

بهبود کارایی یک معماری موتور جستجوی معنایی با استفاده از تشابه معنایی و در نظر گرفتن نزدیکی مفهومی

زهرا جعفری^۱، حمید رستگاری^{۲*}

۱- دانشکده مهندسی کامپیوتر، واحد نجف آباد، دانشگاه آزاد اسلامی، نجف آباد، ایران

۲- دانشکده مهندسی کامپیوتر، واحد نجف آباد، دانشگاه آزاد اسلامی، نجف آباد، ایران.

خلاصه

با پیشرفت وب و گسترش مطالب موجود در آن و استفاده از آن اهمیت موتورهای جستجو روز به روز بیشتر می شود. از طرف دیگر سطح انتظارات و نیازهای کاربران نیز روز به روز در حال گسترش است و کاربران تمایل به استفاده از سرویس هایی دارند که نه تنها درخواست آنها را بررسی کرده و پاسخ دهند، بلکه درخواست آنها را درک و تجزیه و تحلیل کرده و مرتبط ترین موارد را یافت کنند. به همین دلیل بحث استفاده از موتورهای جستجوی معنایی روز به روز پر اهمیت تر می شود. در این مقاله در جهت بهبود کارایی بازیابی اسناد مرتبط با نیازمندی های کاربر روشی پیشنهاد شده است که بر اساس آن علاوه بر در نظر گرفتن کلمات کلیدی و مفاهیم بر اساس اطلاعات زبان شناسی و تشخیص حوزه سند، تشابه معنایی کلمات (کلمات کلیدی مشابه) و نزدیکی مفهومی با استفاده از اطلاعات RDF در نظر گرفته شده است. روش پیشنهادی بر اساس چهارچوب معماری جدید SBIRS است و مؤلفه های این معماری با اضافه شدن دیتاست RDF و مؤلفه تشابه معنایی در راستای بهبود آن در نظر گرفته شده است. برای نشان دادن نزدیکی مفهومی از اطلاعات دیتاست RDF و اسناد استخراج شده مرتبط با حوزه آن و همچنین برای تشابه معنایی کلمات نیز از اطلاعات اسناد استفاده شده است. در تشابه معنایی از ماشین بردار پشتیبان برای یادگیری کلمات مشابه و نامشابه استفاده شده است. آزمایش ها روی قسمتی از اسناد استخراج شده مرتبط با حوزه RDF و جمع آوری شده از اینترنت انجام شده است تا بتوان از کارایی روش پیشنهاد شده اطمینان حاصل کرد و تحلیل درستی از نتایج ارائه داد. آزمایشات نشان می دهد که با استفاده از روش پیشنهادی به دلیل در نظر گرفتن بهتر ارتباط های معنایی و مفهومی بین کلمات در مقایسه با معماری پایه SBIRS معیارهای فراخوانی و معیار F بررسی گردیده و این نتیجه حاصل شد که فراخوانی و معیار F افزایش بیشتری داشته است.

واژه های کلیدی: وب معنایی، جستجو معنایی، استخراج مفاهیم، شباهت معنایی، RDF، شباهت مفهومی

*Corresponding author Islamic Azad University, Najaf Abad Branch, Faculty of Computer Engineering, Zahra Jafari, 09139591368
Email: taranom1225@gmail.com

۱. مقدمه

وجود موتورهای جستجویی که بتوانند بهترین پاسخ را به درخواست کاربر بدهند امری ضروری است با توجه به اینکه تعداد صفحات وب موجود در اینترنت دائماً در حال افزایش می‌باشد و پیدا کردن اسناد مناسب برای کاربران مشکل است. این واقعیت که امروزه نیاز به استفاده از موتورهای جستجوی پیشرفته برای یافتن صفحات وب در بین حجم عظیم صفحات وجود دارد واقعیتی انکارناشدنی است. از طرفی بنا به دلایل متعدد موتورهای جستجوی سنتی نمی‌توانند همواره پاسخ‌های بهینه را ارائه کنند. از آن جمله می‌توان به عدم استفاده کاربر از کلید واژه‌های صحیح در بسیاری از جستجوها و یا عدم تطابق کلمات اشاره کرد. از این رو ضرورت استفاده از موتور جستجوی معنایی که با تجزیه و تحلیل درخواست کاربر و صفحات بتواند مرتبط‌ترین صفحات را ارائه کند، محسوس است.

موتورهای جستجوی معنایی یکی از مواردی هستند که از لغت‌شناسی و امکانات جدیدی که وب در اختیار قرار داده است، استفاده می‌کنند. جستجوی کارآمد، یکی از کاربردهای پیش روی نسل بعدی وب است که تحت عنوان وب معنایی شناخته می‌شود [۱].

روش‌های موجود جستجوی وب، اسناد را به عنوان یک مجموعه کلمه می‌بینند و تنها به تکرار و وجود آنها توجه می‌کنند. موتور جستجویی پیش رو دیگر محدود به تطبیق کلمات کلیدی بین درخواست کاربر و اسناد نیستند و سعی می‌کنند با استفاده از اطلاعات و آنالیز بیشتر جوابی دقیق تر و ساخت یافته تر به دست آورند. در این صورت دیگر تطبیق دقیق درخواست کاربر و کلمات کلیدی نیاز نبوده و موتور جستجو با تحلیل منظور کاربر، خود بهترین جوابها را پیشنهاد خواهد داد [۲]. می‌توان در نظر گرفت که هر صفحه وب حاوی اطلاعات مفهومی است که می‌تواند برای بررسی بیشتر خود صفحه در نظر گرفته شود. این اطلاعات مفهومی که معمولاً تحت عنوان لغت‌شناسی شناخته می‌شوند می‌توانند در تحلیل صفحات وب کمک بسزایی انجام دهند. به عنوان مثال در تشخیص مرتبط بودن موضوع یک صفحه وب با درخواست کاربر می‌توان از این اطلاعات استفاده کرد. ایده اصلی این مدل این است که در خصوص مربوط بودن یک مدل با توجه به مدل فکری کاربر که از اطلاعات به دست آمده است استفاده شود. در این صورت سندی که دقیقاً منطبق بر درخواست کاربر است به عنوان سند هسته در نظر گرفته می‌شود و با بسط مفاهیم و بررسی‌های بیشتر سایر اسناد نیز استخراج می‌شوند [۳].

اما روشی که در اینجا استفاده شده علاوه بر در نظر گرفتن کلمات کلیدی و مفاهیم بر اساس اطلاعات زبان‌شناسی و تشخیص حوزه سند، تشابه معنایی کلمات و نزدیکی مفهومی با استفاده از اطلاعات RDF را در نظر گرفته است. چهارچوب روش بر اساس معماری جدید SBIRS است دیتاست RDF استاندارد و مؤلفه تشابه معنایی در راستای بهبود آن گرفته شده است. برای استخراج نزدیکی مفهومی از اطلاعات دیتاست RDF استفاده شده است و در تشابه معنایی از ماشین بردار پشتیبان برای یادگیری کلمات مشابه و نامشابه استفاده شده است.

۲. معماری SBIRS

معماری SBIRS برای بازیابی اطلاعات از وب معنایی از نمایش‌های مفهومی محتوا در راستای کلیدواژه‌ها به عنوان پایگاه دانش استفاده می‌کند و نمایش‌های مفهومی نیازهای کاربر را ارائه می‌دهد [۴]. این معماری نمایش‌های مفهومی محتوا، بسط‌های پرس‌وجو، مطابقت معنایی، استخراج نتایج مرتبط بر اساس ارتباط مفاهیم را با کمک مؤلفه‌های پیمایشگر، پیش‌پردازشگر، مفسر معنایی، شاخص‌کننده معنایی، مبدل پرس‌وجوی معنایی، بازیاب‌کننده محتوای معنایی و رتبه‌دهنده معنایی انجام می‌دهد. این مؤلفه‌ها در لایه‌های مختلف معماری گروه‌بندی شده‌اند. لایه داده فیزیکی پایگاه داده وب را با مؤلفه‌های پیمایشگر و پیش‌پردازشگر ایجاد می‌کند. لایه مفسر معنایی پایگاه دانش را با مفسر معنایی و اندیس

معنایی ایجاد می کند. لایه مطابقت معنایی وظیفه مطابقت بین محتوای معنایی و جستجوی معنایی را به عهده دارد. لایه بازیابی، نتایج بازیابی شده را رتبه بندی می کند و به لایه کاربرد تحویل می دهد.

۳. شباهت معنایی کلمات

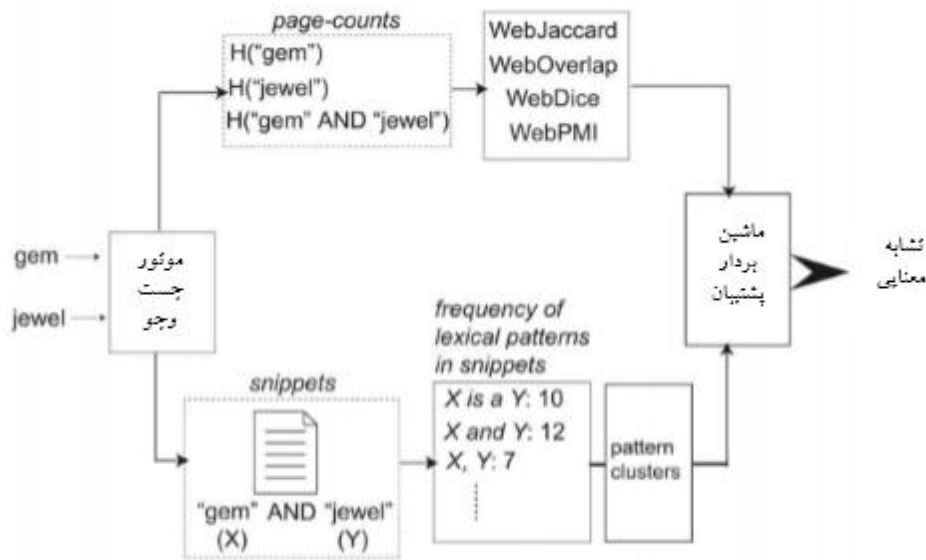
دو کلمه P و Q را در نظر بگیرید. ما مسئله اندازه گیری شباهت معنایی بین P و Q را به صورت ساختن یک تابع $\text{sim}(P, Q)$ در نظر گرفته ایم، که مقداری را در رنج $[0, 1]$ برمی گرداند. اگر P و Q به شدت مشابه باشند (مثلاً مترادف باشند)، ما انتظار داریم $\text{sim}(P, Q)$ نزدیک به ۱ باشد. از طرف دیگر، اگر P و Q از نظر معنایی شبیه نباشند، پس انتظار داریم $\text{sim}(P, Q)$ به صفر نزدیک باشد. ما ویژگی های مختلفی که شباهت بین P و Q را بیان می دارند، با استفاده از شمارش صفحه و اسنیپت بازیابی شده از یک موتور جستجو برای دو کلمه داده شده تعریف می کنیم. با استفاده از این نوع نمایش ویژگی کلمات، ما یک ماشین بردار پشتیبانی دو-دسته ای* را جهت دسته بندی جفت واژه های مترادف و غیر مترادف آموزش می دهیم. سپس تابع $\text{sim}(P, Q)$ با امتیاز اعتماد SVM آموزش داده شده تخمین زده می شود [۵].

شکل ۱ مثالی را جهت محاسبه شباهت معنایی بین دو کلمه gem و jewel نشان می دهد. ابتدا پرس و جوی خود را به موتور جستجو صادر می کنیم و شمارش صفحه را برای دو کلمه و برای عطف آن ها «gem»، «jewel» و «gem AND jewel» بازیابی می کنیم. در بخش بعدی ما چهار امتیاز شباهت با استفاده از شمارش صفحه تعریف می کنیم. امتیازهای شباهت مبتنی بر شمارش صفحه حضور همزمان کلی دو کلمه را در وب در نظر می گیرند. باین وجود، آن ها متن محلی را که کلمات در آن به طور همزمان ظاهر شده اند، مدنظر قرار نمی دهند. از طرف دیگر، اسنیپت های برگردانده شده توسط یک موتور جستجو متن محلی که دو کلمه در آن ظاهر شده اند را ارائه می دهند. در نتیجه، ما میزان تکرار الگوهای معنایی لغوی متعددی را در اسنیپت ها برای پرس و جوی عطف دو کلمه پیدا می کنیم.

الگوهای لغوی به صورت خودکار استخراج می شوند. باین وجود، لازم است به این نکته اشاره کنیم که یک رابطه معنایی می تواند با استفاده از چند الگوی لغوی بیان شود. گروه بندی الگوهای لغوی مختلف که رابطه معنایی یکسانی را حمل می کنند، ما را قادر به نمایش یک رابطه معنایی دقیق بین دو کلمه می سازد. آن ها برای این منظور، یک الگوریتم خوشه بندی الگوی متوالی ارائه نموده اند. هم امتیازات شباهت مبتنی بر شمارش و هم خوشه های الگوی لغوی جهت تعریف ویژگی های مختلفی که رابطه بین دو کلمه را نمایش می دهند، مورد استفاده واقع شده اند. با استفاده از نمایش ویژگی جفت کلمات، یک ماشین بردار پشتیبانی دو کلاسه را آموزش داده اند.

* Two-class

شکل ۱- نمایش کلی روش شباهت معنایی



۱.۳. معیارهای حضور همزمان مبتنی بر شمارش صفحه

شمارش صفحه برای پرس‌وجوی P AND Q را می‌توان به صورت یک تقریب از وقوع همزمان دو کلمه (دو عبارت چندکلمه‌ای) P و Q در وب در نظر گرفت. با این وجود، شمارش صفحه برای پرس‌وجوی P AND Q به تنهایی شباهت معنایی را به صورت دقیق بیان نمی‌دارد. برای مثال، گوگل ۱۱,۳۰۰,۰۰۰ کلمه را به عنوان شمارش صفحه برای پرس‌وجوی "car" AND "apple" برمی‌گرداند. در حالی که برای پرس‌وجوی "car" AND "apple" این مقدار برابر با ۴۹۰۰۰۰۰۰ است. با وجود اینکه واژه automobile از نظر معنایی شباهت بیشتری به واژه car دارد، اما شمارش صفحه‌ای که برای پرس‌وجوی "car" AND "apple" برگردانده شده است، تقریباً ۴ برابر بیشتر از شمارش صفحه پرس‌وجوی "automobile" AND "car" است. باید شمارش صفحه را نه تنها برای پرس‌وجوی P AND Q ، بلکه برای Q و P به صورت جداگانه نیز در نظر گرفت تا بتوان شباهت معنایی بین P و Q را ارزیابی کرد.

بر این اساس چهار معیار وقوع همزمان رایج؛ جاکارد*، همپوشانی†، دایس‡ و اطلاعات دوجانبه نقطه به نقطه§ را جهت محاسبه شباهت معنایی با استفاده از شمارش صفحه در نظر گرفته شده است. در ادامه از $H(P)$ برای نشان دادن شمارش صفحه برای پرس‌وجوی P در یک موتور جستجو استفاده شده است. ضریب جاکارد و** بین دو کلمه (یا دو عبارت چندکلمه‌ای) P و Q ، به صورت زیر تعریف می‌شود:

$$WebJaccard(P, Q) = \begin{cases} 0 & \text{if } H(P \cap Q) \leq C \\ \frac{H(P \cap Q)}{H(P) + H(Q) - H(P \cap Q)} & \text{Otherwise} \end{cases} \quad (1)$$

در رابطه بالا، $P \cap Q$ پرس‌وجوی عطف P AND Q را نشان می‌دهد. با در نظر گرفتن نویز و مقیاس در داده‌های وب، ممکن است دو کلمه با وجود اینکه رابطه‌ای ندارند در برخی صفحات همزمان حضور پیدا کنند. جهت کاهش تأثیرات منفی

*

† Overlap

‡ Dice

§ Pointwise mutual information (PMI)

** WebJaccard

چنین حضورهای همزمانی، در حالی که شمارش صفحه برای پرس و جوی $P \cap Q$ کمتر از یک آستانه c باشد، ضریب WebJaccard را برابر با صفر قرار داده‌ایم. به‌طور مشابه، WebOverlap به‌صورت زیر تعریف شده است:

$$Web\ Overlap(P, Q) = \begin{cases} 0 & \text{if } H(P \cap Q) < c \\ \frac{H(P \cap Q)}{\min(H(P), H(Q))} & \text{Otherwise} \end{cases}$$

WebOverlap یک تغییر طبیعی در ضریب همپوشانی است. ما ضریب WebDice را به‌صورت گونه‌ای از ضریب Dice تعریف کرده‌ایم. WebDice به‌صورت زیر تعریف شده است:

$$Web\ Dice(P, Q) = \begin{cases} 0 & \text{if } H(P \cap Q) < c \\ \frac{2 * H(P \cap Q)}{H(P) + H(Q)} & \text{Otherwise} \end{cases} \quad (3)$$

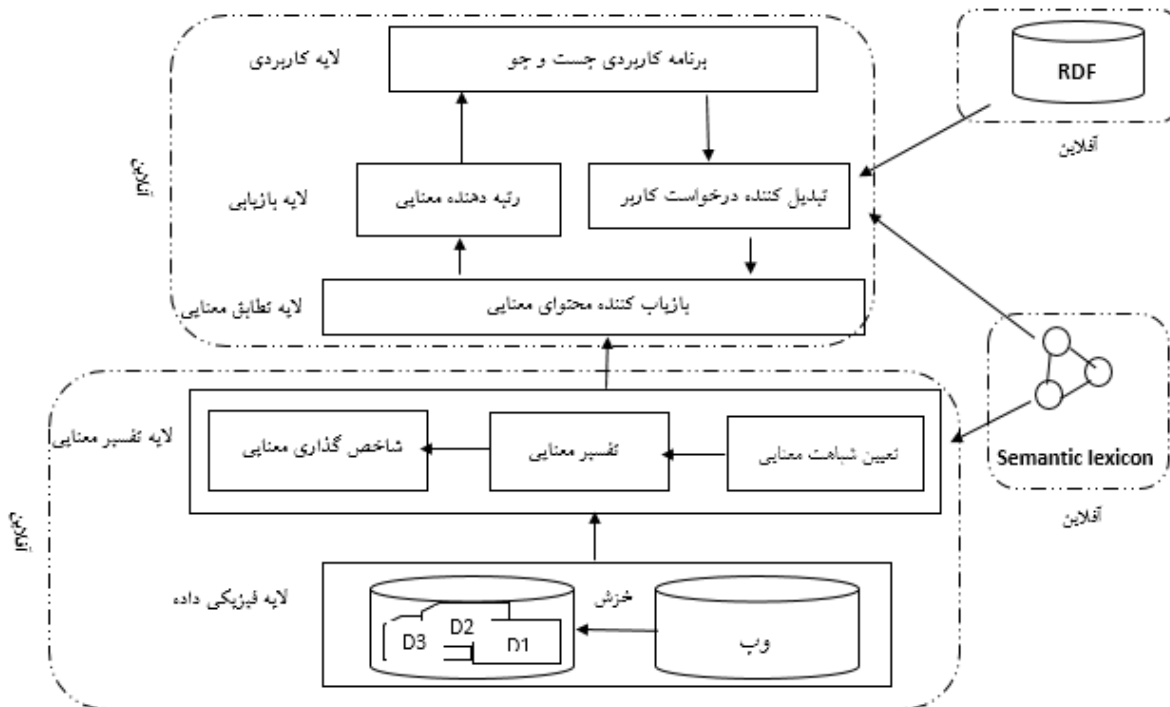
اطلاعات دوجانبه نقطه‌به‌نقطه معیاری است که از تئوری اطلاعات نشأت گرفته است؛ جهت بازتاب وابستگی بین دو رخداد احتمالی ارائه شده است. ما WebPMI را به‌صورت گونه‌ای از معیار اطلاعات دوجانبه نقطه‌به‌نقطه با استفاده از شمارش صفحه تعریف کرده‌ایم:

$$Web\ PMI(P < Q) = \begin{cases} 0 & \text{if } H(P \cap Q) < c \\ \log_2 \left(\frac{H(P \cap Q) \cdot N}{H(P) \cdot H(Q)} \right) & \text{otherwise} \end{cases} \quad (4)$$

۴. معماری موتور جستجوی بهبود داده شده

معماری بهبود داده شده سعی در اکتشاف حوزه اسناد و درخواست کاربر با استفاده از لغتنامه معنایی، استخراج نزدیکی معنایی کلمات از خود کلمات اسناد وب و نزدیکی مفهومی کلمات با استفاده از RDF را دارد. در این معماری پس از کشف نزدیکی معنایی و مفهومی در فرآیند جستجو، کلمات کلیدی استخراج شده را نیز مدنظر قرار می‌دهد تا بتواند دقت و کارایی جستجو را بهبود ببخشد. به‌علاوه با در نظر گرفتن نزدیکی حوزه اسناد و درخواست کاربر سعی می‌شود که اسنادی که به حوزه مدنظر کاربر نزدیک‌تر هستند در رتبه بالاتری نمایش داده شوند. در ادامه، ابتدا به توضیح موتور جستجوی مورد استفاده و معماری آن پرداخته شده است پس از آن روش اکتشاف نزدیکی معنایی کلمات و نزدیکی مفهومی کلمات توضیح داده شده و شیوه به‌کارگیری آن در معماری پیشنهادی توضیح داده شده است. طرح کلی روش پیشنهادی در شکل ۲ نشان داده شده است.

شکل ۲- طرح کلی معماری بهبود داده شده



۴.۱. دریافت اطلاعات

اولین مرحله در هر موتور جستجو، دریافت اطلاعات صفحات وب است. برای دریافت اطلاعات صفحات وب دو راهکار وجود دارد: راه اول پیاده سازی یک نرم افزار پیمایشگر* برای کاوش وب و استخراج اطلاعات هر یک از صفحات وب و راه دوم استفاده از دیتاست های موجود که قبلاً استخراج و پالایش شده هستند، است. به دلیل وجود دیتاست ها، در این پایان نامه نرم افزار پیمایشگر پیاده سازی نشده است و از یک دیتاست RDF استاندارد به همراه اسناد استخراج شده مرتبط با آن استفاده شده است.

۴.۲. پیش پردازش

در این مرحله اطلاعات اسناد پردازش می شوند و داده های اضافی از آنها حذف می شوند. از جمله مواردی که در این مرحله انجام می شوند شامل:

حذف کلمات اضافی مانند am, is, are, has, have, ...

ریشه یابی کلمات مانند حذف ing یا s یا es از انتهای کلمات

* Crawler

۴.۳. تحلیل مفهومی و استخراج کلمات کلیدی

پس از استخراج اسناد نوبت به آنالیز مفهومی و استخراج کلمات کلیدی هر یک از آن ها می رسد. کلمات کلیدی با شناسایی کلمات منحصر به فردی که تکرارشان زیاد است، مشخص می شوند. به عبارت دیگر هر کلمه ای که تعداد تکرارش در سند از حد آستانه ای بیشتر باشد به عنوان کلمه کلیدی آن سند در نظر گرفته می شود. برای آنالیز مفهومی از دو نوع از داده ها استفاده می شود:

لغتنامه معنایی*

قالب توضیح منابع†

لغتنامه معنایی به کلماتی گفته می شود که می توانند حوزه یک متن را مشخص کنند. به عنوان مثال در مورد حوزه بانکداری می توان به چک، سفته، وام، شماره حساب، موجودی حساب اشاره کرد. به این ترتیب در صورتی که یکی از این کلمات در متن یک سند و یا درخواست کاربر موجود باشد می توان حوزه مدنظر را مشخص کرد و به این ترتیب می توان با داشتن این اطلاعات به پالایش بهینه نتایج قابل نمایش برای کاربر پرداخت.

در اینجا برای بهبود موتور جستجو از آنتولوژی استفاده نموده ایم. روش های مختلفی برای آنالیز مفهومی با استفاده از آنتولوژی وجود دارند. در این خصوص کارهای وسیعی انجام شده است که به برخی از آن ها اشاره شده است. در مورد آنتولوژی کارهایی که انجام شده اغلب بر پایه دایره لغت‡ است.

در کاربردی جدیدتر استفاده از لغتنامه معنایی پیشنهاد شده است. به طوری که واژگانی وجود دارند که در صورت استفاده در یک متن خاص می توانند زمینه متن را با دقت خوبی نشان دهند. به عنوان مثال کلمات «وام، چک، سفته، موجودی» می توانند نشان دهند که متن مورد نظر مربوط به موارد بانکی است. به این ترتیب با داشتن این دایره واژگان معنایی می توان ارتباط مفهومی بین درخواست کاربر و اسناد را به طور دقیق تری در نظر گرفت. برای این منظور در این پایان نامه از لغتنامه معنایی مرتبط با اسناد استفاده شده است و در صورتی که درخواست کاربر و اسناد کلمه خاصی از این واژگان را مشترک داشته باشند، با یک وزن بیشتر در امتیازدهی آن سند در بالاتر نشان داده می شود. به این دلیل که تشخیص داده می شود که این سند در حوزه مدنظر کاربر قرار دارد. توجه شود که لزوماً نباید کلمه درخواست کاربر و کلمه موجود در سند یکی باشند، بلکه مربوط به یک حوزه باشند کفایت می کند. به این ترتیب که ممکن است که کاربر کلمه «وام» را در درخواست خود وارد کرده باشد و سند دارای کلمه «سفته» باشد. در این صورت نیز تشخیص داده می شود که حوزه یکسان است و سند مذکور در رتبه بندی بالاتر نشان داده می شود.

از سوی دیگر قالب توضیح منابع هم اطلاعات آنالیز جملات را در خود دارد که با استفاده از آن می توان اطلاعات بیشتر و دقیق تری از سند را در دست داشت.

۴.۴. استخراج شباهت معنایی کلمات

برای استخراج شباهت معنایی کلمات، روش پیشنهادی به این صورت کار می کند که از آنجاکه بررسی دوبه دوی ارتباط تمام کلمات با یکدیگر آن هم در حجم داده های وب بسیار زمان بر است، ما سعی کردیم که با ارائه روشی تنها ارتباط دوبه دوی قسمتی از کلمات را مقایسه کنیم. به این منظور با این پیش فرض که کلماتی که معنای نزدیکی داشته باشند

* Semantic Lexicon

† Resource Description Framework

‡ Dictionary

به احتمال زیاد در یکی از متن های صفحات وب در کنار یکدیگر آمده اند و با یک کلمه عطف و یا فصل جدا شده اند به استخراج تمامی کلماتی که عطف و یا فصلی بین آنها بوده پرداختیم. جدول ۴-۱ نشان دهنده کلمات عطف و فصل مورد استفاده در این تحقیق را نشان می دهد.

پس از به دست آوردن کلمات، یک سری ملاک بین دو کلمه در بین تمام اسناد موجود به دست می آید. این پارامترها عبارت اند از:

تعداد دفعاتی که عطف بین دو کلمه آمده است

تعداد کلماتی که فصل بین دو کلمه آمده است

جدول ۱- کلمات عطف و فصل

لیست	نوع
yet, so, or, because, even, both, and, either	کلمات عطف
Not, Nor	کلمات فصل

ملاک های حضور کلمات باهم در اسناد مختلف شامل ملاک جاکارد، همپوشانی کلمات، دایس و PMI می باشد. در روابط ۱ تا ۴ به ترتیب فرمول ملاک های جاکارد، همپوشانی کلمات، دایس و PMI بیان شده است.

پس از استخراج این پارامترها، آنها را به SVM برای آموزش می دهیم. به این ترتیب که SVM تلاش می کند که ارتباط معنایی این کلمات را با توجه به این شش ویژگی فرابگیرد. برای این منظور پس از استخراج اطلاعات کلمات از اسناد موجود، SVM را با برخی از کلماتی که خود میدانیم هم معنی هستند و برخی کلماتی که میدانیم هم معنی نیستند آموزش داده تا بعد بتواند در خصوص سایر کلمات تصمیم گیری کند. تا این قسمت شباهت معنایی کلمات استخراج شد و نیز سعی شد که حوزه مفهومی درخواست کاربر و اسناد نیز بررسی شود.

۴.۵. شاخص کننده معنایی

در این مرحله وزن مربوط به هر یک از کلمات کلیدی و مفاهیم در هر سند مشخص می شود. این وزن ها idf و tf برای هر یک از کلمات کلیدی و مفاهیم است. tf تعداد تکرار هر کلمه کلیدی یا مفهوم در آن سند است. $*idf$ تعداد تکرار هر کلمه کلیدی یا مفهوم در تمام اسناد است. این اطلاعات برای محاسبه رتبه اسناد به مراحل بعدی ارسال می شوند.

* document frequency

۴.۶. ورود و تبدیل درخواست کاربر

در این مرحله درخواست کاربر به صورت جمله یا کلمات کلیدی از ورودی دریافت می شود. در ادامه باید کلمات کلیدی و مفاهیم و همچنین معانی مشابه و نزدیکی مفهومی استخراج شوند

در مرحله بعد برای بررسی دقیق تر ارتباط مفهومی کلمات از دیتاست RDF استفاده شده است. در این دیتاست اطلاعات به صورت مفاهیم و ارتباط بین آنها است. لذا با استفاده از استخراج کلماتی که در جاهای متعدد دارای ارتباط مفهومی با یکدیگر هستند سعی شده است علاوه بر استخراج شباهت معنایی کلمات و در نظر گرفتن حوزه به استخراج ارتباط مفهومی کلمات و در نظر گرفتن این مورد نیز در جستجو پرداخته شود. به عنوان مثال جمله «من سیب خوردم» را اگر در نظر بگیریم در سه تایی RDF بر اساس قوانین پردازش زبان طبیعی جمله واکاوی شده و به استخراج «سیب»، «خوردن» می پردازد. پس ما می توان فهمید که سیب با خوردن ارتباط دارد و به عنوان مثال اگر کاربر به دنبال خوردنی بود یک سندی که به سیب و چند مورد مشابه دیگر اشاره می کند می تواند مدنظر کاربر باشد، در صورتی که این سند در نگاه اول و توسط روش های سنتی استخراج نمی شود.

۴.۷. تطابق معنایی

در این مرحله جدول درهمی با توجه به کلمات کلیدی و کلمات مشابه و همچنین مفاهیم و نزدیکی مفهومی (مفاهیم مشابه در صورت وجود) به دست آمده از درخواست کاربر در مرحله قبلی، ساخته می شود. تطابق بین درخواست کاربر و اسناد با توجه به کلمات کلیدی و مفاهیم استخراج شده از آنها است. جدول درهم (hash table) برای سند آدم به صورت زیر محاسبه می شود (۴-۱).

$$hash_i = \sum_{j=1}^{NQ} f1(w_{ij}) + \sum_{j=1}^{NQ} \sum_{w \in \text{similarity of } w_{ij}} f2(w) + \sum_{j=1}^{NS} f3(c_{ij}) + \sum_{j=1}^{NS} \sum_{s \in \text{similarity of } c_{ij}} f4(c)$$

(۵)

$$f1(w_{ij}) = \begin{cases} W_{keyword} & \text{if keyword } j \text{ is document } i \\ 0 & \text{otherwise} \end{cases}$$

$$f2(w_{ij}) = \begin{cases} W_{simKeyword} & \text{if keyword } j \text{ is document } i \\ 0 & \text{otherwise} \end{cases}$$

$$f3(c_{ij}) = \begin{cases} W_{semantic} & \text{if semantic } j \text{ is document } i \\ 0 & \text{otherwise} \end{cases}$$

$$f4(w_{ij}) = \begin{cases} W_{simSemantic} & \text{if semantic } j \text{ is document } i \\ 0 & \text{otherwise} \end{cases}$$

که NQ تعداد کلمات کلیدی، NS تعداد مفاهیم، Wkeyword وزن تأثیر کلمات کلیدی مستقیماً وارد شده در توسط کاربر، WsimKeyword وزن تأثیر کلمات مشابه با کلمات کلیدی درخواست کاربر، Wsemantic وزن تأثیر مفاهیم استخراج شده از درخواست کاربر، WsimSemantic وزن تأثیر مفاهیم مشابه با مفاهیم درخواست کاربر می باشند.

در نهایت لیست اسناد قابل نمایش با آستانه گذاری بر روی ضرایب hash به دست می آید (۶) یعنی:

$$lst = \{i \mid hash_i > T, 1 \leq i \leq N\} \quad (6)$$

که T حد آستانه و N تعداد اسناد است. به عبارت دیگر اگر تأثیر کلمات کلیدی و مفاهیم با نزدیکی معنایی و مفاهیمی آن‌ها از حد آستانه‌ای بیشتر باشد، آنگاه آن سند برای رتبه‌بندی به مرحله بعد ارسال می‌شود. به عبارتی دیگر در پایان این مرحله لیست اسناد کاندید برای نمایش مشخص می‌شوند.

۴.۸. رتبه دهنده معنایی

در این مرحله رتبه هر سند (که از مرحله قبلی ارسال شده است) تعیین می‌شود و سپس اسناد برحسب رتبه از زیاد به کم برای کاربر نمایش داده می‌شوند. در ابتدا معکوس فراوانی سند را تعریف می‌کنیم که به صورت (۷) محاسبه می‌شود:

$$idf = \log \frac{N}{df} \quad (7)$$

که df فراوانی سند است. وزن کلمه کلیدی w_j یا مفهوم c_j در سند D_i به صورت (۸) محاسبه می‌شود:

$$wt_{ij}(w_{ij}) = tf(w_{ij}) * idf(w_i) \quad (8)$$

$$wt_{ij}(c_j) = tf(c_{ij}) * idf(c_j)$$

در نهایت ضریب شباهت* بین درخواست کاربر و سند وب Q توسط رابطه (۹) به دست می‌آید:

$$sr(Q, D_i) = \sum_{j=1}^t wt_{qj} * wt_{ij} \quad (9)$$

در این مرحله اسنادی که امتیازشان خیلی کم باشد (پایین‌تر از حد آستانه از پیش تعیین‌شده باشد) از لیست حذف می‌شوند. لیست نهایی با مرتب‌سازی اسناد بر اساس این ضریب شباهت به کاربر نمایش داده می‌شود.

۵. پیاده سازی و نتایج

۵.۱. پارامترهای پیکربندی

پارامترهای استفاده شده برای وزن دهی اسناد و مرتب‌سازی آنها در جدول ۲ خلاصه شده است.

جدول ۲- پارامترهای پیکربندی

پارامتر	مقدار	توصیف
$W_{keyword}$	۱	وزن تأثیر کلمات کلیدی مستقیماً وارد شده در درخواست کاربر
$W_{simKeyword}$	۰.۸	وزن تأثیر کلمات مشابه با کلمات کلیدی درخواست کاربر
$W_{semantic}$	۱	وزن تأثیر مفاهیم استخراج‌شده از درخواست کاربر
$W_{simSemantic}$	۰.۷	وزن تأثیر مفاهیم مشابه با مفاهیم درخواست کاربر
T	۲.۵	حد آستانه حداقل تعلق درخواست کاربر به اسناد
T_2	۱.۸۵	حد آستانه حداقل امتیاز اسناد لیست نهایی
C	۰.۰۵	حد آستانه احتمال حضور همزمان دو کلمه در معیارهای مختلف

* similarity coefficient

۵.۲. نحوه بدست آوردن داده ها

یکی از راهکارهای به دست آوردن اطلاعات صفحات وب برای موتورهای جستجو استفاده از اطلاعات استخراج شده از صفحات وب توسط پیمایشگر است. از آنجا که روش پیشنهادی بر مبنای استفاده از اطلاعات RDF و صفحات وب است، برای اینکه اطلاعات استخراج شده از RDF را جهت بهبود جستجوی صفحات وب معمولی استفاده کرد، به این منظور به صورت دستی به استخراج یک سری از اسناد صفحات وب در این حوزه که تا حدی مرتبط به RDF مورد استفاده در ارزیابی می شد، اقدام شد. بر این اساس از RDF در زمینه music استفاده کرده و به استخراج اسناد وب در زمینه مربوطه و مرتبط با RDF استفاده شده پرداخته شده است و تعدادی سند وب در این حوزه جمع آوری شد که بتوان از دیتاست استاندارد RDF نیز در کنار آن استفاده شود.

۵.۳. معیارهای ارزیابی

برای بررسی تأثیر هر گام، اول باید ملاک های ارزیابی و داده های مورد استفاده مشخص شود. از آنجا که مهم ترین ملاک های ارزیابی موتورهای جستجو (چه سنتی و چه معنایی) سه ملاک دقت^{*}، بازیابی[†] و F[‡] هستند، ما نیز این سه ملاک را برای ارزیابی کارایی استفاده کردیم. ملاک دقت نشان دهنده این موضوع است که چه درصدی از نتایج نشان داده شده توسط موتور جستجو، نتایج صحیح و مدنظر می باشد رابطه (۱۰)، ملاک بازیابی بیانگر این موضوع است که چند درصد از اسناد مرتبط توسط موتور جستجو بازیابی شده و نمایش داده شده اند رابطه (۱۱).

در نهایت ملاک ارزیابی F بیانگر ترکیبی از دو ملاک گفته شده است رابطه (۱۲)

$$\text{Precision} = \frac{\text{No. of relevant documents retrieved}}{\text{Total No. of documents retrieved}} \quad (10)$$

$$\text{Recall} = \frac{\text{No. of relevant documents retrieved}}{\text{Total No. of documents relevant}} \quad (11)$$

$$F = 2 * \frac{\text{Precision} * \text{Recall}}{\text{Precision} + \text{Recall}} \quad (12)$$

که Precision دقت و Recall بازیابی است. از طرف دیگر هر چند مجموعه داده هایی برای بررسی قدرت موتورهای جستجو موجود است ولی در اکثر مقالات این حوزه ذکر شده است که بررسی دقیق کارایی تنها با کمک فرد خبره ممکن است. به نحوی که فرد خبره در خصوص تک تک نتایج ارائه شده توسط الگوریتم نظر بدهد و به این ترتیب کارایی محاسبه شود. به عبارت دیگر برخلاف بسیاری از حوزه های دیگر مجموعه داده ها دارای برچسبی که برای مقایسه و محاسبه کارایی قابل استفاده باشد، نمی باشند. به این ترتیب سعی شده است تا در خصوص هر یک از موارد تست شده به دقت کارایی محاسبه شود.

* Precision

† Recall

‡ F-measure

۵.۳. طرح آزمایش

در بررسی تأثیر گام‌های روش پیشنهادی، سه حالت مختلف در نظر گرفته می‌شود: موتور جستجوی SBIRS، استفاده از شباهت معنایی، استفاده از شباهت مفهومی و شباهت معنایی به این ترتیب تأثیر افزوده شدن هر گام را به مراحل قبلی با استفاده از ملاک‌های ذکر شده بررسی شده است برای ارزیابی از تعدادی درخواست تصادفی استفاده شده است. به علاوه برای دقت بهتر ارزیابی تعدادی از اسناد مورد آزمایش نیز توسط تعدادی فرد خبره بررسی شده و میزان نزدیکی آن‌ها به هر یک از درخواست‌ها گردآوری شده است.

جدول ۳- انطباق بخشی از اسناد مورد آزمایش و درخواست‌های تصادفی کاربر

درخواست کاربر	نتیجه بررسی فرد خبره و میزان ارتباط اسناد مورد آزمایش با درخواست کاربر
۱	Elvis-Aaron-Presley ElvisWorldrestaurantreviews (could be) PublicPerformanceRevised (10% related) group1a_gi-project09 (Semantically and conceptually related)
۲	Christian-hip hop music (40% related) music-assesshb2011-2015 Atton_Popular_Music_Rev3 group1a_gi-project09 (50%) gradek_music_sol_strategies(10% preferably unrelated) JSTOR_Project_MUSE_Titles (30%)
۳	Elvis Elvis Elvis - 2014 Spring Tour Press Release (60%) Elvis The King Tribute Cruise (70%) Elvis-Aaron-Presley Leman Book Review (30%) Presstext (30%) elvis aint dead worksheet (10%) ELVIS PRESLEY (80%) ElvisWorldrestaurantreviews (10%) activity 2 (80%) 0608ElvisGospelService-press (10%)
۴	Atton_Popular_Music_Rev3 gradek_music_sol_strategies magini (40%) Christian-hip hop music (70%) group1a_gi-project09 (20%)
۵	Elvis-Aaron-Presley PublicPerformanceRevised (10% related) group1a_gi-project09 (20% semantically related)
۶	Elvis-Aaron-Presley group1a_gi-project09 (40%) Sacredandprofaneicon-workOulu (50%) activity 2 (30%)
۷	Elvis-Aaron-Presley coyne606-1 (70%) Leman Book Review group1a_gi-project09 (30%)

۵,۴. روش SBIRS

در جستجوی پایه اسناد بر مبنای روش SBIRS مورد مقایسه قرار گرفته است. در این گام نتایج اجرای این روش بر روی مجموعه داده مورد تست نشان داده شده است. جدول ۴ مقادیر دقت و بازیابی نشان داده شده است.

جدول ۴- دقت و بازیابی روش SBIRS

F-measure	Recall	Precision	درخواست کاربر
۰.۳۳	۰.۲۵	۱	۱
۰.۴۹	۰.۳۳	۱	۲
۰.۸۲	۰.۹۱	۰.۷۵	۳
۱	۱	۱	۴
۰.۴۹	۰.۳۳	۱	۵
۰.۳۳	۰.۲	۱	۶
۰.۷۵	۰.۶	۱	۷

۵,۵. افزودن شباهت معنایی

با افزودن گام بررسی شباهت معنایی کلمات، موتور جستجو علاوه بر جستجوی کلمات وارد شده توسط کاربر به جستجوی کلمات مشابه با کلمات کلیدی وارد شده توسط کاربر نیز می پردازد و به این ترتیب شانس یافتن اسناد مورد نظر کاربر افزایش پیدا می کند. به این ترتیب بازیابی موتور جستجو افزایش یافته است. به دلیل بازیابی اسناد بیشتر با توجه به این نکته که سندی اگر دارای کلمه کلیدی مشابه با درخواست کاربر باشد نیز بازیابی می شود، دقت مشابه با روش پایه به دست آمده است. مقادیر دقت و بازیابی این حالت نیز در جدول ۵ نشان داده شده است.

جدول ۵- دقت و بازیابی و F-measure با در نظر گرفتن شباهت معنایی

F-measure	بازیابی Recall	دقت Precision	درخواست
۰.۴۹	۰.۴	۱	۱
۰.۹۰	۰.۸۳	۱	۲
۰.۸۲	۰.۹۱	۰.۷۵	۳
۱	۱	۱	۴
۰.۴۹	۰.۳۳	۱	۵
۰.۵۷	۰.۴	۱	۶
۰.۸۸	۰.۸	۱	۷

۵,۶. در نظر گرفتن شباهت معنایی و شباهت مفهومی

در گام بعد با استفاده از اطلاعات استخراج شده از RDF ها و در نظر گرفتن شباهت مفهومی کلمات سعی در بهبود بیشتر کارایی الگوریتم شده است. در این گام علاوه بر در نظر گرفته شدن حوزه درخواست کاربر و تطبیق با اسناد مرتبط و همچنین در نظر گرفته شدن اسنادی که شامل کلمات کلیدی مشابه هستند تلاش شده است تا نزدیکی مفهومی از درخواست کاربر و اطلاعات اسناد نیز با استفاده از اطلاعات RDF استفاده شده در نظر گرفته شود و با استفاده از این راهکار نیز در برخی از موارد کارایی افزایش پیدا کرده است. نتایج این حالت نیز در جدول ۶ نشان داده شده است.

جدول ۹- دقت و بازیابی و F-measure با در نظر گرفتن شباهت مفهومی و شباهت معنایی

F-measure	بازیابی Recall	دقت Precision	درخواست
۰.۷۵	۰.۶	۱	۱
۰.۹۰	۰.۸۳	۱	۲
۰.۸۲	۰.۹۱	۰.۷۵	۳
۱	۱	۱	۴
۰.۴۹	۰.۳۳	۱	۵
۰.۷۵	۰.۶	۱	۶
۱	۱	۱	۷

۶. نتیجه گیری

در این مقاله، تلاش شد تا روشی بر مبنای در نظر گرفتن شباهت معنایی و نزدیکی مفهومی کلمات و نیز اطلاعات زبان‌شناسی برای جستجوی مفهومی صفحات وب ارائه شود. در این روش ابتدا درخواست کاربر به صورت متنی دریافت می‌شود، سپس این درخواست بسط داده شده و کلمات کلیدی و مفاهیم از آن استخراج می‌شود. همچنین در این مرحله کلمات کلیدی مشابه با درخواست کاربر (شباهت معنایی) و مفاهیم نزدیک به آن (نزدیکی مفهومی) نیز استخراج می‌شوند. کلمات مشابه، کلماتی هستند که بین آن‌ها کلمات عطف یا فصل وجود داشته باشد و شباهت معنایی توسط ذخیره و آموزش ماشین بردار پشتیبان با کمک کلمات مشابه و غیرمشابه به دست می‌آید. نزدیکی مفهومی توسط RDF انجام می‌شود که شامل اطلاعات مفهومی است. این روش سعی در بهبود روش SBIRS برای یافتن کلمات کلیدی مشابه و در نظر گرفتن مفاهیم مشابه در اسناد به منظور بالابردن معیار فراخوانی و یافتن اسناد مرتبط بیشتر با نیازمندی های کاربر دارد.

مراجع

1. E. Motta, and M. Sabou, "Next generation semantic web applications," The Semantic Web–ASWC 2006, pp. 24-29: Springer, 2006.
2. B. Aleman-Meza, B. Arpinar, M. V. Nural, and A. P. Sheth, "Ranking documents semantically using ontological relationships." pp. 299-304.
3. V. Jindal, S. Bawa, and S. Batra, "A review of ranking approaches for semantic search on Web," Information Processing & Management, vol. 50, no. 2, pp. 416-425, 2014.
4. M. Thangaraj, and G. Sujatha, "An architectural design for effective information retrieval in semantic web," Expert Systems with Applications, vol. 41, no. 18, pp. 8225-8233, 2014.
5. Bollegala, D., Y. Matsuo, and M. Ishizuka, "A web search engine-based approach to measure semantic similarity between words". Knowledge and Data Engineering, IEEE Transactions on, Vol 23(7), pp 977-990, 2011.