

موتور جستجوگر چگونه کار می کند؟

وقتی جستجویی در یک موتور جستجوگر انجام و نتایج جستجو ارایه 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 امکان دارد گوگل در کمتر از يك ثانیه تمام سایت های وب را بگردد و میلیون ها صفحه را در نتایج جستجوی خود ارایه کند؟ نه گوگل و نه هیچ موتور جستجوگر دیگری توانایی انجام این کار را ندارند. همه آنها در زمان پاسخ گویی به 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 و آنرا در پایگاه داده اش ذخیره می نماید و به هنگام جستجوی کاربر تنها در همین پایگاه داده می گردد .

بخش های مجزای يك موتور جستجوگر عبارتند از:

- Spider یا عنكبوت
- Crawler یا خزنده
- Indexer یا پایگانی کننده
- Database یا پایگاه داده
- Ranker یا سیستم رتبه بندی

الف) Spider- عنكبوت

اسپایدر یا روبات (Robot)، نرم افزاری است که کار جمع آوری اطلاعات مورد نیاز يك موتور جستجوگر را بر عهده دارد. اسپایدر به صفحات مختلف سر می زند، محتوای آنها را می خواند، اطلاعات مورد نیاز را جمع آوری می کند و آنرا در اختیار سایر بخش های موتور جستجوگر قرار می دهد .

کار يك 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.

اما یک اسپایدر آنرا چگونه می بیند؟
برای این که شما هم بتوانید دنیای وب را از دیدگاه يك اسپایدر ببینید، کافی است که کدهای HTML صفحات را مشاهده کنید.

آیا این دنیای متنی برای شما جذاب است؟

اسپایدر، به هنگام مشاهده صفحات، از خود 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 آنرا داشته باشید، می توانید مشخص کنید که اسپایدر کدام یک از موتورهای جستجوگر صفحات سایت را مورد بازدید قرار داده اند.
یکی از فعالیتهای اصلی که در SEM انجام می شود تحلیل آمار همین دید و بازدیدها می باشد.

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 پست الکترونیکی (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 در پایگاه داده اش ذخیره کند و 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 از صفحات سایت تان دور کنید و اجازه دسترسی به محتوای آن صفحات را به آنها ندهید. تنظیم میزان دسترسی موتورهای جستجوگر به محتوای یک سایت توسط پروتکل Robots انجام می شود که در مقالات دیگر سایت به آن پرداخته شده است. به عمل کراولر، خزش (Crawling) می گویند.

ج) Indexer - بایگانی کننده

تمام اطلاعات جمع آورش 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 کدام است، کلمات چندبار تکرار شده است، 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 صرف نظر می کند. کلماتی نظیر a، an، the، www، و ... 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.

ه) 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 پاسخگویی 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 مرتبیطی 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 دهد. 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
- پایگانی
- سیستم رتبه بندی

تذکر- برای سهولت در بیان مطالب بعدی هر گاه صحبت از پایگانی کردن (شدن) به میان می آید، مقصود این است که صفحه تجزیه و تحلیل شده و به پایگاه داده موتور جستجوگر وارد می شود.

برای آنکه تصور درستی از نحوه کار يك 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 ابیورد، 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 اوربال (نوعی آهو) را از شکاربانی منطقه دریافت می کند .

-کار اسپایدر

او اوربالی 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 همسر او بسته به ذائقه مهمانان برای آنها غذا طبخ می کند. ممکن است عده ای کله پاچه، عده ای آبگوشت، عده ای جگر و ... 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 اورپال 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.
- 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.
- به عنوان آخرین نکته این مقاله یاد آوری می کنم که به شکار اورپالی رعنا آن هم در منطقه حفاظت شده ابیورد، اصلاً فکر نکنید. اما توصیه می شود که حتماً از طبیعت بکر آن دیدن فرمایید (بدون اسلحه.!).