که برای تشخیص المان‌های مشابه و مرتبط علاوه بر اطلاعات ساختاری به محتویات المان‌ها نیز توجه شود. یعنی همانطور که از روی اطلاعات برگرفته از درخت XML می‌توان المان‌های هم‌زمینه که مضامین شبیه و مرتبطی دارند را تشخیص داد، از روی محتویات این المان‌ها نیز می‌توان دقت سنجش این تشابه را افزود.

از این رو و با توجه به توضیحات بالا تعریفی از نزدیکی المان‌ها ارائه خواهد شد که میزان نزدیکی المان‌ها را وابسته به دو عامل می‌داند. این دو عامل عبارتند از: نزدیکی‌های برگرفته از ساختار و نزدیکی برگرفته از محتوای المان‌ها.

فرمول ۱ همین تعریف را بیان می‌کند.

(۱)

$$N(e_i, e_j) = \gamma \times NS(e_i, e_j) + (1 - \gamma) \times NC(e_i, e_j)$$

$NS(e_i, e_j)$ نزدیکی ساختاری، $NC(e_i, e_j)$ نزدیکی محتوایی و e_i المان i را نشان می‌دهد. در ادامه در این خصوص توضیحات بیشتری ارائه خواهد شد.

اما از آنجا که این روش تنها به نحو کلمات توجه می کند روش کاملی نیست. برای روشن تر شدن موضوع به مثال زیر توجه کنید:

در یک سیستم هوشمند انتظار می رود تا اگر محتویات متنی U_i شامل کلمات "clever men" و محتویات متنی U_j شامل کلمات "Bright person" است، شباهتی بین این دو المان در نظر گرفته شود. اما از آنجایی که بردار هر المان بر اساس حضور و یا عدم حضور کلمات شاخص ساخته می شود، شباهت بین این دو متن صفر در نظر گرفته می شود. به همین منظور و به سبب غنی ساختن بردار شباهت های متنی دو المان، در این رهاورد شباهت معنایی جاری بین دو المان نیز بررسی می شود. در ادامه در این خصوص توضیحات بیشتری را ارائه خواهد شد.

۴-۲-۲- مفهوم شباهت معنایی

شباهت معنایی معیاری که میزان نزدیکی معنایی دو کلمه را در نشان می دهد. در این مقاله از شبکه واژگان برای استخراج روابط معنایی بین کلمات موجود در محتویات المان ها استفاده شده است. شبکه واژگان یک فرهنگ واژه ای الکترونیکی است که در آن معنای واژه ها از طریق روابط بین آنها بیان می شود. روابط معنایی که در این مقاله مورد بررسی قرار گرفته است رابطه ابرمعنایی، زیرمعنایی می باشد. کوتاهترین فاصله وزن دار کلمات المان در شبکه واژگان می تواند معیاری برای نزدیکی معنایی کلمات المان و در نتیجه خود المان ها باشد. در این مقاله وزن رابطه هم معنایی ۱ و وزن رابطه ابر معنایی و زیر معنایی ۰٫۸۵ در نظر گرفته شده است.

(۶)

$$Nsem(W_i, W_j) = \frac{1}{distance(W_i, W_j)}$$

$distance(W_i, W_j)$ کوتاهترین فاصله وزن دار دو کلمه W_i و W_j در شبکه واژگان را نشان می دهد که باید از حد آستانه ای که توسط کاربر تعیین می گردد کمتر باشد (در این مقاله حد آستانه پنج در نظر گرفته شده است).

(۷)

$$Nsem(e_i, e_j) = \sum_{i=1}^m \sum_{j=1}^n Nsem(W_i, W_j)$$

m و n تعداد کلمات المان ها را نشان می دهد.

۵- انتخاب کلمات مناسب

پس از آنکه المان های مشابه با المان های مرتبط تشخیص داده شدند، نوبت به انتخاب کلمه از میان محتوای متنی آن ها می رسد. جستجو در محتویات این المان ها سبب می شود تا محتویات همه المان ها برای

۴-۲- مفهوم نزدیکی بر گرفته از محتوا

آنچه که در این رهاورد نزدیکی برگرفته از محتوا نامیده می شود، میزان شباهتی است که بین المان ها به دلیل محتوای آنها وجود دارد. در این تحقیق نزدیکی برگرفته از محتوای المان ها تلفیقی از شباهت های معنایی و لغوی المان ها در نظر گرفته شد. برآورد این نوع شباهت از برآورد شباهت میان مستندات متنی الهام گرفته شده است

(۳)

$$NC(e_i, e_j) = Nsyn(e_i, e_j) \times Nsem(e_i, e_j)$$

۴-۲-۱- مفهوم شباهت لغوی

در این بخش میزان شباهت موجود بین محتویات المان ها از نظر لغوی محک خورده است. به عبارت دیگر در این بخش برای تعیین میزان شباهت محتوایی متنی المان ها تنها روی نحو کلمات تمرکز شده و توجهی به معنای کلمات نشده است.

شباهت کسینوسی معیار مناسبی برای سنجش شباهت بین دو بردار n بعدی است که با محاسبه زاویه بین دو بردار این مهم را انجام می دهد. در این روش ابتدا محتویات متنی مستندات به روش بردار ذخیره خواهد شد. اجزای بردارها، کلمات شاخص برگرفته از محتویات المان ها هستند که طبق فرمول های وزن دهی، وزن دهی خواهند شد. در این رهاورد، کلمات شاخص طبق دو معیار آماری (چگالی کلمه در متن و فرکانس کلمه) و با استفاده از فرمول OKAPI گسترش یافته [13] وزن دهی شده اند.

(۴)

$$Relevance(W_i, U_i) =$$

$$\left(0.3 * \frac{tf}{\left(0.5 + 1.5 * \frac{L_i}{L} \right) + tf} * \log \frac{N - n_j + 0.5}{n_j + 0.9} \right)$$

$$\bar{U}_i = (Relevance(W_1, U_i), Relevance(W_2, U_i), \dots, Relevance(W_j, U_i))$$

U_i محتویات متنی المان ها، tf فرکانس کلمه W_i در U_i ، L_i طول U_i ، n_j تعداد u های شامل کلمه W_j ، N تعداد کل U ها در مجموعه، $\bar{L} = \frac{\sum_{i=1}^N L_i}{N}$ میانگین طول U است.

بعد از محاسبه وزن ها و ساختن بردار هر متن، نوبت به محاسبه شباهت متن ها می رسد.

(۵)

$$Nsyn(e_i, e_j) = Sim_{synthetic} = \frac{\bar{U}_i \cdot \bar{U}_j}{\|\bar{U}_i\| * \|\bar{U}_j\|}$$

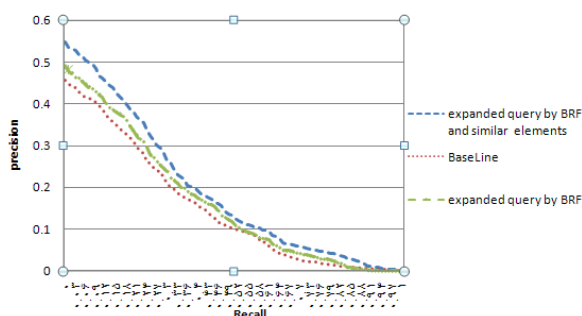
این روش اگر چه در حال حاضر یکی از ساده ترین و در عین حال پر کاربردترین روش های موجود برای ارزیابی شباهت مستندات متنی است.

۶- پیاده سازی و ارزیابی

در این پژوهش، چارچوب پیشنهادی در محیط جاوا و با استفاده از پایگاه داده‌های Sql server روی ماشینی با ۵۱۲ مگابایت RAM و ۱.۶ مگا هرتز Cpu پیاده سازی شده است. همچنین برای ارزیابی روش پیشنهادی از محک INEX ۲۰۰۹ [12] استفاده شده است که شامل مجموعه‌ای از مستندات XML و یک مجموعه از پرس‌وجوهای NEXI به همراه مجموعه ارزیابی نتایج برای هر پرس‌وجو می‌باشد. برای ارزیابی مستندات XML ای از موتور جستجوی Topx استفاده شده است. معیار ارزیابی روش پیشنهادی این مقاله بر مبنای دقت فراخوانی می‌باشد. همانطور که در شکل ۲ مشاهده می‌شود اگر تنها المنت‌های بازخورد کور به‌عنوان بستر کلمات گسترش مورد کاوش قرار گیرند دقت ارزیابی‌ها نسبت به مبنا افزایش می‌یابد اما اگر از المنت‌های مشابه المنت‌های بازخورد کور نیز به‌عنوان بستر کلمات گسترش استفاده شود دقت ارزیابی‌ها نسبت به حالت اول و مبنا بهبود بیشتری می‌یابد.

۷- نتیجه گیری

این مقاله با استفاده مناسب از اطلاعات ساختاری موجود در درخت XML گام موثری در استفاده از بازخورد رابطه کور در گسترش پرس‌وجوها برداشته است. همچنین با کم‌اهمیت کردن کلمات چند معنا در فرایند وزن‌دهی کلمات کاندید گسترش موجب کاهش ابهام در پرس‌وجوهای جدید شده است. نتایج ارزیابی‌ها افزایش دقت ارزیابی‌ها را نشان می‌دهد.



شکل ۲: نمودار دقت-فراخوانی

مراجع

- [1] S. R. Joty and S.Sadid-AlHasan, "Advances in Focused Retrieval: A General Review", 10th international conference on Computer and information technology, 2008. iccit 2007.
- [2] B.Ribeiro, R.Baeza and." Modern Information Retrieval". NewYork: ACM Press, 1999,pp.117-120.

انتخاب کلمات مناسب مورد کاوش قرار نگیرد که این کار سبب صرفه جویی در زمان و منابع می‌گردد.

به منظور انتخاب کلمه مناسب، تمام کلمات موجود در محتویات المان‌های مورد بررسی طبق رابطه (۱۰) وزن‌دهی شده و نهایتاً تعدادی از کلمات که وزن بالاتری دارند برای اضافه شدن به پرس‌وجو انتخاب می‌شوند. رابطه (۱۰) به سه عامل مهم وابسته می‌باشد:

(۱) اهمیت کلمه در محتوای متنی المان که طبق اطلاعات آماری محاسبه می‌گردد.

(۲) میزان ارتباط المان حاوی کلمه مورد نظر نسبت به پرس‌وجوی اولیه

(۸)

$$N(e_j, q) = \frac{\sum_i^m N_s(e_j, e_i)}{m}$$

q پرس‌وجوی اولیه، m تعداد کلمات پرس‌وجوی اولیه بعد از حذف کلمات توقف، e_i المان شامل i امین کلمه پرس‌وجو و $N(e_j, q)$ درجه ارتباط المان e_j با پرس‌وجوی q را نشان می‌دهد.

ذکر این نکته ضروری است که چنانچه چندین المان شامل i امین کلمه‌ی پرس‌وجو باشند از المانی به عنوان e_i در رابطه بالا استفاده می‌گردد که نزدیکترین فاصله را نسبت به المان مورد مقایسه داشته باشد.

(۳) کمیابی معنایی کلمه که در رهاورد با درجه چندمعنایی^۳

کلمه تعریف می‌شود. از آنجایی که افزودن کلماتی به پرس‌وجو که دارای چندین معنا هستند باعث انحراف در ارزیابی‌ها می‌شود، این عامل مورد بررسی قرار گرفته است. رابطه زیر برای تعریف کمیابی معنایی آمده است

(۹)

$$rarity(w_i) = \log \left(\frac{Max-Polysemy}{|senses(w_i)|} \right)$$

Max-Polysemy عددی ثابت است که مقدار آن بیشتر از تعداد معنای کلمه‌ای است که بیشترین معنی را در شبکه واژگان دارد. $|senses(w_i)|$ تعداد معنای کلمه w_i را نشان می‌دهد.

(۱۰)

$$weight(w_i) = Relevance(w_i, u_j) * rarity(w_i) * N(e_j, q)$$

- [3] I.Ruthven and M. Lalmas, "A survey on the use of relevance feedback for information access systems", Knowledge Engineering Review, vol. 18, no.2, pp. 95 - 145 , june 2003.
- [4] M.Mandelbrod, Y.Mass and. "Relevance Feedback for XML Retrieval." In INEX 2004 Workshop, 2004. pp.154-157.
- [5] C.J.Crouch, A.Bellamkonda and A.Mahajan. "Flexible XML retrieval based on the extended vector model." INEX 2004 Workshop, 2004. pp.149-153.
- [6] H.Pan, "Relevance feedback in XML Retrieval." , Ph.D. thesis, Max-Planck-Institut Informatik, Germany .2003.
- [7] M.Theobald, R.Schenkel and. "Relevance Feedback for Structural Query Expansion," in Advances in XML Information Retrieval and Evaluation , Vol. 3977/2006, Heidelberg: Springer, 2006, pp.344-357.
- [8] M.Theobald, R.Schenkel and. "Structural Feedback for Keyword-Based XML Retrieval," in Advances in Information Retrieval. Vol. 3936/2006, Heidelberg: Springer, 2006,pp. 326-337.
- [9] M.Theobald, R.Schenkel and,"Feedback-driven structural query expansion for ranked retrieval of XML data," In 10th International Conference on Extending Database Technologies (EDBT 2006), Munich, Germany, Mar. 2006.
- [10] W. Hsu, M. L. Lee and X. Wu. "Path-augmented keyword search for XML documents," ICTAI 2004, 2004 , pp.526-530.
- [11] M.Boughanem, L.Hlaoua and ,"Towards Contextual and Structural Relevance Feedback in XML Retrieval,"In Workshop on Open Source Web Information Retrieval(OSWIR). 2005.
- [12] Otago university ,.inex., [Online]. vailable :<http://www.inex.otago.ac.nz/> [access:octobr,20, 2008]
- [13] Robertson, S., Walker, S., Jones, S., HancockBeaulieu,M. and Gatford, M. OKAPI at TREC3. 3rd Text Retrieval Conference, 1995

زیر نویس ها

¹ Implicit Semantic

² Lowest Common Ancestor

³ Polysemy