

بِسْمِ اللَّهِ الرَّحْمَنِ الرَّحِيمِ

۱۵۸۴۲۹-۲۰۲۳۵۷۸



دانشگاه اصفهان

دانشکده فنی و مهندسی

گروه مهندسی کامپیوتر

پایان نامه‌ی کارشناسی ارشد رشته‌ی مهندسی کامپیوتر گرایش نرم افزار

گسترش معنایی پرس و جو

استاد راهنما:

دکتر محمد علی نعمت بخش

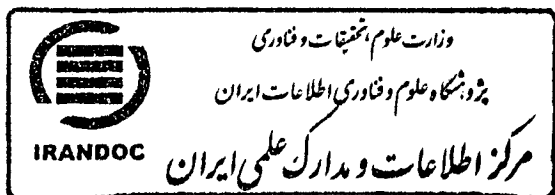
استاد مشاور:

دکتر ناصر نعمت بخش

پژوهشگر:

مژگان شبان زاده حبیب آبادی

مهرماه ۱۳۸۹



۱۵۸۴۲۶

۱۳۹۰/۳/۱۶

کلیه حقوق مادی مترتب بر نتایج مطالعات، ابتکارات  
و نوآوری های ناشی از تحقیق موضوع این پایان نامه  
متعلق به دانشگاه اصفهان است.

شوه نگارش پایان نامه  
رعایت شده است.  
تحصیلات تکمیلی دانشگاه اصفهان



دانشگاه اصفهان

دانشکده فنی و مهندسی

گروه مهندسی کامپیوتر

پایان نامه ی کارشناسی ارشد رشته ی مهندسی کامپیوتر گرایش نرم افزار

خانم مژگان شبان زاده حبیب آبادی تحت عنوان

### گسترش معنایی پرس و جو

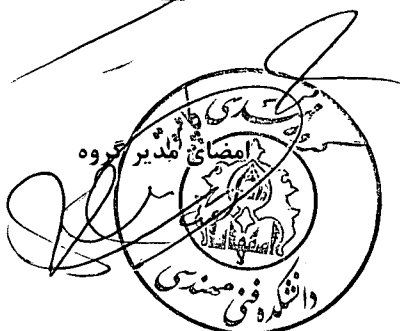
در تاریخ ۸۹/۷/۴ توسط هیأت داوران زیر بررسی و با درجه عالی به تصویب نهایی رسید.

۱- استاد راهنمای پایان نامه دکتر محمد علی نعمت بخش با مرتبه ی علمی دانشیار امضا

۲- استاد مشاور پایان نامه دکتر ناصر نعمت بخش با مرتبه ی علمی استادیار امضا

۳- استاد داور داخل گروه دکتر محمدرضا خیامباشی با مرتبه ی علمی استادیار امضا

۴- استاد داور خارج از گروه دکتر سید رسول موسوی با مرتبه ی علمی استادیار امضا



## سپاسگزاری

بر خود لازم می‌دانم تا از استاد راهنمای ارجمندم جناب آقای دکتر محمدعلی نعمت بخش و همچنین استاد مشاور گرانقدرم جناب آقای دکتر ناصر نعمت بخش کمال تشکر و قدردانی را داشته باشم و از خداوند توفیق ایشان را طلب می‌کنم.

از خانواده مهربانم که همواره مشوق و حامی اینجانب در طول زندگی بوده‌اند، سپاسگزاری کرده، سلامتی و موفقیت آن‌ها را از پروردگار مهربان مسئلت دارم. همچنین از تمامی دوستانی که در این راه یاریگر اینجانب بوده‌اند مراتب تشکر و قدردانی را دارم.

تقدیم بہ

پرو مادرم کہ موفقیت من آرزوی آن ما و سر بلندی آن ما

آرزوی من است.

## چکیده

بازیابی اطلاعات یکی از اصلی‌ترین نیازهای کاربران است؛ روزانه کاربران زیادی به جستجو در وب و دیگر منابع به منظور پاسخ گویی به نیاز اطلاعاتی خود می‌پردازند. مسائل موجود در زبان طبیعی از جمله عدم تطابق لغوی، کلمات چند معنایی، کوتاه و مبهم بودن پرس‌وجو و دانش ناقص کاربران از موضوع مورد نظر در بازیابی اطلاعات منجر به بازیابی نتایج نامرتب و کاهش رضایت کاربران از نتایج بازیابی شده می‌شود. گسترش پرس‌وجو با بررسی پرس‌وجوهای کاربران و افزودن خودکار کلمات مناسب و با ارزش به آن‌ها کمک می‌کند تا اسناد مرتبط با نیاز و منظور کاربر جستجو و بازیابی گردد. اگر گسترش پرس‌وجو به صورت هوشمندانه‌ای انجام نشود، با انحراف پرس‌وجو و فاصله گرفتن آن از منظور کاربر منجر به بازیابی نتایج نامرتب‌تری نسبت به نتایج پرس‌وجوی اولیه خواهد شد. مساله‌ی دیگری که در گسترش پرس‌وجو باید به آن توجه کرد این است که در مورد پرس‌وجوهای مشتمل بر بیش از یک کلمه، انتخاب واژگان گسترشی که تنها با یکی از این کلمات مرتبط باشند باعث رخداد مشکل خروج از تعادل پرس‌وجو خواهد شد.

در این رساله، روش جدیدی برای گسترش معنایی پرس‌وجو به منظور تطبیق دادن پرس‌وجو با منظور کاربر ارائه شده است. روش پیشنهادی با استفاده از یک الگوریتم رفع ابهام مبتنی بر هستی‌شناسی به رفع ابهام از کلمات پرس‌وجو می‌پردازد. سپس، به منظور در نظر گرفتن روابط بین لغات در پرس‌وجوهای چند کلمه‌ای و اجتناب از مشکل خروج از تعادل پرس‌وجو به گروه‌بندی کلمات آن بر مبنای تشابه معنایی بین آن‌ها می‌پردازد. در ادامه با استفاده از روابط موجود در شبکه واژگان، یک شبکه‌ی معنایی از واژگان هر گروه ایجاد شده از کلمات پرس‌وجو و لغات مرتبط با آن‌ها از نظر معنایی می‌سازد. این روش بر طبق روابط و سلسله مراتب شبکه‌ی ساخته شده، مهمترین کلمات برای گسترش پرس‌وجو را مشخص می‌کند. از بین کلمات انتخاب شده، کلماتی که باعث ایجاد ابهام و نویز در پرس‌وجو نشوند، به عنوان کلمات گسترش انتخاب می‌شوند و وزن مناسبی برای آن‌ها محاسبه می‌شود. به این ترتیب پرس‌وجوی گسترش یافته ساخته می‌شود و عملیات جستجو با این پرس‌وجوی جدید انجام می‌شود. این روش با در نظر گرفتن معیارهای فراخوانی و دقت بر روی مجموعه داده‌ی TIME ارزیابی شده است. نتایج ارزیابی نشان دهنده‌ی افزایش نرخ فراخوانی و دقت بازیابی می‌باشد.

**کلمات کلیدی:** بازیابی اطلاعات، موتور جستجو، گسترش پرس‌وجو، روابط معنایی.

## فهرست مطالب

صفحه

عنوان

### فصل اول: مقدمه

- ۱-۱ مقدمه ..... ۱
- ۲-۱ شرح و بیان مسأله‌ی پژوهشی ..... ۲
- ۱-۲-۱ مقایسه‌ی سیستم‌های بازیابی اطلاعات با سیستم‌های بازیابی داده ..... ۳
- ۲-۲-۱ فرآیند بازیابی اطلاعات ..... ۳
- ۳-۲-۱ بازیابی اطلاعات کلیدواژه‌ای ..... ۵
- ۳-۱ هدف تحقیق ..... ۹
- ۴-۱ اهمیت و کاربرد نتایج تحقیق ..... ۱۰
- ۵-۱ ساختار پایان نامه ..... ۱۰

### فصل دوم: پیشینه و زمینه‌ی تحقیق

- ۱-۲ مقدمه ..... ۱۲
- ۲-۲ گسترش پرس‌وجو ..... ۱۲
- ۱-۲-۲ روش‌های محلی گسترش پرس‌وجو ..... ۱۵
- ۲-۲-۲ روش‌های سراسری گسترش پرس‌وجو ..... ۱۶
- ۳-۲ مشکلات موجود در گسترش پرس‌وجو ..... ۱۸
- ۴-۲ پیشینه‌ی تحقیق ..... ۱۸
- ۵-۲ مدل فضای بردار ..... ۲۷
- ۶-۲ هستی‌شناسی ..... ۲۸
- ۱-۶-۲ شبکه‌ی واژگان ..... ۲۹
- ۷-۲ تشابه معنایی ..... ۳۱
- ۸-۲ رفع ابهام ..... ۳۸
- ۱-۸-۲ روش‌های مبتنی بر مجموعه‌ی نوشتار ..... ۳۸



۴۰	۲-۸-۲ روش‌های مبتنی بر پایگاه دانش.....
۴۳	۹-۲ الگوریتم فعال سازی گسترشی.....
۴۶	۱۰-۲ جمع‌بندی.....

### فصل سوم: طراحی و پیاده سازی یک الگوریتم جدید گسترش معنایی پرس‌وجو

۴۷	۱-۳ مقدمه.....
۴۷	۲-۳ تعریف مساله.....
۴۸	۳-۳ اهداف و انگیزه‌ی طراحی روش جدیدی برای گسترش معنایی پرس‌وجو.....
۵۰	۴-۳ معماری سیستم.....
۵۰	۵-۳ مفاهیم اولیه.....
۵۰	۶-۳ ساختار کلی الگوریتم ارائه شده.....
۵۲	۱-۶-۳ پیش‌پردازش پرس‌وجو.....
۵۷	۲-۶-۳ گروه‌بندی کلمات پرس‌وجو.....
۵۸	۳-۶-۳ ساخت گراف معنایی.....
۶۰	۴-۶-۳ یافتن کلمات کاندید برای گسترش پرس‌وجو.....
۶۱	۵-۶-۳ اعمال صافی.....
۶۲	۶-۶-۳ وزن دهی پرس‌وجوی گسترش یافته.....
۶۳	۷-۳ جمع‌بندی.....

### فصل چهارم: ارزیابی و تحلیل عملکرد الگوریتم طراحی شده

۶۵	۱-۴ مقدمه.....
۶۵	۲-۴ مجموعه داده.....
۶۸	۳-۴ معیارهای ارزیابی.....
۶۸	۱-۳-۴ فراخوانی.....
۶۸	۲-۳-۴ دقت.....

۶۹	۴-۴ بستر بازیابی اطلاعات
۶۹	۵-۴ نتایج آزمایشات
۶۹	۱-۵-۴ تعیین مقادیر پارامترهای الگوریتم
۷۴	۲-۵-۴ ارزیابی الگوریتم پیشنهادی
۷۷	۶-۴ مقایسه‌ی الگوریتم پیشنهادی با دیگر الگوریتم‌های گسترش پرس‌وجو و نتیجه‌گیری
۷۸	۷-۴ جمع‌بندی

### فصل پنجم: نتیجه‌گیری و کارهای آینده

۷۹	۱-۵ مقدمه
۷۹	۲-۵ نتایج
۸۰	۳-۵ کارهای آینده
۸۳	واژه‌نامه
۸۶	منابع و مأخذ

## فهرست شکل‌ها

صفحه	عنوان
۴	شکل ۱-۱ فرآیند بازیابی اطلاعات.....
۸	شکل ۲-۱ نتایج پرس‌وجوی "Apple Growing" با موتور جستجوی گوگل.....
۸	شکل ۳-۱ نتایج پرس‌وجوی "Apple Growing Computers" با موتور جستجوی گوگل.....
۱۴	شکل ۱-۲ روش‌های گسترش پرس‌وجو و منابع استخراج واژگان گسترش.....
۱۵	شکل ۲-۲ فرآیند کلی روش‌های محلی گسترش پرس‌وجو.....
۱۷	شکل ۳-۲ فرآیند کلی روش‌های سراسری برای گسترش پرس‌وجو.....
۲۰	شکل ۴-۲ شبه کد الگوریتم پیشنهادی در مرجع [۲۳].....
۲۵	شکل ۵-۲ شبکه‌ی معنایی تولید شده برای اولین معنای bus.....
۲۵	شکل ۶-۲ الگوهای بین سومین معنای mountain و اولین معنای top.....
۳۰	شکل ۷-۲ معنای door در شبکه واژگان.....
۳۱	شکل ۸-۲ روابط زیرمعنایی در شبکه واژگان برای لغت door.....
۴۴	شکل ۹-۲ ورودی، پارامترها و شبه کد الگوریتم فعال سازی گسترشی.....
۴۵	شکل ۱۰-۲ اولین مرحله‌ی الگوریتم فعال سازی گسترشی.....
۴۵	شکل ۱۱-۲ خروجی مرحله‌ی اول الگوریتم فعال سازی گسترشی.....
۴۵	شکل ۱۲-۲ خاتمه‌ی الگوریتم فعال سازی گسترشی.....
۴۸	شکل ۱-۳ فرآیند کلی روش ارائه شده برای گسترش پرس‌وجو.....
۵۱	شکل ۲-۳ مروری بر روش ارائه شده برای گسترش معنایی پرس‌وجو.....
۵۶	شکل ۳-۳ تعدادی از کلمات توقف در زبان انگلیسی.....
۵۸	شکل ۴-۳ روابط بین کلمات در شبکه واژگان.....
۶۲	شکل ۵-۳ برخی از کلمات پرس‌وجوی نمونه‌ی گسترش یافته.....
۶۴	شکل ۶-۳ شبه کد الگوریتم پیشنهادی.....
۶۶	شکل ۱-۴ محتوای مقاله‌ی شماره‌ی ۲۳ مجموعه داده‌ی TIME.....

## عنوان

## صفحه

شکل ۲-۴ میانگین دقت بازیابی نسبت به تعداد معانی هر کلمه ..... ۷۴

شکل ۳-۴ نمودار دقت- فراخوانی بازیابی کلید واژه‌ای و گسترش پرس‌وجو ..... ۷۵

شکل ۴-۴ نمودار دقت- فراخوانی بازیابی کلید واژه‌ای و گسترش وزن دار پرس‌وجو ..... ۷۶

## فهرست جدول‌ها

عنوان	صفحه
جدول ۱-۱ مقایسه‌ی سیستم‌های بازیابی اطلاعات با سیستم‌های بازیابی داده	۳
جدول ۱-۲ مقادیر بهینه‌ی پارامترهای $\alpha$ ، $\beta$ ، $\gamma$ و $\delta$	۲۱
جدول ۲-۲ دسته بندی ماژول‌های تشابه معنایی و نقاط قوت و ضعف آن‌ها	۳۶
جدول ۳-۲ دسته بندی روش‌های رفع ابهام و نقاط ضعف آن‌ها	۴۲
جدول ۱-۳ خروجی نرم افزار WordNet-SenseRelate-AllWords با دریافت پرس‌وجوی نمونه	۵۴
جدول ۲-۳ کلمات غیرتوقف رفع ابهام شده‌ی پرس‌وجوی نمونه	۵۵
جدول ۱-۴ پرس‌وجوهای مجموعه داده‌ی TIME	۶۷
جدول ۲-۴ قسمتی از محتوای فایلی که اسناد مرتبط با هر پرس‌وجو را مشخص می‌کند	۶۷
جدول ۳-۴ میانگین دقت بازیابی نسبت به تشابه معنایی	۷۰
جدول ۴-۴ میانگین دقت بازیابی نسبت به وزن یال‌ها	۷۱
جدول ۵-۴ میانگین دقت بازیابی نسبت به عامل کاهنده	۷۲
جدول ۶-۴ میانگین دقت بازیابی نسبت به حد آستانه‌ی فعال سازی	۷۳

## فصل اول

### مقدمه

#### ۱-۱ مقدمه

فرآیند تشخیص ارتباط بین اسناد و پرس و جوی کاربر و بازگرداندن اسنادی که احتمالاً نیاز کاربر را تامین می‌کنند، بازیابی اطلاعات نامیده می‌شود. بارزترین کاربرد سیستم‌های بازیابی اطلاعات موتورهای جستجو و کتابخانه‌های دیجیتالی می‌باشد. اگرچه با وجود موتورهای جستجو و کتابخانه‌های دیجیتالی دسترسی به اطلاعات آسان‌تر شده است، اما از آن جایی که حجم اطلاعات وب و کاربران آن روز به روز در حال افزایش می‌باشند، دسترسی به اطلاعات مورد نیاز به چالشی برای کاربران تبدیل شده است. علاوه بر این، موتورهای جستجوی کلیدواژه‌ای در پاسخ به نیاز کاربران ضعیف عمل می‌کنند و بسیاری از کاربران نمی‌دانند چگونه پرس و جوی مناسب ایجاد کنند. حتی وقتی کاربران باتجربه در حوزه‌ی ناشناخته‌ای به جستجو می‌پردازند، قادر به ارائه‌ی پرس و جوی مناسب و بیان دقیق نیاز اطلاعاتی خویش نمی‌باشند [۱].

گسترش پرس و جو با افزودن کلمات به پرس و جو نتایج بازیابی اطلاعات را بهبود می‌بخشد. در این پایان نامه با ارائه‌ی روش جدیدی برای گسترش پرس و جو سعی می‌شود، بازیابی اطلاعات بهبود بخشیده شود.

در این فصل به معرفی سیستم‌های بازیابی اطلاعات، مشکلاتی که سیستم‌های بازیابی اطلاعات کلید واژه‌ای با آن مواجه می‌باشند و شرح مسأله‌ی پژوهشی در این پایان‌نامه پرداخته می‌شود. سپس اهمیت و کاربرد موضوع بررسی

می‌شود و همچنین اهداف پایان‌نامه تبیین خواهد شد.

## ۲-۱ شرح و بیان مسأله‌ی پژوهشی

بازیابی اطلاعات با مسائلی از قبیل نمایش، ذخیره، سازمان دهی و دسترسی اطلاعات مواجه می‌باشد [۲]. نمایش و سازمان دهی اطلاعات باید به گونه‌ای باشد که دسترسی آسان به اطلاعات دلخواه را برای کاربران فراهم کند. بیان دقیق و توصیف صفات اختصاصی نیاز اطلاعاتی<sup>۱</sup> کاربران مسأله‌ی ساده‌ای نیست. برای مثال نیاز اطلاعاتی زیر در وب را در نظر بگیرید.

چرا تابع آزادسازی حافظه<sup>۲</sup> در زبان برنامه‌نویسی ++C در محیط ویژوال بر روی سیستم عامل لینوکس با خطای دسترسی مواجه می‌شود؟

بدیهی است که این توضیح کامل از نیاز اطلاعاتی را نمی‌توان مستقیماً برای درخواست اطلاعات با استفاده از واسط‌های کاربری موتورهای جستجوی امروزی به کار برد. بنابراین کاربر باید در ابتدا نیاز اطلاعاتی خود را به پرس‌وجویی قابل پردازش توسط موتورهای جستجو تبدیل کند.

معمولاً خلاصه‌ای از نیاز اطلاعاتی کاربر توسط پرس‌وجویی مشتمل بر مجموعه‌ای از کلمات کلیدی بیان می‌شود. نیاز اطلاعاتی مطرح شده در مثال بالا برای یک موتور جستجوی انگلیسی زبان ممکن است به صورت یکی از پرس‌وجوهای زیر مطرح شود:

- free error access violation
- error at runtime and break
- "free function" access violation
- "free function" \_CrtIsValidHeapPointer(pUserData)
- linux free error
- free function runtime error

با در اختیار داشتن پرس‌وجوی کاربر وظیفه‌ی اصلی سیستم‌های بازیابی اطلاعات، بازیابی اطلاعاتی است که مرتبط با نیاز اطلاعاتی کاربر باشند. ذکر این نکته ضروری است که سیستم‌های بازیابی اطلاعات متفاوت از سیستم‌های بازیابی داده می‌باشند.

<sup>۱</sup> Information Need

<sup>۲</sup> Free Function

### ۱-۲-۱ مقایسه‌ی سیستم‌های بازیابی اطلاعات با سیستم‌های بازیابی داده

در جدول ۱-۱ سیستم‌های بازیابی اطلاعات و سیستم‌های بازیابی داده با هم مقایسه شده‌اند. در سیستم‌های بازیابی داده مانند پایگاه داده‌های رابطه‌ای<sup>۱</sup> داده‌ها ساختار یافته‌اند و از پرس‌وجوهای ساختار یافته برای بازیابی داده استفاده می‌شود. در صورت بازیابی پاسخ غلط در این سیستم‌ها، کل سیستم دچار خطا می‌باشد [۲]. در حالی که در سیستم‌های بازیابی اطلاعات، داده‌ها ساختار نیافته می‌باشند به این معنا که داده‌ها متنی می‌باشند و ساختار واضح و آسانی برای کامپیوترها ندارند. علاوه بر این در این سیستم‌ها پرس‌وجوها نیز غیر ساختار یافته و اغلب به زبان طبیعی می‌باشند. در نتیجه ممکن است نتایجی بازیابی شود که با پرس‌وجوی کاربر ارتباطی ندارند و یا ارتباط کمی دارند. چون پرس‌وجوها و داده‌ها به زبان طبیعی می‌باشند، زبان طبیعی ساختار نیافته است و ممکن است ابهام داشته باشند [۲].

جدول ۱-۱ مقایسه‌ی سیستم‌های بازیابی اطلاعات با سیستم‌های بازیابی داده

سیستم‌های بازیابی داده	سیستم‌های بازیابی اطلاعات
داده‌های ساختار یافته	داده‌های ساختار نیافته
پرس‌وجوهای ساختار یافته	پرس‌وجوها به زبان طبیعی
عدم بازیابی نتایج غلط	امکان بازیابی نتایج نامرتب

برای پاسخگویی به نیاز اطلاعاتی کاربران، سیستم‌های بازیابی اطلاعات باید محتوای اطلاعات را تفسیر و بر طبق میزان مرتبط بودن با پرس‌وجوی کاربر رتبه بندی<sup>۲</sup> کنند. این تفسیر شامل استخراج اطلاعات از اسناد به منظور انطباق آن‌ها با نیاز کاربران است. سختی کار نه تنها دانستن چگونگی استخراج اطلاعات می‌باشد بلکه دانستن چگونگی تصمیم‌گیری در مورد مرتبط بودن اسناد نیز می‌باشد. هدف اصلی بازیابی اطلاعات، بازیابی تمامی اسنادی است که با پرس‌وجوی کاربر مرتبط می‌باشند و امکان بازیابی تعداد کمی اسناد نامرتب نیز وجود دارد.

### ۱-۲-۲ فرآیند بازیابی اطلاعات

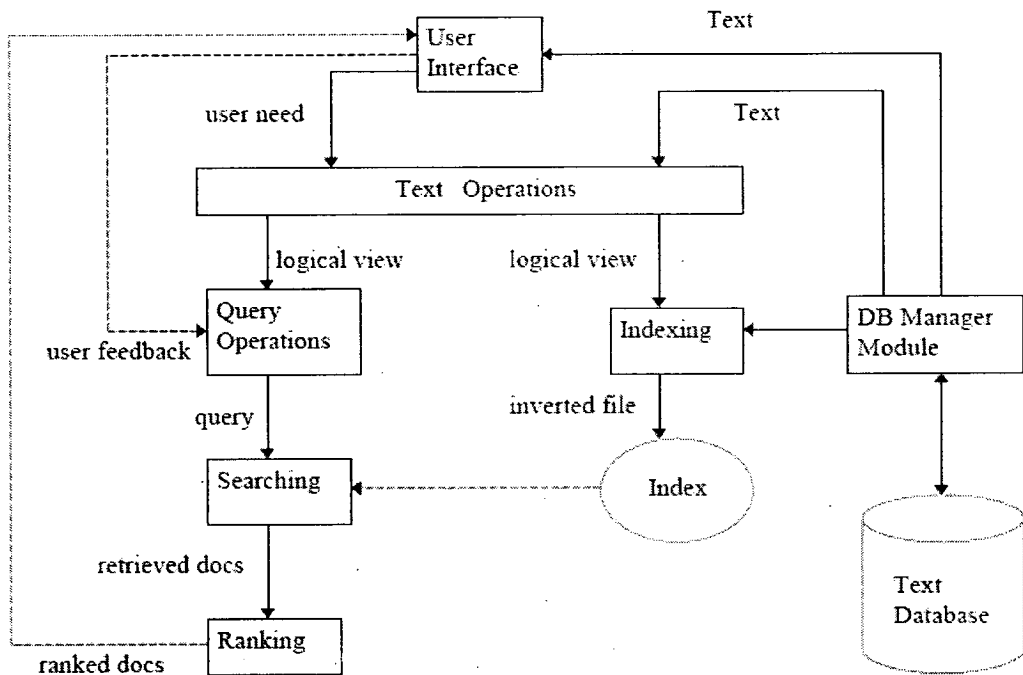
در شکل ۱-۱ فرآیند بازیابی اطلاعات نشان داده شده است. به منظور فراهم سازی امکان بازیابی، لازم است که پایگاه متون تعریف شود. این کار با استفاده از مدیر پایگاه داده انجام می‌شود، مدیر پایگاه داده موارد زیر را

<sup>1</sup> Relational Database

<sup>2</sup> Rank



مشخص می‌کند: الف - اسنادی که استفاده می‌شوند. ب - اعمالی که روی متون انجام می‌شود. ج - مدل متن (ساختار متن و عناصری که می‌توانند بازیابی شوند). اعمال متنی اسناد را تغییر شکل می‌دهند و دیدی منطقی از آن‌ها تولید می‌کنند به این ترتیب که هر سند را با مجموعه‌ای از کلید واژه‌ها نشان می‌دهند. کلید واژه‌ها یا به صورت خودکار از اسناد استخراج می‌شوند و یا به صورت دستی برای هر سند مشخص می‌شوند. برای استخراج کلید واژه‌ها، کلمات توقف<sup>۱</sup> از متن حذف می‌شوند. کلمات توقف کلماتی هستند که تقریباً در تمامی اسناد رخ می‌دهند و قادر به نفکیک اسناد مرتبط از اسناد نامرتبط نمی‌باشند. بقیه‌ی واژه‌ها ریشه‌یابی می‌شوند تا تعداد کلیدواژه‌ها کاهش پیدا کند، به این پردازش‌ها اعمال متنی گفته می‌شود.



شکل ۱-۱ فرآیند بازیابی اطلاعات [۲]

با ساخته شدن دید منطقی از اسناد، مدیر پایگاه داده، ایندکسی از متون می‌سازد. شاخص ساختار مهمی است و امکان جستجوی سریع روی حجم عظیمی از داده‌ها را فراهم می‌کند. معروف‌ترین ساختار شاخص فایل معکوس است. در فایل معکوس هر واژه به اسنادی اشاره می‌کند که این واژه در آن‌ها رخ داده است. اگرچه تولید شاخص

<sup>۱</sup> Stopwords

و پایگاه داده هزینه‌ی زمانی و حافظه‌ای در بر دارد اما چون از سیستم‌های بازیابی اطلاعات بسیار زیاد استفاده می‌شود، این هزینه قابل چشم‌پوشی است.

پس از ایندکس گذاری اسناد پایگاه داده می‌توان عمل جستجو را انجام داد. ابتدا کاربر نیاز اطلاعاتی خود را در قالب یک پرس‌وجو مشخص می‌کند. نیاز اطلاعاتی کاربر توسط همان اعمال متنی که به اسناد اعمال شده‌اند، مورد پردازش قرار می‌گیرد و تجزیه می‌شود. بعد از آن اعمال پرس‌وجو انجام می‌شود و نیاز کاربر به فرم قابل ارائه به سیستم تبدیل می‌شود. سپس پرس‌وجو برای استخراج اسناد پردازش می‌شود. با وجود شاخص، پرس‌وجو به سرعت پردازش می‌شود.

قبل از ارسال نتایج به کاربر اسناد بازیابی شده بر اساس میزان مرتبط بودن با نیاز اطلاعاتی کاربر رتبه بندی می‌شوند. سپس کاربر اسناد رتبه بندی شده را به منظور پیدا کردن اطلاعات مفید بررسی می‌کند. در این مرحله ممکن است، کاربر زیرمجموعه‌ای از اسنادی که مشاهده کرده است را به عنوان اسناد دلخواه مشخص کند و چرخه‌ی بازخورد کاربر را آغاز کند. در این چرخه، سیستم با استفاده از اسنادی که کاربر انتخاب کرده است، اطلاعات پرس‌وجو را تغییر می‌دهد و به احتمال زیاد این پرس‌وجوی تغییر یافته نمایش بهتری از نیاز واقعی کاربر فراهم می‌کند.

### ۳-۲-۱ بازیابی اطلاعات کلیدواژه‌ای<sup>۱</sup>

در سیستم‌های بازیابی اطلاعات کلید واژه‌ای معیار میزان مرتبط بودن اسناد با پرس‌وجوی کاربر و رتبه بندی اسناد، تعداد رخداد واژه‌های پرس‌وجو در آن‌ها می‌باشد. در این سیستم‌ها هر سند و هر پرس‌وجو به عنوان مجموعه‌ای از واژه‌ها در نظر گرفته می‌شود و سندی که بیشترین تعداد رخداد واژه‌های پرس‌وجو را داشته باشد، به عنوان مرتبط‌ترین منبع برای کاربر محسوب می‌شود.

اگرچه موثر بودن بازیابی کلیدواژه‌ای ثابت شده است اما این نوع بازیابی اطلاعات با مشکلاتی نیز مواجه است، این مشکلات را می‌توان به سه دسته‌ی زیر تقسیم‌بندی کرد:

۱- مشکلاتی که از ذات زبان طبیعی ناشی می‌شوند:

<sup>۱</sup> Keyword

• عدم تطابق<sup>۱</sup> واژگان. مشکل عدم تطابق واژگان به این علت رخ می‌دهد که مفاهیم ممکن است با واژگان متفاوتی در پرس‌وجوی کاربر و اسناد بیان شوند [۳]. در هر زبانی ممکن است که واژه‌ای یک یا چندین واژه‌ی هم معنا<sup>۲</sup> داشته باشد. بنابراین، این احتمال وجود دارد که کاربران برای توصیف مفاهیم در پرس‌وجوهای خود از واژگانی متفاوت با واژگانی که نویسندگان برای توصیف همان مفهوم به کار می‌برند، استفاده کنند. علاوه بر این نویسندگان متفاوت از واژگان متفاوتی برای بیان مفاهیم یکسان استفاده می‌کنند [۴]. برطبق [۵] کاربران با نیاز اطلاعاتی یکسان تنها در ۲۰٪ موارد پرس‌وجوی یکسان به کار می‌برند، بنابراین با وجود نیاز یکسان نتایج متفاوتی دریافت می‌کنند. در این موتورهای جستجو مراجعی که هر چند برای کاربر بسیار مفیدند ولی به دلیل عدم وجود واژگان پرس‌وجوی کاربر در آن‌ها و رخداد واژگان هم‌معنا در نتایج جستجو دیده نمی‌شوند. اگر مشکل عدم تطابق در سیستم‌های بازبایی اطلاعات به طور مناسبی حل نشود، کارایی آن‌ها کاهش پیدا خواهد کرد.

• واژگان چندین معنایی<sup>۳</sup>. علاوه بر وجود واژگان هم‌معنا، واژگانی نیز وجود دارند که چندین معنا دارند و معنای آن‌ها با توجه به متنی که در آن به کار برده شده‌اند، مشخص می‌شود. واژگان چند معنایی مشکلاتی را برای موتورهای جستجو به وجود می‌آورند؛ وقتی کاربر در پرس و جوی خود واژه‌ای را به کار می‌برد، موتورهای جستجویی که تنها بر اساس انطباق واژگان عمل جستجو را انجام می‌دهند، تنها مراجعی را به کاربر ارائه می‌دهند که واژه‌ی مذکور در آن‌ها به کار رفته باشد. پس ممکن است نتایجی ارائه شود که واژه‌ی مذکور در معنایی غیر از معنای مورد نظر کاربر در آن‌ها استفاده شده باشد [۶].

بنابراین با وجود واژگان هم‌معنا و چندمعنا، در بازبایی اطلاعات کلیدواژه‌ای نتایج جستجو بسیار حساس به واژه‌های به کار رفته در پرس‌وجوی کاربر است، به این ترتیب که ممکن است جستجو با واژه‌ای نتیجه‌ای در بر نداشته باشد در حالی که با واژه‌ی مترادف آن به نتیجه برسد. علاوه بر این صحت نتایج جستجو منوط به این است که واژه‌هایی که هم در پرس‌وجوی کاربر و هم در اسناد بازبایی شده وجود دارند، به یک معنا باشند.

<sup>1</sup> Mismatch

<sup>2</sup> Synonym

<sup>3</sup> Polysemy

## ۲- کاربران

- پرس وجوهای کوتاه. کاربران در اغلب موارد اطلاعات مورد نیاز خود را توسط چند واژه بیان می کنند، بنابراین معمولا پرس وجوهای کاربران کوتاه می باشند. در [۷] بیان شده است که وقتی تعداد کلمات پرس و جو زیاد باشد، گسترش دادن پرس وجو کارایی زیادی به دنبال نخواهد داشت ولی پرس وجوهای کوتاه معمولا از گسترش پرس وجو سود می برند. در واقع علت این مساله ابهام کم تر پرس وجوهای طولانی و احتمال بیشتر رخداد واژه های پرس وجوی طولانی در اسناد مرتبط با آنها می باشد. بنابراین مساله ی عدم تطابق در مورد پرس وجوهای کوتاه نسبت به پرس وجوهای طولانی شدیدتر است. ذکر این نکته ضروری است که معمولا پرس وجوهای کاربران کوتاه می باشند، به طوری که متوسط طول پرس وجو های کاربران دو کلمه ذکر شده است [۸].

- دانش ناقص کاربران. کاربران در اغلب موارد، دانش کاملی در مورد موضوع مورد نظر خود ندارند، بنابراین پرس وجوهای کاربران معمولا مبهم اند و بیان ضعیفی از نیاز اطلاعاتی کاربر می باشند. پرس وجوهای ضعیف باز یابی ضعیف اطلاعات را به دنبال خواهند داشت.

۳- سیستم های باز یابی اطلاعات. علاوه بر مسائلی که تاکنون مطرح شد، در باز یابی اطلاعات کلید واژه ای تعداد نتایج زیاد است، در بین آنها موارد نامرتب و تکراری زیاد دیده می شود و نتایج مرتبط زیادی نیز باز یابی نمی شوند. بنابراین بعد از ارائه ی نتایج، کاربر باید آنها را بررسی کند و از بین آنها موارد مرتبط را پیدا کند [۹، ۱۰]. علاوه بر این، تعداد زیاد نتایج جستجو از یک طرف باعث سر در گمی کاربران می شود و از سوی دیگر باعث می شود کاربران تنها چند نتیجه ی اولیه ی باز یابی را بررسی کنند. در نتیجه ممکن است کاربر به سراغ سند مرتبطی که جزء آخرین نتایج باز یابی است، نرود.

برای روشن تر شدن مساله در این جا به ارائه ی یک مثال پرداخته می شود. پرس وجوی "Apple Growing" را در نظر بگیرید، همان طور که در شکل ۱-۲ نیز مشاهده می شود، در نتایجی که با موتور جستجوی گوگل باز یابی شده است، "Apple" به معنای سیب در نظر گرفته شده است. بنابراین اگر نیاز اطلاعاتی کاربر کامپیوترهای اپل باشد با این پرس وجو به نتیجه ای نمی رسد و همانطور که در شکل ۱-۳ نیز مشاهده می شود کاربر مجبور است پرس وجوی خود را تغییر دهد. با این مثال مشخص می شود که نتایج جستجوی کلید واژه ای حساس به واژه های به کار رفته در پرس وجوی کاربر است.