

رسالة محمد

دانشگاه یزد
دانشکده مهندسی برق و کامپیوتر
گروه مهندسی کامپیوتر

پایان نامه
برای دریافت درجه کارشناسی ارشد
مهندسی فناوری اطلاعات - شبکه‌های کامپیوتری

خزش و رتبه‌بندی کارا مبتنی بر ویژگی‌های گراف وب

استاد راهنما:
دکتر علی محمد زارع‌بیدی

استاد مشاور:
دکتر ولی درهمی

پژوهش و نگارش:
محمدامین گلشنی

اسفندماه ۱۳۹۰

**این پایان نامه با حمایت های مالی
مرکز تحقیقات مخابرات ایران
به انجام رسیده است.**

تقدیم به:

کوهران پاک‌پانی که روشنایی بخش دیدگانم در تمام زندگی ام بودند

پدرم و مادرم

آنها که گرمای وجودشان پشتیبان من

در سردترین روزگار ان بودند

برادران فداکارم

سپاسگزاری

این ناخیز مجموعه ای است که با الطاف و عنایات حضرت احدیت به سرانجام رسیده که لازم میدانم دست شکر به آستانش کشوده، رحمت بی کرانش را بستیم تا مزید ابر رحمانی او را به نظاره. نشینم و بر سبیل ((من لم یشکر المخلوق لم یشکر الخالق)) از زحمات بی دریغ استاد را بنمایم دکتر علی محمد زارع بیدکی که در طی این مدت بارها بنیانی های ارزنده و صبر و شکیلی خود، مرا به مسیر افتخار علم و پژوهش هدایت نموده اند کمال شکر را داشته باشم و نیز از اساتید بزرگوار جناب دکتر ولی دربی (استاد مشاور)، جناب دکتر محمد قاسم زاده (داور داخلی) و جناب دکتر سید ابوالفضل شاهزاده فاضلی (داور خارجی) که در نیل بدین سویاریم نمودند شکر و قدر دانی نمایم. در پایان از کلمه عزیزانی که مراد انجام این پایان نامه یاری نمودند صمیمانه شکر می نمایم. امید است در سایه توجهات ساحت اقدس الهی سهامی دلمان نوری از منبع لایزال الهی گیرد.

چکیده

امروزه وب جهان گستر به عنوان بهترین محیط برای تولید اطلاعات، انتشار و دسترسی به دانش مورد نیاز کاربران تبدیل شده است. نیاز به اشتراک گذاری داده‌ها و هزینه پایین تولید در دسترس قرار دادن صفحات وب، منجر به پویایی وب شده به طوری که رشد بی‌رویه ساختار و پوشش متنوع موضوعات خود را نمایان می‌سازد. نکته حائز اهمیت در این محیط روش یافتن اطلاعات مورد نیاز کاربر می‌باشد. آمارهای منتشر شده حاکی از آن است که اکثر کاربران اینترنتی، از موتورهای جستجو به عنوان درگاه دسترسی و جستجو در این حجم عظیم از اطلاعات استفاده می‌کنند. با توجه به وابستگی بالای کاربران به موتور جستجو، این ابزار به جزئی لاینفک از وب تبدیل شده و در حال حاضر بر اساس فناوری‌های موجود می‌توان ادعا کرد جایگزین مناسب‌تری جهت استخراج اطلاعات از این محیط وسیع و غنی وجود ندارد.

در این پایان نامه جهت پوشش مناسب صفحات وب (پیدا کردن سریع صفحه‌ها مهم)، الگوریتم خزشی به نام ¹IECA مبتنی بر چندین ویژگی نظیر درجات خروجی، درجات ورودی، فاصله لگاریتمی و خاصیت ساختاری گراف وب ارائه گردید. جهت ارزیابی دقت الگوریتم IECA از چهار گراف مختلف وب ایران، ایتالیا، انگلستان و دانشگاه برکلی استفاده می‌گردد. نتایج آزمایشات حاکی از کاراتر بودن الگوریتم IECA نسبت به سایر الگوریتم‌های بررسی شده می‌باشد. همچنین یک روش رتبه‌بندی مبتنی بر انتشار² امتیاز محبوبیت³ صفحات وب و اطلاعات موجود در آدرس صفحات وب (تعداد اسلش‌های، "/"، موجود در آدرس صفحه) ارائه و جهت بررسی کارایی الگوریتم ارائه شده مجموعه داده LETOR 3 مورد استفاده قرار گرفت. بر اساس نتایج آزمایشات استفاده از آدرس صفحات در عمل انتشار اهمیت صفحات در افزایش دقت فرآیند رتبه‌بندی موثر می‌باشد.

کلمات کلیدی: موتور جستجو، خزشگر، گراف وب، رتبه‌بندی، انتشار وابستگی، تعداد اسلش‌های موجود در آدرس

¹ Intelligent Effective Crawling Algorithm

² Propagation

³ Popularity

فهرست مطالب

عنوان	صفحه
۱ مقدمه.....	۱
۱-۱ شرح پایان نامه.....	۲
۲-۱ مقدمه.....	۲
۳-۱ ساختار وب.....	۵
۴-۱ موتور جستجو در یک نگاه.....	۶
۵-۱ رفتار خزشگر.....	۱۱
۶-۱ کارکرد خزشگر.....	۱۲
۷-۱ چالش‌های خزش.....	۱۳
۸-۱ اهداف پایان نامه.....	۱۵
۹-۱ ساختار پایان نامه.....	۱۶
۲ مروری بر کارهای انجام شده.....	۱۷
۱-۲ مقدمه.....	۱۸
۲-۲ رتبه‌بندی.....	۱۹
۱-۲-۲ رتبه‌بندی بر مبنای متن.....	۲۰
۱-۱-۲-۲ مدل فضای برداری.....	۲۰
۲-۱-۲-۲ مدل احتمالی.....	۲۱
۲-۲-۲ رتبه‌بندی مبتنی بر اتصال.....	۲۳
۱-۲-۲-۲ رتبه‌بندی مستقل از پرس‌وجو.....	۲۳
۲-۲-۲-۲ رتبه‌بندی وابسته به پرس‌وجو.....	۲۷
۳-۲-۲ رتبه‌بندی ترکیبی.....	۲۹
۳-۲ خزش.....	۳۱
۱-۳-۲ خزش دسته‌ای.....	۳۱

۳۲ ۱-۱-۳-۲ عرض اول
۳۲ ۲-۱-۳-۲ درجه ورودی
۳۳ Partial PageRank الگوریتم ۳-۱-۳-۲
۳۶ OPIC الگوریتم ۴-۱-۳-۲
۳۷ FICA الگوریتم ۵-۱-۳-۲
۳۹ IECA الگوریتم ۶-۱-۳-۲
۳۹ ۲-۳-۲ خزشگرهای موضوعی
۴۰ ۳-۳-۲ خزشگرهای موازی
۴۰ ۱-۳-۳-۲ چالش‌های خزشگرهای موازی
۴۰ ۲-۳-۳-۲ مزایای خزشگرهای موازی
۴۰ ۳-۳-۳-۲ معماری خزشگرهای توزیع شده
۴۱ ۴-۳-۳-۲ روش‌های خزش در تخصیص استاتیک
۴۲ ۵-۳-۳-۲ روش‌های بخش‌بندی وب
۴۲ ۴-۳-۲ صفحات نامطلوب
۴۴ ۵-۳-۲ خزش عمیق وب
۴۶ ۲ ارائه و ارزیابی الگوریتم IECA
۴۷ ۱-۳ مقدمه
۴۷ ۲-۳ الگوریتم IECA
۴۹ ۱-۲-۳ ویژگی ساختاری وب
۵۵ ۳-۳ نتایج آزمایشات
۶۱ ۴-۳ بررسی صحت الگوریتم
۶۱ ۵-۳ پیچیدگی IECA
۶۳ ۴ ارائه و ارزیابی الگوریتم رتبه‌بندی انتشاری
۶۴ ۱-۴ مقدمه

۶۴ مدل رتبه‌بندی مبتنی بر انتشار	۲-۴
۶۵ Hyperlink-based Slash-Score propagation method (HSS)	۱-۲-۴
۶۷Hyperlink-based Slash-Term propagation method (HST)	۲-۲-۴
۶۷ نتایج ارزیابی	۳-۴
۶۷ مجموعه داده آزمایشی استفاده شده	۱-۳-۴
۶۷ نحوه پیاده‌سازی	۲-۳-۴
۶۹ نتایج آزمایشات	۳-۳-۴
۶۹ معیارهای ارزیابی استفاده شده	۴-۳-۴
۷۵ نتیجه‌گیری و کارهای آینده	۵
۷۶ نگاهی کلی به مطالب ارائه شده	۱-۵
۷۷ کارهای آینده	۲-۵
۷۸ پیوست‌ها	۶
۷۹ پیوست الف. معرفی یک الگوریتم جدید خزش صفحات وب (FICA+)	۱-۶
۸۳ پیوست ب. لیست مقالات ارائه شده	۲-۶

فهرست شکل‌ها

عنوان	صفحه
شکل ۱-۱: ساختار پایونی وب [۳-۵].....	۶
شکل ۲-۱: مثالی از نمایه سازی معکوس [۳].....	۷
شکل ۳-۱: مثالی از جستجوی موفق [۴].....	۹
شکل ۴-۱: مثالی از جستجوی ناموفق [۴].....	۱۰
شکل ۵-۱: منظر موتور جستجو از دید کاربر مانند وضعیت ستارگان از دید ما می‌باشد [۱۴]..	۱۱
شکل ۶-۱: معماری خزشگر [۴].....	۱۳
شکل ۱-۲: ساختار فصل دوم.....	۱۹
شکل ۲-۲: مثالی از PageRank [۱۰].....	۲۴
شکل ۳-۲: مثالی از گراف پایه که از گره‌های {۱و۲و۳و۴} شروع شده است [۱۱].....	۲۷
شکل ۴-۲: مثالی از hub و authority [۳].....	۲۸
شکل ۵-۲: رتبه بندی گراف بر اساس درجه ورودی.....	۳۳
شکل ۶-۲: تاثیر مرتب سازی صف بر کیفیت صفحات بارگذاری شده [۳۲].....	۳۴
شکل ۷-۲: مقایسه الگوریتم‌های خزش [۱۶].....	۳۷
شکل ۱-۳: فاصله لگاریتمی در گراف خزش [۱۶].....	۴۸
شکل ۲-۳: میانگین امتیاز PageRank صفحات در الگوریتم عرض اول [۵۳].....	۵۰
شکل ۳-۳: نمونه ای از گراف وب.....	۵۱
شکل ۴-۳: رابطه ی بین δ_i و صفحات خزش شده.....	۵۳
شکل ۵-۳: گراف وب دانشگاه برکلی_۲۰۰۸ (چهار میلیون صفحه وب).....	۵۸
شکل ۶-۳: گراف وب ایتالیا_۲۰۰۴ (ده میلیون صفحه وب).....	۵۹
شکل ۷-۳: گراف وب انگلستان_۲۰۰۵ (شش میلیون و دویست و پنجاه هزار صفحه وب).....	۵۹
شکل ۸-۳: گراف وب انگلستان_۲۰۰۵ (دوازده میلیون صفحه وب).....	۶۰
شکل ۹-۳: گراف وب ایران_۲۰۱۰ (سه میلیون صفحه وب).....	۶۰

- شکل ۴-۱: مثالی از گراف وب..... ۶۶
- شکل ۴-۲: فلوچارت ساخت مجموعه کاری..... ۶۸
- شکل ۴-۳: مقایسه بین الگوریتم‌ها بر اساس P@n..... ۷۲
- شکل ۴-۴: مقایسه بین الگوریتم‌ها بر اساس NDCG@n..... ۷۳
- شکل ۴-۵: مقایسه بین الگوریتم‌های بر اساس P@n..... ۷۳
- شکل ۴-۶: مقایسه بین الگوریتم‌های بر اساس NDCG@n..... ۷۴

فهرست جداول

عنوان	صفحه
جدول ۱-۱: سهم موتور جستجوهای اصلی در پاسخ به پرس‌جوهای کاربران در امریکا [۲].	۲
جدول ۱-۲: متغیرهای استفاده شده در روش BM25	۲۲
جدول ۲-۲: انواع مدل‌های ترکیبی موجود [۲۹].	۳۱
جدول ۳-۲: دسته بندی الگوریتم‌ها بر اساس هدف آنها	۳۹
جدول ۱-۳: مقایسه بین الگوریتم‌های خزش بر روی گراف وب انگلستان	۵۷
جدول ۲-۳: ضریب کندال الگوریتم‌های خزش در مقایسه با الگوریتم پایه	۶۱
جدول ۱-۴: لیست الگوریتم‌های انتشاری به همراه مخفف‌ها.	۶۹
جدول ۲-۴: بالاترین کارایی در دو مدل HS و HT در TREC-2003	۷۰
جدول ۳-۴: بالاترین کارایی در دو مدل HS و HT در TREC-2004	۷۱
جدول ۴-۴: بالاترین کارایی در مدل‌های HS، HT، HSS و HST در TREC-2003	۷۱
جدول ۵-۴: بالاترین کارایی در مدل‌های HS، HT، HSS و HST در TREC-2004	۷۱
جدول ۶-۴: مقایسه بین الگوریتم‌ها بر اساس NDCG در TREC-2003	۷۲
جدول ۷-۴: مقایسه بین الگوریتم‌ها بر اساس NDCG در TREC-2004	۷۲

فهرست علائم اختصاری

عبارت انگلیسی	مخفف عبارت
Average Precision	AP
Hypertext-Induced Topic Search	HITS
Hyperlink-based Score Propagation	HS
Hyperlink-based Slash-Score Propagation	HSS
Hyperlink-based Slash-Term Propagation	HST
Hyperlink-based Term Propagation	HT
Inverse Document Frequency	IDF
Term Frequency	TF
Topic-Sensitive PageRank	TSPR
Uniform Out-link	UO
Weighted In-link	WI
Weighted Out-link	WO

فصل اول

مقدمه

۱-۱ شرح پایان نامه

موضوع پایان نامه به خزش و رتبه بندی صفحات وب اختصاص داده شده است. خزش و رتبه بندی از بخش های مهم در موتورهای جستجو می باشند که به ترتیب وظیفه ی بارگذاری صفحات با کیفیت و رتبه بندی نتایج بازگشتی به کاربر را بر عهده دارند. تا کنون الگوریتم های خزش متنوعی پیشنهاد شده اند، اما یا بازدهی مناسبی نداشته اند یا از پیچیدگی بالایی برخوردار هستند. لذا الگوریتم خزش کارایی به نام IECA (Intelligent Effective Crawling Algorithm) پیشنهاد می گردد که نسبت به الگوریتم های فعلی از کارایی بالاتر و پیچیدگی پایین تری برخوردار می باشد. در روش پیشنهادی اهمیت صفحات بر اساس چندین ویژگی نظیر درجات ورودی، فاصله لگاریتمی و خاصیت ساختاری گراف وب (بالا بودن درجه ورودی در صفحات با کیفیت) تعیین می گردد. نتایج آزمایشات بر روی چهار گراف مختلف وب (ایران، انگلستان، ایتالیا و گراف وب دانشگاه برکلی) حاکی از کارا تر بودن الگوریتم پیشنهادی نسبت به سایر الگوریتم های خزش می باشد. همچنین بخشی از پایان نامه به مبحث رتبه بندی اختصاص گرفته و مدلی مبتنی بر انتشار شامل دو الگوریتم به نام های Hyperlink-based Slash-Score propagation (HSS) و Hyperlink-based Slash- Term propagation (HST) جهت تعیین رتبه ی صفحات پیشنهاد گردید. جهت ارزیابی مدل رتبه بندی پیشنهادی LETOR 3 مورد استفاده قرار گرفت. بر اساس آزمایشات دو الگوریتم پیشنهادی نسبت به سایر الگوریتم های رتبه بندی انتشاری از دقت بهتری برخوردار هستند.

۲-۱ مقدمه

امروزه با توجه به رشد روز افزون اطلاعات و محتوای موجود در وب و همچنین تغییرات زیاد در اطلاعات موجود، موتورهای جستجو نقش مهمی در بازیابی اطلاعات از اینترنت ایفاء می نمایند. قابل ذکر است که حدود ۸۰٪ از افراد از طریق موتورهای جستجو به سایت ها و اطلاعات مورد نظرشان دسترسی پیدا می کنند [۱].

در روزهای اولیه ی پیدایش وب به دلیل کم بودن اسناد، تعیین محل دسترسی آنها بدون نیاز به موتور جستجو (به روش دایرکتوری مثل ساختار فعلی موتور جستجوی Yahoo) امکان پذیر بود اما امروزه با توجه به رشد چشمگیر وب، این کار امکان پذیر نیست. لذا موتورهای جستجو تبدیل به ابزاری مهم جهت بازیابی اطلاعات از این محیط پهناور شده اند به طوری که ماهانه حدود ۹/۶ میلیارد جستجو توسط موتورهای جستجوی اصلی در ایالت متحده امریکا انجام می شود [۲]. جدول ۱-۱ سهم هر موتور جستجو در انجام پرس و جوهای کاربران و تعداد پرس و جوهای دریافتی از کاربران را نشان می دهد.

جدول ۱-۱: سهم موتور جستجوهای اصلی در پاسخ به پرس و جوهای کاربران در امریکا [۲].

موتور جستجوهای اصلی	سهمی که هر موتور جستجو در پرس و جوهای
---------------------	---------------------------------------

دریافتی داشته (درصد)			
مقایسه بین دو ماه	دسامبر ۲۰۰۷	نوامبر ۲۰۰۷	
			کل پرس وجوهای انجام شده
٪۰/۰	٪۱۰۰	٪۱۰۰	موتور جستجوی Google
-٪۰/۲	٪۵۴/۴	٪۵۸/۶	موتور جستجوی Yahoo
٪۰/۵	٪۲۲/۹	٪۲۲/۴	موتور جستجوی Microsoft sites
٪۰/۰	٪۹/۸	٪۹/۸	موتور جستجوی Time Warner Network
٪۰/۱	٪۴/۶	٪۴/۵	موتور جستجوی ASK
-٪۰/۳	٪۴/۳	٪۴/۶	

هدف اصلی موتور جستجو ارائه‌ی نتایج مرتبط و با کیفیت با توجه به پرس‌وجوی ارسالی توسط کاربر است. برای یک موتور جستجو مشکلاتی وجود دارد که باید از عهده آنها برآید، در ادامه به این مشکلات اشاره شده است [۳].

۱) اولین مشکل مشخص کردن صفت "مرتبط" می‌باشد که باعث شده موتورهای جستجو اهمیت زیادی به الگوریتم‌های رتبه‌بندی بدهند. بازبایی صفحه‌ی مرتبط با پرس‌وجوی کاربر در نتایج نخست و بخصوص در صفحه‌ی اول بسیار پر اهمیت است چرا که معمولاً کاربران کمی عجول بوده و فقط ۱۰ نتیجه‌ی اول را بررسی می‌کنند.

۲) مسئله‌ی دیگری که مطرح می‌شود مسئله‌ی پهنش رتبه^۱ می‌باشد. با توجه به راه‌اندازی تجارت الکترونیک در محیط وب برای شرکت‌های تجاری مهم است که در نتایجی که موتور جستجو به کاربر نشان می‌دهد جزء نتایج اولیه باشند. لذا برای رسیدن به این هدف به پهنش رتبه رو می‌آورند، پهنش رتبه به دو صورت تغییر محتوا و یا افزایش تعداد پیوندها انجام می‌گیرد. در روش تغییر محتوا دسته‌ای از کلمات خاص مدام در صفحه تکرار می‌شوند که ممکن است رنگ فونت اختصاص داده شده به این دسته از کلمات هم‌رنگ با رنگ پس زمینه باشد یا این دسته از کلمات دارای فونت کوچکی باشند که کاربر متوجه آنها نشود.

۳) مسئله‌ی دیگر که در رابطه با الگوریتم‌های مبتنی بر پیوند مطرح می‌شود، مسئله‌ی غنی‌تر شدن اغنیا^۲ می‌باشد. یعنی اینکه صفحات محبوب در صدر نتایج قرار می‌گیرند و صفحات محبوب تازه متولد شده که فعلاً کسی آنها را ندیده و اشاره‌ای به آنها نمی‌شود در دید کاربران قرار نمی‌گیرند، در نتیجه پیوند به صفحات محبوب بیشتر شده و محبوب‌تر می‌شوند.

¹ Spam

² Rich get richer

۴) مسئله دیگری که در رابطه با موتور جستجو مطرح می‌شود، مسئله‌ی پوشش^۱ و تازگی^۲ است. این دو معیار در جهت مخالف یکدیگر عمل می‌کنند چرا که اگر خزشگر هر چه صفحات بیشتری بارگذاری کند کمتر می‌تواند به تازگی صفحات بارگذاری شده بپردازد. بزرگترین موتور جستجو در سال ۱۹۹۹ یعنی Northern Light در حدود ۱۶٪ از وب را تحت پوشش قرار می‌داد و در سال ۲۰۱۰ موتور جستجوی گوگل در حدود ۷۰٪ از وب را تحت پوشش قرار داده است.

جهت حل مسئله‌ی پوشش می‌توان از فراجویشگرها^۳ استفاده کرد به این صورت که پرس‌وجوی کاربر را به چندین موتور جستجو ارسال می‌کنند و نتایج بازگشتی را با یکدیگر ادغام می‌کنند اما چون روش‌های رتبه‌بندی موتور جستجوهای مختلف با یکدیگر متفاوت است ادغام نتایج بازگشتی از موتورهای جستجو کار راحتی نبوده و از آنجا که موتورهای جستجو دارای سرعت متفاوتی هستند، ممکن است عمل جستجو کند گردد.

۵) مشکل کلیدی دیگری که می‌توان به آن اشاره کرد انتظارات افراد مختلف از یک پرس‌وجو است مثلاً هدف یک محقق از پرس‌وجوی موتور جستجو می‌تواند مقالات علمی باشد در حالی که هدف یک دانش‌آموز از یک پرس‌وجو می‌تواند یافتن اطلاعاتی ساده در مورد موتور جستجو باشد.

۶) مسئله‌ی دیگر مبهم بودن یکسری کلمات استفاده شده در پرس‌وجو می‌باشد، در مواقعی که کاربر نتواند درخواست خود را با چندین عبارت بیان کند. به عنوان مثال کلمه Python را در نظر بگیرید که دارای سه معنی مختلف، غیبگو، افعی و زبان برنامه‌نویسی است.

مطابق با مباحث ذکر شده، یکی از بخش‌های مهم موتور جستجو، بخش رتبه‌بندی می‌باشد. رتبه‌بندی فرآیندی است که کیفیت یک صفحه از لحاظ ارتباط با پرس‌وجوی کاربر توسط موتور جستجو تخمین زده می‌شود. با توجه به اینکه به ازای هر پرس‌وجوی کاربر معمولاً هزاران صفحه‌ی مرتبط وجود دارد، لازم است آنها را اولویت بندی کرده و بهترین‌ها به کاربر نشان داده شود. مسائلی مانند حوصله کم کاربران (مشاهده‌ی ۱۰ تا ۲۰ نتیجه‌ی اول)، کم بودن کلمات پرس‌وجو، حجم زیاد اطلاعات، فرآیند رتبه‌بندی را با مشکلاتی جدی روبه‌رو می‌کند [۳]. در حال حاضر سه روش مهم رتبه‌بندی مبتنی بر محتوا (استفاده شده در روش‌های سنتی بازیابی اطلاعات)، مبتنی بر ساختار (پیوندهای بین صفحات وب) و ترکیبی از محتوا و ساختار وجود دارد. در روش‌های استفاده شده در بخش سنتی بازیابی اطلاعات مدل‌هایی مانند بولی^۴، احتمالی و فضای برداری^۵ جهت رتبه‌بندی اسناد بر مبنای محتوای آنها ارائه شده است. در فصل دوم مروری بر الگوریتم‌های رتبه‌بندی انجام خواهد گرفت.

¹ Coverage

² Freshness

³ Meta-Search engines

⁴ Boolean

⁵ Polysemy

۳-۱ ساختار وب

می‌توان وب را به صورت یک گراف در نظر گرفت که در آن صفحات به منزله‌ی گره‌های گراف و پیوند بین صفحات به منزله‌ی یال‌های موجود در گراف می‌باشند. بر اساس تحقیقات انجام شده بیشتر از ۹۰٪ وب دارای ساختاری مانند شکل ۱-۱ است که از چهار قسمت اصلی تشکیل شده است [۳-۵].

(۱) هسته‌ی مرکزی^۱. شامل اسنادی است که به هم متصل‌اند و تمام اسناد موجود در این بخش قابل دسترسی‌اند.

(۲) بخش IN. این بخش شامل اسنادی است که از بخش هسته‌ی مرکزی به آنها پیوندی وجود ندارد اما از این بخش به هسته‌ی مرکزی پیوند وجود دارد. می‌توان گفت که در این بخش اسنادی قرار می‌گیرند که یا تازه متولد شده‌اند، یا ناشناخته‌اند یا اسنادی‌اند که دارای اهمیت زیادی نمی‌باشند. اما این نکته نباید از دید پنهان باشد که این اسناد، نقاط خوبی برای شروع عمل خزش^۲ می‌باشند.

(۳) بخش OUT. این بخش متشکل از اسنادی است که از هسته‌ی مرکزی به آنها پیوند وجود دارد اما از اسناد موجود در این بخش هیچ پیوندی به اسناد موجود در هسته‌ی مرکزی وجود ندارد. به عبارتی می‌توان گفت اسناد قرار گرفته در این بخش، اسنادی‌اند که شامل پیوند به یکدیگر می‌باشند. مثلاً یک وب سایت که از صفحات دیگر به آنها اشاره شده است اما آنها به صفحات دیگری اشاره نکرده‌اند و فقط به یکدیگر پیوند دارند.

(۴) بخش دیگر پیچک‌ها^۳ می‌باشند که جزء بخش IN یا بخش OUT می‌باشند اما به هسته‌ی مرکزی پیوندی ندارند. قسمت دیگر لوله یا تونل^۴ می‌باشد که به اسنادی گفته می‌شود که در بخش IN بوده و به بخش OUT متصل شده‌اند و پیوندی به بخش هسته‌ی مرکزی ندارند و در نهایت آخرین قسمت جزایرند^۵ که توسط هیچ یک از اسنادی که در این ساختار پاپیونی شکل وجود دارد، پیوندی به آنها وجود ندارد.

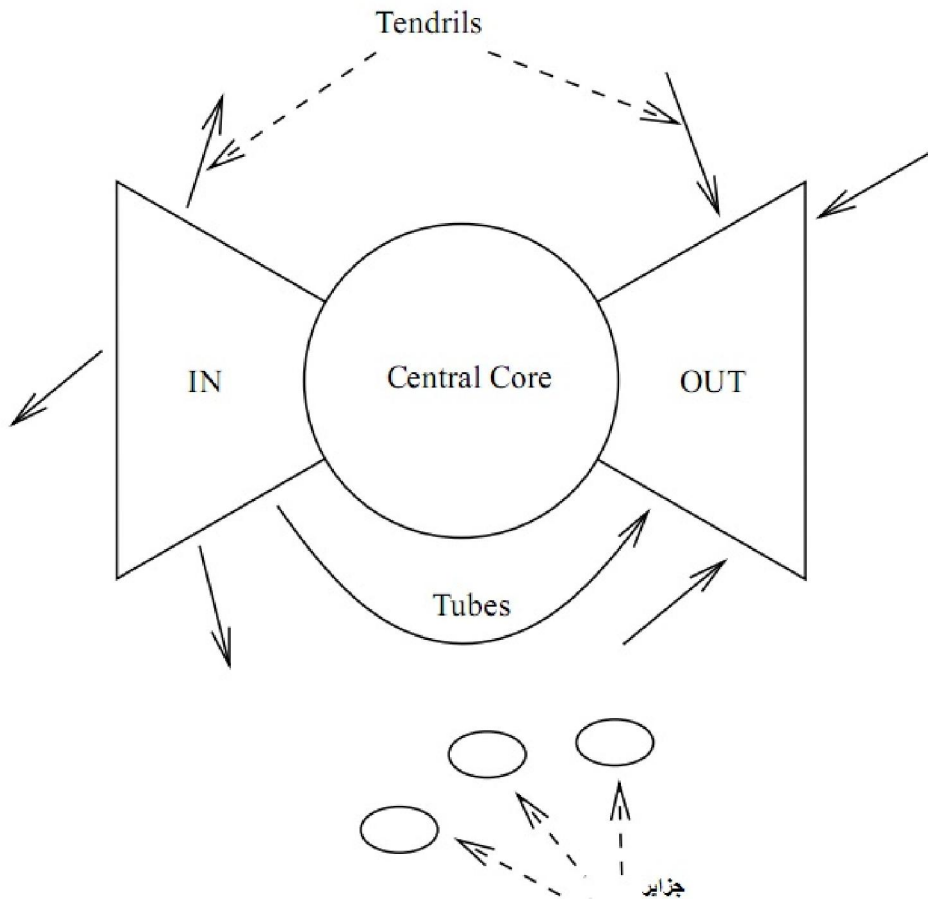
¹Central Core

²Seed

³Tendrils

⁴ Tube

⁵ Island



شکل ۱-۱: ساختار پایونی وب [۳-۵]

۴-۱ موتور جستجو در یک نگاه

موتور جستجو از چهار قسمت اصلی تشکیل شده است (که ممکن است این بخش‌ها با یکدیگر ادغام شوند یا به بخش‌های بیشتری تقسیم شوند) [۳].

- خزشگر^۱ - ربات^۲، عنکبوت^۳، راه‌رونده^۴ یا سرگردان^۵

برنامه‌ای است که وظیفه‌ی بازیابی صفحات از وب و ذخیره کردن آنها را در مخزن بر عهده دارد [۳, ۴, ۶, ۷]. خزشگر دو هدف تازه‌سازی و پوشش بالارا به دنبال دارد که با

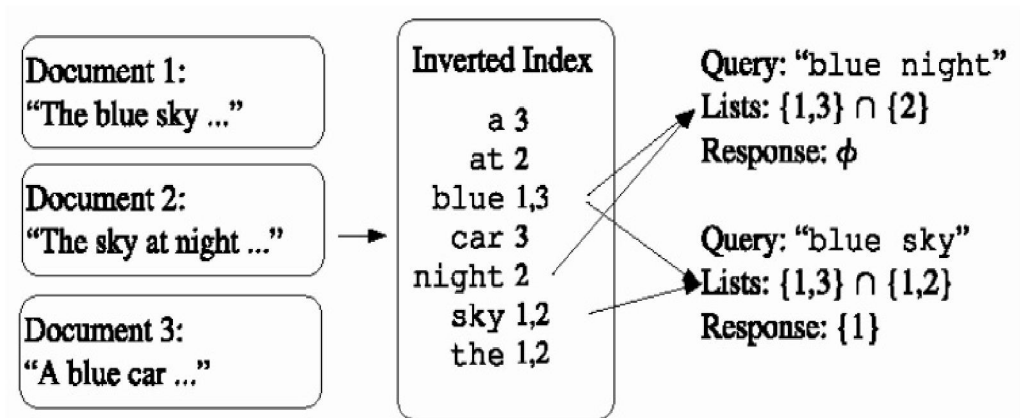
¹Crawler
²Robots
³Spider
⁴Walker
⁵wanderer

یکدیگر رابطه عکس دارند، بعد از اتمام عمل پیمایش علاوه بر صفحات، گراف وب نیز قابل استخراج می‌باشد.

• نمایه‌ساز

این واحد اسناد ذخیره شده در مخزن‌ها را پردازش کرده و نمایه سازی^۱ می‌کند. جهت بالا بردن سرعت دسترسی از نمایه‌سازی معکوس استفاده می‌شود. در نمایه سازی معکوس به جای اینکه مشخص شود که یک سند شامل چه کلماتی می‌باشد، مشخص می‌گردد که یک کلمه در چه سندهایی ظاهر شده است. نمایه‌ساز معکوس از دو بخش واژگان، که شامل تمام کلمات موجود در همه‌ی اسناد و لیست مکانی آنها شامل تمام اسنادی که آن کلمه در آنها ظاهر شده است می‌باشد.

شکل زیر نمایه‌ساز معکوس را نشان می‌دهد که شامل حروف اضافه نیز می‌باشد.



شکل ۱-۲: مثالی از نمایه سازی معکوس [۳].

• موتور بازیابی^۲

موتور بازیابی رابطه‌ی مستقیم با دو بخش نمایه‌سازی و واسط کاربر دارد. وظیفه‌ی اصلی این واحد رتبه‌بندی جواب‌ها با استفاده از نتایج تهیه شده توسط واحد نمایه‌سازی و گراف تهیه شده توسط خزشگر می‌باشد. با توجه به اینکه کاربر پرس‌وجوی را به زبان طبیعی بیان می‌کند از کارهای مهم واحد پردازش پرس‌وجو، عملیات پردازش زبانی پرس‌وجو مانند نرمال سازی^۳ و بسط پرس‌وجو می‌باشد. از وظایف دیگر این واحد، انجام عملیاتی مانند مدل کردن کاربر، و نیز اعمال بازخورد^۴ خواهد بود. بعد از پردازش پرس‌وجو و بدست

¹ index

² Retrieval Engine

³ Normalization

⁴ Relevance Feedback