

فهرست مطالب :

۱	مقدمه
۳	پیشینه پژوهش
۵	روش شناسی
۹	یافته ها
۱۰	امتیاز تنظیم بر اساس ارتباط
۱۱	جستجو در اینترنت
۱۲	انواع موتورهای جستجو
۱۳	نحوه کار موتورهای جستجو
۱۵	عوامل مهم در انتخاب موتور جستجو
۱۷	دسته بندی موتورهای جستجو
۱۸	بررسی یک موتور جستجوی پیمایشی
۲۰	مکان و تکرار
۲۱	عوامل خارج از صفحه
۲۱	نتیجه گیری ها
۲۴	مهمترین موتورهای جستجو
۲۵	منابع و پی نوشت ها

مقدمه

رشد اینترنت شگفت‌آور شده است. با توجه به تحقیق میدانی در سال ۱۹۹۶ این رشد تصاعدي همچنان ادامه دارد. و تخمین زده شده که شبکه از نظر اندازه و حجم هر ۱۲ تا ۱۵ ماه دوبرابر می‌شود. بطور تقریبی ۱۰۰/۰۰۰ وب‌گاه در اگوست ۱۹۹۵ وجود داشته و این تعداد در اگوست ۱۹۹۶ به ۵۳۶/۰۴۱ رسیده است. از آنجائی که هر پایگاه می‌تواند بسیاری از صفحات وب را در خود داشته باشد این باعث می‌شود که تعداد بیشتری از صفحات وب بوجود آید. در حالیکه کیفیت بسیاری از صفحات ممکن است مورد سؤال باشد و نگهداری بسیاری از صفحات پراکنده است. اما صفحات معتبری هم وجود دارد که اطلاعات با ارزشی در مورد بسیاری از موضوعات ارائه می‌دهد. آنچه استفاده‌کنندگان به آن نیاز دارند يك سیستم جداکننده است که مو از ماست بکشد.

تلاشهای زیادی صورت گرفته که این وظیفه را آسان کند: کتابداران و متخصصان موضوعی راهنماهای موضوعی را گردآوری کرده‌اند. «اخبار کتابخانه‌های تحقیقاتی و دانشکده‌ای» ۱ بطور دوره‌ای راهنماهای منابع اینترنتی را در موضوعات انتخابی منتشر می‌کند. علاوه بر آن را (که يك فهرست آزمایشی برای اینترنت است) بوجود آورده‌اند. Intercat فهرست نویسان را برای نمایه‌سازی منابع اینترنت با سرعنوانهای موضوعی Netfirst نیز پایگاه OCLC کتابخانه‌کنگره و اعداد طرح دهدهی دیویی ایجاد کرده است. کتابداران مرجع يك پایگاه وب را بوجود آورده‌اند که مرور منابع اینترنتی را به اشتراك بگذارند و کتابداران Infofilter به نام رده‌بندی موضوعی را برای سازماندهی منابع اینترنت مورد آزمایش قرار داده‌اند. اما چیزی که بیشترین توجه استفاده‌کنندگان اینترنت را به خود معطوف داشته است، احتمالاً عنکبوتها و رباتهایی هستند که خدمات جستجو را انتخاب می‌کنند. برای بسیاری از جستجوگران اینترنت، این موتورها با راه دادن آنها به فضای اطلاعاتی عظیم کمک موقتی ارائه می‌دهند. کاربران اینترنت بزودی فهمیدند که این موتورها کامل و مناسب نیستند زیرا آنها منطقه جغرافیایی خاصی را پوشش می‌دهند، بصورت متفاوت نمایه‌سازی می‌شوند و منابع را با کلیدواژه‌ها بازیابی می‌کنند. استفاده‌کنندگان هرگز نمی‌توانند اعتماد کنند که يك جستجو جامع یا قطعی باشد. با وجود اینکه نتایج بازیابی ظاهراً بوسیله میزان ارتباط مرتب شده است و استفاده‌کنندگان مبارزه

با ریزش کاذب مواد تکراری و نامربوط را ادامه می‌دهند. در کل پیشرفت خوبی در کمک به استفاده‌کنندگان برای جهت‌یابی در اینترنت بوجود آمده، اما این ابزارها آنقدر زیاد هستند که انتخاب صحیح یکی از آنها کار سختی است.

علاوه بر آن انتخاب موتورهای جستجوی مناسب ممکن است برای استفاده‌کنندگان عمومی و نیز کتابداران، به علت رویه‌های پیچیده، مهمتر باشد. بعنوان مثال، اینفوسیک سرویس رایگان دارد، اما سرویسهای هزینه بر، یعنی متخصصان اینفوسیک، نمایه بزرگتر و قابلیت‌های جستجوی قویتر ارائه می‌دهند.

وب کروکر یک درخواست حق عضویت دارد که زمانی آن را تحمیل می‌کند. اگر این عمل در مقابل هزینه برای خدمات یک رویه شود، لازم خواهد بود برای استفاده‌کنندگان و بخصوص کتابداران که بدانند کدامیک از موتورهای جستجو را باید به خدمت بگیرند.

این تحقیق تلاش کرده که کارآیی موتورهای کاوش را در آدرس دهی نیازهای اطلاعاتی ارزیابی کند. آیا آنها می‌توانند جوابهایی برای سؤالات مرجع واقعی بازیابی کنند؟ آیا آنها منابع خوبی برای سؤالات موضوعی ارائه می‌دهند؟ آنها تا چه حد نتایج جستجو را براساس میزان ارتباط مرتب می‌کنند؟ کدامیک از موتورهای جستجو بهتر عمل می‌کنند؟ جواب این سؤالات به ما کمک خواهد کرد که نقاط ضعف و قوت موتورهای کاوش را بهتر بفهمیم و ما را قادر می‌سازد که برای برطرف کردن نیازهای اطلاعاتی موتور کاوش مناسب را انتخاب کنیم.

پیشینه پژوهش

Netsearch موتورهای کاوش بسیاری موجود هستند و براحتی در دسترس قرار می‌گیرند مربوط به کتابخانه کنگره، هر دو با هم `www by subject or keyword` متعلق به نت اسکپ و موتورهای کاوش اصلی را گرد آورده‌اند. همچنین ابرموتورهایی وجود دارد که به جستجوگران قادر Savy search اجازه می‌دهد که سریعاً به چندین موتور کاوش دسترسی پیدا کنند، اما فقط گروهی دیگر از ابزارهایی را "All-in-one" به جستجوی بیش از ۵ پایگاه در یک زمان است. " که ارائه دهنده نمونه‌های جستجو بسیاری از موتورهای کاوش در یک پایگاه وب برای آسان‌سازی عمل جستجو هستند را نشان می‌دهد.

کار عنکبوتها، روباتها و دیگر برنامه‌های خودکار بوسیله پروسیس خلاصه شده است. (۱۹۹۵)

و محدودیت‌های این ابزارها بوسیله کاستر تجزیه و تحلیل شده‌اند. (۱۹۹۵) چندین مقاله ادعا کرده‌اند که ابزارهای جستجوی اینترنت را ارزیابی کرده‌اند، اما اغلب آنها توصیف‌هایی از شبکه جهانی وب یا موتورهای کاوش ارائه داده‌اند. برینکلی و بیورک (۱۹۹۵) هایتلنت، آرچی، گوگر، و ایزو شبکه جهانی وب را شرح داده‌اند. کورتولیس، بیرواستارک (۱۹۹۵) پرسش‌هایی بکار برده‌اند برای اینکه موتورهای کاوش و نمایه‌های وب را آزمایش کنند. اما گزارش آنها بیشتر توصیفی است. کایمل (۱۹۹۶) تاریخچه‌ای از پایگاه‌های تولید شده بوسیله رباتها را ارائه داده و به جستجوگران مبتدی راهنمایی‌های خوبی در مورد موتورهای کاوش پیشنهاد کرده است. وندیتو (۱۹۹۶) هفت موتور کاوش را آزمایش کرده و ویژگی‌های جستجوی آنها را شرح داده است. گزارش‌های مشابه در مورد اینترنت فراوان است و بسیاری از آنها را می‌توان در لیت کمپبل پیدا کرد. نمونه‌های خوب بسیار کم هستند. مثل گزارش‌های لین (۱۹۹۵)، وین شیب (۱۹۹۵) و بستر و پانول (۱۹۹۵).

مطالعات ارزشیابی نسبتاً کمی وجود داشته که شامل تحقیق میدانی در مورد موتورهای جستجو باشد. دسای (۱۹۹۵) قدرت بازیابی سیزده ابزار جستجو را با یک سؤال آزمایش کرده است. او قادر بود با جستجوی نامش تعیین کند که چگونه بسیاری از اسناد وبی او بازیابی شده است. اینفوسیک و لایکاس با بازیابی هفت سند از ۲۴ سند بهتر عمل کردند. در حالیکه دیگر موتورها و نمایه‌ها نظیر وب کرولر و یاهو ضعیف عمل کردند. لیتون (۱۹۹۵) عملکرد اینفوسیک، لایکاس، وب کرولر و ورد واید وب وارم را با استفاده از ۸ سؤال مقایسه کرده است. او این آزمایش را با ۴ معیار انجام داد-نسبت تکراری بودن، دقت، دقت کامل و حداکثر ۱۰ دقت- و نتیجه گرفت که لایکاس و اینفوسیک بهتر از بقیه عمل کردند. پکروتوماپولو ۲ سؤال مرجع را در آلتاویستا، ماژلان، اینفوسیک، لایکاس و یونیت جستجو کردند. امتیاز دقت آنها مبتنی بود بر ۱۰ نتیجه اول، آنها فهمیدند که آلتاویستا بهترین عملکرد را داشت بعد از آن اینفوسیک، لایکاس، ماژلان و پوینت. مقاباب (۱۹۹۵) ۵ سؤال را برای امتحان کردن ۵ موتور کاوش با اجرای جستجوهای اصلی و اصلاح شده در هر موتور بکار برد. او دقت نتیجه اول را تغییر داد. با استفاده از ۲۵ نتیجه بعنوان پایه و مبنا. او یاهو را بعنوان بهترین عمل کننده شناخت. یافته‌های این مطالعات بطور قطعی بهترین موتورهای کاوش را مشخص نکرد بعلاوه سؤالات مختلف، تعداد متفاوت سؤالات و مقیاس‌های مختلف بکار رفته برای ارزشیابی. با این وجود، این مطالعات شیوه‌های مختلف

ارزیابی موتورهای کاوش را ثابت کرد و معیارهای جدید و منطقی برای اجرای جستجو پیشنهاد کرد.

روش‌شناسی

این مطالعه هشت موتور کاوش را که مشهور هستند و برای عموم رایگانند ارزیابی می‌کند. این موتورها عبارتند از: آلتاویستا، اکسایت، اینفوسیک، گاید، لیکاس، مازلان، این تکست، وب کرولر و ورد واید وب وارم.

عوامل بسیاری ممکن است برای موفقیت یک جستجو مشارکت کنند. درک صحیح از تقاضای جستجو، استراتژی جستجو، پایگاه اطلاعاتی، موتور کاوش، و قضاوت میزان ارتباط بوسیله جستجوگران. در این تحقیق تعدادی از این معیارها کنترل شده بودند بنحوی که تفاوت موتورهای کاوش قابل مشاهده بود. موتورهای کاوش از نظر اندازه، محتوای پایگاه اطلاعاتی‌شان، خطمشی‌های نمایه‌سازی‌شان، کنترل کیفیت، شیوه‌های بازیابی و ارائه نتایج جستجو متفاوتند.

فرض بر این شد که سوالات آزمایشی با پیچیدگی و وضع متفاوت در تعیین بهترین موتور جستجو با ارزش‌تر خواهد بود. با این وجود، مطالعات پیشین مشخص کردند که امکان اینکه یک موتور جستجو در جوابگویی همه نوع سوالات بهتر از همه باشد، وجود ندارد. ما تقاضاهای جستجو را با استفاده از ۲۰ سؤال مطرح شده در میز مرجع استاندارد کردیم. و ۵ سؤال موضوعی که در حوزه‌هایی که منابع اینترنتی بسیاری داشت، بوجود آمده بود-سرگرمی، تجارت، سیاست اقتصاد و بهداشت. سوالات مرجع گردآوری شده شامل سوالات تخصصی و پرسشهای موضوعی وسیع بود و از نظر اینکه توانایی موتورهای کاوش را در جوابگویی به سوالات مرجع واقعی مورد آزمایش قرار داد با ارزش بودند. ۵ سؤال موضوعی ساختگی بودند، اما این طراحی ما را قادر می‌ساخت که تجزیه و تحلیل معنی‌دار بیشتری انجام دهیم. با این همه، هیچ ارزیابی خیلی خردمندانه نبود اگر ما از سوالاتی که برای آن هیچ چیز قابل بازیابی نبود، استفاده می‌کردیم.

مجموعه داده‌ها ۲ به هشت موتور کاوش چهار جستجو اختصاص یافته بود برای اطمینان از اینکه

هر سؤال دو بار در يك موتور جستجو شده است. به جستجوگران آموزش داده شد که از مرورگر نت اسکپ برای دسترسی به اینترنت استفاده کنند و سؤالات داده شده را در موتورهای کاوش تعیین شده جستجو کنند و نسخه‌های چاپی از نتایج جستجو ارائه دهند. جستجوگران کار را در آوریل شروع کردند و در ژوئن ۱۹۹۶ نتیجه‌گیری انجام شد. جمله‌بندی سؤالات مرجع برای جستجوهای اینترنتی کمی تغییر داده شد. برای اطمینان یافتن از تکنیکهای مشابه بکار رفته، کلیدواژه‌ها تعیین شده بودند و پارامترهای اساسی در مورد اینکه در هر موتور کاوش جستجو به چه صورت انجام شود ارائه شدند. به جستجوگران گفته شد که بهترین قضاوتشان را در ارزیابی ارتباط منابع بازیابی شده بکار گیرند.

یکی از سؤالات مرجع دو بخش داشت، بنابراین سؤالات مرجع به ۲۱ تبدیل شد. هر کدامیک از سؤالات ۲۱ گانه مرجع و ۵ سؤال موضوعی دو بار در هر موتور، جستجو شده بود. اما در اکسایت ۴ بار جستجو شد، زیرا در این موتور، جستجوگر قادر به جستجوی کلیدواژه‌های و نیز جستجوی مفهومی می‌باشد. در میانه راه با مازلان جستجوگران گزینه‌ای از جستجو را در بخش خاصی از پایگاههای اطلاعاتی یا کل پایگاهها به منظور هماهنگی و یکدستی همه جستجوگران در کل پایگاهها انجام دادند. در کل ۴۶۸ جستجو انجام شده بود.

متغیرهای وابسته. چهار متغیر برای این تحقیق اندازه‌گیری شدند. «دقت» ۳ که بطور سنتی تعریف شده بود: تعداد منابع مرتبط بازیابی شده تقسیم بر تعداد منابع بازیابی شده و یک معیار استاندارد برای سیستمهای بازیابی اطلاعات بوده است. از آنجایی که ارزیابی ارتباط تعداد زیاد صفحات بازیابی شده بوسیله موتورهای کاوش غیرممکن بود. این متغیر در این تحقیق بطور عملیاتی چنین تعریف شده:

دقت: تعداد منابع مرتبط در ۱۰ گزینه اول

شیوه استفاده از ده گزینه اول قابل توصیه است، زیرا این گزینه‌ها بیشتر امکان دارد که به وسیله جستجوگران دیده شود. این معیار را لیتون، پیکروتومایولو بکار برده است. اما برخلاف تحقیق لیتون، این تحقیق پایگاههای ارجاعی و تکراری را در معیار دقت در نظر می‌گیرد. زیرا آنها بالقوه مفید بودند (در صورتی که گزینه‌های تکراری مرتبط باشد) و حذف آنها باعث می‌شد که پایه مقایسه (که ده تا بود) کوچکتر شود.

«تکراری بودن» ۴: در همان اوایل جستجو در موتورهای کاوش گزارشهای حکایت گونه‌ای از تکراریها بوجود آمد. در نتیجه این معیار در ارزیابی‌های ما وارد شد. تکراری بودن بطور عملیاتی «تعداد گزینه‌هایی که تکرار شدند و قبل از آن نیز ارائه شده بودند» تعریف شده بود. پایگاههای ارجاعی هم جزء تکراریها به حساب آمدند. در این مطالعه ما تعداد گزینه‌های تکراری را بر اساس ده نتیجه اول در نظر گرفتیم.

«امتیاز مرتب‌ترین گزینه» ۵: همه موتورهای کاوش انتخابی، نتایج ارزیابی را با استفاده از الگوریتم متفاوت مرتب می‌کنند و بهترین تطبیقها را اول ارائه می‌دهند. اما تنظیم همیشه مفید نبوده است. این متغیر برای امتحان کردن توانایی درجه‌بندی موتورهای کاوش طراحی شده بود. که بر این فرضیه مبتنی است که شیوه درجه‌بندی مؤثر، مرتب‌ترین گزینه‌ها را در بالاترین لیست نتایج جستجو قرار می‌دهد. جستجوگران، بطور عملیاتی مرتب‌ترین گزینه از بین ده گزینه تعریف کردند و به آن بخاطر جایگاهش يك امتیاز دادند. اگر این گزینه در اولین، دومین یا سومین گزینه بود این موتور امتیازی بین ۱ یا ۲ یا سه می‌گرفت. اگر این گزینه جای دیگر ظاهر می‌شد به آن امتیاز ۶ داده می‌شد. عدد ۶ به این علت انتخاب شده بود که نشان می‌داد این گزینه در خارج از اولین نیمه لیست ده تایی قرار گرفته است. امتیاز پائین در مورد مرتب‌ترین گزینه نشان دهنده این بود که آن موتور بهترین درجه‌بندی گزینه‌های مرتبط را داشته است. «امتیاز درجه‌بندی میزان ارتباط» ۶: این متغیر نیز درجه‌بندی میزان ارتباط در موتورهای کاوش را ارزیابی کرد اما به شیوه‌ای متفاوت. این متغیر به عنوان درصد گزینه‌های مرتبط که در اولین نیمه لیست ده گزینه‌ای ظاهر شدند تعریف شد. این تعریف مبتنی بود بر این فرضیه که میزان ارتباط گزینه‌ها کاهش خواهد یافت هر چقدر که به گزینه‌های پائین‌تر می‌رسیم. جستجوگران تعداد گزینه‌های مرتبط را در هر نیمه از لیست ده گزینه‌ای ثبت کردند و مأموران تحقیق این تعداد را با توجه به فرمول زیر برای رسیدن به امتیاز درجه‌بندی براساس میزان ارتباط تبدیل کردند:

تعداد گزینه‌های مرتبط در اولین لیست

کل تعداد گزینه‌های مرتبط در لیست ده‌تایی

«بازیابی» ۷: یکی دیگر از معیارهای استاندارد برای بازیابی اطلاعات است و چنین تعریف

شده: تعداد گزینه‌های مرتبط بازیابی شده تقسیم بر کل تعداد گزینه‌های مرتبط در يك فایل

اطلاعاتی. این معیار برای استفاده دشوار بود زیرا جستجوگران می‌بایست همه گزینه‌های مرتبط را در کل یک پایگاه یا فهرست شناسایی کنند. این مشکل در شبکه جهانی وب خیلی شدیدتر است. با هزاران هزار صفحه وب نمایه شده به وسیله موتورهای جستجو انتخابی غیرممکن بود که همه صفحات وب مرتبط با موضوع جستجو شناسایی شود. به این ترتیب بازیابی در این مطالعه استفاده نشد.

«تجزیه و تحلیل داده‌ها» ۸: از ۶۸ جستجوی انجام شده ۴ معیار برای هر جستجو در موتور کاوش ثبت شد. بسامد و میانگین این معیارها برای هر موتور جستجو با نوع سؤالات حساب شده بودند.

یافته‌ها

دقت

سؤالات مرجع عمومی متنوع بود بطوری که شاید یک کتابدار مرجع برای یافتن پاسخ آنها از اینترنت استفاده نمی‌کرد. با این وجود همه سؤالات در موتورهای کاوش جستجو شده بودند که توانایی‌شان را در پاسخگویی به سؤالات مرجع ارزیابی کنند. موتورهای کاوش این کار را خوبی انجام ندادند. میانگین امتیاز دقت خیلی پائین بود. بین ۰/۳۱ و ۲/۹۳. این تکست بالاترین تعداد گزینه‌های مرتبط را بازیابی کرد. بعد از آن آلتاویستا و اینفوسیک و سپس لایکاس با اختلاف کم چهارم شد. برای نشان دادن جنبه دیگری از این جستجوها، اطلاعاتی در مورد نقاط کور هر موتور در جدول ۱ قرار گرفتند که نشان داد اکسایت پائین‌ترین تعداد نقاط کور را داشت و بعد از آن این تکست و لایکاس. بر رویهم رفته، این دو مجموعه از داده‌ها این تکست را بعنوان بهترین موتور در برخورد با سؤالات مرجع معرفی کرد. این موتور صفحات وب را برای این سؤالات بازیابی کرد و نتایج جستجویش بالاترین امتیاز میزان دقت را داشت. در این پژوهش موتورهای جستجو با سؤالات موضوعی ساختگی بهتر عمل کردند. سؤالات موضوعی پائین‌ترین میانگین امتیاز دقت (۳/۲) نسبت به بالاترین امتیاز میزان دقت (۲/۹۳) در سؤالات مرجع واقعی بالاتر بود. اینفوسیک بهتر عمل کرد بعد از آن مازلان و این تکست و باز هم لایکاس با اختلاف کم چهارم شد. از آنجایی که سؤالات مرجع برای حوزه‌هایی طراحی شده

بودند که در مورد آن اطلاعات بیشتری در وب موجود باشد، مشکل نقاط کور در این سؤالات خیلی جدی نبود. در سؤالات مرجع کیفیت گزینه‌های بازیابی شده، «دقت خاص» تعداد جستجو‌هایی که بیش از ۵ گزینه مرتبط را بازیابی کردند در نظر گرفته شد که در جدول ۲ آمده است.

اینفوسیک باز هم برنده ظاهر شد، بعد از آن ماژلان و این تکست و لایکاس و وب کرولر هر سه بطور مساوی در جایگاه سوم قرار گرفتند. این اطلاعات نشان داد که اینفوسیک در برخورد با سؤالات بهترین بود. این موتور بیش از ۵ گزینه مرتبط را برای اغلب پرسش‌های موضوعی بازیابی کرد و نتایج جستجویش بالاترین امتیاز میزان دقت را داشت. تکراری بودن

تکراریها در بازیابی زمان جستجوگران را تلف می‌کنند و باعث سردرگمی می‌شوند. علاوه بر قصه شکایت در مورد تکراریها، به هر حال، این مسأله به نظر می‌رسد که در بیشتر موتورهای کاوش مطرح بوده است. میانگین تعداد موارد تکراری برای هر دو سؤالات مرجع و سؤالات موضوعی در هر موتور کاوش ناچیز است (کمتر از یک). اما سؤالات موضوعی شانسشان برای داشتن موارد تکراری بیشتر بود. این اطلاعات نشان می‌دهد که نمایش داده‌های تکراری حتی زمانی که گزینه‌های مرتبط زیادی بازیابی شده بود اهمیت چندانی نداشتند. امتیاز مرتبط‌ترین گزینه

این امتیاز توانایی هر موتور کاوش را برای نشان دادن اولین گزینهء مرتبط اندازه‌گیری کرد. بخاطر اینکه امتیازی به محل گزینه‌ها اختصاص یافته بود، پائین‌ترین امتیازها عملکردهای بهتر را نشان می‌داد. برای سؤالات مرجع امتیاز موتورهای جستجو بین ۳/۳ و ۵/۳ قرار داشت، اول این تکست بعنوان برنده بعد از آن اکسایت و آلتاویستا قرار گرفتند. موتورهای کاوش با سؤالات موضوعی خوب عمل نکردند. امتیازات آنها بین ۲/۵ تا ۴/۲ قرار داشت. این تکست و بعد از آن اینفوسیک و وب کرولر بهترین عملکرد را داشتند. این تکست در ارائه مرتب‌ترین گزینه همیشه بهترین بود.

امتیاز تنظیم براساس ارتباط

این امتیاز قدرت موتورهای کاوش را اندازه‌گیری کرد برای ارائه گزینه‌های مرتبط در اولین نیمه نتایج جستجو. برای سؤالات مرجع، امتیازات موتورهای کاوش بین ۱۰/۵% تا ۴۵/۱% و

با پیشتازی این تکست و بعد از آن اینفوسیک و اکسایت قرار داشت. برای سؤالات موضوعی امتیازاتشان بین ۲۳٪ تا ۵۲/۸٪ قرار گرفت. اینفوسیک بعنوان بهترین عمل کننده لایکاس در جایگاه دوم و اکسایت به عنوان سومین جایگاه.

عملکرد جامع

چهار معیار جنبه‌های قدرت بازیابی موتورهای کاوش را اندازه‌گیری کردند. نمودار ۲ دقت، تکراری بودن و امتیاز مرتب‌ترین گزینه‌ها را برای سؤالات مرجع خلاصه کرده است. امتیاز رتبه‌بندی براساس میزان ارتباط در آن وارد نشد، زیرا دامنه آنها خیلی بالاتر بود و نمی‌توانست بطور کامل در این نمودار وارد شود. بهترین موتور جستجو بالاترین دقت، پائین‌ترین موارد تکراری، پائین‌ترین امتیاز مرتبط‌ترین گزینه و بهترین امتیاز تنظیم براساس دقت را دارد. این نتایج در نمودار ۲ روشن است، به هر حال موتورهای کاوش چنین عمل کردند: این تکست بالاترین مانعیت و پائین‌ترین امتیاز مرتب‌ترین گزینه را داشت. اما اکسایت و اینفوسیک پایین‌ترین تعداد موارد تکراری را داشتند. از این ۴ معیار، این تکست بهترین امتیاز را از بین آنها داشت و توانست بهترین عمل کننده برای این نوع سؤالات باشد. رتبه دوم مشخص نبود چون این موتورها فقط در یک یا دو معیار ممتاز بودند. با این وجود این امکان وجود داشت که آنها را به دو گروه تقسیم کنیم: آلتاویستا، اکسایت، اینفوسیک و لایکاس نسبتاً بهتر از ماژلان، وب کرولر و وردواید وب وارم عمل کردند.

نمودار ۳ شباهت دشواری را در تعیین برنده برای سؤالات موضوعی نشان می‌دهد. اطلاعات موجود بر روی نمودار ۳ اینفوسیک را بهترین عملگر می‌داند و امتیاز تنظیم براساس ارتباط آنرا تقویت می‌کند. بقیه موارد برای اعلام کردن خیلی مشکل بود.

جستجو در اینترنت

در سال ۲۰۰۰ حدود یکصد میلیون پایگاه وب بر روی شبکه جهانی اینترنت وجود دارد و پیش بینی می‌شود که تنها پس از گذشت ۲ سال، در سال ۲۰۰۲ به ۲۵۰ میلیون پایگاه برسد. با رشد تصاعدی حجم اطلاعات، یافتن اطلاعات موردنظر در این دریای پهناور کار مشکلی است و بکارگیری ابزارهای جستجوی مناسب یکی از ضروریات کار باشبکه می‌باشد.

موتورهاي جستجو از سال ۱۹۹۴ مورد استفاده قرار گرفتند. در ابتدا فعاليت آنها فقط جستجو در وب بود ولي با گذشت زمان ، خدمات ديگري از جمله فروش کالا، اجاره فضاي وب و پست الكترونيك ، تحليل سايتها و... به فعاليت آنها اضافه شد.

انواع موتور جستجو

موتورهاي عمومي كه در كليه پايگاهها فارغ از نوع آن جستجو مي كنند.

موتورهاي عمومي معروف عبارتند از:

altavista.com, google.com, go.com, hotbot.com .

"deja.com"، تجارت "news.com" موتورهاي تخصصي در يك رشته خاص مانند اخبار "

"، مقالات و انتشارات www.whowhere.com ۲- "، افراد www.yellowpages.com شركتها "

"[infojump.com](http://www.infojump.com)"

موتورهاي تخصصي در خدمات اينترنت [deja.com](http://www.deja.com) و [magellan.com](http://www.magellan.com) بهترين پايگاه

3-مانند گروههاي خبري و مباحثه

۴- موتورهاي كلان ؛ اين موتورها عبارت مورد جستجو را همزمان به چند موتور جستجو

داده و پاسخها را اولويت بندي کرده و با ذكر نام

موتور جستجو نمايش مي دهند.مانند: [mamma.com](http://www.mamma.com), [savvysearch.com](http://www.savvysearch.com) -

۵-نقطه شروع ؛ اينگونه سايتها موتورهاي جستجو را برحسب موضوعات مختلف معرفي مي

كنند. در صورتي كه موتورهاي تخصصي و عمومي رانمي شناسيد از اين پايگاهها شروع كنيد.

مانند:

۶-موتور اختصاصي پايگاهها؛ بعضي از پايگاههاي بزرگ مانند مايكروسافت ، جنرال

جديدا مي توان از موتورهاي جستجو الكتريك ،... از داخل پايگاه خود موتور جستجو دارند.

...، در يك پايگاه جهت جستجو در آن استفاده كرد. [hotbot](http://www.hotbot.com), [altavista](http://www.altavista.com) عمومي مانند

تقسيم بندي ديگري كه از موتور جستجو مي توان كرد، موتورهاي جستجوگر، و دوم موتورهايي

گفته می web directory کرده اند و به آنها ۷- است که پایگاهها را دسته بندی موضوعی است . در حال حاضر اکثر موتورهای جستجو webcrawear,yahoo میشود. معروفترین آنها دسته بندی موضوعی نیز دارند و در هر دسته و یا زیرشاخه های بعدی می توان جستجو کرد.

نحوه کار موتورهای جستجو

در این مقاله فقط نحوه کار موتورهای عمومی بررسی می شود. موتورهای برای یافتن و spider و یا crawler,robot عمومی از برنامه هایی معروف به مرور صفحات وب استفاده می کنند. نحوه کار این برنامه ها بدین صورت است که با یافتن يك صفحه کلمات مورداستفاده در آن را شناسایی کرده و به جداول فهرست بانک اطلاعاتی خود اضافه می کنند درواقع موتورها صفحات وب را در بانک اطلاعاتی نگهداری نمی کنند بلکه دربانك اطلاعاتی فهرستی از کلمات و آدرس صفحات مشمول این کلمات می باشد.

کار دیگر این برنامه ها این است که به صفحات فهرست شده قبلی مراجعه کرده و در صورت به روز شدن صفحات ، مجددا آنها را فهرست بندی می کنند. ممکن است پایگاه موردنظر موجود نبوده و یا آدرس آن عوض شده باشد.

عوامل مهم در فهرست کردن يك صفحه وب ، تعداد وقوع کلمه در صفحه ، محل قرارگیری آن ، نوع فایلهاي مورداستفاده در صفحه ، درجه اهمیت کلمه در صفحه با توجه به کلید واژه های تعیین شده توسط مالك صفحه و توضیحات آمده در بخش در شناسنامه صفحه می باشد. موتورهای جستجو باتوجه به حجم بانک meta دستورات اطلاعاتی و برنامه هایشان به پایگاههای جدید مراجعه می کنند ولي مطمئن

تمام پایگاههای وب را شامل نمی شوند. بزرگترین موتور جستجوی عمومی حدود می شود. در صورتی که می خواهید پایگاه وب شما ۵۰۰ میلیون صفحه وب را شامل به بانک اطلاعاتی يك موتور جستجو اضافه شود پایگاه وب خود را به آن موتور

جستجو معرفی کنید در صفحات وب ، در بخش دستورات شناسنامه ای صفحه ، کلید واژه های مورد نظر خود را معرفی کنید. موتورهای جستجو عمومی به دو روش کلمات را فهرست بندی می کنند.

اغلب موتورهای جستجو بر اساس کلمات فهرست بندی می کنند. در واقع در جستجو کلمات هم

keyword indexing معنی را تشخیص نمی دهند

excite.com معروفترین موتور می که بر اساس مفهوم جستجو می کند می باشد.

نکته دیگری که در فهرست بندی باید بدانید این است که هر موتور چه بخشهایی از یک صفحه ها و یا چند خط `hyperlink, heading, title` را فهرست می کند. برای مثال ممکن است فقط می توان تعیین `opentext` اول صفحه را فهرست کند. در برخی از موتورهای جستجو مانند نکته `heading, title` کرد که کلمه مورد جستجو در کجای صفحه باشد. برای مثال در...

هستند ما `stop word` دیگر کلمات معروف به

و... بعضی از موتورها این کلمات را در نظر نمی `web, and, or, the, is, an, a`

گیرند. مانند

عوامل مهم در انتخاب موتور جستجو.

عوامل زیر در انتخاب موتور جستجو مهم هستند:

حجم بانک اطلاعاتی موتور جستجو و تعداد صفحات مرور شده توسط آن

به روز بودن بانک اطلاعاتی

تعداد صفحات مرور شده در روز

سرعت برگرداندن نتایج جستجو

تعداد سرویس دهنده های آن در شبکه اینترنت جهت کاهش ترافیک و افزایش سرعت

نحوه نمایش نتایج جستجو و کنترل آن توسط کاربر

نحوه اولویت بندی نتایج حاصله و ارتباط آنها با یکدیگر و کنترل آن توسط کاربر

راحتی استفاده

صفحات معرفی شده به آن طی چند روز در فهرست قرار می گیرند

امنیت در پایگاهها `imagemap,frame` پیشتیبانی ,

قابلیت جستجو در نتایج

stop words پیشتیبانی

حساس به حروف بزرگ و کوچک

پشتیبانی عبارت

عدم محدودیت در تعداد حروف عبارت جستجو

دسته بندی موضوعی و امکان جستجو در هر دسته

جستجو در خدمات اینترنت شامل وب ، گروه های خبری و مباحثه ، ...

پشتیبانی عملگرهای جستجو

پشتیبانی زبانهای مختلف و تبدیل زبانها به یکدیگر

جستجو بر اساس تاریخ

" banner ارائه امکانات بهتر از جمله ارسال نتایج به آدرس پست الکترونیک ، ارائه تبلیغات "

مرتبط با عبارت .

جستجو بر اساس نوع فایل مانند تصویر

جستجو در مکان خاصی از صفحه وب

" صفحه وب tag جستجو در دستورات "

ارائه کلید واژه ها و نتایج مشابه

دسته بندی موتور های جستجو

موتور های جستجو به دو دسته کلی تقسیم می شوند. موتور های جستجوی پیمایشی و فهرستهای

تکمیل دستی. هر کدام از آنها برای تکمیل فهرست خود از روشهای متفاوتی استفاده میکنند که

هر يك را بطور جداگانه مورد بررسی قرار می دهیم:

موتورهای جستجوی پیمایشی یا Crawler-Based Search Engine

Engine

لیست خود را بصورت خودکار تشکیل می‌دهند. Google موتورهای جستجوی پیمایشی مانند آنها وب را پیمایش کرده و سپس کاربران آنچه را که می‌خواهند از میانشان جستجو می‌کنند. اگر شما در صفحه وب خود تغییراتی را اعمال نمایید، موتورهای جستجوی پیمایشی آنها را به خودی خود می‌یابند و سپس این تغییرات لیست خواهند شد. عنوان، متن و دیگر عناصر صفحه، همگی شامل این لیست خواهند بود.

فهرستهای تکمیل دستی یا Human-Powered Directories

وابسته به کاربرانی است که Dmoz مثل Open Directory فهرست تکمیل دستی مانند یک آنرا تکمیل می‌کنند. شما صفحه مورد نظر را به همراه توضیح مختصر در فهرست ثبت می‌کنید یا این کار توسط ویراستارهایی که برای آن فهرست در نظر گرفته شده انجام می‌شود. عمل جستجو در این حالت تنها بر روی توضیحات ثبت شده صورت می‌گیرد و در صورت تغییر روی صفحه وب، روی فهرست تغییری بوجود نخواهد آورد. چیزهایی که برای بهبود یک فهرست بندی در یک موتور جستجو مفید هستند، تأثیری بر بهبود فهرست بندی یک دایرکتوری ندارند. تنها استثناء این است که یک سایت خوب با پایگاه داده‌ای با محتوای خوب شانس بیشتری به نسبت یک سایت با پایگاه داده ضعیف دارد.

موتورهای جستجوی ترکیبی با نتایج مختلط

به موتورهای اطلاق می‌شود که هر دو حالت را در کنار هم نمایش می‌دهند. غالباً، یک موتور جستجوی ترکیبی در صورت نمایش نتیجه جستجو از هر یک از دسته‌های فوق، نتایج حاصل از بیشتر نتایج حاصل از MSN دسته دیگر را هم مورد توجه قرار می‌دهد. مثلاً موتور جستجوی فهرستهای تکمیل دستی را نشان می‌دهد اما در کنار آن نیم‌نگاهی هم به نتایج حاصل از جستجوی پیمایشی دارد.

بررسی يك موتور جستجوی پیمایشی

موتورهای جستجوی پیمایشی شامل سه عنصر اصلی هستند. اولی در اصطلاح عنکبوت است که پیمایشگر می‌شود. پیمایشگر همینکه به يك صفحه می‌رسد، آنرا می‌خواند و سپس لینکهای آن به صفحات دیگر را دنبال می‌نماید. این چیزیست که برای يك سایت پیمایش شده

اتفاق افتاده است. پیمایشگر با يك روال منظم، مثلاً يك یا دو بار در ماه به سایت Crawled مراجعه می‌کند تا تغییرات موجود در آنرا بیابد. هر چیزی که پیمایشگر بیابد به عنصر دوم يك موتور جستجو یعنی فهرست انتقال پیدا می‌کند. فهرست اغلب به کاتالوگی بزرگ اطلاق می‌شود که شامل لیستی از آنچه است که پیمایشگر یافته است. مانند کتاب عظیمی که فهرستی را از آنچه که پیمایشگرها از صفحات وب یافته‌اند، شامل شده است. هرگاه سایتی دچار تغییر شود، این فهرست نیز به روز خواهد شد.

از زمانی که تغییری در صفحه‌ای از سایت ایجاد شده تا هنگامیکه آن تغییر در فهرست موتور جستجو ثبت شود مدت زمانی طول خواهد کشید. پس ممکن است که يك سایت پیمایش شده باشد اما فهرست شده نباشد. تا زمانی که این فهرست‌بندی برای آن تغییر ثبت نشده باشد، نمی‌توان انتظار داشت که در نتایج جستجو آن تغییر را ببینیم. نرم‌افزار موتور جستجو، سومین عنصر يك موتور جستجو است و به برنامه‌ای اطلاق می‌شود که بصورت هوشمندانه‌ای داده‌های موجود در فهرست را دسته‌بندی کرده و آنها را بر اساس اهمیت طبقه‌بندی می‌کند تا نتیجه جستجو با کلمه‌های درخواست شده هر چه بیشتر منطبق و مربوط باشد.

چگونه موتورهای جستجو صفحات وب را رتبه‌بندی می‌کنند؟

وقتی شما از موتورهای جستجوی پیمایشی چیزی را برای جستجو درخواست می‌نمایید، تقریباً بلافاصله این جستجو از میان میلیونها صفحه صورت گرفته و مرتب می‌شود بطوریکه مربوطترین آنها نسبت به موضوع مورد درخواست شما رتبه بالاتری را احراز نماید.

البته باید در نظر داشته باشید که موتورهای جستجو همواره نتایج درستی را به شما ارائه نخواهند داد و مسلماً صفحات نامربوطی را هم در نتیجه جستجو دریافت می‌کنید و گاهی اوقات مجبور هستید که جستجوی دقیقتری را برای آنچه که می‌خواهید انجام دهید اما موتورهای جستجو کار حیرت‌انگیز دیگری نیز انجام می‌دهند.

فرض کنید که شما به يك كتابدار مراجعه مي كنيد و از وي درباره «سفر» كتابي مي خواهيد. او براي اينكه جواب درستي به شما بدهد و كتاب مفيدي را به شما ارائه نمايد با پرسیدن سوالاتي از شما و با استفاده از تجارب خود كتاب مورد نظرتان را به شما تحويل خواهد داد. موتورهاي جستجو همچنين توانايي ندارند اما به نوعي آنها را شبیه‌سازي مي‌کنند. پس موتورهاي جستجوي پيمائشي چگونه به جواب مورد نظرتان از ميان ميليونها صفحه وب مي‌رسند؟ آنها يك مجموعه از قوانين را دارند که الگوريتم ناميده مي‌شود. الگوريتمهاي مورد نظر براي هر موتور جستجويي خاص و تقريباً سري هستند اما به هر حال از قوانين زير پيروي مي‌کنند:

مکان و تکرار

يکي از قوانين اصلي در الگوريتمهاي رتبه‌بندي موقعيت و تعداد تکرار کلماتي است که در (Location/Frequency Methode) صفحه مورد استفاده قرار گرفته‌اند که ناميده مي‌شود بطور خلاصه روش مکان-تکرار (Methode)

كتابدار مذکور را به خاطر مي‌آورد لازم است که او کتابهاي در رابطه با کلمه «سفر» را طبق درخواست شما بيابد. او در حله اول احساس مي‌کند که شما به دنبال کتابهايي هستيد که در نامشان کلمه «سفر» را شامل شوند. موتورهاي جستجو هم دقيقاً همان کار را انجام مي‌دهند. حاوي HTML موجود در کد Title آنها هم صفحاتي را براي تان ليست مي‌کنند که در برچسب کلمه «سفر» باشند.

موتورهاي جستجو همچنين به دنبال کلمه مورد نظر در بالاي صفحات و يا در ابتداي پاراگرافها هستند. آنها فرض مي‌کنند که صفحاتي که حاوي آن کلمه در بالاي خود و يا ابتداي پاراگرافها و عناوين باشند به نتيجه مورد نظر شما مربوطتر هستند.

عامل بزرگ و مهم ديگري است که موتورهاي جستجو از طريق آن Frequency تکرار يا صفحات مربوط را شناسايي مي‌نمايند. موتورهاي جستجو صفحات را تجزيه کرده و با توجه به تکرار کلمه‌اي در صفحه متوجه مي‌شوند که آن کلمه نسبت به ديگر کلمات اهميت بيشتري در آن صفحه دارد و آن صفحه را در درجه بالاتري نسبت به صفحات ديگر قرار مي‌دهند.

عوامل خارج از صفحه

موتورهای جستجوی پیمایشی اکنون تجربه فراوانی در رابطه با وب مسترهایی دارند که صفحات خود را برای کسب رتبه بهتر مرتباً بازنویسی می‌کنند. بعضی از وب مسترهای خبره حتی ممکن است به سمت روشهایی مانند مهندسی معکوس برای کشف چگونگی روشهای مکان-تکرار بروند. به همین دلیل، تمامی موتورهای جستجوی معروف از روشهای امتیازبندی «خارج از صفحه» استفاده می‌کنند. عوامل خارج از صفحه عواملی هستند که از تیررس وبمسترها خارجند و آنها نمی‌توانند در آن دخالت کنند و مساله مهم در آن تحلیل ارتباطات و لینکهاست. بوسیله تجزیه صفحات، موتورهای جستجو لینکها را بررسی کرده و از محبوبیت آنها می‌فهمند که آن صفحات مهم بوده و شایسته ترفیع رتبه هستند. بعلاوه تکنیکهای پیشرفته به گونه‌ای است که از ایجاد لینکهای مصنوعی توسط وبمسترها برای فریب موتورهای جستجو جلوگیری می‌نماید. علاوه بر آن موتورهای جستجو بررسی می‌کنند که کدام صفحه توسط یک کاربر که کلمه‌ای را جستجو کرده انتخاب می‌شود و سپس با توجه به تعداد انتخابها، رتبه صفحه مورد نظر را تعیین کرده و مقام آنرا در نتیجه جستجو جابجا می‌نمایند.

نتیجه‌گیری‌ها

این پژوهش ۸ موتور کاوش اصلی را با دوبار جستجوی ۲۶ سؤال در هر کدام از آنها (۴ بار در اکسایت) ارزشیابی کرد. اطلاعات نشان داد که موتورهای کاوش انتخابی نمی‌توانند نتایج خوبی برای سؤالات مرجع واقعی ارائه دهند. اما در مورد سؤالات موضوعی ساختگی خوب عمل کردند. این نکته نیز فهمیده شد که موتورهای کاوش برای دو نوع سؤال بطور متفاوت عمل کردند: اینفوسیک در سؤالات موضوعی بهتر عمل کرد؛ در حالیکه این تکست در سؤالات مرجع بهترین بود. از این پژوهش فهمیده شد که موارد تکراری یک مشکل نمی‌تواند باشد. با تعریف متغیر تنظیم براساس میزان ارتباط در موتورهای کاوش می‌تواند ارزشیابی شود. با ارائه اطلاعاتی در مورد این ۴ متغیر این پژوهش چندین جنبه از عملکرد موتورهای کاوش را روشن کرد.

این پژوهش بدون محدودیت نبود. اول، داده‌های آن لحظات ناپایدار را در اینترنت ثبت کرد. آنها

عکس‌هایی ارائه دادند از اینکه چگونه موتورهای کاوش از آوریل تا ژوئن ۱۹۹۶ کار کردند. و این تصاویر ممکن است کاملاً نهایی نباشد که قبلاً بود زیرا اینترنت سریعاً در حال گسترش است. با این وجود، مشابه مطالعه لیتون، این پژوهش دریافت که اینفوسیک یکی از بهترین موتورهای کاوش است. اگر این تحقیق تکرار شده و همان یافته‌ها که بدست آمده. به هر حال یکی از تحقیقات قادر خواهد بود که اطمینان بیشتری در مورد این یافته‌ها بدهد. دوم، سؤالات مرجع از يك كتابخانه دانشگاهي جمع‌آوری شده بود و سخت بود تعیین کردن اینکه آنها نمونه‌ای از سؤالات مرجع بودند. باز هم، تکرار این پژوهش اعتبار یافته‌ها را افزایش می‌داد. سوم، تعداد سؤالات آزمایش احتمالاً می‌توانست زیاد باشد اگرچه این پژوهش سؤالات بیشتری نسبت به اغلب پژوهش‌های دیگری بکار برد.

گذشته از محدودیتها، این پژوهش متغیرهای جدید برای ارزشیابی تنظیم براساس ارتباط تولید کرد و يك طرح تحقيقي براي مقایسه عملکرد موتورهای کاوش برای دو نوع سؤال بکار برد، درك عملکرد موتورهای کاوش را افزایش داد، توصیه‌هایی در مورد اینکه چگونه طراحان سیستم می‌توانند سیستم‌های خود را بهبود بخشند ارائه داد و اشاره کرد که چگونه کتابداران می‌توانند خودشان و مردم را برای جستجو در اینترنت آماده کنند.

جدول ۱: میانگین مانعیت و نقاط کور برای سؤالات مرجع در موتورهای کاوش

موتورها	دقت	نقاط کور	درصد
آلتاویستا	05/2	21	50%
اکسایت	75/1	12	29%
اینفوسیک	95/1	19	45%
لایکاس	93/1	16	38%
ماژلان	33/1	27	64%
اپن تکست	93/2	15	36%
وب کرولر	10/1	24	57%
ورلد واید وب	31/0	32	76%

وارم			
------	--	--	--

جدول ۲- میانگین مانعیت و تعداد نقاط کور موتورهای کاوش برای سوالات موضوعی

موتورها	دقت	نقاط کور (فراوانی)	دقت خاص (فراوانی)
آلتاویستا	4/5	0	3
اکسایت	2/4	1	3
اینفوسیک	3/7	1	8
لایکاس	3/6	0	6
ماژلان	7/6	1	7
اپن تکست	5/6	2	6
وب کرولر	3/5	0	6
ورلد واید وب	2/3	5	2
وارم			

دقت خاص به تعداد جستجوهایی که بیش از ۵ گزینه مرتبط را بازیابی کردند اشاره دارد. 1.

مهم ترین موتور های جستجو

www.google.com

<http://www.yahoo.com/>

<http://www.37.com/>

<http://www.infoseek.com/>

<http://www.aliweb.com/>

<http://www.ask.com/>

<http://www.mamma.com/>

<http://www.altavista.com/>

منابع:

[www.irandoc.ac.ir/ETELA-ART/ 18/18 3 4 10.htm](http://www.irandoc.ac.ir/ETELA-ART/18/18_3_4_10.htm)

<http://www.tarighat-e.com/information/ShowArticle.asp?ID=341>

پی نوشت ها:

1. College & Research Libraries News.
2. Data collection
3. Precision
4. Duplicate
5. Most-relevant-item score (MRI)
6. Relevancy score.
7. Recall
8. Data analysis

مؤلفان
های جستجو

نویسنده: ابراهیم آرام

مؤلف: هادی جستنجو

این کتاب به صورت رایگان ارائه میشود

مؤلفان مقامات جسته‌جسته