

جستوجوی اطلاعات و ساز و کارهای بهینه سازی آن

علی گزنسی^۱

تاریخ دریافت: ۷۹/۳/۱۷

چکیده: تحقیق در زمینه بازیابی اطلاعات راهگشای طراحی، کاربرد، و استفاده مؤثر از ابزارهای کنترل اطلاعات بوده، و کنترل بر روی حجم گسترده‌ای از اطلاعات، محبوطه‌های اینترنت و بانک‌های صوتی - تصویری را به همراه دارد. مقاله حاضر به صورت مختصر به تشریح یک نظام بازیابی اطلاعات (نرم افزار بازیابی اطلاعات) پرداخته و سپس بر روی فرمول جستجو در این نظام‌ها متمرکز می‌شود، و آنگاه به بررسی امکانات و قابلیت‌های موجود در یک نظام بازیابی اطلاعات برای شکل‌گیری فرمول جستجو می‌پردازد. امکانات و قابلیت‌های مورد بحث در این مقاله شامل عملگرهای منطقی، عملگرهای مجاورت، جستجوی املائی، کوئنه‌سازی، تعیین محدوده در داده‌های کمی، گسترش دامنه مفاهیم و اصطلاحات، ضرب و زنی، سیاهه واژه‌ها، مجموعه‌های بولی نمره‌گذاری شده، و آنالیز: گری مشابهت‌هاست.

کلید واژه‌ها: بازیابی اطلاعات، نرم افزار، اندازه‌گیری، سیاهه واژگان

مقدمه

مسئله بازیابی اطلاعات از زمانی آغاز شد که بشر سعی کرد محیط پیرامون خود را کنترل یا حداقل از فشارهای خارجی که باعث نابودی او می‌شدند جلوگیری کند. بشر برای ایجاد محیطی مطلوب جهت ادامه زندگی نیاز به اتخاذ تصمیمات سریع، صحیح، و دقیق داشت. کیفیت این تصمیمات به توانایی تصمیم‌گیرنده در حل مسائل وابسته بود ولی قبل از آن به میزان ارتباط و کیفیت اطلاعاتی وابسته بود که تصمیم‌گیرنده برای حل مشکل فراهم می‌آورد. به تدریج و در طول تاریخ، جمع‌آوری، سازماندهی و نگهداری اطلاعات امری متداول و مرسوم شد.

جنگ جهانی دوم را می‌توانیم نقطه تحرکی برای نظام‌های بازیابی اطلاعات^۱ در نظر بگیریم، زیرا حجم انتشارات خصوصاً در زمینه‌های علمی، به صورت تصاعدی در حال افزایش بود. دامنه تحقیقات و پژوهش‌ها روزبه روز گسترده‌تر می‌گردید و این افزایش تدریجی، کنترل اطلاعات با استفاده از ابزارها و شیوه‌های قبلی را غیرممکن می‌ساخت.

واژه بازیابی اطلاعات در سال ۱۹۵۱ برای اولین بار توسط تئوریسینی بنام موئرز^۲ به کار رفت. وی معتقد بود که یک نظام بازیابی اطلاعات باید امکان نمایه سازی و جستجوی اطلاعات موجود در کتابخانه را فراهم آورد، او با استفاده از روش‌های نظری و بررسی ساختار تلگراف توانست برای کدگذاری اطلاعات از راهی مؤثر و کم خطأ استفاده کند. اولین کتاب نظام‌های بازیابی اطلاعات در سال ۱۹۶۱ توسط ویکری^۳ ارائه شد. در این کتاب نحوه ایجاد نمایه‌ها و توصیف موضوعی مدارک به صورت گسترده‌تری مورد بررسی قرار گرفت. انگیزه اصلی ویکری از تألیف این کتاب، رشد سریع انتشارات و بررسی روش‌های نمایه‌سازی در کتابداری بود.

(لنکستر^۴، ۱۹۷۹) بازیابی اطلاعات را عبارت از فرآیند جستجو در میان مجموعه‌ای از مدارک با هدف تعیین آن دسته از مدارک که در حیطه موضوعی درخواست شده، باشند می‌داند. (سات‌کلیف^۵، ۱۹۹۲) مشخصه‌های اصلی نظام‌های بازیابی اطلاعات را به شرح زیر برمی‌شمارد:

۱. مجموعه‌ای از مدارک
۲. نوع سازماندهی و ساختار اطلاعات (شكل منطقی و فیزیکی اطلاعات ذخیره شده)
۳. روش بازیابی
۴. فرآیند بازیابی
۵. عملکرد بازیابی

به طور کلی نظام‌های بازیابی اطلاعات را می‌توانیم این‌گونه تعریف کنیم: نظام‌هایی که به منظور پردازش و بازیابی داده‌های ساختار نیافته طراحی شده و به لحاظ نوع سازماندهی اطلاعات، ساختار پایگاه‌ها، روش، بازیابی، گروههای خدمت گیرنده و فرآیندی که در طی آن به درخواست‌های اطلاعاتی کاربران پاسخ می‌دهند، از دیگر نظام‌های موجود متمایز هستند.

مجموعه داده‌ها

ارائه یک شمای کلی از نظام‌های بازیابی اطلاعات به صورت‌های مختلف امکان‌پذیر است.

1. Information Retrieval System

2. Mooers

3. Vickery

4. Lancaster

5. Sutcliffe

یک راه آن است که این نظام‌ها را به عنوان مجموعه‌ای از رکوردها و فیلدهای اطلاعاتی در نظر بگیریم، بدین ترتیب که هر بانک اطلاعاتی دارای ساختار خاصی، متشكل از فیلدهای مختلف است و هر بانک اطلاعاتی می‌تواند بین یک تا N رکورد اطلاعاتی را در خود جای دهد است. گردآوری و سازماندهی اطلاعات هر بانک اطلاعاتی براساس خط مشی تولید کننده بانک مربوط صورت می‌گیرد، بنابراین هر بانک اطلاعاتی دارای مشخصاتی مخصوص به خود است (به عنوان نمونه حوزه یا حوزه‌های موضوعی هر بانک). عناصر اطلاعاتی بانک‌های اطلاعاتی متفاوت بوده و شکل استانداری ندارد (مهراد و مقدسی، ۱۳۸۰). اطلاعات موجود در هر بانک مورد تجزیه و تحلیل قرار گرفته و به کوچکترین واحدهای اطلاعاتی معنادار تجزیه می‌شوند، آنگاه روابط لازم بین این واحدهای اطلاعاتی برقرار می‌شود (از جمله مشخص کردن روابط اعم و اخص، گروه‌بندی واحدهای اطلاعاتی و) و در نهایت واحدهای اطلاعاتی بدست آمده به همراه روابط میان آنها در یک فایل جدید ذخیره و نگهداری می‌شوند.

چرخه دستورات

مطالعات انجام شده در زمینه جست و جوی اطلاعات نشان دادند که، جست و جوگران بانک‌های متنی و کتابشناختی از یک چرخه دستورات به صورت مشترک استفاده می‌کنند. در این مطالعات تأکید بر روی دستورات پراستفاده‌تر بود تا اینکه بخواهند سلسه مراتب استفاده از این دستورات را مورد توجه قرار دهند. این مطالعات نشان دادند که بر مبنای دستورات انتخابی کاربران می‌توان پیش‌بینی کرد که جست و جوگران در چه مرحله‌ای از جست و جو قرار دارند و هدف آنها چیست؟ (پنای من^۱؛ چاپمن^۲، ۱۹۷۸؛ ۱۹۸۱).

هر دستور ممکن است با شاخص‌های مختلفی فراخوانی شود، برای مثال ترکیب فرمول جست و جو با روشنی که برای بازبینی اطلاعات انتخاب می‌شود، از یک کاربر تا یک کاربر دیگر متفاوت است. در یک بررسی دیگر که توسط (چاپمن، ۱۹۸۱) صورت گرفت، وی توانست با گروه‌بندی دستورات، ترتیب استفاده از آنها و شاخص‌های استفاده شده به دسته‌بندی کاربران اطلاعاتی بپردازند. در ادامه به بررسی دستورات طبقه‌بندی شده و مورد استفاده در این چرخه می‌پردازیم:

۱. انتخاب یک بانک اطلاعاتی برای انجام جست و جو
۲. جست و جو برای واژه‌های مورد نظر در بانک واژگان
۳. ایجاد فرمول جست و جو و انجام جست و جو

۴. نمایش و بازبینی رکوردها
 ۵. سفارش مدارک
 ۶. درخواست برای اطلاعاتی درباره نظام بازیابی اطلاعات
 ۷. برقراری شاخص‌های نمایش و ارتباطی
- ارائه این الگوها و تکرار آنها خصوصاً برای اهدافی همانند طراحی نظام‌های بازیابی اطلاعات و تعلیم جست‌وجوگران بسیار حائز اهمیت است و براساس اطلاعات باز خوردهای حاصل از تکرار این الگوها می‌توان به طراحی نظام‌های بازیابی اطلاعات بهینه پرداخت.

انتخاب یک بانک اطلاعاتی برای انجام جست‌وجو
 بعد از انتخاب بانک اطلاعاتی توسط کاربر، خلاصه‌ای از اطلاعات، همانند محدوده تاریخی رکوردها، تعداد رکوردها، قیمت و... را در اختیار کاربر قرار می‌دهد.

جست‌وجو برای واژه‌های مورد نظر در بانک واژگان

در نظام‌های بازیابی اطلاعات بعد از سازماندهی و تجزیه و تحلیل اطلاعات براساس داده‌های موجود در بانک‌های اطلاعاتی نظام، تعدادی فایل کمکی به منظور کمک به عملیات بازیابی اطلاعات ایجاد می‌شود. این فایل‌ها را می‌توان به عنوان سیاهه‌ای از کل واژه‌های موجود در بانک اطلاعاتی دانست که با نظم الفبایی در کنار هم قرار گرفته‌اند. این فایل‌ها از یک نظام به نظام دیگر می‌تواند ساختار بسیار ساده یا پیچیده‌ای داشته باشد. در شکل ساده، این فایل‌ها شامل ریشه واژه‌ها، پسوندها، پیشوندها، و تعداد تکرار آنها و در شکل پیچیده علاوه بر شکل ساده حاوی ارجاعات، سلسله مراتب‌ها، عبارات موجود در بانک اطلاعاتی و غیره است.

این بخش می‌تواند از اهمیت زیادی برخوردار باشد زیرا اطلاعاتی که بعداً توسط کاربر مورد بازیابی قرار می‌گیرد، بالقوه وابسته به اطلاعات صحیحی است که در این قسمت انتخاب می‌شود. حداقل استفاده‌ای که از بانک واژگان می‌شود این است که کاربر خواهد فهمید که آیا واژه‌های مورد نظر او در بانک اطلاعاتی وجود دارند؟ شکل صحیح آنها به چه صورت است؟ واژه در چند رکورد، در کدام فیلدیها و به چه میزان تکرار شده است؟ اگر واژه درخواستی کاربر بد تایپ شده باشد یا در بانک واژگان وجود نداشته باشد، یک مقدار صفر برای آن نمایش داده می‌شود. کاربر با استفاده از این فایل‌ها می‌تواند به واژه‌های مرتبط با واژه‌های مورد نظر خود دست یابد و به اخص‌تر یا اعم‌تر کردن فرمول جست‌وجو بپردازد. یک نمونه بسیار ساده از این بانک‌ها در جدول ۱ قابل مشاهده می‌باشد.

جدول ۱. بانک واژگان براساس موجودی نظام بازیابی اطلاعات

تعداد تکرار	واژه
۱۲۳۴	کتاب
۹۹۹	کتاب
۶۵۷	کتابخوانی
۳۰۰۱	کتابخانه
۱۰۰	کتابخانه‌ها
۲۵۳۷	کتابخانه‌های
۱۲۰۰	کتابدار
۲۰۰	کتابداری
۱۲۹	کتابشناختی
۱۹۰۰	کتابشناسی

فرمول جست و جو

کاربر به وسیله یک عبارت یا مجموعه‌ای از عبارات، درخواست اطلاعاتی خود را به نظام اعلام می‌کند که به آن فرمول جست و جو یا به بیان دیگر پرسش گفته می‌شود. پرسش‌ها بر حسب قابلیت نظام‌های بازیابی اطلاعات می‌توانند دارای گزینه‌هایی برای تعیین شکل، قالب خروجی و نتایج بازیابی نیز باشند. ارائه پرسش‌ها به نظام‌های بازیابی اطلاعات عمده‌تاً به دو شیوه صورت می‌گیرد: زبان طبیعی و زبان رایانه. زبان طبیعی علاوه بر آنکه نیاز به یادگیری ندارد به لحاظ ماهیت از قابلیت بهتری برای بیان درخواست اطلاعاتی برخوردار است، اما باید گفت که تجزیه و تحلیل زبان طبیعی برای نظام مشکل تر بوده و نیاز به توانایی منطقی بالائی دارد. از سوی دیگر درخواست اطلاعات به زبان رایانه نیاز به یادگیری، کسب مهارت و تجربه دارد، اما تجزیه و تحلیل اطلاعاتی که با این شیوه طرح می‌شوند برای نظام بسیار آسان است.

ترکیب فرمول جست و جو با استفاده از عملگرهای موجود

بعد از آنکه واژه‌های مورد نظر برای بیان درخواست اطلاعات تعیین شد، با استفاده از امکانات، ابزار، و روشی که نظام در اختیار ما قرار می‌دهد، به گونه‌ای این واژه‌ها را در یک فرمول جست و جو با هم ترکیب می‌کنیم تا نتیجه دلخواه از جست و جو حاصل گردد. در ادامه به توضیح این قابلیت‌ها، ابزارها و روش‌ها می‌پردازیم که شامل عملگرهای منطق بولی، عملگرهای مجاورت، جست و جوی املائی، کوتاه‌سازی، تعیین محدوده در داده‌های کمی، گسترش دامنه مفاهیم و اصطلاحات، ضریب وزنی، سیاهه واژه‌ها، مجموعه بولی نمره‌گذاری شده، و اندازه‌گیری مشابهت‌هاست.

الف. عملگرهای بولی^۱

با استفاده از این عملگرها (AND, OR, NOT, XOR) کاربر می‌تواند میان تصورات و واژه‌هایی که برای بیان درخواست اطلاعاتی خود برگزیده است ارتباطی منطقی برقرار کند (داورپناه، ۱۳۷۶). عملگر منطقی نمادی است که برای مشخص کردن رابطه منطقی بین دو مقدار یا مفهوم مورد استفاده قرار می‌گیرد. با استفاده از عملگر منطقی OR می‌توان محدوده رکوردهای مورد نظر را افزایش و با استفاده از عملگر منطقی AND, NOT نیز محدوده آن را کاهش داد. در جدول ۲ این عملگرها مشاهده می‌شود.

جدول ۲. عملگرهای منطق بولی در نظام‌های بازیابی اطلاعات

عملگر	عملکرد	مثال	تشریح مثال
AND	در صورتی که به عنوان عملگر در کنار یک واحد اطلاعاتی ظاهر گردد، یک مجموعه فقط در صورتی که حاوی آن اطلاعات باشد در نتایج بازیابی ظاهر خواهد شد (اشتراک).	موضوع = "رایانه" و تاریخ = ۱۳۷۸	رکوردهای اطلاعاتی مورد بازیابی فرار می‌گیرد که فیلد موضع آنها "رایانه" و همین‌طور فیلد تاریخ آنها برابر ۱۳۷۸ باشد.
OR	در صورتی که به همراه تعدادی از واحدهای اطلاعاتی به کار برود، در صورت وجود این واحدها در رکوردهای اطلاعاتی، آن رکوردهای مورد بازیابی فرار خواهد گرفت (احتماع).	موضوع = "کتابخانه" یا "اطلاع‌رسانی" باشد.	رکوردهای اطلاعاتی مورد بازیابی فرار می‌گیرد که فیلد موضع آنها "کتابخانه" یا "اطلاع‌رسانی" باشد.
NOT	در صورتی که به همراه یک تعداد واحد اطلاعاتی به کار برود، هر یک از واحدها در رکوردهای اطلاعاتی مشاهده شوند رکوردهای مربوط مورد بازیابی فرار نخواهد گرفت.	موضوع = "اطلاع‌رسانی" و تاریخ نامساوی ۱۳۷۸	رکوردهای اطلاعاتی که موضع آنها "اطلاع‌رسانی" باشد و تاریخ آنها برابر هر سالی غیر از سال ۱۳۷۸ باشد.
XOR	به دلیل درک دشوار آن توسط کاربران کمتر مورد استفاده فرار می‌گیرد ولی اگر به همراه دو واحد اطلاعاتی به کار برود به این مفهوم می‌باشد که هر یک از دو واحد اطلاعاتی که به صورت منفرد مورد مشاهده فرار گیرند مورد بازیابی فرار می‌گیرد ولی اگر با هم در یک مجموعه مشاهده شوند آن مجموعه مورد بازیابی فرار ننمی‌گیرد.	موضوع = "رایانه" یا "اطلاع‌رسانی" باشد ولی نباید با هم در یک رکورد تکرار شده باشند.	رکوردهای اطلاعاتی که موضع آنها "رایانه" یا "اطلاع‌رسانی" باشد ولی نباید با هم در یک رکورد تکرار شده باشند.

ب. عملگرهای مجاورت^۲

واحدهای اطلاعاتی در موقعیت‌های مکانی مشخصی از رکوردهای اطلاعاتی قرار دارند. فوائل این واحدهای اطلاعاتی از یکدیگر می‌تواند نشان دهنده ارتباط یا عدم ارتباط معنائی

این واحدها با یکدیگر باشد. برای تشخیص این ارتباط معنایی و بالا بردن ضریب دقت در بازیابی اطلاعات می‌توان از عملگرهای مجاورت استفاده کرد. به بیان دیگر با ترکیب عملگرهای منطق بولی و عملگرهای مجاورت می‌توانیم رکوردهای اطلاعاتی را مورد بازیابی قرار دهیم که حاوی واحدهای اطلاعاتی در فواصل مکانی معینی نسبت به هم باشند. به عنوان نمونه اگر دو واحد اطلاعاتی "کتابخانه" و "ساختمان" در یک بند اطلاعاتی^۱ در فواصل نزدیک به هم قرار بگیرند احتمال اینکه در رابطه با "ساختمان کتابخانه" بحث شده باشد خیلی بالاست. شکل اصلی فرمول مجاورت به این صورت است:

واحد اطلاعاتی اول در M فاصله بند اطلاعاتی از واحد اطلاعاتی دوم قرار گرفته باشد.

مثال ۱: "ساختمان" در یک پاراگراف با "کتابخانه" باشند (ساختمان و کتابخانه باید در یک پاراگراف قرار داشته باشند)

مثال ۲: "ساختمان" در یک جمله با "کتابخانه" باشند (ساختمان و کتابخانه باید در یک جمله قرار داشته باشند)

یک نمونه دیگر از عملگرهای مجاورت، "ADJ" نامیده می‌شود و در صورتی که بین دو یا چند واحد اطلاعاتی به کار رود به شرطی که این واحدها با رعایت تقدم و بلا فاصله بعد از هم قرار بگیرند، مورد بازیابی قرار خواهد گرفت، به عنوان مثال فرمول "استاندارد" ADJ "کتابخانه" رکوردهای اطلاعاتی را مورد بازیابی قرار خواهد داد که عبارت "استاندارد کتابخانه" را در خود جای داد باشند نه "کتابخانه استاندارد". نمونه‌ای از عملگرهای مجاورت در جدول ۳ مشاهده می‌شود.

جدول ۳. مهم‌ترین نمونه‌های عملگرهای مجاورت

عملگر	توصیف
W	با رعایت تقدم بین دو واژه تعداد صفر تا N واژه قرار خواهد گرفت.
N	بدون رعایت تقدم بین دو واژه تعداد صفر تا N واژه قرار خواهد گرفت.
P	دو واژه با هم در یک پاراگراف تکرار شده باشند.
S	دو واژه با هم در یک جمله یا یک فبلد فرعی تکرار شده باشند.
F	دو واژه با هم در یک فبلد فرعی داشته باشد، که می‌توان نام فبلد را نیز به فرمول جست و جو اضافه کرد.

ج. جست و جوی املائی (هجایی)^۲

اشتباهات تایپی حاصل از غلط‌های املائی واردکنندگان اطلاعات، کاربران نظام و همچنین

۱. یک بند اطلاعاتی می‌تواند از یک واژه تا چند واژه، یک جمله، یک پاراگراف، یک فبلد یا یک فبلد فرعی تشکیل شده باشد.

2. Spelling Search

استفاده از روش‌های جدید برای درون داد اطلاعات همانند استفاده از پویشگر^۱ و نرم‌افزار تشخیص نوری حروف، سبب شده تا از شیوهٔ جست‌وجوی املائی برای کاستن تأثیر این خطاهای در بازیابی اطلاعات استفاده شود^۲.

این روش هم در جست‌وجوی واژه‌ها و هم در تجزیه و تحلیل ابتدایی رکوردهای اطلاعاتی به منظور ایجاد فایل‌های کمکی در نظام کاربرد دارد. به عنوان مثال کاربر واژه COMPUTER را برای جست‌جو وارد می‌کند و آنگاه واژه‌های COMPUTER, COMMPITER, COMPUTTER, COMPUTER نیز مورد جست‌جو قرار خواهد گرفت. استفاده از این روش باعث افزایش ضریب بازیافت و کاهش ضریب دقت می‌شود.

د. کوتاه‌سازی^۳

ریشه هر واژه براساس پیشوند یا پسوندی که به آن اضافه می‌شود، اشکال مختلفی به خود می‌گیرد. بنابراین ریشه با ملحقاتی که به آن اضافه می‌شود به شکل‌های متفاوتی نوشته می‌شود. از سوی دیگر بعضی واحدهای اطلاعاتی با وجود یکسان بودن به صورت‌های مختلفی نوشته می‌شوند (همانند اسم پدید آورندگان). اهمیت کوتاه‌سازی بخصوص در نظام‌هایی که ریشه‌یابی واژه‌ها (در مراحله تجزیه و تحلیل داده‌ها) و یکدست سازی اطلاعات بخوبی صورت نمی‌گیرد بسیار پراهمیت است.

کوتاه‌سازی این امکان را فراهم می‌آورد که هنگام جست‌جو بدون در نظر گرفتن حروف یا حرفهایی از یک واژه یا یک واحد اطلاعاتی به بازیابی پردازم. به این دلیل که هنگام تجزیه و تحلیل داده‌ها بعضی حروف و کاراکترهای خاص حذف می‌شوند. برای نمایش شکل کوتاه‌سازی در دستورات از این کاراکترها استفاده می‌شود (به عنوان نمونه #,*). دو الگوی رایج کوتاه‌سازی در جدول ۴ قابل مشاهده می‌باشند. کوتاه‌سازی باعث افزایش ضریب بازیافت و کاهش ضریب دقت می‌شود.

1. Scanner

۲. یکی از شیوه‌های درون داد اطلاعات در مراکز کتابخانه‌ای و اطلاع‌رسانی استفاده از اسکنر و نرم‌افزار تشخیص نوری حروف است. در این شیوه یک صفحه چاپی بدون نیاز به تایپ به نظام وارد می‌شود به نحوی که قابل ویرایش از نظر واژه‌پردازی باشد. در یک صفحه چاپی باکیفیت و واضح خوب می‌توان بین ۹۰ تا ۹۹ درصد بر روی صحت داده درون داد شده از جهت واژه‌پردازی حساب کرد.

3. Truncation

جدول ۴. نمونه هایی از عملگرهای کوتاه سازی

نحوه	نتایج حاصل از اجرای نمونه	شرح
کتابخانه ای کتابخانه ها کتابخانه	کتابخانه ??	کوتاه سازی برابر یک یا چند حرف با تعداد مشخص در مکان های مشخصی از واژه یا واحد اطلاعاتی
کتابخانه کتابخوانی کتابدار کتابشناسی کتابشناسختی	*	کوتاه سازی برای یک مجموعه از حروف با تعداد نامشخص از واژه با واحد اطلاعاتی
LABOUR LABOR	LABO#R	کوتاه سازی به منظور نادیده گرفتن وجود یا عدم وجود یک یا چند حرف در واژه با واحد اطلاعاتی

ه. تعیین محدوده در داده های کمی

داده های کمی می توانند بخشی از اطلاعات موجود در بانک های اطلاعاتی را تشکیل دهند. از نمونه های بارز این گونه داده ها می توانیم از مقادیر عددی و تاریخی نام ببریم. بنابر اهمیت این نوع داده ها نیاز خواهد بود که علاوه بر یک مورد خاص، در محدوده های خاصی نیز جست و جو صورت گیرد. عملگرهای مورد استفاده در تعیین محدوده ها عبارتند از: کوچکترین مساوی (=<)، بزرگتر مساوی (=>)، بزرگتر (=)، مساوی (=)، و نامساوی (=!). استفاده از این عملگرها باعث افزایش ضریب دقت و کاهش ضریب بازیافت می شوند.

و. گسترش مفاهیم و اصطلاحات

نظام های بازیابی اطلاعات با استفاده از ابزارها و شیوه های مختلف به گسترش مفاهیم موجود در درخواست اطلاعاتی کاربر می پردازند و برای این منظور با استفاده از اصطلاحات نامه ها و واژه نامه های موجود یا ایجاد آنها بر پایه اطلاعات موجود در نظام می پردازنند. ایجاد این اصطلاحات نامه ها و واژه نامه ها بر اساس داده های موجود در بانک اطلاعاتی مبحث گسترده ای است که در مقاله حاضر به بحث درباره آن نمی پردازیم.

ز. کاربرد ضریب وزنی ^۱

اگر بخواهیم هر عضو فرمول جست و جو در بازیابی از دقت بالاتری برخوردار باشد، می توانیم از روش ضریب وزنی استفاده کیم. منظور از ضریب وزنی این است که کاربر به هر

عضو دلخواه از فرمول جستجو یک عدد صحیح که مشخص کننده اهمیت عضو در میان اعضا است اختصاص می‌دهد. برای مثال به پرسش زیر توجه کنید:

(TENNIS(.8)OR GOLF (.4)AND CHAMPION(.6)

در چنین حالتی، نظام بازیابی اطلاعات بدون در نظر گرفتن ضرایب وزنی تعیین شده به جستجوی رکوردها می‌پردازد. سپس با استفاده از ضرایب وزنی تعیین شده توسط کاربر به نمره‌گذاری رکوردها می‌پردازد. جدول ۵. بیانگر موضوع است.

جدول ۵. محاسبه ضرایب وزنی

نمره‌ای که به رکورد اختصاص داده می‌شود	محتوای رکورد
.8+.6=1.4	اگر یک رکورد حاوی CHAMPION و TENNIS باشد
.4+.6=1.0	اگر یک رکورد حاوی GOLF و CHAMPION باشد
.8+.4+.6=1.8	اگر یک رکورد حاوی هر سه واژه باشد

به ازای هر بار تکرار هر یک از واژه‌های عضو بیش از یک بار در رکوردهای اطلاعاتی، ضریب وزنی واژه مربوط، به کل نمره رکورد بازیابی شده اضافه می‌گردد. بعد از انجام جستجو و انجام نمره‌گذاری توسط نظام، رکوردها بر حسب بالاترین نمرات مرتب می‌شوند و کاربر می‌تواند تا N رکورد بازیابی شده را که دارای نمره بالاتری هستند، انتخاب کند. با وجود مزایای این روش کاربران باید یاد بگیرند که چگونه باید ضرایب وزنی را به عضوهای فرمول اختصاص دهند و در کنار آن باید متنزکر شد که نتایج بدون ضریب وزنی از ارزش بالایی برخوردار نخواهند بود.

ح. سیاهه واژه‌ها

برای آنکه یک روش آسان‌تر نسبت به روش تعیین ضرایب وزنی برای کاربران تدارک بینیم می‌توانیم از روش فرمول نویسی "سیاهه واژه‌ها" استفاده کنیم. این شیوه همان‌کار ضرایب وزنی را با یک روش دیگر انجام می‌دهد. به فرمول زیر توجه کنید:

TENNIS,GOLF,CHAMPION,CUT,TOURNAMENT,WINNER

در سیاهه فوق بین هر جفت واژه عملگر OR فرض می‌شود و هر کدام از آنها دارای ضرایب وزنی یکسان هستند. معادل این فرمول با استفاده از عملگرهای بولی به صورت زیر خواهد بود:

TENNIS.OR.GOLF.OR.CHAMPION.OR.CUT.OR.TOURNAMENT.OR.WINNER

در نهایت، رکوردها براساس بالاترین تعداد واژه‌های فوق که در رکوردها ظاهر می‌شوند و همچنین براساس تعداد تکرار آنها در رکوردها نمره‌گذاری می‌شوند و براساس همین نمره‌گذاری

مرتب می‌شوند. در این شیوه از عملگر AND نمی‌توان استفاده کرد و به همین دلیل ضریب دقت پایین است. در این روش می‌توان برای هر واژه یک ضریب وزنی نیز مشخص کرد.

ت. مجموعه‌های بولی نمره‌گذاری شده

همانند دو روش قبل در این روش نیز براساس تعداد تکرار اعضای فرمول جستجو در رکوردهای بازیابی شده، به هر کدام از رکوردها یک نمره اختصاص می‌یابد و در نهایت رکوردها بر حسب همین نمره در خروجی مرتب می‌شوند. این نمره نشان دهنده درصد ارتباط رکورد بازیابی شده حاضر با در خواست اطلاعاتی کاربر است. به عنوان مثال بعد از اجرای فرمول Relevance صورت زیر نمره‌گذاری شده و بازگردانده شد. این نمره‌گذاری در جلو عبارت Infoseek (Information And Retrieval And System) به صورت درصد بیان شده است.

نتایج جستجو و نمره‌گذاری شده سایت Infoseek

1. Food Safety Information from NC State

<<http://www.ces.ncsu.edu/depts/foodsci/agentinfo/>>

Food Safety resources for consumers, educators, researchers and absolutely anyone that wants to know more about food safety prevention and issues. Developed by The NC Cooperative Extension Service and...

Relevance: 100% Date: 11 May 1999, Size 9.0K,

<http://www.ces.ncsu.edu/depts/foodsci/agentinfo/> find similar pages

2. CIIR Demo Page <<http://ciir.cs.umass.edu/>>

Sorry, this site is geared for browsers that support frames.

Relevance: 98% Date: 28 Sep 1999, Size 1.9K, <http://ciir.cs.umass.edu/find> similar pages

3. California Highway conditions <http://www.amdahl.com/internet/general/travel/ca_hiqhway.html>

Access highway condition information by region of highway type.

Relevance: 85% Date: 8 Oct 1999, Size 25.6K,

http://www.amdahl.com/internet/general/travel/ca_hiqhway.html find similar pages

4. CAIRSS <<http://imr.utsa.edu/CAIRSS.html>>

Institute for Music Research ACIRSS for Music (Computer-Assisted Information Retrieval Service System) Search CAIRSS CAIRSS is a bibliographic database of music research literature in music education, music psychology, music...

Relevance: 81% **Date:** 28 Aug 1998, **Size:** 3.7K, <http://imr.utsa.edu/CAIRSS.html> Find similar pages.

۵. اندازه‌گیری مشابهت

دسترسی آسان کاربر به حداقل یک رکورد که بالاترین میزان ارتباط را با درخواست اطلاعاتی او داشته باشد بسیار حائز اهمیت است. اولین رکوردي که به این صورت توسط نظام مورد بازبایی قرار گیرد می‌تواند اساسی برای جستجوهای بعدی کاربر باشد. اگر کاربر چنین رکوردها را بازبایی، و به جهت مرتبط بودن با درخواست اطلاعاتیش آن را تائید کرد، آنگاه نظام به تجزیه و تحلیل رکورد می‌پردازد و با معیارهای ارزشیابی متن و دادن ضریب وزنی به هر یک از واژه‌ها اقدام به انجام یک جستجوی دیگر می‌کند و رکوردهای مرتبط را مورد بازبایی قرار می‌دهد. در این روش کاربر نیاز نخواهد داشت که به بازنویسی دوباره فرمول بپردازد.

انجام جستجو

پس از آنکه درخواست اطلاعاتی توسط کاربر به نظام وارد شد، نظام درخواست اطلاعاتی را مورد تجزیه و تحلیل قرار داده و آنگاه با انجام عملیات جستجو مجموعه‌هایی را که با آن منطبق هستند، به صورت خروجی به کاربر نشان می‌دهد که آن را نتیجه جستجو می‌نامند. منظور از تجزیه و تحلیل، ترجمه و تفسیر درخواست اطلاعاتی مطابق با زبان رایانه است. با بررسی نتیجه جستجو کاربر در می‌یابد که آیا جستجو رضایت‌بخش بوده است یا باید از فرمول جدیدی استفاده کند، و یا اینکه باید بعد از انجام تغییرات لازم دوباره فرمول جستجو به اجرا بگذارد. در مجموعه رکوردهای اطلاعاتی بازبایی شده رکوردهایی وجود دارند که با درخواست اطلاعاتی کاربر منطبق نیستند که آنها را ریزش کاذب اطلاعاتی می‌نامند. ریزش کاذب اطلاعات می‌تواند از عوامل مختلفی ناشو، شود، از جمله عدم سازماندهی بهینه اطلاعات، محدودیت‌های زبان نرم‌افزاری برای بیان پرسش، بازخورد ضعیف نظام در تعامل با کاربر، تشریح بد یاناقص درخواست اطلاعاتی.

نمایش و بازبینی رکوردها

بازخورد خوب نظام در این قسمت می‌تواند نقش مهمی در هدایت کاربر برای رسیدن به

اطلاعات مورد نظرش داشته باشد. در این قسمت کاربر باید بتواند به سوالات زیر پاسخ دهد:
چه کارهایی با نیاز اطلاعاتی کاربر مطابقت دارند؟

آیا اطلاعاتی وجود دارد که بتوان واژه‌هایی که حوزه‌های موضوعی وسیعی را پوشش داده‌اند تشخیص داد؟

آیا اطلاعاتی وجود دارد تا بتوان واژه‌هایی که می‌تواند موجب اخص‌تر شدن دامنه جست و جو شوند را تشخیص داد؟

آیا مجموعه رکوردهای بازیابی شده با نیازهای اطلاعاتی کاربر مطابقت دارد؟
آیا اطلاعاتی وجود دارد که از طریق آن کاربر بتواند به تصحیح منطق بولی فرمول جست و جو پردازد؟

سفارش مدارک

در بعضی از نظام‌ها فقط اطلاعات کتابشناختی نگهداری می‌شود، در این‌گونه نظام‌ها میان زمان جست‌وجو و دریافت اطلاعات مدت زمانی تلف خواهد شد. اما در نظام‌های تمام‌متن کاربر قادر خواهد بود، مستقیماً مدارک را مورد مطالعه قرار داده و در مورد انتخاب یا عدم انتخاب آنها تصمیم‌گیری کند، در این نظام‌ها کاربر مستقیماً مدارک را دریافت می‌کند و نیازی به سفارش آنها نیست.

درخواست برای اطلاعاتی درباره نظام بازیابی اطلاعات

بیشتر نظام‌های بازیابی اطلاعات از راهنمای پیوسته برخوردارند. و اگر کاربر در مورد دستورات یا دیگر ترکیبات نظام نیاز به راهنمایی داشت، اطلاعات لازم در اختیار او قرار می‌گیرد. نظام‌های بازیابی اطلاعات حرفه‌ای همیشه سعی دارند به سوالات کاربران نهایی پاسخ دهند، همانند (در اینجا چه کاری در حال انجام است؟ در اینجا چگونه باید عمل کنم؟) برای پاسخ دادن به بعضی از این سوالات نظام باید از اهداف کاربر آگاه باشد و یک چارچوب برای ارزیابی کاربر در اختیار داشته باشد. پاسخ دادن به بسیاری از سوالات دیگر به راحتی و با تفسیر دستورات یا گزینه‌های انتخاب شده توسط کاربر صورت می‌گیرد. کاربرانی که از تجربه کمتری برخوردار هستند نیاز به راهنمایی بیشتری در این زمینه خواهند داشت؟

برقراری شاخص‌های نمایشی و ارتباطی

شاید یک از علل مهم ناموفق بودن جست‌وجو این باشد که کاربران نمی‌توانند با خدمات

اطلاعاتی میزبان ارتباط برقرار کنند، زیرا کلمات عبور را اشتباه وارد می‌کنند یا شماره تلفن مربوط به شبکه ارتباطی را درست وارد نمی‌کنند. کاربران در این گونه موقع نمی‌دانند که چگونه باید از صحبت کار خود مطمئن شوند. در این گونه موقع هنوز ارتباط با نظام بازیابی اطلاعات برقرار نشده است که نظام بتواند از خود عکس العمل نشان دهد.

نتیجه‌گیری

در مجموعه چرخه دستورات باید چهار مقوله زیر را به عنوان مهم‌ترین اهداف خود مدنظر قرار دهیم:

- بهینه سازی انتخاب واژگان جست‌وجو، توسط کاربر؛
- بهینه سازی فرمول جست‌وجوی کاربر؛
- بهینه سازی و تنظیم تعداد رکوردهای بازیابی شده؛
- بهینه سازی ضریب دقت و بازیابی یا بهینه سازی کل نظام بازیابی اطلاعات

بعد از انجام این مراحل در صورت رضایت جست‌وجوگر، جست‌وجو به مرحله بعد منتقل خواهد شد در غیر اینصورت جست‌وجو دوباره تکرار خواهد شد. معنی تکرار انجام دوباره یک کار است. در جست‌وجو تکرار به معنی تکرار پردازش فرمول جست‌وجو و اصلاح آن تا رسیدن به نتیجه دلخواه است. هدف ممکن است رسیدن به یک اندازه دلخواه یا اطلاعاتی خاص یا رکوردهای خاصی باشد، که دارای بالاترین ارزش اطلاعاتی هستند. اطلاعات به دست آمده امکان تصمیم‌گیری در مورد این پرسش که آیا جست‌وجو باید دوباره انجام شود را در اختیار ما قرار می‌دهد، و این قضاوتی است که از باز خورد حاصل از نتایج جست‌وجو بدست می‌آید. در این مورد بعضی عوامل به صورت مرسوم و متعارف چنانچه قبلًا در چرخه دستورات به آن اشاره شد وجود دارند، ولی قالب‌های باز خورد به صورت کامل در جست‌وجو مشخص نشده‌اند.

بعضی از کاربران ترجیح می‌دهند که با توجه به محتوای رکوردهای بازیابی شده، درباره نحوه و چگونگی ایجاد و اجرای دوباره فرمول جست‌وجو تصمیم‌گیری کنند، نه براساس تعداد رکوردهای بازیابی شده، اما در موقعي که تعداد رکوردها از حد معینی بالاتر باشد، انجام این بررسی مشکل به نظر می‌رسد. حالت معمول این قضیه به این صورت است که اگر بعد از اجرای فرمول جست‌وجو تعداد رکوردهای زیادی مورد بازیابی قرار نگرفته، یک فرمول دیگر ایجاد و اجرا می‌شود، به عنوان مثال اگر نتیجه یک بازیابی اطلاعات ۳۵۰۰ رکورد باشد، حتی اگر رکوردهای مورد نظر کاربر در میان این رکوردها وجود داشته باشند، بررسی آنها بسیار مشکل است و باید فرمول جست‌وجو اصلاح شود تا تعداد رکوردهای کمتری مورد بازیابی قرار بگیرند.

از سوی دیگر فرض کنید برای ایجاد یک فرمول جست و جو تعداد رکوردهای بازیابی شده ۲۳ عدد باشد این اشتباه است اگر از این تعداد خوشحال شده و فکر کنیم جست و جوی خوبی انجام شده است یا برعکس فکر کنیم جست و جوی خوبی صورت نگرفته است، زیرا هر چند تعداد این رکوردها کم است ولی احتمالاً اطلاعات مفیدی از این تعداد رکورد بدست خواهد آمد. در این حالت تعداد رکوردها به آن اندازه کوچک است که کاربر می‌تواند آن را مورد بازبینی قرار داده و درباره فرمول جست و جوی بعدی تصمیم‌گیری کند. بازخورد به اصلاح فرمول جست و جو کمک می‌کند، و این خود باعث اجرای دوباره جست و جو و بازیابی تعداد زیاد یا کمتری از رکوردها می‌شود. در این هنگام مسئله قابلیت و توانایی‌های به تصویر کشیدن اطلاعات نمایان می‌شود. جست و جوگران با تجربه و متخصصان بازیابی باید به این نکته توجه داشته باشند که انجام بازیابی مستلزم صرف وقت و هزینه است و در این مسیر باید با اصلاح و بازنویسی و تکرار چرخه دستورات به نتایج دلخواه خود دست یابیم.

مأخذ

۱. داورینا، محمد رضا (۱۳۷۶). "قواعد استنتاج و فرمولبندی پرسشن در جست و جوی کامپیوتري". پیام کتابخانه، دیپرخانه هیأت امنای کتابخانه های عمومی کشور، دوره ۷ (شماره ۲).
۲. مهراد، جعفر؛ مقدس، جلیل (۱۳۷۸). "جست و جوی رایانه ای پایگاه های اطلاعاتی علم شبیه". پیام کتابخانه، دیپرخانه هیأت امنای کتابخانه های عمومی کشور، دوره ۹ (شماره ۳).
3. Chapman, J.L.(1980). "A State Transition Analysis of On- line Information - Seeking Behavior", *Journal of the American Society for Information Science* 32(2): 107-116.
4. Lancaster, F.W.(1979). *Information Retrieval Systems: Characteristics, Testing, and Evaluation*. 2nd ed. New York: Wiley.
5. Penniman, W.D. *A Stochastic Process Analysis of On- line Behavior*, Proceedings of the American Society for Information Science 38th Annual Meeting.
6. Tague - Sutcliffe, J.(1992). "The Pragmatics of Information Retrieval Experimentation, Revisited," *Information Processing & Management*, 28: 467-490.