



پورتال مقالات کامپیوتر و فناوری اطلاعات

مجموعه مقالات

موتورهای جستجوگر

شماره ۲



پورتال مقالات کامپیوتر و فناوری اطلاعات

فهرست مقالات

روش کار موتور جستجو

موتورهای جستجو چگونه کار می کنند؟

موتور های جستجو

آداب جستجو در اینترنت

چه نوع موتور جستجو (Search Engine) یا دایرکتوری (Directory) باید استفاده کرد؟

جایگاه موتور های جست و جو

بررسی کمی و کیفی موتور های جستجو

وب پنهان، اطلاعاتی که موتور جستجوگر بدان راهی ندارد!

جویشگر چیست ؟

مفاهیم و اصطلاحات دنیای جستجو و موتورهای جستجوگر



پورتال مقالات کامپیوتر و فناوری اطلاعات

روش کار موتور جستجو

وقتی جستجویی در یک موتور جستجوگر انجام و نتایج جستجو ارائه می شود. کاربران در واقع نتیجه کار بخش های متفاوت موتور جستجوگر را می بینند. موتور جستجوگر قبلاً پایگاه داده اش را آماده کرده است و این گونه نیست که درست در همان لحظه جستجو. تمام وب را بگردد. بسیاری از خود می پرسند که چگونه ممکن است گوگل در کمتر از یک ثانیه تمام سایت های وب را بگردد و میلیون ها صفحه را در نتایج جستجوی خود ارائه کند؟

گوگل و هیچ موتور جستجوگر دیگری توانایی انجام این کار را ندارند. همه آنها در زمان پاسخ گویی به جستجوهای کاربران. تنها در پایگاه داده ای که در اختیار دارند به جستجو می پردازند و نه در وب! موتور جستجوگر به کمک بخش های متفاوت خود. اطلاعات مورد نیاز را قبلاً جمع آوری. تجزیه و تحلیل می کند. آنها در پایگاه داده اش ذخیره می نماید و به هنگام جستجوی کاربر تنها در همین پایگاه داده می گردد. بخش های مجزای یک موتور جستجوگر عبارتند از:

* Spider یا عنکبوت

* Crawler یا خزنده

* Indexer یا بایگانی کننده

* Database یا پایگاه داده

* Ranker یا سیستم رتبه بندی

الف) Spider - (عنکبوت)

اسپایدر یا روبات (Robot)، نرم افزاری است که کار جمع آوری اطلاعات مورد نیاز یک موتور جستجوگر را بر عهده دارد. اسپایدر به صفحات مختلف سر می زند. محتوای آنها را می خواند. لینکها را دنبال می کند. اطلاعات مورد نیاز را جمع آوری می کند و آنها در اختیار سایر بخش های موتور جستجوگر قرار می دهد. کار یک اسپایدر بسیار شبیه کار کاربران وب است. همانطور که کاربران. صفحات مختلف را بازدید می کنند. اسپایدر هم درست این کار را انجام می دهد با این تفاوت که اسپایدر کدهای HTML صفحات را می بیند اما کاربران نتیجه حاصل از کنار هم قرار گرفتن این کدها را. index.html صفحه ای است که کاربران آنها می بینند:

اما یک اسپایدر آنها چگونه می بیند؟

برای این که شما هم بتوانید دنیای وب را از دیدگاه یک اسپایدر ببینید. کافی است که کدهای HTML صفحات را مشاهده کنید.

آیا این دنیای متنی برای شما جذاب است؟

اسپایدر. به هنگام مشاهده صفحات. بر روی سرورها رد پا برجای می گذارد. شما اگر اجازه دسترسی به آمار دید و بازدیدهای صورت گرفته از یک سایت و اتفاقات انجام شده در آن را داشته باشید. می توانید مشخص کنید که اسپایدر کدام یک از موتورهای جستجوگر صفحات سایت را مورد بازدید قرار داده است. یکی از فعالیتهای اصلی که در SEM انجام می شود تحلیل آمار همین دید و بازدیدها است.

اسپایدرها کاربردهای دیگری نیز دارند. به عنوان مثال عده ای از آنها به سایت های مختلف مراجعه می کنند و فقط به بررسی فعال بودن لینک های آنها می پردازند و یا به دنبال آدرس ایمیل (Email) می گردند.

ب) Crawler (خزنده)

کراولر. نرم افزاری است که به عنوان یک فرمانده برای اسپایدر عمل می کند. آن مشخص می کند که اسپایدر کدام صفحات را مورد بازدید قرار دهد. در واقع کراولر تصمیم می گیرد که کدام یک از لینک های صفحه ای که اسپایدر در حال حاضر در آن قرار دارد. دنبال شود. ممکن است همه آنها را دنبال کند. بعضی ها را دنبال کند و یا هیچ کدام را دنبال نکند.

کراولر. ممکن است قبلاً برنامه ریزی شده باشد که آدرس های خاصی را طبق برنامه. در اختیار اسپایدر قرار دهد تا از آنها دیدن کند. دنبال کردن لینک های یک صفحه به این بستگی دارد که موتور جستجوگر چه حجمی از اطلاعات یک سایت را می تواند (می خواهد) در پایگاه داده اش



پورتال مقالات کامپیوتر و فناوری اطلاعات

ذخیره کند. همچنین ممکن است اجازه دسترسی به بعضی از صفحات به موتورهای جستجوگر داده نشده باشد.

شما به عنوان دارنده سایت، همان طور که دوست دارید موتورهای جستجوگر اطلاعات سایت شما را با خود ببرند، می توانید آنها را از بعضی صفحات سایت تان دور کنید و اجازه دسترسی به محتوای آن صفحات را به آنها ندهید. موتور جستجو اگر مودب باشد قبل از ورود به هر سایتی ابتدا قوانین دسترسی به محتوای سایت را (در صورت وجود) در فایلی خاص بررسی می کند و از حقوق دسترسی خود اطلاع می یابد. تنظیم میزان دسترسی موتورهای جستجوگر به محتوای یک سایت توسط پروتکل Robots انجام می شود. به عمل کراولر، خزش (Crawling) می گویند.

ج) Indexer (بایگانی کننده)

تمام اطلاعات جمع آوری شده توسط اسپایدر در اختیار ایندکسر قرار می گیرد. در این بخش اطلاعات ارسالی مورد تجزیه و تحلیل قرار می گیرند و به بخش های متفاوتی تقسیم می شوند. تجزیه و تحلیل بدین معنی است که مشخص می شود اطلاعات از کدام صفحه ارسال شده است، چه حجمی دارد، کلمات موجود در آن کدامند، کلمات چندبار تکرار شده اند، کلمات در کجای صفحه قرار دارند و ...

در حقیقت ایندکسر، صفحه را به پارامترهای آن خرد می کند و تمام این پارامترها را به یک مقیاس عددی تبدیل می کند تا سیستم رتبه بندی بتواند پارامترهای صفحات مختلف را با هم مقایسه کند. در زمان تجزیه و تحلیل اطلاعات، ایندکسر برای کاهش حجم داده ها از بعضی کلمات که بسیار رایج هستند صرفنظر می کند. کلماتی نظیر a, an, the, www, is از این گونه کلمات هستند.

د) DataBase (پایگاه داده)

تمام داده های تجزیه و تحلیل شده در ایندکسر، به پایگاه داده ارسال می گردد. در این بخش داده ها گروه بندی، کدگذاری و ذخیره می شود. همچنین داده ها قبل از آنکه ذخیره شوند، طبق تکنیکهای خاصی فشرده می شوند تا حجم کمی از پایگاه داده را اشغال کنند. یک موتور جستجوگر باید پایگاه داده عظیمی داشته باشد و به طور مداوم حجم محتوای آنرا گسترش دهد و البته اطلاعات قدیمی را هم به روز رسانی نماید. بزرگی و به روز بودن پایگاه داده یک موتور جستجوگر برای آن امتیاز محسوب می گردد. یکی از تفاوت های اصلی موتورهای جستجوگر در حجم پایگاه داده آنها و همچنین روش ذخیره سازی داده ها در پایگاه داده است.

و) Ranker (سیستم رتبه بندی)

بعد از آنکه تمام مراحل قبل انجام شد، موتور جستجوگر آماده پاسخ گویی به سوالات کاربران است. کاربران چند کلمه را در جعبه جستجوی (Search Box) آن وارد می کنند و سپس با فشردن Enter منتظر پاسخ می مانند. برای پاسخ گویی به درخواست کاربر، ابتدا تمام صفحات موجود در پایگاه داده که به موضوع جستجو شده، مرتبط هستند، مشخص می شوند. پس از آن سیستم رتبه بندی وارد عمل شده، آنها را از بیشترین ارتباط تا کمترین ارتباط مرتب می کند و به عنوان نتایج جستجو به کاربر نمایش می دهد.

حتی اگر موتور جستجوگر بهترین و کامل ترین پایگاه داده را داشته باشد اما نتواند پاسخ های مرتبطی را ارائه کند، یک موتور جستجوگر ضعیف خواهد بود. در حقیقت سیستم رتبه بندی قلب تپنده یک موتور جستجوگر است و تفاوت اصلی موتورهای جستجوگر در این بخش قرار دارد. سیستم رتبه بندی برای پاسخ گویی به سوالات کاربران، پارامترهای بسیاری را در نظر می گیرد تا بتواند بهترین پاسخ ها را در اختیار آنها قرار دارد.

حرفه ای های دنیای SEM به طور خلاصه از آن به Algo (الگوریتم) یاد می کنند. الگوریتم، مجموعه ای از دستورالعمل ها است که موتور جستجوگر با اعمال آنها بر پارامترهای صفحات موجود در پایگاه داده اش، تصمیم می گیرد که صفحات مرتبط را چگونه در نتایج جستجو مرتب کند. در حال حاضر قدرتمندترین سیستم رتبه بندی را گوگل در اختیار دارد.

می توان با ادغام کردن اسپایدر با کراولر و همچنین ایندکسر با پایگاه داده، موتور جستجوگر را شامل سه بخش زیر دانست که این گونه تقسیم بندی هم درست می باشد:

* کراولر

* بایگانی

* سیستم رتبه بندی

تذکر- برای سهولت در بیان مطالب بعدی هر گاه صحبت از بایگانی کردن (شدن) به میان می آید، مقصود این است که صفحه تجزیه و تحلیل



پورتال مقالات کامپیوتر و فناوری اطلاعات

شده و به پایگاه داده موتور جستجوگر وارد می شود.

برای آنکه تصور درستی از نحوه کار یک موتور جستجوگر داشته باشید داستان نامتعارف زیر را با هم بررسی می کنیم. داستان ما یک شکارچی دارد. او تصمیم به شکار می گیرد:

- کار کراولر:

او قصد دارد برای شکار به منطقه حفاظت شده آبیورد. واقع در شهرستان درگز (شمالی ترین شهر خراسان بزرگ) برود.

- پروتکل Robots :

ابتدا تمام محدودیت های موجود برای شکار در این منطقه را بررسی می کند:

* آیا در این منطقه می توان به شکار پرداخت؟

* کدام حیوانات را می توان شکار کرد؟

* حداکثر تعداد شکار چه میزانی است؟

* و

فرض می کنیم او مجوز شکار یک اوریل (نوعی آهو) را از شکاربانی منطقه دریافت می کند.

- کار اسپایدر

او اوریلی رعنا را شکار می کند و سپس آنرا با خود به منزل می برد.

- کار ایندکسر

شکار را تکه تکه کرده، گوشت، استخوان، دل و قلوه، کله پاچه و ... آنرا بسته بندی می کند و بخش های زاید شکار را دور می ریزد.

- کار پایگاه داده

بسته های حاصل را درون فریزر قرار داده، ذخیره می کند.

- کار سیستم رتبه بندی

مهمانان سراغ او می آیند و همسرش بسته به ذائقه مهمانان برای آنها غذا طبخ می کند. ممکن است عده ای کله پاچه، عده ای آبگوشت، عده ای ... دوست داشته باشند. پخت غذا طبق سلیقه مهمانان کار سختی است. ممکن است همه آنها آبگوشت بخواهند اما آنها مسلماً "بامزه ترین آبگوشت را می خواهند!"

+ نکته ها:

* شکارچی می توانست برای شکار کبک یا اوریل و یا هر دو به آن منطقه برود همانطور که موتور جستجوگر می تواند از سرور سایت شما انواع فایل (عکس، فایل متنی، فایل اجرایی و ...) درخواست کند.

* شکارچی می تواند شب به شکار برود یا روز. موتور جستجوگر هم ممکن است شب به سایت شما مراجعه کند یا روز. بنابراین همواره مطمئن باشید که سایت شما آپ است و موتور جستجوگر می تواند در آن به شکار فایلها بپردازد.

* غذای خوشمزه را می توانید با نتایج جستجوی دقیق و مرتبط مقایسه کنید. اگر شکارچی بهترین شکار را با خود به منزل ببرد اما غذایی خوشمزه و مطابق سلیقه مهمانان طبخ نگردد، تمام زحمات هدر رفته است.

منبع : مرکز کامپیوتر ستارگان



پورتال مقالات کامپیوتر و فناوری اطلاعات

موتورهای جستجو چگونه کار می کنند؟

موتورهای جستجو به دو دسته کلی تقسیم می شوند. موتورهای جستجوی پیمایشی و فهرستهای تکمیل دستی. هر کدام از آنها برای تکمیل فهرست خود از روشهای متفاوتی استفاده می کنند که هر یک را بطور جداگانه مورد بررسی قرار می دهیم:

موتورهای جستجوی پیمایشی یا Crawler-Based Search Engines

موتورهای جستجوی پیمایشی مانند Google لیست خود را بصورت خودکار تشکیل می دهند. آنها وب را پیمایش کرده و سپس کاربران آنچه را که می خواهند از میانشان جستجو می کنند. اگر شما در صفحه وب خود تغییراتی را اعمال نمایید، موتورهای جستجوی پیمایشی آنها را به خودی خود می یابند 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

یک فهرست تکمیل دستی مانند یک Open Directory مثل Dmoz وابسته به کاربرانی است که آنرا تکمیل می کنند. شما صفحه مورد نظر را به همراه توضیحی مختصر در فهرست ثبت می کنید یا این کار توسط ویراستارهایی که برای آن فهرست در نظر گرفته شده انجام می شود. عمل جستجو در این حالت تنها بر روی توضیحات ثبت شده صورت می گیرد 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، تاثیری بر بهبود فهرست بندی یک دایرکتوری ندارند. تنها استثناء این است که یک سایت خوب با پایگاه داده ای با محتوای خوب شانس بیشتری به نسبت یک سایت با پایگاه داده ضعیف دارد.

موتورهای جستجوی ترکیبی با نتایج مختلط

به موتورهایی اطلاق می شود که هر دو حالت را در کنار هم نمایش می دهند. غالباً، یک موتور جستجوی ترکیبی در صورت نمایش نتیجه جستجو از هر یک از دسته های فوق، نتایج حاصل از دسته دیگر را هم مورد توجه قرار می دهد. مثلاً موتور جستجوی MSN بیشتر نتایج حاصل از فهرستهای تکمیل دستی را نشان می دهد اما در کنار آن نیم نگاهی هم به نتایج حاصل از جستجوی پیمایشی دارد.

بررسی یک موتور جستجوی پیمایشی

موتورهای جستجوی پیمایشی شامل سه عنصر اصلی هستند. اولی در اصطلاح عنکبوت (Spider) است که پیمایشگر (Crawler) هم نامیده می شود. پیمایشگر همین که به یک صفحه می رسد، آنرا می خواند و سپس لینکهای آن به صفحات دیگر را دنبال می نماید. این چیز است که برای یک سایت پیمایش شده (Crawled) اتفاق افتاده است. پیمایشگر با یک روال منظم، مثلاً یک یا دو بار در ماه به سایت مراجعه می کند تا تغییرات موجود در آنرا بیابد. هر چیزی که پیمایشگر بیابد به عنصر دوم یک موتور جستجو یعنی فهرست انتقال پیدا می کند. فهرست اغلب به کاتالوگی بزرگ اطلاق می شود که شامل لیستی از آنچه است که پیمایشگر یافته است. مانند کتاب عظیمی که فهرستی را از آنچه که پیمایشگرها از صفحات وب یافته اند، شامل شده است. هرگاه سایتی دچار تغییر شود، این فهرست نیز به روز خواهد شد.

از زمانی که تغییری در صفحه ای از سایت ایجاد شده تا هنگامی که آن تغییر در فهرست موتور جستجو ثبت شود مدت زمانی طول خواهد کشید. پس ممکن است که یک سایت پیمایش شده باشد اما فهرست نشده باشد. تا زمانی که این فهرست بندی برای آن تغییر ثبت نشده باشد، نمی توان انتظار داشت که در نتایج جستجو آن تغییر را ببینیم. نرم افزار موتور جستجو، سومین عنصر یک موتور جستجو است و به برنامه ای اطلاق می شود که بصورت هوشمندانه ای داده های موجود در فهرست را دسته بندی کرده و آنها را بر اساس اهمیت طبقه بندی می کند تا نتیجه جستجو با کلمه های درخواست شده هر چه بیشتر منطبق و مربوط باشد.

چگونه موتورهای جستجو صفحات وب را رتبه بندی می کنند؟

وقتی شما از موتورهای جستجوی پیمایشی چیزی را برای جستجو درخواست می نمایید، تقریباً بلافاصله این جستجو از میان میلیونها صفحه صورت گرفته و مرتب می شود 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



پورتال مقالات کامپیوتر و فناوری اطلاعات

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 یک کتابدار مراجعه می‌کنید و از وی درباره «سفر» کتابی می‌خواهید. او برای اینکه جواب درستی به شما بدهد و کتاب مفیدی را به شما ارائه نماید با پرسیدن سؤالاتی از شما و با استفاده از تجارب خود کتاب مورد نظران را به شما تحویل خواهد داد. موتورهای جستجو همچنین توانایی ندارند اما به نوعی آنها را شبیه‌سازی می‌کنند. پس موتورهای جستجوی پیمایشی چگونه به جواب مورد نظران از میان میلیونها صفحه وب می‌رسند؟ آنها یک مجموعه از قوانین را دارند که الگوریتم نامیده می‌شود. الگوریتمهای مورد نظر برای هر موتور جستجویی خاص و تقریباً سری هستند اما به هر حال از قوانین زیر پیروی می‌کنند:

مکان و تکرار

یکی از قوانین اصلی در الگوریتمهای رتبه‌بندی موقعیت و تعداد تکرار کلماتی است که در صفحه مورد استفاده قرار گرفته‌اند که بطور خلاصه روش مکان - تکرار (Location/Frequency Methode) نامیده می‌شود.

کتابدار مذکور را به خاطر می‌آورد؟ لازم است که او کتابهای در رابطه با کلمه «سفر» را طبق درخواست شما بیابد. او در وحله اول احساس می‌کند که شما به دنبال کتابهایی هستید که در نامشان کلمه «سفر» را شامل شوند. موتورهای جستجو هم دقیقاً همان کار را انجام می‌دهند. آنها هم صفحاتی را برایتان لیست می‌کنند که در برچسب Title موجود در کد HTML حاوی کلمه «سفر» باشند. موتورهای جستجو همچنین به دنبال کلمه مورد نظر در بالای صفحات و یا در ابتدای پاراگرافها هستند. آنها فرض می‌کنند که صفحاتی که حاوی آن کلمه در بالای خود و یا در ابتدای پاراگرافها و عناوین باشند به نتیجه مورد نظر شما مربوط تر هستند. تکرار یا Frequency عامل بزرگ و مهم دیگری است که موتورهای جستجو از طریق آن صفحات مربوط را شناسایی می‌نمایند. موتورهای جستجو صفحات را تجزیه کرده و با توجه به تکرار کلمه‌ای در صفحه متوجه می‌شوند که آن کلمه نسبت به دیگر کلمات اهمیت بیشتری در آن صفحه دارد و آن صفحه را در درجه بالاتری نسبت به صفحات دیگر قرار می‌دهند.

دستور آشپزی

خب آشپزی چه ربطی به موتورهای جستجو دارد؟ رابطه در اینجا است. همانطور که آشپزهای حرفه‌ای دستورات آشپزی خود را در لفافه نگه می‌دارند و مقدار و چگونگی ادویه‌های افزودنی به غذاهای خود را افشا نمی‌کنند. چگونگی کارکرد دقیق موتورهای جستجو درباره روشهایی از قبیل مکان-تکرار فاش نمی‌شود و هر موتور جستجویی روش خود را دنبال می‌کند. به همین دلیل است که وقتی شما کلمات واحدی را در موتورهای متفاوت جستجو می‌کنید، به نتایج متفاوتی می‌رسید. برخی موتورهای جستجو نسبت به برخی دیگر صفحات بیشتری را فهرست کرده‌اند. نتیجه این خواهد شد که هیچ موتور جستجویی نتیجه جستجوی مشترکی با موتور دیگر نخواهد داشت و شما نتایج متفاوتی را از آنها دریافت می‌کنید. موتورهای جستجو همچنین ممکن است که برخی از صفحات را از فهرست خود حذف کنند البته به شرطی که آن صفحات با Spam شدن سعی در گول زدن موتورهای جستجوگر داشته باشند. Spamming روشی است که برخی از صفحات برای احراز رتبه بالاتر در موتورهای جستجو در پیش می‌گیرند و آن به این صورت است که با تکرار بیش از حد کلمات بطور عمدی سعی در برهم زدن تعادل و در نتیجه فریب موتورهای جستجو دارند. آنها سعی دارند که با افزایش عامل تکرار، در رتبه بالاتری قرار بگیرند. موتورهای جستجو راههای متنوعی برای جلوگیری از Spamming دارند و در این راه از گزارشهای کاربران خود نیز بهره می‌برند.

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.

منبع : جنوبی‌ها



پورتال مقالات کامپیوتر و فناوری اطلاعات

موتور های جستجو

امروزه بر روی اینترنت و مهمترین سرویس آن (وب)، صدها میلیون صفحه حاوی اطلاعات وجود دارد. کاربران اینترنت با آگاهی از آدرس یک سایت، قادر به اتصال به سایت مورد نظر و استفاده از منابع اطلاعاتی موجود بر روی سایت خواهند بود. ما با دریائی از اطلاعات مواجه هستیم، در صورتی که قصد یافتن اطلاعاتی خاص را داشته باشیم، از چه امکاناتی در این زمینه می توان استفاده کرد؟ برای جستجو و یافتن اطلاعات مورد نیاز از مراکز جستجوی اطلاعات در اینترنت استفاده می گردد. به مراکز فوق Search engines نیز می گویند.

مراکز جستجو در اینترنت، نوع خاصی از سایت های موجود در وب بوده که با هدف کمک برای یافتن اطلاعات، ایجاد شده اند. مراکز جستجو در اینترنت بمنظور پاسخگویی به کاربران متقاضی و جستجوکنندگان اطلاعات از سه روش متفاوت استفاده می نمایند. نحوه عملکرد سه روش با یکدیگر مشابه بوده و تنها تفاوت موجود میدان انتخاب شده برای عملیات جستجو است.

* اینترنت و یا بخشی از آن بر اساس کلمات مهم، جستجو می گردد.

* از کلمات پیدا شده یک ایندکس به همراه محل مربوط به هر یک، ایجاد می نمایند.

* به کاربران امکان جستجو برای کلمات خاص و یا ترکیبی از آنها که در فایل ایندکس موجود می باشند، داده می شود.

مراکز جستجوی اولیه در اینترنت، صرفاً اطلاعات مربوط به چندین هزار صفحه وب را ایندکس و روزانه دو تا سه هزار کاربر متقاضی به آنها مراجعه می کردند. مراکز جستجوی فعلی در اینترنت اطلاعات مربوط به صدها میلیون صفحه را ایندکس نموده و روزانه به بیش از دهها میلیون متقاضی پاسخ می دهند.

وب

اغلب مردم زمانی که از مراکز جستجو در اینترنت سخن می گویند، منظور آنها مراکز جستجوی وب است. قبل از مطرح شدن وب (مشهورترین بخش اینترنت)، از مراکز جستجوی اطلاعات برای کمک به کاربران برای یافتن اطلاعات استفاده می گردید. برنامه هائی نظیر: « gopher » و « Archie » از فایل های ذخیره شده بر روی سرویس دهنده های متصل به اینترنت، یک ایندکس ایجاد می کردند. بدین ترتیب جستجو و دسترسی به اطلاعات و مستندات مورد نظر در اسرع وقت انجام می گردید. در اواخر سال ۱۹۸۰ اکثر کاربران مستلزم دارا بودن دانش کافی در رابطه با استفاده از gopher, Archie و Veronica بودند. امروزه اکثر کاربران اینترنت دامنه جستجوی خود را محدود به وب نموده اند.

قبل از اینکه یک مرکز جستجو قادر به ارائه آدرس و محل فایل مورد نظر باشد، می بایست فایل مورد نظر پیدا شود. بمنظور یافتن اطلاعات مربوط به صدها میلیون صفحه وب موجود، مراکز جستجو می بایست از یک نرم افزار خاص با نام Spider (عنکبوت) برای ایجاد لیست های شامل کلمات موجود در هر یک از صفحات وب، استفاده نمایند. فرآیند ایجاد لیست های مربوطه توسط Spider، اصطلاحاً « web crawling » نامیده می شود. برای ایجاد و نگهداری یک لیست مفید از کلمات، Spider های مراکز جستجو می بایست تعداد زیادی از صفحات وب را بررسی و مشاهده نمایند. نحوه حرکت Spider در وب به چه صورت است؟ نقاط شروع، لیستی از سرویس دهندگان با ترافیک و اطلاعات بالا و صفحات وب متداول است. Spider از یک سایت رایج عملیات خود را آغاز و پس از ایندکس نمودن کلمات موجود در صفحات وب، هر یک از لینک های موجود در صفحات را برای ادامه حرکت خود انتخاب خواهد کرد. بدین ترتیب سیستم مبتنی بر Spider بسرعت حرکت خود در طول وب را آغاز خواهد کرد.

Google یکی از مراکز جستجوی دانشگاهی و معتبر است. در سیستم فوق از چندین Spider (معمولاً سه Spider در هر لحظه) برای ایجاد مقادیر اولیه برای سیستم، استفاده می گردد. هر Spider قادر به نگهداری ارتباط خود با بیش از ۳۰۰ صفحه وب در یک لحظه است. با استفاده از چهار spider، سیستم فوق قادر به جستجوی ۱۰۰ صفحه در ثانیه و تولید ۶۰۰ کیلوبایت اطلاعات در هر ثانیه است. اطلاعات مورد نیاز هر یک از spider ها می بایست بسرعت در اختیار آنان گذاشته شود. سیستم اولیه Google، دارای یک سرویس دهنده اختصاصی بمنظور تغذیه آدرس های URL مورد نیاز برای هر یک از Spider ها بود. بمنظور افزایش سرعت عملیات، Google از یک سیستم DNS اختصاصی استفاده می کرد. در سایر موارد از DNS مربوط به ISP استفاده می گردد. زمانیکه Spider به یک صفحه وب شامل تگ های Html برخورد می نماید، دو آیتیم در رابطه با آن را یادداشت خواهد کرد:

کلمات موجود در صفحه

محلی که کلمات پیدا شده اند.

از کلمات موجود در عنوان (title)، زیرعناوین (Subtitles)، تگ های متا و سایر مکانهای مهم یادداشت برداشته شده تا در آینده با توجه به خواسته کاربر، امکان پاسخگویی مناسب به آنها فراهم گردد. Spider مربوط به Google، از کلمات موجود در هر یک از صفحات وب ایندکس ایجاد و



پورتال مقالات کامپیوتر و فناوری اطلاعات

کلماتی نظیر: «an», «a» و «the» را حذف می نمایند. سایر Spider ها هر یک دارای رویکردهای خاص خود در این زمینه می باشند.

سیاست های استفاده شده در رابطه با نحوه ایندکس نمودن اطلاعات صفحات وب، مستقیماً بر سرعت عملکرد spider ها تأثیر گذاشته و به کاربران امکان جستجوی قدرتمندتر و کارآ را خواهد داد. مثلاً برخی از Spider ها، علاوه بر نگهداری اطلاعات مربوط به عناوین و لینک ها، یکصد کلمه با فرکانس تکرار بیشتر در صفحه وب و کلمات موجود در بیست خط اولیه را نیز نگهداری خواهند کرد. مرکز جستجوی Lycos از رویکرد فوق استفاده می نماید.

سیستم های دیگر نظیر «Altavista»، از روش خاص خود در این زمینه استفاده می نمایند. در سیستم فوق برای هر یک از کلمات موجود در صفحه شامل «an»، «a» و «the» و سایر کلمات مشابه نیز ایندکس ایجاد می گردد.

تگ های متا

با استفاده از تگ های متا، ایجاد کنندگان صفحات وب می توانند کلمات کلیدی موجود در صفحه و نحوه ایندکس نمودن آن را مشخص نمایند. روش فوق در مواردی که یک کلمه دارای بیش از یک معنی می باشد بسیار مفید و کارساز خواهد بود. بدین ترتیب تگ های فوق به مراکز جستجو راهنمایی لازم در خصوص انتخاب معنی مربوط به کلمات مورد نظر را خواهند داد. در این راستا ممکن است برخی از راهنمایی های انجام شده نیز اغفال کننده بوده و نتایج مثبتی را برای مراکز جستجو بدنبال نداشته باشد. بمنظور پیشگیری از راهنمایی های اغفال کننده توسط تگ های متا، برنامه های Spider عملیات بازبینی محتویات یک صفحه وب را بمنظور تطبیق با اطلاعات ارائه شده توسط تگ های متا، انجام می دهند. اطلاعات نادرست ارائه شده بوسیله تگ های متا، توسط Spider ها نادیده گرفته می شود.

تمام موارد فوق مفروض به حالتی است که ایجاد کننده صفحه وب قصد معرفی صفحه ایجاد شده خود را به مراکز جستجو دارد. در برخی موارد ممکن است تمایلی به انجام این کار وجود نداشته باشد.

ایجاد ایندکس

پس از اینکه عملیات Spider ها در رابطه با یافتن اطلاعات به اتمام رسید، (عملیات فوق در عمل با توجه به ماهیت وب و استقرار صفحات وب جدید هرگز به پایان نخواهد رسید، بنابراین همواره عملیات جستجو و یافتن اطلاعات توسط Spider ها انجام می گیرد) مراکز جستجو می بایست اطلاعات مورد نظر را بگونه ای ذخیره نمایند که قابل استفاده باشند. دو روش عمده در این راستا وجود دارد:

* اطلاعات به همراه داده ذخیره گردند.

* با استفاده از روشی اطلاعات ایندکس گردند.

در ساده ترین حالت، یک مرکز جستجو می تواند صرفاً کلمه و آدرس URL آن را ذخیره نماید. روش فوق در بازیابی اطلاعات و جستجو توسط کاربران ایجاد محدودیت خواهد کرد. با استفاده از روش فوق می توان جایگاه و وزن یک کلمه در یک صفحه وب را مشخص نمود. مثلاً می توان تشخیص داد که کلمه مورد نظر چند مرتبه در صفحه تکرار شده و یا لینک های موجود در صفحه نیز شامل کلمه مورد نظر می باشند یا خیر. بدین ترتیب امکان ارائه یک لیست از صفحات وب که شامل کلمه مورد نظر بر اساس میزان تکرار می باشند، وجود نخواهد داشت.

بمنظور ارائه نتایج مفیدتر توسط مراکز جستجو، اکثر مراکز جستجو صرفاً کلمه و آدرس URL را ذخیره نمی نمایند. در این حالت مواردی نظیر تعداد تکرار کلمه در صفحه نیز ذخیره خواهد شد. مراکز جستجو همچنین به هر entry یک وزن را نسبت خواهند داد. وزن نسبت داده شده، نشان دهنده جایگاه کلمه در صفحه است (ابتدای صفحه، در لینک ها، در تگ های متا و یا در عنوان صفحه) هر یک از مراکز جستجو برای اختصاص یک وزن مناسب به کلمه مورد نظر از یک فورمول استفاده می نمایند. موضوع فوق یکی از دلایلی است که جستجو یک کلمه توسط دو مرکز جستجو، نتایج مشابه ای را بدنبال خواهد داشت.

مراکز جستجو بدلیل استفاده بهینه از فضای ذخیره سازی، اطلاعات مورد نظر را بصورت رمز شده ذخیره می نمایند. مثلاً در نسخه اولیه سایت Google از دو بایت بمنظور ذخیره سازی اطلاعات مربوط به کلمات در یک صفحه استفاده می کردند. کلمات بصورت حروف بزرگ به همراه اندازه فونت، وزن و موقعیت آن ذخیره می گردید. هر یک از فاکتورهای فوق دو و یا سه بیت از دو بایت اشاره شده را به خود اختصاص می دادند. بدین ترتیب اطلاعات گسترده ای بصورت فشرده ذخیره و سپس عملیات ایجاد ایندکس انجام می گیرد.

ایندکس دارای صرفاً یک هدف است: امکان یافتن اطلاعات با سرعت بالا. برای ایجاد ایندکس از روش های متعددی استفاده می گردد. یکی از بهترین روش های موجود، ایجاد یک جدول Hash است. در روش hashing، از یک فرمول برای اختصاص یک عدد به یک کلمه استفاده می گردد. توزیع عددی با توزیع کلمات بصورت الفبائی با یکدیگر متفاوت بوده و همین امر، موثر بودن جدول hash را بدنبال خواهد داشت.



پورتال مقالات کامپیوتر و فناوری اطلاعات

در زبان انگلیسی حروفی وجود دارد که کلمات زیادی با آنان آغاز می‌گردد. مثلاً «بخش حرف M». در یک دیکشنری به‌راتب قطورتر از حرف «X» است. بدین ترتیب جستجو و یافتن کلماتی که با حرف M شروع می‌گردند زمانی به مراتب بیشتری نسبت به یافتن کلماتی که با حرف «X» آغاز می‌گردند، را طلب می‌کند. در روش hashing، با آگاهی از موارد فوق، بگونه‌ای رفتار می‌گردد که متوسط زمان بازبایی هر کلمه کاهش یابد. در روش فوق ایندکس از داده واقعی متمایز می‌گردد. جدول hash، شامل شماره hash به همراه اشاره‌گری است که به داده مورد نظر اشاره می‌نماید. با ایجاد یک سیستم ایندکس مناسب و ذخیره‌سازی مطلوب اطلاعات، امکان ارائه نتایج مفید برای کاربران را فراهم خواهد کرد.

جستجو

کاربران برای جستجوی اطلاعات مورد نیاز، پس از ورود به سایت مرکز جستجو، یک query را ایجاد می‌نمایند. query می‌تواند ساده و صرفاً شامل یک کلمه و یا پیچیده و استفاده از عملگرهای منطقی باشد. برخی از عملگرهای منطقی عبارتند از:

* AND. تمام کلماتی که توسط AND یکدیگر ملحق شده‌اند، می‌بایست در صفحه موجود باشند. در برخی از مراکز جستجو از عملگر «+» بعنوان عملگر جایگزین AND نیز استفاده می‌شود.

* OR. حداقل یکی از کلماتی که توسط OR یکدیگر ملحق شده‌اند، می‌بایست در صفحه موجود باشد.

* NOT. کلمه بعد از NOT نباید در صفحه موجود باشد. برخی از مراکز جستجو از عملگر «-» نیز استفاده می‌نمایند.

* Followed by. یکی از کلمات می‌بایست مستقیماً پس از کلمه دیگر وجود داشته باشد.

آینده مراکز جستجو

جستجوهای که توسط عملگرهای منطقی تعریف می‌گردند از نوع جستجوهای literal می‌باشند. مراکز جستجو بمنظور ارائه نتایج مورد نظر کاربر، دقیقاً کلمه و یا کلمات مشخص شده توسط کاربر در بانک اطلاعاتی جستجو می‌نمایند. روش فوق در مواردی که یک کلمه دارای بیش از یک معنی باشد، نتایج مثبتی را بدنبال نخواهد داشت. برای اخذ نتایج دلخواه، کاربران اینترنت می‌توانند با استفاده از عملگرهای منطقی محدودیت هائی را ایجاد نمایند، روش مناسب این است که محدودیت فوق از طریق مرکز جستجو اعمال گردد.

یکی از مواردی که اخیراً توسط محققین متفاوتی در مراکز جستجو دنبال می‌گردد، جستجو بر اساس مفهوم است. در روش فوق با استفاده از آنالیزهای آماری بر روی صفحات شامل کلمات سعی در ارائه نتایج مطلوبتری وجود دارد. در برخی موارد دیگر استفاده از زبانهای طبیعی برای جستجو دنبال می‌گردد. در روش فوق برای طرح سوال خود از یک مرکز جستجو از روشی که توسط انسان برای طرح سوالات مربوطه استفاده می‌گیرد، استفاده خواهد شد. در این راستا ضرورتی به استفاده از عملگرهای منطقی و یا query های پیچیده نخواهد بود.

منبع: جنوبی‌ها



پورتال مقالات کامپیوتر و فناوری اطلاعات

آداب جستجو در اینترنت

اغلب کسانی که با شبکه اینترنت کار می کنند، زمان زیادی را به جستجو بر روی شبکه می گذرانند. یافتن یک فرد، یک محصول، یک شرکت و بسیاری موارد دیگر می توانند هدف کاربر از جستجو بر روی شبکه باشند. شما نیز اگر تا به حال به جستجو بر روی اینترنت پرداخته باشید، متوجه شده اید که با توجه به حجم بسیار زیاد اطلاعات موجود بر روی شبکه یافتن اطلاعات مفید و مناسب کار ساده ای نیست. در واقع هرگاه به دنبال اطلاعاتی می گردید با دو سؤال مهم روبه رو هستید:

* چگونه جستجو کنیم ؟

* کجا جستجو کنیم ؟

در صورتی که بتوانید به این دو سؤال پاسخ مناسب و روشنی بدهید، به راحتی خواهید توانست اطلاعات مورد نظر خود را بر روی شبکه بیابید.

چگونه جستجو کنیم ؟

جستجو کردن نیز مانند هر کار دیگری آداب و رسوم خاص خود را دارد. در واقع اکثر سایت هایی که امکان جستجو را برای شما فراهم می آورند، الگوی یکسان و استانداردی را برای دریافت عبارت پرس وجو از شما دارند. شما می توانید با به کار بردن این الگوها و استفاده از نحوه نگارش صحیح عبارات جستجو، اطلاعات صحیح را به دست آورید و اپراتورهای زیر، اپراتورهای استاندارد هستند که در اکثر سایت هایی که امکان جستجو را فراهم آورده اند، قابل استفاده می باشند.

* AND : هنگامی که شما به دنبال صفحاتی می باشید که حاوی چند کلمه یا عبارت معین باشند از این اپراتور استفاده کنید. به طور مثال زمانی که به دنبال اطلاعاتی در مورد اینترنت و امنیت هستند عبارت Internet And Security را مورد استفاده قرار دهید.

* OR : زمانی که وجود حداقل یکی از چند کلمه یا عبارات معین در صفحات، مورد نظر شما باشد، اپراتور OR را در بین این کلمات به کار گیرید. به طور مثال Internet Or Security، صفحات سایت هایی را در اختیار شما قرار خواهد داد که دارای اطلاعاتی در مورد اینترنت، امنیت یا هر دو باشند.

* + : این اپراتور به عنوان پیشوند قبل از کلمات مورد نظر به کار می رود. در واقع کلماتی که پیش از آنها این اپراتور گذارده شدند، لزوماً در نتایج جستجو وجود خواهند داشت. Internet+Security صفحات و سایت هایی را در نتایج حاصله از جستجو بر می گرداند که لزوماً دارای لغت امنیت هستند، اما اینترنت می تواند در آنها وجود داشته یا نداشته باشد.

* - : این اپراتور که همانند اپراتور قبلی به صورت پیشوند به کار می رود، عدم وجود کلمه پسین خود را از سایت جستجو می خواهد. یعنی نتایج حاصل از جستجو لزوماً کلمه پس از این اپراتور را نخواهد داشت. به طور مثال Internet-Security صفحاتی را در اختیار شما قرار می دهد که دارای لغت اینترنت باشد، اما نامی از امنیت در آن صفحات ذکر نشده باشد.

* «» : هرگاه به دنبال یک عبارت هستید و می خواهید کلمات به همان شکل و ترتیب در متن سایت های نتایج جستجو یافت شوند، از این اپراتور استفاده کنید. به طور مثال زمانی که به دنبال اطلاعاتی در مورد امنیت شبکه های کامپیوتری هستید باید عبارت "Network Security" را وارد نمایید. در صورتی که به جای این عبارت از Network+Security استفاده نمایید، اگر در صفحه ای کلمه Security در ابتدای آن و Network در جای دیگری ذکر شده باشد، این صفحه نیز به عنوان نتیجه جستجو در اختیار شما قرار خواهد گرفت در حالی که حاوی اطلاعات مورد نظر شما نیست.

این اپراتور به عنوان پسوندی به معنای هر حرف یا مجموعه از حروف به کار می رود به طور مثال در صورتی که به دنبال Inter* بگردید، تمام صفحاتی که شامل کلمه های Internet, Internal, International و هر کلمه دیگری که با Inter شروع می شود و دارای پسوندی در ادامه است، به عنوان نتیجه به شما برگردانده خواهد شد.

* () : در صورتی که بخواهید عبارات جستجوی پیچیده تری را ایجاد نمایید، می توانید از پرانتز استفاده کنید. به طور مثال اگر به دنبال مشکلات نرم افزاری یا سخت افزاری هستید، می توانید از عبارت (Software Or Hardware) Problems استفاده نمایید.

تعداد این اپراتورها در موتورهای جستجو بیشتر است اما به دلیل کاربرد کم سایر موارد و طولانی شدن مطلب، از ذکر آنها صرف نظر می کنیم. دوستان علاقه مند می توانند در منابع آموزشی اینترنت اطلاعات کاملی بدست آورند.



پورتال مقالات کامپیوتر و فناوری اطلاعات

کجا جستجو کنیم ؟

حال که توانستیم عبارت مناسب برای جستجو را بسازیم، باید این عبارت را در یک سایت جستجو نماییم تا به نتایج مورد نظر دست یابیم. یکی از سایت های بسیار متداول برای جستجو Google است. اگر با استفاده از این سایت به نتایج دلخواه خود نرسیدید، استفاده از سایت هایی مانند Yahoo و Altavista توصیه می شود.

در صورتی که به دنبال موضوعی کمیاب هستید و با استفاده از دو مجموعه فوق نتوانستید به جواب مناسب دست یابید، از سایت هایی مانند Mamma و Meta Crawler استفاده نمایید. این دو سایت از مجموعه ای هستند که اصطلاحاً به آنها موتور فرا جستجو (Meta Search Engine) گفته می شود و زمانی که شما بر روی یکی از اعضای این گروه سایت ها جستجو می کنید، پرس و جو به سایت های جستجوی دیگر فرستاده می شود و پاسخ های گرفته شده از مجموعه لیست ها به صورت یک لیست در اختیار شما قرار می گیرد.

منبع : www.aftab.ir



پورتال مقالات کامپیوتر و فناوری اطلاعات

چه نوع موتور جستجو (Search Engine) یا دایرکتوری (Directory) باید استفاده کرد؟

بیشتر مردم به ابزارهایی علاقمندند که اطلاعات را در شبکه برایشان پیدا کند. اصولاً دو نوع سرویس جستجو در وب وجود دارد:

* موتورهای (Engines)

* دایرکتوریها (Directories)

سرویس جستجوی دایرکتوریها (Directories)

دایرکتوریها یک سلسله بانکهای اطلاعاتی هستند با لینکهایی به وب سایت مرجع. که این وب سایتها بوسیله اشخاص در زمان حال بوجود آمده است و بر اساس قانون مختص به سرویس جستجو طبقه بندی می شود.

Yahoo مادر همه دایرکتوریها است.

Look smart هم کاملاً عمومی است و این دایرکتوری را می توانید در سایتهایی مانند MSN پیدا کنید.

دایرکتوریها بسیار مفید هستند زمانیکه شما بیش از یک مفهوم یا توصیف برای آنچه که جستجو می کنید ندارید.

صفحه اول معمولاً مقالات مورد نظر را به شما می دهد. اکنون از میان مقاله های موجود می توانید مقاله مورد نظر خود را بیابید و آدرس مقاله یا صفحه یا سایتی را که برای شما جالبتر است انتخاب کنید و شروع به خواندن کنید اگر شما این شکل جستجو را برای پیدا کردن اطلاعات در زمینه مورد علاقه خود انتخاب کردید و به این روش سایتی را کشف کردید باید به خاطر داشته باشید که شما واقعاً آن متنی را که پیدا کرده اید جستجو نکرده اید بلکه متنی را که شامل تیترو سایت و توصیف آن است جستجو کرده اید. اینها (تیترو توصیف آن) بوسیله تنظیم کننده دایرکتوری و اغلب مبنی بر پیشنهاد صاحبان سایتها ساخته شده است. البته بعضی از دایرکتوریها اگر نتوانند در خواست شما را پاسخگو باشند اطلاعاتی را از موتورهای جستجو برای شما پیدا می کنند.

موتورهای جستجو (Search Engines)

موتورهای جستجو موتورهای روباتهایی هستند که صفحات وب (Web pages) را می پیمایند و صفحات جدید را پیدا می کنند. این روباتها صفحات وب را می خوانند و تمام یا قسمتی از متن را در یک بانک اطلاعاتی یا فهرستی که شما امکان دستیابی به آن را داشته باشید قرار می دهند هیچکدام از آنها تمام شبکه را پوشش نمی دهند ولی بعضی از آنها واقعاً بزرگ هستند.

بازیکنان اصلی در این زمینه عبارتند از : Alta vista , Google, Fast All the Web , Askjeeves , Inktomi

«Inktomi» در اصل یک سایت جستجو نیست ولی MSN , Hot bot را تغذیه اطلاعاتی می کند. Fast بیشتر پرتالهایی مانند lycos را تغذیه می کند.

موتورهای جستجو باید اولین انتخاب شما باشد زمانیکه دقیقاً می دانید که دنبال چه چیزی می گردید. آنها همچنین قسمت خیلی بزرگتری از وب را در مقایسه با دایرکتوریها می پوشانند.

هر چند فرق میان دایرکتوریها و موتورهای جستجو آنقدر که در گذشته مشخص بود اکنون نیست همه دایرکتوریهای جستجو نتایج جستجو را ابتدا از موتورهای جستجو پیدا می کنند و اگر در آن جا جستجو موفق نبود به سراغ دایرکتوری خودشان می روند مثلاً «Yahoo» از موتور جستجوی «Google» برای این منظور استفاده می کند.

از طرف دیگر بعضی موتورهای جستجو اطلاعات را از دایرکتوریهای جستجو تهیه می کنند قبل از اینکه از بانک اطلاعاتی موتورهای جستجو اطلاعاتی را به شما بدهند.

Ask Jeeves سعی می کند جوابها را به سئوالات شما از زبان طبیعی بانک اطلاعاتی خودش بیرون بکشد. یعنی یک سری پاسخها به به پرسشهای متداول را در خود دارد.

Lycos , Google , Aol به شما قابلیت دسترسی به دایرکتوریای باز را می دهند MSN , Alta Vista یک سری سلسله دایرکتوری دارند مشابه Yahoo اینها بر اساس Look smart directory ساخته شده اند اگر آنچه را که جستجو می کنید نتوانستید پیدا کنید نتایج را در سایتهای مشابه برای شما پیدا خواهد کرد.



پورتال مقالات کامپیوتر و فناوری اطلاعات

Metasearch engines

همچنین سرویسهای Metasearch وجود دارند مانند search.com و یا Metacrawler که مربوط به GO۲NET و یا Meta search وب سایت Pandia

Meta search engine ها در یک زمان در چندین موتور جستجو و دایرکتوری دنبال پاسخ شما می گردند و سعی می کنند مناسب ترین پاسخ را از بین همه بیرون بکشند.

ممکن است برای شما استفاده از یکی از آنها مفید واقع نشود فقط برای اینکه یک دیدگاه کلی از آنچه در خارج از یک موتور جستجو به تنهایی بدست می آورید.

مشکل بیان جمله و کلمات در موتورهای جستجو است به هر حال این ممکن است از یک موتور جستجو تا موتور جستجوی دیگر فرق کند و این به این معنی است که فرا موتور جستجو (Metaserch engine) سعی خواهد کرد سؤال شما را به زبانی که هر موتور جستجو می فهمد ترجمه کند اما اغلب این گونه عمل نمی کند.

برای جستجوهای پیچیده شما مستقیماً باید موتور جستجوی مناسب برای آن را پیدا کنید همچنین توجه داشته باشید که Metasearch engine به شما قسمت کوچکی از نتایج هر یک از موتورهای جستجو را می دهد.

search utilities

مشابه Metasearch engine ها هستند با این تفاوت که برنامه هایی هستند که از روی کامپیوتر شما اجرا می شوند می توانند در اینترنت به جستجو پردازند که اصطلاحاً به آنها search utility می گویند.

بیشتر آنها مانند فرا موتورهای جستجو (Metasearch engines) کار می کنند. آنها از چندین موتور جستجو نتایج را جمع آوری می کنند و این زمانی است که شما درخواست خود را مطرح می کنید.

برای این منظور از «copernic» برای PC و از «sherlock» برای مکینتاش می توانید استفاده کنید. این برنامه ها برای یک جستجوی ساده خیلی مفید هستند اما اگر شما می خواهید جستجوی پیشرفته تری داشته باشید باید در هر بار جستجوی خود را به یک سرویس جستجو محدود کنید.

بهترین سرویس جستجو

احتمالاً بیش از ۷ میلیارد پرونده (متن و ..) در وب موجود می باشد در حال حاضر سریعترین سرویسها Fast All The Web و Google و inktomi هستند که مجموعاً بیش از ۲ میلیارد صفحه وب را فهرست می کنند. حقیقتی که در این میان باقی می ماند این است که موتورهای جستجو همه یک قسمت را پوشش نمی دهند و این دلیلی است که شما را وادار می دارد تا از بیش از یک موتور جستجو استفاده کنید.

بر اساس یک سری تحقیقات Google و Alta Vista و Fast All The Web بهترین موتورهای جستجو هستند در حالیکه به نظر می رسد Yahoo بهترین دایرکتوری باشد برای Metasearch Ma Visimo و Inquich را پیشنهاد می کنیم.

منبع : www.yadbegir.com



پورتال مقالات کامپیوتر و فناوری اطلاعات

جایگاه موتور های جست و جو

موتورهای جست و جو (search engines) اکنون در فعالیت های اینترنتی، نقش غیرقابل انکاری پیدا کرده اند. بسیاری از تحقیقات نشان داده اند که موتور های جست و جو ابزاری مهم و محبوب برای یافتن اطلاعات مورد نیاز کاربران است. در آگوست ۲۰۰۴، مطالعه ای توسط موسسه Pew و ComScore در زمینه جایگاه موتور های جست و جو انجام شده است که خلاصه ای از نتایج آن را در این نوشته می خوانید. این تحقیق، با استفاده از پرسش تلفنی از ۱۳۹۹ کاربر اینترنتی در سراسر آمریکا انجام گرفته است.

مطالعه مورد اشاره نشان داده است که استفاده از موتور های جست و جو یکی از محبوب ترین فعالیت های اینترنتی است. اکثر اوقات، استفاده از موتور های جست و جو، پس از پست الکترونیک در جایگاه دوم قرار دارد. البته گاهی که اتفاق مهمی رخ می دهد (مثل جنگ در عراق) دریافت اخبار از طریق اینترنت بر استفاده از موتور های جست و جو پیشی می گیرد. از نظر آماری ۸۴ درصد کاربران اینترنتی بیان داشته اند که از موتور های جست و جو استفاده می کنند. حجم جست و جو نیز قابل توجه است. هر کاربر اینترنتی به طور میانگین ۳۳ جست و جو انجام می دهد. به این ترتیب، فقط در آمریکا، با استفاده از ۲۵ موتور جست و جوی متداول، ۹/۳ میلیارد جست و جو در ماه انجام می شود. همچنین در حالی که ۶۴ میلیون نفر از شهروندان آمریکا هر روز به اینترنت متصل می شوند، بیش از ۳۸ میلیون نفر از آنان از موتور های جست و جو استفاده می کنند.

در مورد میزان سودمندی موتور های جست و جو نیز بیشتر کاربران، اعتماد و رضایت خود را هنگام کار با موتور های جست و جو نشان داده اند. ۸۷ درصد کاربران ابراز کرده اند که بیشتر اوقات به آنچه جست و جو کرده اند، دست یافته اند. از این گروه ۲۰ درصد همیشه از نتایج جست و جو رضایت داشته اند.

همچنین بسیاری از کاربران مهارت خود را در استفاده از موتور های جست و جو افزایش داده اند. ۹۲ درصد جست و جو گران اعلام کرده اند که به مهارت خود در یافتن نتایج اطمینان دارند. بیش از نیمی از آنان گفته اند که کاملاً مطمئن هستند و می توانند آنچه را در نظر دارند با استفاده از جست و جوی اینترنتی پیدا کنند.

در دسترس بودن موتور های جست و جوی مطمئن و ساده، نحوه دستیابی کاربران به اطلاعات را تغییر داده است. بسیاری از کاربران به میزان زیادی، برای دستیابی به اطلاعاتی که برای آنان بسیار حیاتی است وابسته به موتور های جست و جو هستند: ۴۴ درصد جست و جوگران بیان داشته اند اطلاعاتی که از طریق موتور های جست و جو می یابند، برایشان کاملاً ضروری است. اهمیت این آمار آنگاه دوچندان می شود که بدانیم یک سوم جست و جو کنندگان از طریق اینترنت بیان کرده اند که بدون موتور های جست و جو نمی توانند به فعالیت های خود ادامه دهند. اهمیت موتور های جست و جو برای نیمی از کاربران کمتر است. آنان ابراز داشته اند که موتور های جست و جو علاقه دارند ولی در صورت لزوم می توانند اطلاعات مورد نظر خود را از روش های دیگر به دست آورند. در میان موتور های جست و جوی متداول میزان رضایت از گوگل (Google) بسیار زیاد است. در این مطالعه، ۴۷ درصد پرسش شوندگان اعلام کرده اند گوگل انتخاب اول آنان برای جست و جو است.

ياهو (Yahoo) با ۲۶ درصد در مقام دوم قرار دارد. پس از این دو، MSN، AOL، Askjeeves و Netscape با ۷، ۵، ۲ و ۱ درصد قرار گرفته اند. به این ترتیب نتایج این مطالعه برای گوگل و ياهو بسیار خوشحال کننده است. میزان اعتماد کاربران به موتور های جست و جو نیز جالب است. ۴۴ درصد کاربران گفته اند تنها از یک موتور جست و جو برای یافتن نتایج مورد نظرشان استفاده می کنند. ۴۸ درصد نیز ۲ یا ۳ موتور جست و جو را تجربه می کنند.

به عنوان بخشی از این مطالعه، موسسه ComScore برای سنجش میزان رضایت کاربران از آنان سؤال کرد کدام ویژگی موتور های جست و جو برای آنان بیشتر مهم بوده است و موتور جست و جوی مورد استفاده آنان از این نظر در چه جایگاهی قرار دارد.

نتیجه این بخش از مطالعه آن است که گرچه میزان رضایت کاربران از موتور های جست و جو به طور کلی بسیار زیاد است، ولی فاصله نسبتاً زیادی میان توقع آنان از نظر ویژگی های مورد نظرشان با آنچه از موتورهای جست و جو کسب کرده اند وجود دارد. برای مثال، ۹۱ درصد از پرسش شوندگان، عبارت «نتایجی برمی گرداند که با نیازهای شما تطبیق دارد» را به عنوان مهمترین ویژگی مد نظر اعلام کرده اند. در حالی که تنها ۶۴ درصد آنان گفته اند که موتور جست و جوی مورد علاقه شان توقع آنان را از این نظر برآورده می سازد.

همچنین ۸۵ درصد کاربران محفوظ نگه داشتن اطلاعات مورد جست و جو را مهم دانسته اند، اما تنها ۵۵ درصد اعلام کرده اند که موتور جست و جوی مورد علاقه شان این موضوع را رعایت می کند. از دیگر پارامترهای مهم مورد نظر کاربران می توان به این موارد اشاره کرد: ۹۰ درصد کاربران، راحتی استفاده را عامل مهم می دانند. از نظر ۸۹ درصد آنان، سرعت بازگرداندن اطلاعات بسیار مهم است.

جالب است که پیشنهاد کلمه جست و جوی جدید هنگامی که کلمه خوبی از سوی کاربر جست و جو نشده است از ویژگی های مورد علاقه ۷۰ درصد کاربران بوده است.

منبع : مرکز توسعه و تبادل دانش فناوری اطلاعات



پورتال مقالات کامپیوتر و فناوری اطلاعات

بررسی کمی و کیفی موتور های جستجو

گوگل با ادعای ایندکس کردن ۸ میلیارد صفحه سهم خود را در مبارزه بین اندازه موتور های جستجو گر افزایش داده است.

مایکروسافت هم که قصد دارد عنوان بزرگترین موتور جستجوگر جهان را کسب کند امروز ۵ میلیارد صفحه را ایندکس کرده است. این مسئله موجب می شود تا از ۴.۲ میلیارد صفحه به گزارش خود گوگل در سال گذشته ایندکس شده است پیشی بگیرد. ما پیشتر از این نیز در جریان جنگ اندازه قرار داشتیم.

این نزاع وقتی بوجود می آید که یکی از موتور های جستجو در تلاش برای بدست آوردن مشخصه مثبتی باشد که آن را نسبت به سایرین در جایگاه بالاتری قرار می دهد. البته باید توجه داشت که بزرگی اندازه دلیلی برای اثبات برتری نیست.

یک موتور جستجو به همراه شمار زیادی از صفحات، اگر به موقع بروز رسانی نشود و یا ارتباطات صحیح وجود نداشته باشد ممکن است حقیقتا ضعیف تر از موتور جستجویی باشد که صفحات کمتری را در بردارد. بطور خلاصه در مورد میزان کمیت و کیفیت موتور های جستجو دو معیار میزان صفحات ایندکس شده و میزان عمق صفحات بسیار مهم می باشند. در ادامه این دو معیار را به تفصیل بررسی می کنیم.

اندازه گزارش شده Reported Size

این عدد در حقیقت آن چیزی که موتور های جستجو ادعا می کنند. در مورد گوگل این عدد برخی اوقات شامل صفحات ایندکس شده نسبی «partially-indexed» و یا به معنای روشنتر صفحات فقط لینک-link-only است. اینها صفحاتی بوده اند که گوگل منحصر از طریق لینک می شناسد و خود این صفحات در واقع ایندکس نشده اند.

نوعا، موتور های جستجو قادر به شناسایی صفحات تکراری، صفحات اسپم و نظیر آن نیستند. ما اینجا به این مسئله نمی پردازیم و فرض می کنیم چنین قابلیتی دارند.

اما یاهو سعی می کند از بازی اندازه صفحات خودش را بیرون بکشد. اوایل امسال وقتی تکنولوژی جدید جستجوی خود را پیاده سازی می کرد، از اعلام میزان صفحات ایندکس شده خودداری کرد و به گفتن این جمله که «موتور جستجوی ما قابل مقایسه با دیگران است» اکتفا کرد. سخنگوی یاهو در این باره می گوید: «همانند گذشته یاهو سائز صفحات ایندکس شده خود را به دلایل رقابتی ذکر نمی کند. ما اعتقاد داریم ایندکس ما کاملا قابل رقابت با دیگران است.

کیفیت جستجو شامل فاکتور های متفاوتی نظیر بروز رسانی، میزان وابستگی و ... است. ما تلاش می کنیم تا نتایج با کیفیت بالاتری برای مصرف کنندگان برسیم تا به آنها اطمینان دهیم که قادر خواهند بود بهترین نتیجه را مورد نظرشان را از طریق موتور جستجوی ما بدست بیاورند. این رویه یاهو هم خوشایند است هم ناخوشایند. خوشایند از آن جهت که گفتن این عدد مانع از وارد شدنشان به بحث بیهوده اندازه صفحات ایندکس شده می شود. نا خوشایند از آن که گفتن این عدد موجب می شود تا آمار دقیقی از عملکرد آنها نداشته باشیم. به نظر می رسد همه موتور های جستجوگر باید آنرا بدون ریاکاری ذکر کنند. با این حساب می توان میزان صفحات ایدکس شده یاهو را با توجه به ادعای قابل رقابت بودن موتور جستجوی آنها ۴ و ۲ میلیارد صفحه فرض کرد. این رقم معادل میزان صفحاتی است که چند ماه پیش گوگل است.

عمق صفحات Page Depth Amount

بررسی عمق صفحات ایندکس شده بسیار جالبتر خواهد بود. فرض کنید که شما میلیون صفحه را می خواهید ایندکس کنید. آیا حقیقتا کل متن صفحات را ایندکس می کنید یا بخشی از آنرا؟ با استفاده از این معیار می توان فهمید که عملکرد یک موتور جستجو چگونه است. در مورد گوگل باید تقریبا می توان گفت که بعضی از صفحات را تنها بصورت نسبی و بخشی از آنرا ایندکس می کند. در گذشته اگر اندازه صفحه ای بیشتر از ۱۰۱K بود تنها ۱۰۱K اول متن توسط گوگل ایندکس می شد و بقیه متن کنار گذاشته می شد. بر طبق فرضیات دقیق من گوگل هنوز هم به همین شکل عمل می کند.

میزان عمق صفحات ایندکس شده MSN را می توان از گفته هایی که در آگوست سال گذشته به هنگام جلسه ملاقات کرولر ها «Crawlers» در همایش موتور های جستجو در سن خوزه بیان شد فهمید. البته در مورد نسخه جدید جستجوگر MSN شاید شرایط فرق کرده باشد. میزان این معیار در یاهو نیز در همان جلسه ای ذکرش رفت بیان شد. اما مدیران Ask Jeeves از گفتن میزان عمق صفحاتشان در جلسه مذکور خودداری کردند و به جمله «ما در حد دیگران هستیم» قناعت کردند. لذا می توان عمق صفحات آنها را چیزی شبیه به گوگل دانست. بررسی پیدا کردن میزان این معیار در موتور های جستجو بسیار آسان است. برای این کار کافی است یک صفحه حجیم که طول زیادی دارد را پیدا کنید و به



پورتال مقالات کامپیوتر و فناوری اطلاعات

جستجویی مطلبی که در انتهای آن صفحه قرار دارد دارد پردازید. نتایج بدست آمده را مقایسه کنید تا ببینید کدام جستجوگر عمق پایین تر را بهتر سرچ می کند. در بررسی نتیجه یک تحقیق مشخص شد که یاهو حقیقتا قادر است تا عمق ۸۰۰K صفحات را نیز جستجو کند.

در پایان باید گفت هر چند بیشتر بودن معیار میزان صفحات ایندکس شده یک موتور جستجوگر مفید است، اما هرگز به معنای مرتبط بودن و درست بودن صفحات ایندکس شده نیست. هر چند گوگل تقریبا دو برابر یاهو صفحه ایندکس کرده است اما به این معنی نیست که دو برابر یاهو بهتر و کارآمدتر است.

منبع : انجمن علمی دانشگاه شیخ بهایی



پورتال مقالات کامپیوتر و فناوری اطلاعات

وب پنهان، اطلاعاتی که موتور جستجوگر بدان راهی ندارد!

واقعیت آن است که چالش عمده ما در حال حاضر، نبود اطلاعات نیست، بلکه دسترسی به اطلاعات مهمتر شده است. آنهم دسترسی به اطلاعات دقیق و معتبر و در زمان مورد نیاز. گفته می شود وب منبع بزرگ اطلاعاتی عصر حاضر است و تقریباً درباره هر موضوعی می توان در آن اطلاعاتی یافت. در آن می توانیم درباره موضوعاتی از "پرورش لاک پشت" تا "طراحی موشک" اطلاعاتی بیابیم. اما کجا؟ به عبارت دیگر در کدام سایت؟

وب راهنمایی دارد که به کاربران برای یافتن اطلاعات کمک کند. سایتهایی وجود دارند که کاربران وب با مراجعه به آنها پاسخ سوالات خود را می یابند. ما اینگونه سایتها را با عنوان "موتورهای جستجوگر" می شناسیم. در حقیقت موتور جستجوگر سایتی است که کاربر وب با مراجعه به آن و نوشتن چند کلمه می تواند هزاران پاسخ برای سوال خود بیابد. علاوه بر مراجعه به موتورهای جستجوگر یکی از راههای دیگر جستجوی اطلاعات، استفاده از "وب پنهان" است. به راستی وب پنهان چیست؟

انواع اطلاعات موجود در اینترنت را می توان به سه دسته زیر تقسیم بندی کرد:

(۱) اطلاعات رایگان و پیدا

(۲) اطلاعات رایگان و ناپیدا

(۳) اطلاعات جاری

اطلاعات رایگان و پیدا اطلاعاتی هستند که در دسترس همگان قرار داده شده اند و با جستجو در موتورهای جستجوگر می توانیم آنها را بیابیم. موتور جستجوگر هر چقدر هم از پایگاه داده بزرگی برخوردار باشد نمی تواند تمام اطلاعات وب را در خود داشته باشد.

اطلاعات جاری اطلاعاتی هستند که برای استفاده از آن باید مبلغی پرداخت شود. و در آخر اطلاعات رایگان و ناپیدا اطلاعاتی اند که نمی توانیم از طریق موتورهای جستجوگر به آنها دسترسی داشته باشیم.

وب پنهان چیست؟

در حقیقت بخش اعظم وب از دسترسی موتورهای جستجوگر دور است که به آن وب پنهان گفته می شود. در مقابل می توانید وب نمایان را بخشی از وب بدانید که موتورهای جستجوگر می توانند به آن دسترسی داشته باشند و در نتایج جستجو به مراجعه کنندگان خود نمایش دهند.

پایگاههای داده قابل جستجو

بخش اعظم وب پنهان همین پایگاههای داده هستند. برای استفاده از اطلاعات موجود در آنها ابتدا باید کاربر فرمی را پر کند. چون موتور جستجوگر توانایی انجام این کار را ندارد بنابراین نمی تواند به اطلاعات آن دسترسی داشته باشد. در این پایگاههای داده متناسب با نیاز کاربر صفحه ساخته می شود و با توجه به حجم بالای اطلاعات عملانی توان تمام حالتها را مورد نیاز کاربر را شناسایی کرد و از قبل برای آن صفحه ای ساخت. اگر کاربری لینک مستقیمی به یکی از این صفحات تولید شده ایجاد کند آنگاه موتور جستجوگر شاید بتواند به آن اطلاعات دسترسی پیدا کند.

صفحات منفک شده

پاره ای از صفحات نیز به دلایلی از تیررس موتورهای جستجوگر دور نگاه داشته شده اند. سیاست کاری صاحبان سایتها و ضعف طراحان سایتها مهم ترین این دلایل هستند. فایلها پنهان، اسناد نیازمند رمز عبور برای خواندن و پایگاههای داده جاری مواردی (استفاده از اطلاعات به شرط پرداخت حق عضویت در سایت) را می توان در این گروه قرار داد.

آخرین مطالعه آکادمیک صورت گرفته نشان می دهد که وب نمایان ۵/۱۱ میلیارد سند دارد و موتورهای جستجوگر ۸۵ درصد آنرا می شناسند. این مطالعه همچنین حجم "وب پنهان" را ۵۰۰ میلیارد سند برآورد کرده است. در این مطالعه گوگل با بایگانی کردن ۸/۱ میلیارد سند رتبه اول را دارد (۶۹/۱ درصد) و پس از آن یاهو با ۶/۴ میلیارد صفحه (۵۷/۴ درصد) در جایگاه دوم قرار گرفته است. با این که مطالعه انجام شده بر اساس تخمینهای بسیار بوده است اما در نوع خود یکی از آخرین تلاشها برای برآورد حجم واقعی وب و میزان پوشش آن توسط موتورهای جستجوگر است.



پورتال مقالات کامپیوتر و فناوری اطلاعات

وب "پنهان" را چگونه "نمایان" کنیم؟

برای استفاده از وب پنهان ابتدا باید آدرس یکی از آنها را به کمک موتورهای جستجوگر بیابید. به عنوان مثال اگر موضوع پزشکی مد نظر شماست، کافیهست که در گوگل اینگونه جستجو کنید: پایگاه داده پزشکی یا پایگاه داده علوم پزشکی. در هر صورت کلمه "پایگاه داده"-Database- یک کلمه کلیدی است.

فراموش نکنید که "وب پنهان" به هر حال وجود دارد و صرف کمی وقت برای استفاده از آنها می تواند مکمل مناسبی باشد برای آنچه که از طریق موتورهای جستجوگری نظیر یاهو و گوگل می یابید. برخی از آنها عبارتند از:

Librarians Index

AcademicInfo

Infomine

با استفاده از www.invisible-web.net می توانید موارد بسیار دیگری از وب پنهان متناسب با نیازتان را مشخص کنید.

اینترنت منبع مهمی برای دستیابی به اطلاعات معتبر و موثق است. مهم آن است که کاربر تکنیکهای جستجو و ارزیابی اطلاعات را بداند تا بتواند بهتر و سریعتر به آنچه که می خواهد دست یابد. از سوی دیگر تمام اطلاعات از طریق موتورهای جستجوگر قابل دسترسی نیست.

ماهیت جاری بسیاری از سایتها و همین طور محدودیتهای تکنولوژیکی موتور جستجوگر را از دسترسی به تمام اطلاعات وب دور می کند. در این مورد کاربر باید بتواند از اطلاعات موجود در "وب پنهان" بهره گیرد.

منبع : www.irandevolvers.com



پورتال مقالات کامپیوتر و فناوری اطلاعات

جویشر چیست ؟

جویشر یا موتور جستجو (به انگلیسی: Search Engine)، در فرهنگ رایانه، به طور عمومی به برنامه‌ای گفته می‌شود که کلمات کلیدی را در یک سند یا بانک اطلاعاتی جستجو می‌کند. در اینترنت به برنامه‌ای گفته می‌شود که کلمات کلیدی موجود در فایل‌ها و سندهای وب جهانی، گروه‌های خبری، منوهای گوگر و آرشیوهای FTP را جستجو می‌کند.

برخی از جویشرها برای تنها یک وب‌گاه (پایگاه وب) اینترنت به کار برده می‌شوند و در اصل جویشری اختصاصی آن وب‌گاه هستند و تنها محتویات همان وب‌گاه را جستجو می‌کنند.

برخی دیگر نیز ممکن است با استفاده از SPIDERها محتویات وب‌گاه‌های زیادی را پیمایش کرده و چکیده‌ای از آن را در یک پایگاه اطلاعاتی به شکل شاخص‌گذاری شده نگهداری می‌کنند. کاربران سپس می‌توانند با جستجو کردن در این پایگاه داده به پایگاه وبی که اطلاعات موردنظر آن‌ها را در خود دارد پی ببرند.

انواع جویشرها در اینترنت

جویشرها به دو دسته کلی تقسیم می‌شوند. جویشرهای پیمایشی (خودکار) و فهرست‌های تکمیل‌دستی (غیر خودکار). هر کدام از آن‌ها برای تکمیل فهرست خود از روش‌های متفاوتی استفاده می‌کنند البته لازم به ذکر است که گونه‌ای جدید از جویشرها تحت عنوان «ابرجویشر» (Meta Search Engines) نیز وجود دارد که در ادامه به توضیح هر یک از این موارد خواهیم پرداخت :

جویشرهای پیمایشی

جویشرهای پیمایشی (Crawler-Based Search Engines) مانند گوگل فهرست خود را بصورت خودکار تشکیل می‌دهند. آنها وب را پیمایش کرده، اطلاعاتی را ذخیره می‌کنند. سپس کاربران از میان این اطلاعات ذخیره شده، آنچه را که می‌خواهند جستجو می‌کنند. اگر شما در صفحه وب خود تغییری را اعمال نمایید، جویشرهای پیمایشی آن‌ها را به طور خودکار می‌یابند و سپس این تغییرات در فهرست‌ها اعمال خواهد شد. عنوان، متن و دیگر عناصر صفحه، همگی در این فهرست قرار خواهند گرفت. وجه مشخصه این گروه از جویشرها وجود نرم افزار موسوم به SPIDER در آن‌هاست. این شبه نرم‌افزار کوچک بصورت خودکار به کاوش در شبکه جهانی پرداخته و از پایگاه‌های وب یادداشت‌برداری و فهرست‌برداری می‌کند سپس این اطلاعات را برای تجزیه و تحلیل و طبقه‌بندی به بانک اطلاعاتی جویشر تحویل می‌دهد.

فهرست‌های دست‌نویس شده

فهرست‌های دست‌نویس‌شده یا (Human-Powered Directories) مانند فهرست بازی (Open Directory) مانند Dmoz وابسته به کاربرانی است که آن را تکمیل می‌کنند. شما صفحه مورد نظر را به همراه توضیحی کوتاه 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 به وجود نخواهد آورد. چیزهایی که برای بهبود یک فهرست‌بندی در یک جویشر مفید هستند، تأثیری بر بهبود فهرست‌بندی یک دایرکتوری ندارند. تنها استثناء این است که یک سایت خوب با پایگاه داده‌ای با محتوای خوب شناس بیشتری نسبت به یک سایت با پایگاه داده ضعیف دارد. البته در مورد جویشرهای مشهور مانند گوگل و یاهو، یک مولفه دیگر هم برای بهبود فهرست‌بندی وجود دارد که کمک مالی (یا به اصطلاح اسپانسر) است. یعنی وب‌گاه‌هایی که مایل به بهبود مکان وب‌گاه خود در فهرست بندی هستند، می‌توانند با پرداخت پول به این جویشرها به هدف خویش برسند.

جویشرهای ترکیبی با نتایج مختلف

به موتورهایی گفته می‌شود که هر دو حالت را در کنار هم نمایش می‌دهند. غالباً یک جویشر ترکیبی در صورت نمایش نتیجه جستجو از هر یک از دسته‌های فوق، نتایج حاصل از دسته دیگر را هم مورد توجه قرار می‌دهد. مثلاً جویشر ام.اس.ان (MSN) بیشتر نتایج حاصل از فهرست‌های تکمیل‌دستی را نشان می‌دهد اما در کنار آن نیم نگاهی هم به نتایج حاصل از جستجوی پیمایشی دارد.



پورتال مقالات کامپیوتر و فناوری اطلاعات

ابرجویشگرها

این گونه جدید از جویشگرها که قدمت چندانی نیز ندارند. بصورت هم‌زمان از چندین جویشگر برای کاوش در شبکه برای کلید واژه مورد نظر استفاده می‌کنند. بدین معنی که این جویشگر عبارت مورد نظر شما را در چندین جویشگر دیگر جستجو کرده و نتایج آنها را با هم ترکیب کرده و يك نتیجه کلی به شما ارائه می‌دهد. به‌عنوان مثال جویشگر داگ پایل از نتایج حاصل از موتورهای Google - Yahoo - MSN و ASK استفاده کرده و نتیجه حاصله را به شما ارائه می‌دهد. لازم به ذکر است که روش و یا راهکار مشخص و یکسانی برای ترکیب نتایج حاصله از موتورهای پایه - موتورهای که به عنوان جویشگر استفاده می‌شوند مانند Yahoo که یک موتور پایه برای dogpile می‌باشد - وجود ندارد. اما dogpile قابلیت جستجو به همه زبانها را ندارد و ظاهراً فقط کلمات انگلیسی را پیدا می‌کند.

نوجویشگرها

این گونه از جویشگرها، نسل جدید و متفاوتی از جویشگرهای گذشته هستند. امکان ثبت جستجو و مدل‌سازی فعالیت‌های کاربر و ارائه نتایج جدید به کاربر، به‌صورت متفاوت و تفکیک شده، از امکانات نوجویشگرها است.

بررسی یک جویشگر پیمایشی

جویشگرهای پیمایشی شامل سه عنصر اصلی هستند. اولی در اصطلاح عنکبوت (Spider) است که پیمایشگر (Crawler) هم نامیده می‌شود. پیمایشگر همین که به یک صفحه می‌رسد، آن را می‌خواند و سپس پیوندهای آن به صفحات دیگر را دنبال می‌نماید. این چیز است که برای یک سایت پیمایش شده (Crawled) اتفاق افتاده است. پیمایشگر با یک روال منظم، مثلاً یک یا دو بار در ماه به سایت مراجعه می‌کند تا تغییرات موجود در آن را بیابد. هر چیزی که پیمایشگر بیابد به عنصر دوم یک جویشگر 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 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 هنگامی که آن تغییر در فهرست جویشگر ثبت شود مدت زمانی طول خواهد کشید. پس ممکن است که یک 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 جویشگر، سومین عنصر یک جویشگر است و به برنامه‌ای اطلاق می‌شود که به صورت 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.

رتبه‌بندی صفحات وب توسط جویشگرها

وقتی شما از جویشگرهای 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 جویشگرها 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 می‌خواهید انجام دهید اما جویشگرها کار حیرت‌انگیز دیگری نیز انجام می‌دهند. فرض کنید که شما به یک کتابدار مراجعه می‌کنید و از وی درباره «سفر» کتابی می‌خواهید. او برای این که جواب درستی به شما بدهد و کتاب مفیدی را به شما ارائه نماید با پرسیدن سؤالاتی از شما و با استفاده از جُارب خود کتاب مورد نظرتان را به شما تحویل خواهد داد. جویشگرها همچنین توانایی ندارند اما به نوعی آنها را شبیه‌سازی می‌کنند. پس جویشگرهای پیمایشی چگونه به پاسخ مورد نظران از میان میلیون‌ها صفحه وب می‌رسند؟ آنها یک مجموعه از قوانین را دارند که الگوریتم نامیده می‌شود. الگوریتم‌های مورد نظر برای هر جویشگری خاص و تقریباً سری هستند اما به هر حال از قوانین زیر پیروی می‌کنند:

مکان و بسامد

یکی از قوانین اصلی در الگوریتم‌های رتبه‌بندی موقعیت و بسامد (تعداد تکرار) واژه‌هایی است که در صفحه مورد استفاده قرار گرفته‌اند که بطور خلاصه روش مکان-بسامد (Location/Frequency Methode) نامیده می‌شود. کتابدار مذکور را به خاطر می‌آورید؟ لازم است که او کتاب‌های در رابطه با واژه «سفر» را طبق درخواست شما بیابد. او در وحله اول احساس می‌کند که شما به دنبال کتاب‌هایی هستید که در نامشان کلمه «سفر» را شامل شوند.



پورتال مقالات کامپیوتر و فناوری اطلاعات

جوبیشگرها هم دقیقاً همان کار را انجام می‌دهند. آنها هم صفحاتی را برایتان فهرست می‌کنند که در برجسب عنوان (Title) موجود در کد زبان نشانه‌گذاری آبرمتنی (زنگام) (HTML) حاوی واژه «سفر» باشند. جوبیشگرها همچنین به دنبال واژه مورد نظر در بالای صفحات و یا در آغاز بندها (پاراگرافها) هستند. آنها فرض می‌کنند که صفحاتی که حاوی آن واژه در بالای خود و یا در آغاز بندها و عناوین باشند به نتیجه مورد نظر شما مربوطتر هستند. بسامد عامل بزرگ و مهم دیگری است که جوبیشگرها از طریق آن صفحات مربوط را شناسایی می‌نمایند. جوبیشگرها صفحات را تجزیه کرده و با توجه به تکرار واژه‌ای در صفحه متوجه می‌شوند که آن واژه نسبت به دیگر واژه‌ها اهمیت بیش‌تری در آن صفحه دارد و آن صفحه را در درجه بالاتری نسبت به صفحات دیگر قرار می‌دهند.

چگونگی کارکرد دقیق جوبیشگرها درباره روش‌هایی از قبیل مکان-تکرار فاش نمی‌شود و هر جوبیشگری روش ویژه خود را دنبال می‌کند. به همین دلیل است که وقتی شما واژه‌های همانندی را در موتورهای متفاوت جستجو می‌کنید. به نتایج متفاوتی می‌رسید. الگوریتم‌های اولیه جوبیشگرهای معتبر و بزرگ همچنان محرمانه نگهداری می‌شوند. برخی جوبیشگرها نسبت به برخی دیگر صفحات بیشتری را فهرست کرده‌اند. نتیجه این خواهد شد که هیچ جوبیشگری نتیجه جستجوی مشترکی با موتور دیگر نخواهد داشت و شما نتایج متفاوتی را از آن‌ها دریافت می‌کنید. جوبیشگرها همچنین ممکن است که برخی از صفحات را از فهرست خود حذف کنند البته به شرطی که آن صفحات با هرزنامه (Spam) شدن سعی در گول زدن جوبیشگرها داشته باشند. فرستادن هرزنامه (Spamming) روشی است که برخی از صفحات برای احراز رتبه بالاتر در جوبیشگرها در پیش می‌گیرند و آن به این صورت است که با تکرار بیش از حد واژه‌ها و یا بزرگ نوشتن یا بسیار ریز نوشتن متن‌ها بطور عمدی کوشش در برهم زدن تعادل و در نتیجه فریب جوبیشگرها دارند. آنها سعی دارند که با افزایش عامل تکرار، در رتبه بالاتری قرار بگیرند. البته آنگونه که گفته شد تعداد تکرارها اگر از حد و اندازه خاصی فراتر رود نتیجه معکوس می‌دهد. جوبیشگرها راه‌های متنوعی برای جلوگیری از فرستادن هرزنامه دارند و در این راه از گزارش‌های کاربران خود نیز بهره می‌برند. امروزه بهینه‌سازی سایت‌های اینترنت برای جوبیشگرها یکی از مهم‌ترین روش‌های جلب بازدیدکننده به سایت است.

عوامل خارج از صفحه

جوبیشگرهای گردشی اکنون تجربه فراوانی در رابطه با وب‌دارهایی دارند که صفحات خود را برای کسب رتبه بهتر مرتباً بازنویسی می‌کنند. بعضی از وب‌دارها (وب‌مسترها)ی خبره حتی ممکن است به سمت روش‌هایی مانند مهندسی معکوس برای کشف چگونگی روش‌های مکان-تکرار بروند. به همین دلیل، تمامی جوبیشگرهای معروف از روش‌های امتیازبندی «خارج از صفحه» استفاده می‌کنند. عوامل خارج از صفحه عواملی هستند که از تیررس وب‌دارها خارجند و آنها نمی‌توانند در آن دخالت کنند و مسأله مهم در آن تحلیل ارتباطات و پیوندهاست. به وسیله تجزیه صفحات، جوبیشگرها پیوندها را بررسی کرده و از محبوبیت آنها می‌فهمند که آن صفحات مهم بوده و شایسته ترفیع رتبه هستند. به علاوه تکنیک‌های پیشرفته به گونه‌ای است که از ایجاد پیوندهای مصنوعی توسط وب‌دارها برای فریب جوبیشگرها جلوگیری می‌نماید. علاوه بر آن جوبیشگرها بررسی می‌کنند که کدام صفحه توسط یک کاربر که واژه‌ای را جستجو کرده انتخاب می‌شود و سپس با توجه به تعداد انتخاب‌ها، رتبه صفحه مورد نظر را تعیین کرده و مقام آن را در نتیجه جستجو جابه‌جا می‌نمایند.

سرفصل‌های بهینه سازی

* تدوین استراتژی

- * بازنویسی محتوای سایت با توجه به هدف و با مساعدت شما
- * تحقیق و انتخاب کلمات کلیدی مرتبط با فعالیت و هدف سایت
- * معرفی کامل وب سایت به موتورهای جستجوگر مشهور Google , Yahoo , Msn و ...
- * انتخاب توضیحات متناسب با صفحات سایت
- * بررسی و نحوه تعیین استراتژی ساختار لینک ها
- * طراحی مجدد صفحات سایت با توجه به تنوع مطالب
- * افزایش اهمیت صفحات سایت
- * قرار دادن توضیحات به صورت متنی در قالب جزء و کل
- * ایندکس صفحات سایت
- * افزایش بازدیدکننده هدفمند بر اساس کلمات مرتبط با فعالیت سایت
- * مشاوره و ارائه راه کارهای مناسب با توجه به فرایند انجام کار به صورت ماه به ماه



پورتال مقالات کامپیوتر و فناوری اطلاعات

مفاهیم و اصطلاحات دنیای جستجو و موتورهای جستجوگر

قبل از شروع گفتگو درباره هر موضوعی نیاز به آن است که مفاهیم اولیه و اصطلاحات رایج در آن موضوع بیان شود تا طرفین گفتگو راحت تر به منظور یکدیگر پی ببرند. برخی از مفاهیم و اصطلاحات حوزه SEO در این مقاله شرح داده شده است.

Spider, Crawler, Robot

نرم افزاری است که کار جمع آوری اطلاعات از صفحات سایتهای مختلف را بر عهده دارد.

Directory

فهرست. نوعی از موتورهای جستجوگر که پایگاه داده آن توسط ویراستاران تکمیل می گردد. در آنها سایتهای در گروههایی موضوعی دسته بندی می شوند.

Keyword

به واژه های مهم (کلیدی) هر صفحه گفته می شود. اما غالباً منظور کلماتی است که دوست داریم با آنها رتبه های مناسبی کسب کنیم.

Keyword Density

چگالی کلمه. منظور تعداد دفعات تکرار واژه های کلیدی در مقایسه با سایر کلمات متن است.

Keyword Staffing

تکرار یک کلمه به دفعات و پشت سر هم به منظور بالا بردن چگالی کلمه. این کار تقلب محسوب می شود.

Tinny Text

نوشتن متن با اندازه های بسیار کوچک و ریز به گونه ای که کلمات بسیاری بدین ترتیب در یک خط قرار داده می شود و به سختی نیز در صفحه قابل رویت هستند. نوشتن مطالب به این صورت، تقلب محسوب است.

Invisible Text

متن نامرئی. منظور استفاده از متن های هم رنگ با پس زمینه صفحه است. متن هایی که از دید کاربران مخفی می ماند. به عنوان مثال اگر پس زمینه یک صفحه سیاه است، متن صفحه نیز با رنگ سیاه نوشته می شود تا دیده نشود. این نوع متن ها از مصادیق تقلب می باشند.

Spam

تقلب. به تمام تلاش هایی گفته می شود که به کمک آن سعی می شود از راه های غیر معمول، رتبه های بالایی کسب شود. یا در اختیار گذاردن اطلاعاتی که موتورهای جستجوگر آن را دوست ندارند (اطلاعات ناخواسته) مانند تکرار یک کلمه به دفعات و پشت سر هم، استفاده از متن های هم رنگ زمینه و ...

ALT tag

محتوای این شناسه، متنی است که یک عکس را توضیح می دهد.

Deep Crawl

به معنای این است که موتور جستجوگر، می تواند صفحات زیادی از یک سایت را در پایگاه داده اش قرار دهد. موتور جستجوگر هر چه پایگاه داده اش بزرگتر باشد، صفحات بیشتری از یک سایت را می تواند در پایگاه داده اش قرار دهد. همه موتورهای جستجوگر دارای این ویژگی نمی باشند.

Robots.txt

با این فایل متنی و ساده، میزان دسترسی موتور جستجوگر به محتوای یک «سایت» را می توان کنترل کرد.



پورتال مقالات کامپیوتر و فناوری اطلاعات

META robots tag

به کمک این شناسه میزان دسترسی موتور جستجوگر به محتوای یک «صفحه» را می توان کنترل کرد.

Link

پیوند. در واقع پلی بین دو صفحه است. به کمک آن می توان از یک صفحه به صفحه دیگر رفت.

Link Popularity

مقصود این است که چه تعداد از سایت های دیگر به سایتی مشخص لینک کرده اند یا اینکه از چند سایت دیگر می توان به کمک پیوندها به سایتی مشخص رفت.

Link Reputation

اشاره به این دارد که سایر سایتها درباره سایتی که بدان لینک داده اند. چه می گویند. عموماً در این موارد عنوان. متن لینک و کلمات اطراف لینک در سایت مقصد. بررسی می شوند.

Learn Frequency

بعضی از موتورهای جستجوگر می توانند تشخیص دهند که محتوای صفحات پس از چه مدتی تغییر می کند (به روز می گردد) و بعد از آن مدت به آن صفحات مراجعه می کنند.

URL-Uniform Resource Locator

به آدرس منحصر به فرد هر منبع موجود در اینترنت گفته می شود. این منبع می تواند یک صفحه وب. یک فایل متنی و... باشد

Stop Word

به کلماتی گفته می شود که در کل اینترنت از آنها بسیار استفاده شده است. کلماتی نظیر the, a, an, web www, home page و ...

Meta tags

به کمک این شناسه ها. اطلاعاتی از صفحه در اختیار بینندگان (موتور جستجوگر. مرورگرها و ...) قرار داده می شود.

META Keywords

به کمک آن. کلمات کلیدی صفحه در اختیار موتورهای جستجوگر قرار داده می شود.

META Description

به کمک آن. توضیحی مختصر از صفحه در اختیار موتورهای جستجوگر قرار داده می شود.



پورتال مقالات کامپیوتر و فناوری اطلاعات

Stemming

به معنای این است که موتور جستجوگر می تواند صورت های مختلف یک کلمه را جستجو کند. به عنوان مثال با جستجوی swim موتور جستجوگر به دنبال swimming , swimmer نیز می گردد. همه موتورهای جستجوگر دارای این ویژگی نمی باشند.

Rank

رتبه یک صفحه در نتایج جستجو است زمانی که جستجویی مرتبط با محتوای آن صفحه انجام می شود.

Spamdexing

مختصر شده spam indexing است. منظور طراحی و معرفی صفحاتی به موتورهای جستجوگر است که کیفیت نتایج جستجو را پایین می آورند. موتورهای جستجوگر تمایل دارند که کاربران بارها و بارها به آنها مراجعه کنند و کیفیت بالای نتایج می تواند این روند را تضمین کند. لذا آنها هر کدام به نوعی سعی در تشخیص صفحاتی دارند که کیفیت نتایج جستجو را پایین می آورد. برخی از این موارد عبارتند از: ساختن صفحاتی که همگی دارای محتوای یکسانی اند. تکرار یک کلمه بیش از حد و ...

Comment

توضیحاتی است که طراحان سایت در لا به لای کدهای HTML می گنجانند تا برای فهمیدن وظیفه بخش های متفاوت کدهای HTML در مراجعات آتی نیازی به صرف وقت بسیار نداشته باشند.

منبع : <http://www.iranseo.com>