

نمایه سازی

تهیه کننده:
حیدر قدیم خانی

جلسه اول و دوم مفهوم داده , اطلاعات , دانش

داده data: واقعیات عینی از رویدادها

اطلاعات information: اضافه کردن تفسیر به داده ها و ارتباط آنها به هم اطلاعات را شکل می دهد اطلاعات باید حاوی عناصر آگاهی دهنده و تغییر دهنده باشد

دانش: مفهومی عمیق تر و وسیع تر از اطلاعات دارد. مخلوطی از اطلاعات موجود, تجربه ها, ارزشها و نگرشها.

Knowledge: دانش کمک می کند که اطلاعات و داده های ناخواسته حذف شده و فهم و درک حاصل شود.

چرخه اطلاعات: داده ها و اطلاعات در پایگاه های اطلاعاتی ذخیره می شوند

مجله, مقاله, کتاب (...). چاپی دیجیتالی

پایگاه اطلاعاتی چیست؟ مجموعه ای از فایل‌های مربوط به هم در یک پایگاه اطلاعاتی که در مورد یک موجودیت است مجموعه ای از فایل‌های مورد نیاز و متناسب با کارکردی که از پایگاه اطلاعاتی انتظار می‌رود ایجاد می‌شود. شامل:

فایل ورود اطلاعات: عمل ذخیره سازی داده و تایید صحت آنها در این فایل انجام می‌شود

فایل اصلی: پس از کسب اطمینان از صحت داده ها ، برای نگهداری دائمی وارد فایل اصلی می‌شوند .
رکوردها پشت سر هم و به ترتیب شماره رکورد وارد می‌شود

بایگانی واژه ها و اصطلاحات: مجموعه ای از واژه ها یا عبارات قابل جستجوی در نظر گرفته شده در آن قرار می‌گیرد. معکوس شده فایل اصلی است و به همین دلیل به آن فایل مقلوب می‌گویند. هر واژه یکبار در نظم الفبایی قرار گرفته و در مقابل آن شماره رکوردهای مربوط قرار دارد . نمایه انتهای کتاب

فایل‌های جانبی: برای هر پایگاه اطلاعاتی نوع خاصی ساخته می‌شود. مانند فایل‌های مستند سازی , فایل های تصاویر و

مهمترین عملیات در پایگاه اطلاعاتی

ورود داده ها

اصلاح داده ها

بازیابی داده ها

ضرورت استفاده از فایل مقلوب: جستجو در فایل اصلی مشکل و وقت گیر است زیرا ترتیب شماره مدارک یا رکوردها می‌باشد و برای جستجو باید به ترتیب ورود

داده ها اقدام نمود این جستجو جستجوی ترتیبی نام دارد. مراجعه به فایل مقلوب سرعت بیشتری را در دسترسی به اطلاعات فراهم می‌نماید . (استفاده از اصطلاحات و واژه ها و مراجعه به رکورد مورد نظر. فایل مقلوب با مهارت‌های ویژه نمایه سازی ساخته می‌شود که نقش مهمی در بازیابی اطلاعات صحیح و حذف اطلاعات ناخواسته دارد. در فایل مقلوب پاره هایی از اطلاعات که قابل جستجو بر اساس نیاز واقعی کاربران و قابلیت‌های پردازش

در نرم افزار و سخت افزار قرار می گیرد. پاره های اطلاعاتی قرار گرفته در يك فایل مقلوب می تواند از يك عدد , يك كلمه تا يك عبارت باشد. يك رکورد که يك بار در رایانه ذخیره می شود , می توان مدخلهای زیادی برای جستجوی آن مشخص کرده و در نمایه مقلوب قرار داد. نمایه سازی از شیوه های بازیابی اطلاعات است. کلید واژه: کلمه , عبارت یا بخشی از محتوای يك پیشینه که دارای بار اطلاعاتی است و بعنوان شناسه یا مدخل در نمایه مقلوب ذخیره می شود تا در موقع جستجو بعنوان کلید اصلی دستیابی به آن پیشینه باشد.

افزایش دادن نقاط دسترسی یا بازیابی

ایجاد نمایه های بیشتر

فایلهای بیشتر

هزینه بیشتر

افزایش محل ذخیره سازی

افزایش هزینه روز آمد سازی و نظم دهی فایلها

بنابراین بازدهی مطلوب ندارد.

راه حل چیست ؟

- ترکیب پرسشها و استفاده از نمایه های "ترکیب کننده فایلها"
- در پایگاه های اطلاعاتی از عملگرهای منطقی جستجو استفاده می شود:
- در پایگاه های اطلاعاتی از عملگرهای منطقی جستجو برای بازیابی مطلوب استفاده می شود.

And: ترکیب (اشتراك) – خاصتر

Or : اجتماع – عامتر (بیشتر زمانی بکار می رود که برون داد بسیار کم باشد , همچنین زمانی که جستجو وسیعتر مد نظر قرار گیرد)

Not: جستجو را محدود می نماید .

ل جبر بول	AN	ل	OF	ل	NO	ل
And =	T and T =		T or T =		Not T =	
Or =	T and F =		T or F =		Not F =	
Not =	F and T =		F or T =			
	F and F =		F or F =			

در يك پایگاه اطلاعاتي باید استفاده از این عملگرها را تعریف نمود. ترتیب و تقدم و تاخر عملگر ها را نیز می توان تعریف کرد.

• سایر امکانات جستجو

• مجاورت کلمات History \$ \$ iran

• کلمات هم ریشه (محمود ؟ - محمودیان , محمود زاده, محمود پور و) ...

• عملگرهاي رابطه اي- + < >

نمایه معادل واژه لاتین INDEX به معنای نشان دادن و خاطر نشان کردن ،ریشه واژه لاتین آن INDIC به معنی دلالت کردن است کلمه نمایه از مصدر نمودن یا نماییدن فارسی گرفته شده است.

• هر مجموعه اطلاعاتي نیازمند نمایه است.

نمایه سازي (INDEXING SERVICE): یکی از رایج ترین شیوه های پذیرفته در سازماندهی اطلاعات در سطح جهان است , به معنای عمل توصیف یا شناسایی محتوای موضوعی يك مدرک به قصد سهولت بازیابی در نمایه سازي سنتی به شکل ساده اساسی ترین کلمات متن یا محتوای اطلاعاتي را جمع آوری کرده , بعد الفبایی نموده و آدرس مکان را می دهد در بین انواع نمایه ها , نمایه موضوعی به خاطر مفید بودن و کارسازتر بودن در تنظیم اطلاعات مقبولیت دارد. در این روش می توان محتوای اطلاعاتي مثلا “ يك كتاب خاص , يك موضوع خاص مانند علم کلام , مجموعه آثار یکاندیشمند مانند شهید مطهری را به گونه ای تنظیم کرد که بازیابی به

سهولت و بر اساس خواسته کاربرد در سریعترین زمان انجام یافت. Representation هدف از نمایه سازی باز نمود اطلاعات منتشر شده در قالبی مناسب در پایگاه اطلاعاتی می باشد این پایگاه باز نموده می تواند در قالب چاپی مثل نمایه نامه مهندسی یا در قالب کاردی مثل برگه دان یک کتابخانه سنتی ارائه شود. یا در قالب الکترونیکی مانند اصطلاحنامه های آنلاین یا موتورهای جستجو و راهنماهای موضوعی

تعریف نمایه:

طبق استاندارد نمایه سازی بریتانیا: نمایه سیاهه نظام یافته مدخلهایی است که به منظور کمک به استفاده کنندگان در جابجایی اطلاعات یک مدرک ساخته می شود. انجمن نمایه سازان امریکا: نمایه عبارت است از یک رشته محللهای دستیابی سازمان یافته که استفاده کننده را از اطلاعات معلوم به اطلاعات اضافی ناشناخته، راهنمایی می کند. مجموعه شناسه های الفبایی یا نظام یافته ای که کاربر را به جایگاه اطلاعات درون مدرک هدایت می کند. دانشنامه کتابداری و اطلاع رسانی نمایه چنین تعریف شده است: صورتی از موضوعها و واژه های مهم و دیگر مطالب یک یا چند کتاب با ارجاع به صفحاتی که این مطالب در آنها آمده است.

نیاز به نمایه

هدف اساسی حفظ و انتقال اطلاعات و دانش از نسلی به نسل دیگر بوده است. این ساز و کار رویکردی مشترک دارد یعنی: جمع آوری، ذخیره اطلاعات، بازیابی اطلاعات، استفاده به محض درخواستدهنده 1940 عصر افزایش سریع اطلاعات راهنماهای کلی دیگر کافی نبود. ما از داشتن میلیون ها سند راضی نبودیم، بلکه به تحلیل موثر و سریع محتوای آن اسناد نیاز داشتیم. به موازات تغییر نیازهای اطلاعاتی، تغییر فناوری هم صورت گرفت. یکی از تحولات آرام و بسیار اساسی در نیمه دوم قرن نوزدهم: تغییر بازیابی اطلاعات از اعم به اخصیکی از راههای سازماندهی اطلاعات نمایه سازی است

اهداف نمایه سازی

- نمایه دارای دو هدف کلی است زمان و تلاش برای یافتن اطلاعات را به حداقل میزان کاهش دهد و موفقیت جستجوی استفاده کننده را به حداکثر برساند. (انتخاب واژه مناسب در زمان کمتر و با تلاش کمتر)
- نمایه ها و چکیده ها به ارزش سند می افزایند.
- نمایه یک فرا داده است (داده ای در مورد داده ها)

- فضاي نمايه داراي دو گزاره مربوط به سند نمايه شده است . شمول و عدم شمول يعني چه چيزي در نمايه وجود دارد و چه چيزي در نمايه وجود ندارد.

اهداف نمايه سازي

به طوركلي:

- برقراري ارتباط ميان مفاهيم
- تنظيم شناسه ها به ترتيب نظام مند و موثر
- شناسايي سريع مدارك در يك مجموعه
- سازماندهي اطلاعات به قصد بازيابي سريع و آسان مدارك

تاريخچه نمايه سازي

تاريخ نمايه سازي و چكیده نویسی با تاریخ نوشتن و مجموعه پیشینه های اطلاعات ارتباط نزدیکي دارد . نمايه هاي اوليه به اسامي اشخاص و يا ترتيب آمدن كلمات در متن محدود بود. از جمله اولين نمايه ها : كشف اللغات (concordances) كهين ترين پیشینه اي كه از نمايه سازماندار در تاريخ سراغ داريم، سياهه لوح ههاي گلين كتابخانه هاي اكد و سوم ر است. امانخستين گامهاي جدي در راه تدوين نمايه در خصوص كتب مقدس ديني و در سده چهاردهم ميلادي در چند صومعه صورت گرفت. در مشرق زمين، خاصه در قلمرو تمدن اسلامي كارهايي در زمينه هاي رده بندي، فهرست نویسی و نمايه سازي انجام گرديد. فهرست ابن ندیم یکی از این نمونه هاست . استفاده از نمايه براي طومار پاپيروس مستلزم باز و بسته كردن مكرر طومار بود. و نمايه آن به شكل صفحه كوچكي به طومار متصل مي شده است و با به وجود آمدن كدكس كه شكل جديد كتاب بود گامي اساسي در نمايه شد. و نمايه كتاب براي اولين بار عملي شد . اختراع چاپ، رشد سريع چاپ و توليد كتاب، نياز به نمايه هاي كتاب را افزايش داد . افزودن آستر بدرقه به كتاب ها باعث خواننده گان مطالب و موضوعات مهم را ياد داشتن كنند به عنوان مثال: وكلا از آن براي فهرست الفبائي قوانين مورد علاقه خود استفاده كردند و كاتبان براي نوشتن ارجاعات به كتب مقدس....

- در قرن 16 نمایه های با کیفیت برای کتاب ایجاد شد
- در قرن 17 نوع جدید ابزار اطلاعاتی (نشریات ادواری) به وجود آمد و با رشد سریع آن نمایه ها ضرورت یافتند.
- همزمان با رشد مجلات علمی مجلات چکیده نیز توسعه یافت
- قرن 19 روند حرکت به سمت انتشارات تخصصی، مجلات چکیده رشد سریع خود را آغاز کردند
- انفجار اطلاعات از دهه 1940
- در دهه 1850 دبلیو اف. پول نمایه ای منتشر کرد که مجلات بسیاری را شامل می شد و مفهوم جدید انتشار نمایه مشترک برای شماره های متعدد نشریات گوناگون را آغاز کرد. که اهمیت فراوانی در رشد نمایه سازی داشته است.
- پس از جنگ جهانی دوم، دامنه تحقیقات و پژوهشها روز به روز گسترده تر گردید و این افزایش تدریجی کنترل اطلاعاتی با استفاده از ابزارها و شیوه های قبلی را غیرممکن م یساخت. در سال 1961 اولین کتاب تئوریک در زمینه نظام های بازیابی رشد سریع انتشارات و بررسی روشهای نمایه سازی بود. به «ویکری» اطلاعاتی ارائه شد. انگیزه اصلی تألیف این کتاب از جانب تدریج بر حجم این تلاش ها افزوده شد تا این که در دهه 1970 تحقیقات گسترده ای به روی بانکهای اطلاعاتی و سیستم های خودکار پیوسته صورت گرفت
- . در طول دهه 1980 استفاده از نظام های بازیابی اطلاعات در دو جنبه گسترش پیدا کرد :
- جنبه اول گسترش نظامهای تمام متن بود (تا قبل از این دهه در بازیابی اطلاعات فقط از نمایه ها و چکیده ها استفاده می شد)
- جنبه دوم گسترش نظامهای پیوسته مورد استفاده غیرمتخصصان بود. در دهه 1990 حجم منابع ذخیره اطلاعات رایانه ای افزایش یافت، تصاویر به همراه اطلاعات ارائه شده و راه نمایش اطلاعات به گونه ای متفاوت گردید. تحول مهمی که در این دهه به وقوع پیوست ظهور اینترنت بود. تا قبل از این تحول تعداد ناشران اطلاعاتی بسیار محدود بود ولی در حال حاضر افراد زیادی در سرتاسر جهان اطلاعات مربوط به خود را تولید و دسته بندی می کنند و از طریق صفحات خانگی، آنها را در اختیار دیگران قرار می دهند، امروزه حجم انبوه اطلاعات به خصوص اطلاعات موجود بر روی شبکه جهانی وب موجب ظهور جنبه های جدیدی در بازیابی اطلاعات گردیده است.

نمایه سازان چه کاری انجام می دهند؟

- تحلیل سند و تخصیص واژه های موضوعی که فکر می کند استفاده کننده براساس آن واژه ها جستجو خواهد کرد.

نمایه ساز به چه مواردی توجه می کند؟

- توجه به خط مشی نمایه سازی
- نیاز استفاده کننده
- نقش عامل انسانی در نمایه سازی منابع بسیار مهم است.
- متن را جامع و کامل مطالعه تحلیل کند
- تعیین موضوع و رشته
- تعیین بار اطلاعات و ارزش نمایه سازی آن
- چه میزان امکان عمیق شدن در موضوع وجود دارد
- اصطلاحات را انتخاب کند (گاهی فهم موضوع آسان ولی انتخاب اصطلاح اعم یا اخص یا مترادف مشکل است)
- آشنایی با فن نمایه سازی و اصطلاح نامه نویسی

ویژگی نمایه سازان

- نظم و انضباط و علاقه به جزئیات
- حافظه خوب یک سرمایه است
- سریع خواندن یک ویژگی مثبت است
- نمایه سازان سانسور کننده نیز هستند
- ذهن کنجکاو و خلاق
- نمایه ساز خوب خود را به جای استفاده کننده از اطلاعات قرار می دهد
- با وجود تحولات بسیار در حوزه کامپیوتر و نرم افزارها هنوز هم نمایه سازی نیازمند تحلیل موضوعی است که فرآیندهای ذهنی و انسان مدار است.

مراحل نمایه سازی

با تعارف گفته شده فرآیند نمایه سازی دو مرحله را شامل می شود:

1 - تحلیل مفهومی : یا تشخیص و درک جوهره

2- ترجمه : ارائه جوهره به صورتی که قابلیت پیش بینی داشته و حفظ امانت شود.

1- تحلیل مفهومی :

- تصمیم گیری درباره محتواست
- نمایه سازی برای برآوردن نیاز مخاطبان و بهره گیران از نظام به کار می رود بنابراین یک نمایه کار آمد نه تنها درباره محتوای مدرک باشد بلکه باید مشخص شود موضوع یا محتوای آن برای چه گروهی از بهره گیران تهیه می شود.
- به عبارت دیگر برای بهره گیران می توان با اهداف مختلف از یک مدرک واحد , نمایه سازی های متفاوت انجام داد.

نمایه ساز درباره یک مدرک باید سوالات زیر را مطرح کند:

- مدرک درباره چه چیزی است ؟
- چرا این مدرک به مجموعه اضافه شده است ؟
- بهره گیران به چه جنبه ای از مدرک علاقمندند؟

مثال : گزارش از یک ماموریت فضایی (اصطلاحاتی که این گزارش در هر شرکت زیر آنها نمایه سازی می شود)

- شرکت لاستیک سازی
- شرکت صنایع فلزی
- لباسهای فضایی
- ابزارهای فضایی
- جدید
- جوشکاری
- مصنوعی
- فنون

• لاستیک فلزات

ترکیبات

2 – ترجمه

- گام دوم در نمایه سازی است
- تبدیل تحلیل مفهومی یک مدرک به مجموعه ای از اصطلاحات نمایه ای
- در نمایه سازی هر نوع زبانی که برای رسیدن به قابلیت پیش بینی و حفظ امانت به کار برده شود ، زبان نمایه گفته می شود.
- هرچه میزان حفظ امانت و قابلیت پیش بینی در یک زبان نمایه بالا باشد ، نظم افزایش می یابد.
- در نمایه سازی باید بین اصطلاحات با معانی مختلف تفاوت قائل شد و آن را روشن نمود . (حفظ امانت)

سم گیاهی (گیاهان سمی یا سمی که بر علیه گیاهان به کار می رود)

اگر اصطلاحات به کار برده شده از زبان طبیعی (خود مدرک یا ذهن نمایه ساز) مشتق شده باشد زبان نمایه ، واژگان طبیعی گفته می شود . در غیر این صورت از یک اصطلاحنامه استفاده می شود.

از این جنبه می توان دو نوع نمایه سازی متفاوت را در نظر گرفت:

1- نمایه سازی استخراجی

2- نمایه سازی تخصیصی

1- نمایه سازی استخراجی: تعریف: واژگان و عباراتی را که واقعا در مدرک وجود دارد برای بیان محتوای موضوعی همان مدرک مورد استفاده قرار داد

زبان های نمایی سازی – واژه گزینی: مثال عنوان (تصویر کلی) : نظر سنجی افکار عمومی درباره نگرش آمریکاییان نسبت به خاورمیانه

چکیده (جزئیات بیشتر از عنوان) : یک بررسی تلفنی که در سال 1985 انجام شد مسائل زیر را مورد بررسی قرار داده است:

کمک ایالات متحده آمریکا به اسرائیل و مصر ؛ آیا ایالات متحده آمریکا باید با اسرائیل ، ملیتهای عرب یا سایر کشورها طرف شود ؛ آیا سازمان آزادی بخش فلسطین باید در یک کنفرانس صلح شرکت کند ؛ و آیا ایجاد یک کشور مستقل فلسطینی پیش نیاز صلح است.

اصطلاحات نمایی سازی استخراجی : افکار عمومی ، تحقیقات تلفنی ، ایالات متحده آمریکا ، نگرش ، خاورمیانه ، اسرائیل ، مصر ، کمک ، صلح

2- نمایی سازی تخصیصی تعریف : اختصاص اصطلاحات به یک مدرک از منبعی غیر از خود آن مدرک مثلاً“ از ذهن نمایی ساز یا از یک اصطلاحنامه مدون یا واژگان کنترل شده در حوزه موضوعی.

اصطلاحات نمایی سازی تخصیصی سند قبل:

- کمک خارجی
- ارتباطات خارجی آمریکا
- زبان های نمایی سازی – واژه گذاری

سوال

به دلیل مشکلات و سختی انواع نمایی ، چرا بدون دخالت انسان ، متون اصلی را ذخیره نکنیم و به این ترتیب همه لغات متن را

مورد جستجو قرار ندهیم؟ (ذخیره سازی متن کامل)

دلایل طرفداران ذخیره سازی متن کامل

- هیچ چیز از متن حذف نمی شود
- هیچ چیز به آن اضافه نمی شود
- هر لغت يك توصیفگر است
- لازم نیست استفاده کننده قواعد یچیده زبان مصنوعی را بیاموزد

مخالفان

- این تفکر ساده اندیشانه است
- در بازیابی مشکلاتی ایجاد می شود
- باعث افت اطلاعات می شود
- ممکن است لغت بکاربرده شده در متن دقیقاً به ذهن جستجوگر خطور نکند . مثلاً "اسید سولفوریک به جای جوهر نمک"

- ممکن است بیش از يك بیان لغوی برای مفهوم در متن نباشد . مثلاً "پرخاشگری نوجوانان ، خشونت در نوجوانان , بزهکاری ,

نوجوانان

- ممکن است در متن مفهوم به صورت غیر لغوی و تفسیری بیان شده باشد
- در نظام ذخیره سازی متن کامل اطلاعات غیر متنی مثل جدولها , تصاویر و ... قابل دسترسی نیستند.
- در ذخیره سازی متن کامل تعداد واژه ها زیاد و جم ذخیره سازی افزایش می یابد
- اتلاف وقت زیادی وجود دارد
- در نظام متن کامل مسئله ابهام در بیان حل نشده است . مثلاً "اصطلاح کاج نقره ای (معلوم نیست کاجی از جنس نقره یا نام

یک نوع گیاه است)

- در مراحل نمایه سازی , مرحله ترجمه وجود دارد در ذخیره سازی متن کامل این مرحله حذف می شود.

- در مرحله انتخاب جوهره متن هم صورت نمی گیرد.

تصاصات روش	نیره ساري متن كامل	اياه ساري كنترل شده
اسايي معنا و رفع ابهام		
ليت پيش بيني جوهره		(-
ليت پيش بيني ارائه جوهره و خاب توصيفگر		(-
كردن خلاها		(-
وي ساري تفسيروي		(
ليت ارائه اطلاعات غير متني		(-
ترل مترادفات		
ص بودن كليد واژه ها	(+	(+
بج بودن مفاهيم		
كان نحو در زبان نمايه		

		کان جمع آوری توصیفگرها
	(تترسی به اصطلاحات جا افتاده در ایه سازی
		مان و حجم لازم برای ذخیره سازی

راهنمای جدول:

-نقطه ضعف

+نقطه قوت

(+) تا حدودی قوت

(-) تا حدودی ضعف

(+/-) دارا بودن حالت بین قوت و ضعف

نمایه به دو دسته کلی نمایه توصیفی و نمایه موضوعی تقسیم بندی می شود:

نمایه توصیفی **Descriptive Index**: پرداختن به ویژگیهای اطلاعاتی غیرموضوعی یک مدرک که معمولاً به اطلاعات کتابشناختی موسوم می باشند رami گویند. « توصیفی نمایه»

نمایه موضوعی **Subject Index**: ساخت و پرداخت زمینه های اطلاعاتی مدرک که بار موضوعی دارد و دسترسی بر مدارک را از طریق موضوع امکان پذیر می

خوانده می شود. « نمایه موضوعی » سازد

زبان های نمایه سازی

در نظام های اطلاع رسانی به زبان ساختگی و قراردادی اطلاق می شود که برای مقاصد نمایه سازی به ویژه قابلیت بازیابی اطلاعات و مدارک به کار گرفته می شود. به طور کلی زبان نمایه سازی، استانداردی را مهیا می کند که هم نمایه ساز و هم جستجوگر می توانند از آن استفاده کنند. پس به طور ساده تر: مجموعه ای از روشهای از پیش تعیین شده برای سازمان دهی، بازیابی و اشاعه اطلاعات

زبان های نمایه سازی :

به دو دسته کلی تقسیم می شوند:

1 - نظام های اصطلاح تعیین شده **Assigned term** (دارای ابزارهای کنترل واژگان) نمایه سازی کنترل شده
Controlled language

2 - نظام های اصطلاح مشتق **Derived Term** توصیفگر ها از متن گرفته می شود لذا آن را متن آزاد و طبیعی هم می گویند (نمایه سازی آزاد و طبیعی)

Free language

مزایای زبان طبیعی

زبان متخصصان هر رشته است و برای ارتباط با سایر متخصصان بسیار مهیا تر می باشد

- ساختار طبیعی زبان حفظ می شود
- رو زآمدی در قیاس با زبان ساختگی
- ناهمگونی های اصطلاحات را تحت سیاست واحد قرار می دهد
- امر نمایه سازی برای نمایه ساز با سهولت و راحتی بیشتری صورت می پذیرد زیرا نیاز چندانی به تحلیل موضوعی مدرک نیست
- جامعیت را افزایش می دهد و از پراکندگی اطلاعات جلوگیری می کند
- سرعت نمایه سازی افزایش می یابد

معایب زبان طبیعی

- 1- از آنجا که کلیدواژه ها مطابق با واژگان مؤلف انتخاب شده کاربر مجبور است جهت بازیابی مدرک تمامی کلیدواژه های احتمالی را حدس بزند
- 2- به دلیل وجود مترادف ها و واژه های مشابه احتمال بازیابی مدارک نامرتبب افزایش می یابد
- 3- معمولا وقت و زمان زیادی از کاربر برای جستجو و بازیابی مدرک مورد نظرش گرفته می شود

مثال : زبان طبیعی

عنوان مقاله: اهمیت و نقش کتابخانه های آموزشگاهی در آموزش

چکیده

اهمیت مطالعه در دوران کودکی و نوجوانی و قابلیت گسترش آن در سنین بالا بر کسی پوشیده نیست. جامعه دانش آموزی برای به دست آوردن روحیه پژوهش و مطالعه فردی از نخستین سال دبستان تا آخرین سال دبستان به کتابخانه مجهز نیاز دارد. در این مقاله ابتدا کتابخانه های آموزشگاهی در مقاطع تحصیلی تعریف و سپس اهداف، وظایف، و تشکیلات آن بیان می شود. در پایان برای بهبود وضعیت موجود کتابخانه های آموزشگاهی پیشنهادهایی ارائه شده است.

کلیدواژه ها: کتابخانه های آموزشگاهی، مطالعه، پژوهش، آموزش

مثال : زبان کنترل شده

عنوان مقاله: اهمیت و نقش کتابخانه های آموزشگاهی در آموزش

چکیده

اهمیت مطالعه در دوران کودکی و نوجوانی و قابلیت گسترش آن در سنین بالا بر کسی پوشیده نیست. جامعه دانش آموزی برای به دست آوردن روحیه پژوهش و مطالعه فردی از نخستین سال دبستان تا آخرین سال دبستان به کتابخانه مجهز نیاز دارد. در این مقاله ابتدا کتابخانه های آموزشگاهی در مقاطع تحصیلی تعریف و سپس

اهداف، وظایف، و تشکیلات آن بیان می شود. در پایان برای بهبود وضعیت موجود کتابخانه های آموزشی پیشنهادهایی ارائه شده است.

کلیدواژه ها: کتابخانه های آموزشی، مطالعه، پژوهش، آموزش، ارتقاء کیفیت آموزش، پژوهش محوری
مزایای زبان ساختگی

- بصورت يك بسته آماده به کامپیوتر داده می شود و بین مدرک و زبان کنترل شده يك پیوند برقرار می کند.
- به مراتب توانمند تر در دادن اطلاعات است هر چند مانعیت کمتری دارد.
- وقتی ایجاد سیاست می کنیم (مثل قوانین جدا نویسی ...) خود نوعی کنترل است و زبان ساختگی می شود.
- زبان طبیعی نیاز به ارجاع دارد بنابراین از حالت طبیعی خارج می شود.

کارکرد نمایه سازی (Indexing Function)

- نمایه سازی عموماً سه کارکرد عمده دارد:
- محتوای اطلاعاتی مدارک را فشرده م یسازد.
- به عنوان واسطه برای تطبیق و یکسان سازی زبان مدرک و زبان کاوش به کار می رود.
- به عنوان ابزاری کارا، بر شیوه تدوین راهبردی کاوش در جستجویی اطلاعات نظارت دارد.

گام های ایجاد يك نمایه

- تعیین خط مشی نمایه سازی
- خواندن سریع متن بدون یادداشت و علامت (عنوان , چکیده , خلاصه , نتیجه گیری , تصاویر , فهرست مندرجات که باعث
- درک کلیت موضوع می شوند)
- مشخص کردن کلید واژه ها: باید به علایق جامعه و نوع سازمان توجه کرد

- انتقال هر کلید واژه (مدخل) و ذکر محل بر هر برگه چا پی یا به شکل درون خطی، ثبت روی خود مدرک، ثبت روی نوار صوتی
- الفبایی کردن مدخل ها و ادغام مدخل های تکراری
- ایجاد پیوند مفهومی بین مدخل ها
- استفاده از اصطلاح نامه ها برای ایجاد روابط
- ویژگی های خطی مشی: سیاست نمایه سازی
- تعیین تعداد متوسط کلید واژه ها
- تعیین نوع نمایه
- زبان نمایه سازی (آزاد، کنترل شده، ترکیبی)
- نوع بازیابی (راهنمای استفاده)
- روشن ساختن بیانگر (جمع یا مفرد، تکلیف وضعیت کلمات خارجی)
- سیاست مرجع یا نامرجح کردن بیانگرها
- ارجاعات و جای نماها
- بیان انواع نمایه موجود (موضوع، نویسندگان،).
- استاندارد سازی (ابزارها)
- خودکار سازی نمایه (زبان برنامه نویسی، نرم افزار، چگونگی بازیابی)
- روزآمد سازی

سیاست نمایه سازی

سیاست سازمان متولی نمایه سازی که تصمیم می گیرد نمایه جامع باشد یا جزنگر، تفصیلی یا اجماعی، کم عمق یا با عمق بیشتر؟ نکته موثردیگر در تعیین سیاست نمایه سازی، سطح علم و فهم شخص نمایه ساز است، باید تخصصی علمی او سنجیت داشته و توان تحلیل موضوع و استفاده از اصطلاح نامه را داشته باشد.

نمایه سازی پیش همارا: Pre Coordination Indexing

هرگاه بین دو یا چند جزء به طور تصنعی و ساختگی، پیوند برقرار کنیم (ایجاد نحو) به این عمل و نوع آن، نمایه پیش همارا گویند. به عبارت دیگر در نمایه سازی پیش همارا، ترکیب یا همارایی عناصر تشکیل دهنده موضوع مورد جستجو در هنگام نمایه سازی و به عبارتی پیش از بازیابی صورت می گیرد: مانند : سرعنوان های موضوعی

• اصفهان- تاریخ- قبل از اسلام نمایه سازی پیش همارا بیشتر در نمایه های چاپی به کار می رود.

- برگه دان کتابخانه معمولی ترین نظام بازیابی اطلاعات پیش هماراست
- ترکیب اصطلاحات در هنگام جستجو ساده نیست
- اولین اصطلاح باید مهمترین اصطلاح باشد
- نشان دادن روابط در این نظام مشکل است

نمایه سازی پس همارا: post coordination Index ing

شیوه ای از نمایه سازی که در آن نمایه ساز، سرشناسه ها را از مفاهیم بسیار ساده انتخاب می کند و تعدادی شناسه زیر هر یک اضافه می نماید و نیز تدابیری برای پیوستن آنها به یکدیگر به دست می دهد تا به وسیله آنها جوینده بتواند موضوع مدرک مورد نظر خود را بیابد. ترکیب یا همارایی عناصر تشکی لدهنده موضوع مورد جستجو در هنگام بازیابی انجام می شود.

مانند : اصطلاح نامه ها

تاریخ * قبل از اسلام * اصفهان

- اصطلاحات در فایل مقلوب و اطلاعات در فایل اصلی ذخیره می شوند
- یک نظام بازیابی اطلاعات که به کاوشگر امکان می دهد که به هر طریق اصطلاحات را با هم ترکیب کند نظام پس همارا گویند.
- بیشتر در نمایه های ماشینی کاربرد دارد

- چند بعدي بودن ارتباط بين اصطلاحات حفظ مي شود
- همه اصطلاحات اختصاص یافته به مدرک وزن یکسانی دارند
- مي توان به مدرک مورد نظر با جستجو از طريق يك يا چند اصطلاح به مدرک دست يافت.
- از عملگرهاي منطقي در آن استفاده مي شود

مثلاً“ تاريخ 10

تاريخ * قبل از اسلام 6

تاريخ * قبل از اسلام * اصفهان 2

نمایه بر اساس حجم و سیاست گذاری به دو دسته کلی تقسیم میشود:

- نمایه تک متني
- نمایه مجموعه اي

نمایه تک متني: نمایه اي که متن واحدي را در بر مي گیرد و ناظر بر يك متن واحد است. مانند نمایه کتاب
Collection Index نمایه مجموعه اي نمایه مجموعه يا نمایه نامه شامل محتويات مجموعه اي از آثار هم نوع
و اغلب هم موضوع است. مانند نمایه نشریات ادواري، نمایه نامه نشریات شيمي و نمایه نامه روزنامه ها که
در واقع ناظر بر مجموعه اي از متون و منابع است، که بنا به هدف و سياست ويژه اي گردآوري و کنترل
واژگان در آن اجتناب ناپذير است. واحد اصلي نمایه مدخل است. از مجموعه مدخلها يك نمایه تشكيل ميشود.
يك مدخل متشکل از عناصر زیر مي باشد:

- الف. شناسه

- ب. بیانگر Modifier
 - ج. جایما یا ارجاع صفحه Locator/Refrence
 - د. ارجاع دو سویه Cross-refrence
- شناسه: جزء محوري مدخل را شناسه مي گویند که محل دسترسي به اطلاعات را در نظام تعريف شده و معمولاً الفبایي نشان میدهد.

- شناسه عبارت است از واژه یا واژه ها، نشانه یا نشانه هایی که به عنوان اصطلاح های نمایه ای از يك متن انتخاب میشوند و به صورت الفبایي یا نظام انتخابي دیگر در نمایه مرتب میشوند.

انواع شناسه از نظر شکل:

- شناسه تک واژه ای یا ساده (بسیط): این نوع شناسه تنها از يك واژه تشکیل میشود و این واژه ها می تواند به صورت مفرد یا جمع باشد، مانند موسیقي، هنر، گلهاء،

قلب

- شناسه عبارتي: این شناسه از ترکیب دو یا چند واژه همراه با حرف ربط یا حرف اضافه میان آنها تشکیل میشود . در این نوع شناسه ها به طور معمول دو واژه که لازم است رابطه آنها به یکدیگر نشان داده شود آورده میشود مثل: آموزش و پرورش " و " تکنولوژی و تمدن " که به طور معمول از يك خانواده یا زمینه موضوعي هستند و بیشتر اوقات باهم مطرح میشوند. مانند: کار و کارگر شناسه های متشکل از صفت و موصوف و مضاف و مضاف الیه

این شناسه ها به طور معمول از اسم و صفت و مضاف و مضاف الیه تشکیل میشود مانند:

- روان شناسي رشد
- جامعه اطلاعاتي
- شناسه اسم با توضیحگر

گاهی برای تعیین حدود و معنی و برطرف کردن ابهام و اخص کردن موضوع و گاه برای تمیز قائل شدن میان دو یا چند شناسه مشابه، واژه توضیحگر در مقابل آن و در مقابل پُرانتزافزوده میشوند، و در واقع واژه توضیحگر معنی از "نظر" می دهد:

مانند:

حرکت از نظر فلسفه

حرکت از نظر فیزیک

حرکت (فلسفه)

حرکت (فیزیک)

شناسه مقلوب یا معکوس

این شناسه در واقع همان شناسه های دو واژه ای صفت و موصوف یا مضاف و مضاف الیه هستند که به منظور قرار گرفتن واژه

مهمتر در محل الفبایی خود به شکل مقلوب یا معکوس در می آیند از حالت زبان طبیعی خارج میشوند و میان دو جزء شناسه

نشانه و برگول می نشینند مانند:

اطلاعات، آلودگی

اطلاعات، اقتصاد

کودک، روانشناسی

کاربرد این نوع شناسه های موضوعی در زبان فارسی زیاد رایج نیست.

انتخاب شکل شناسه ها به نوع نمایه و موضوع بستگی دارد.

بیانگر Modifier/ Subheading :

در نمایه برای اخص کردن شناسه و نشان دادن جنبه های مختلف آن از بیانگر استفاده میشود. به این مفهوم که شناسه تک واژه ای یا ساده (بسیط) را میتوان را با افزودن واژه های مناسب خاص تر کرد . بیانگر اساسا به منظور اخص کردن شناسه و نیز برای حفظ جنبه های لاینفک يك شناسه موضوعي با خود شناسه و محدود کردن دامنه معنایی آن بوجود می آید و ممکن است در يك مدخل شناسه، بیانگر نداشته باشد و مدخل تنها متشکل از دو عنصر شناسه و جایما باشد، مانند:

130 ، ارتباطات، 125

انواع بیانگر

-بیانگر موضوعي یا موضوع زیر موضوع

- روان شناسي
- 42 ، 76 ، 97 ، كودك
- اطلاعات
- 11 ، 17 ، 48 ، آلودگي

-بیانگر تاریخي یا زمانی

- اوضاع اجتماعي
- 12 ، 14 ، 17 ، 29 قرن

- 85، 73، 45 قرن 13
- 110، 104، 99 قرن 14

-بیانگر مکانی یا جغرافیایی

- آموزش و پرورش
- 55، 18، ایران 15
- 104، 75، فرانسه 43

جایما یا ارجاع صفحه:

عدد یا نشانه ای که محل دستیابی به اطلاعات را مشخص می سازد به آن جایما یا ارجاع صفحه گفته میشود (.

(Page refrance

- جایما ممکن است شماره صفحه، پاراگراف، بخش یا فصل یا آدرس اینترنتی باشد.
- جایما تا آنجا که امکان دارد باید کامل و دقیق باشد.
- میان شناسه و نخستین جایما و بیانگر، ویرگول و یک فاصله خالی گذاشته میشود.

ارجاع

ارجاع: یعنی راهنمایی از یک شناسه به شناسه دیگر یا رجوع از بیانگرهای یک شناسه به خود شناسه و یا شناسه های دیگر.

الف - - ارجاع ممکن است یک سویه باشد : شناسه نامرجح به شناسه مرجح

معنای آن این است که برای کلیه اطلاعات مورد جستجو باید به اصطلاح مرجح (گزیده نگاه کرد (Perferred .

یا terms)

مانند:

بیوشیمی نگاه کنید به زیست شناسی

ب- ارجاع ممکن است دو سویه باشد: معنای آن این است که برای اطلاعات مورد جستجو به جای دیگر نیز رجوع شود.

مانند: اعتیاد نیز نگاه کنید به مواد مخدر ===== کار و کارفرما " نیز نگاه کنید به " کار و سرمایه مطالب مورد ارجاع به طور معمول اخص تر از مطالب اصلی است.

پرندگان نیز نگاه کنید به قمری ها، کبوترها، ماکیان

ارجاع های دو سویه، جستجوگر را به اطلاعات مرتبط و وابسته راهنمایی می کند.

Blind entry ارجاع کور :

ارجاع از يك شناسه به شناسه ای که در نمایه وجود ندارد و یا ارجاع جایز به صفحه هایی از متن که فاقد اطلاعات و

مطالب ذکر شده باشد، ارجاع کور یا بی فرجام گفته میشود، و در برخی از منابع با عنوان مدخل کور از آن یاد شده است.

ارجاع دور Circuit reference

اشتباه دیگری که در نمایه سازی رخ می دهد ارجاع از واژه ای به واژه دیگر است که به طور معمول به صورت دور تسلسل

انجام می گیرد که به آن ارجاع چرخه ای یا دور گفته میشود.

در نمایه نام ها که بحث مستند سازی نام اشهر مطرح میشود گاهی ارجاع دور پیش می آید.

مانند:

بوعلی نگاه کنید به شیخ الرئیس

شیخ الرئيس نگاه کنید به بوعلی

نمایه کتاب از لحاظ شیوه چاپ و ساختار مدخل

الف – نمایه ساده

ب – نمایه درون بافتی

ج- نمایه برون بافتی

نمایه ساده: به نمایه ای گفته میشود که مدخل یعنی مجموعه شناسه، بیانگر و جایما در يك سطر مي آیند بیانگر به صورت طبیعی در ادامه شناسه ذکر می شود و جایماها نیز به دنبال آن می آیند. در نمایه ساده در مواردی که شناسه به تنهایی خود مدخل قرار گرفته است برای جلوگیری از زیاد شدن جایماها باید يك توضیح گر در برابر آن بیاید و یا در مقدمه نمایه ذکر شود. در نظام نمایه ساده به طور معمول شبکه مفهومی، بین مفاهیم برقرار نمی شود.

نمایه درون بافتی:

به این نوع نمایه نمایه پاراگرافی، پیوسته یا پیوسته سطری نیز گفته میشود. زیرا ترتیب نوشتن بیانگر يك شناسه به صورت نحوه ساختار اصطلاحها در همان بافت سطری و پیوسته است. در این نوع نمایه به لحاظ ساختار نحوی به زبان طبیعی است.

ترتیب نوشتن بیانگر در نمایه درون بافتی، سطری است یعنی تا جایی که سطر اجازه می دهد و از سطر دوم به بعد تو رفتگی دارد. توالی بیانگر ها تابع توالی حضور آنها در متن اصلی است. ترتیب پیدایش بیانگرها در متن است. بنابراین الفبایی نیستند. بیشتر مناسب متون علوم انسانی، تاریخی و سرگذشتنامه ای می باشد

مانند:

جاویدان، داریوش

؛ 20، 15، محل تولد، 12؛ دوران کودکی، 13

35، ؛تحصیلات و ،33-24، دوران نوجوانی، 22

40- آثار علمی، 35

نمایه برون بافتی:

به این نوع نمایه، نمایه خطی یا جدا یا سطری یا گسسته نیز گفته میشود. در این نمایه بیانگرها جدا از شناسه و به صورت الفبایی یا تو رفتگی مرتب می شوندیا به عبارت دیگر بیانگرها به صورت منفرد و هر کدام در سطری مستقل همراه با جایما های مربوط است. در این نمایه توالی بیانگرها تابع نظم الفبایی است.

مانند:

روان شناسی

1950، 184، بالینی، 179

232، 218، رشد، 150

215، 143، کودک، 129

تفاوت های نمایه برون بافتی و درون بافتی

نمایه درون بافتی صورت نحوی دارد.مثلا میگوییم رفتار با کودکان ولی دربرون بافتی با را بکار نمی بریم و نمی دانیم آیا رفتار کودکان است یا رفتار با کودکان .

• نمایه درون بافتی بیشتر به حوزه علوم انسانی تعلق دارد.

- نمایه درون بافتی کل نگر هستند.
- نمایه برون بافتی بیشتر به حوزه های علوم و فنون مرتبط است و بیشتر بن مایه اصطلاحنامه قرار می گیرد.

نمایه از لحاظ پوشش و محتوی:

الف – نمایه نام ها

- نام پدید آور
- نمایه عناوین
- نمایه نام های جغرافیایی
- نمایه نام سازمان و موسسات

ب – نمایه موضوعی

نمایه ها بر حسب روش تنظیم

الفبایی: بیشتر نمایه های کنونی به ترتیب الفبایی است و شامل نمایه الفبایی اسامی اشخاص، سازمان ها و نیز موضوعات می باشد

زمانی: به ترتیب زمان از قدیم به حال می باشد و بیشتر برای نمایه های تاریخی به کار می رود

رده ای یا موضوعی: که براساس رده ها یا سرعنوان موضوعی نظام مند مرتب می گردند. بیشتر در نمایه های علمی به کار می رود که البته می تواند یک نمایه الفبایی رده ای باشد.

تکاملی: برای نمایه های زمین شناسی

ملاك انتخاب اصطلاح مرجح

- رواج داشتن . مثال: جامعه شناسی به جای علم الاجتماع

- بومی بودن . مثال : آرمان گرایی به جای ایده آلیسم
- جدید و امروزی بودن . مثال: هواپیما به جای طیاره
- علمی بودن. مثال: اسید سولفوریک به جای جوهر نمک
- سرنام یا اختصار مشهور. مثال: یونسکو
- پرهیز از بکار بردن کلمات ترکیبی به جز در مواردی که اصطلاح قابل جدا سازی نیست مثل : آبله مرغان
- جامعیت و مانعیت (Recall, Precision): نمایه سازی مدارک مختلف که توسط افراد متفاوت، در مکان های متفاوت و در زمان های متفاوت و احتمالاً با واژگان متفاوت است بنابراین نوعی همگونی و هماهنگی باید پدید آورد که در عین جامعیت دستیابی، موارد زائد بازیابی نشود. پس برای ایجاد همگونی 2 نکته مهم است :
- جامعیت و مانعیت یعنی: از نظر منطق جامع افراد یا افراد همگون زیر یک چتر بروند و مانعیت یعنی اغیار زیر این چتر نروند. تلاش می شود سطح هر دو (جامعیت و مانعیت) بالا برود و توازن آن حفظ شود

سایر انواع نمایه

- نمایه های الفبایی یا واژه ای
- ✓ اساسی ترین مانع در استفاده از آن ها مشکل مربوط به مترادف ها و پراکندگی مدخل هاست
- ✓ نام ها یا مدخل ها در محل الفبایی خود درج شد هاند
- ✓ شکل کلی و غالب نمایه های امروزی است
- نمایه های مؤلف
- ✓ برای ایجاد آنها می توان از رویکرد موضوعی غیرمستقیم بهره برد یعنی با نویسندگان پیشین یا مشهور کار خوشه بندی اسناد و منابع را انجام داد
- ✓ نقاط مدخل آنها شامل اسامی افراد، سازمان ها، پدیدآورندگان تنالگانی، و... م یباشند
- نمایه های کتاب
- ✓ هدف از طراحی آن تسریع در بازیابی اطلاعات موجود در کتاب است و به هیچ عنوان جایگزین اطلاعات موجود در کتاب نیست
- ✓ در انتهای اغلب کتاب ها به چشم می خورد
- نمایه های استنادی

✓ مزیت اصلی آن این است که برخلاف نمایه های سنتی حرکت رو به جلو نیز دارد یعنی مشخص م یکنند که چه مقدار از

✓ ایده های پیشین در ایجاد ایده های جدید نقش داشته اند

✓ برای خوشه بندی اسناد و منابع به منظور قراردادن آن ها در گروه های موضوعی مشابه بسیار مؤثر است

✓ شامل فهرستی از مقالات و همچنین فهرستی فرعی از مقالات منتشر شده ای است که به آن مقالات استناد کرده اند

• نمایه های رد های

✓ هر چند استفاده از آن ها برای متخصصان ساده و کاربردی تر است اما کاربران عادی در استفاده از این نمایه

ها دچار مشکل م یشوند زیرا اغلب با نمایه های الفبایی راحت تر هستند

✓ مدخل های آن بر اساس رده ها یا سرعنوانهای موضوعی مرتب می شود

• نمایه های همارا

✓ پایه و اساس بسیاری از نظام های بازیابی اطلاعات در حال حاضر همین نمایه و نوع "پس همارا" می باشد

✓ دو نوع عمده آن عبارتند از "پیش همارا" و "پس همارا"

✓ نوعی نمایه که به زمان و نحوه ترکیب اصطلاحات نمایه ای اشاره دارد

• نمایه های درهمکرد

✓ بیشتر برای مجلات علمی و کارهای بزرگ و مهم بکار می روند

✓ ترکیبی از چند نمایه که در طول زمان خاصی منتشر می شوند

• نمایه های فرارسانه ای

✓ غالباً در مورد پایگاه های اطلاعاتی و منابع الکترونیکی غیرآنلاین کاربرد دارد

✓ این نمایه به کاربر اجازه می دهد تا از طریق برچسب های الکترونیکی مسیر خود را به سمت اطلاعات مورد

نیاز بیابد

• نمایه های چندرسانه ای

✓ این نمایه مواد متنی، صوتی و تصویری را با هم ترکیب می کند

• نمایه های نشریات ادواری

- ✓ اساس و مبنای کار در آن ها مانند نمایه کتاب است با این تفاوت که در تدوین آن چندصد مؤلف نقش دارند و اساساً کار پردردسرتري است
- ✓ دو نوع م یباشد: نمایه انفرادي براي مجلات انفرادي و نمایه هاي کلي براي گروهی از مجلات
 - نمایه هاي گردشې عنوان
- ✓ ایرادهایی اساسي دارند مانند: پراکندگی مترادفها و اصطلاحات، عدم شفافیت لازم عنوان براي تبیین محتوا، محدودیت
 - ✓ اصطلاحات موجود در عنوان براي بیان صحیح موضوع، و ...
 - ✓ ساخت و ایجاد چنین نمایه اي بسیار سریع و فاقد دردسر مي باشد
 - ✓ در این نمایه فرض بر این است که عنوان به قدر کافي گویاي محتوای سند مي باشد
 - با ورود رایانه به عرصه نمایه سازی دو نوع متداول نمایه گردشې عنوان شکل گرفت
- ✓ (Keyword Out of Context) 2 کواک (-) (Keyword In Context) 1 کوئیک
 - عصر ارتباطات و اطلاعات
 - عصر انفجار اطلاعات
 - عصري که گردآوری، ذخیره ، بازیابی و اشاعه اطلاعات به شدت تغییر یافته است.
 - از تاثیرات این تغییرات :
 - ✓ کاهش زمان جهت بازیابی اطلاعات مورد نظر
 - ✓ اهمیت دستیابی به مطالب روزآمدتر
 - ✓ نقش مهم تحقیق در پیشرفت و توسعه جوامع بشري
 - ✓ رشد پرشتاب دستاوردهای علمي و فني
 - ✓ افزایش متون ومدارك اطلاعاتي
- در حوزه اطلاع رسانی ، نمایه سازی مورد توجه قرارگرفت .
- نمایه سازی حلقه ارتباطي میان تولید کننده و مصرف کننده اطلاعات ، در جهت دستیابی آسان و سریع به دانش موردنیاز کامپیوتر روش های نمایه سازی را دچار تحول کرده در حال حاضر نمایه سازی از طریق کامپیوتر و نرم افزارهای توسعه یافته انجام مي شود.
- به طور کلي کاربرد کامپیوتر درنمایه سازی را به دو دسته تقسیم مي کنند:
 - ✓ نمایه سازی رایانه اي
 - ✓ نمایه سازی با کمک رایانه.

✓ نمایه سازی با کمک رایانه :

در این نظام نخست نمایه ساز با دقت ، شناسه های موضوعی را انتخاب و علامت گذاری می کند و نمایه دستنویس تهیه می کند. آن گاه اپراتور مدخلها را بر اساس نمایه دستنویس از طریق صفحه کلید وارد رایانه می کند بعد از ورود ، رایانه مدخلها را بر اساس نظم خاصی مرتب می کند، توصیفگرها را زیر شناسه ها می برد، جاینامه های هر مدخل را با هم ادغام می کند و بعد یک نمایه الفبایی ارائه می دهد. (بر اساس برنامه های از پیش طراحی شده)

این نوع نمایه سازی با دوروش متفاوت انجام می شود:

1) نمایه سازی در محیط گرافیکی (نمایه سازی درون کاشتی)

(Graphic User Interface (GUI) . امروزه بیشتر نرم افزارها با رابط گرافیکی کار می کنندکار با این

روش به صورت زیر است:

نمایه ساز کلیدواژه ها یا کلمات و عبارتهای مهم را انتخاب می کند. این عبارت انتخاب شده از نظر سیستم یک کلیدواژه محسوب می گردد. نمایه ساز همزمان با انتخاب کردن این کلید واژه ها باید یک پنجره واژه پرداز مستقل دیگر را باز می کند و مدخلهای ارجاعی را به طور دستی

تحریر می نماید. نتیجه این کار یک فایل فرعی نیز به دست می آید. نمایه ساز فایل اصلی مدخلها و فایل فرعی را ادغام و توسط یک برنامه ترتیب بندی مناسب الفبایی می کند.

2) نمایه سازی در محیط دستوری (نمایه سازی فرمانی)

Command Line Interface : این نرم افزار در محیط سی. ال. آی یعنی با رابط دستوری (فرمانی) کار

می کنند کار با این روش چنین است:

نمایه ساز صفحه ها را در صفحه نمایش کامپیوتر مطالعه می کند و کلیدواژه ها را با دونماد کاملاً استثنایی چاپ و راست که مشابه آن در متن وجود نداشته باشد، علامت گذاری می کند (مثلاً کلیدواژه درون سه گوشه قرار می گیردیا...). سپس نمایه ساز در همان فایل ارجاعات مورد نظر تحریر می کند. پس از اتمام علامت گذاری با گوشه ها و ارجاعات و تبدیلهای یک ابزار نرم افزار می تواند کل نمایه را تهیه و الفبایی کند. اگر نرم افزار دارای امکانات نمایه سازی نباشد ویراستار می تواند یک برنامه بسیار ساده (مثلاً به زبان بیسیک بنویسد) که فهرست کلیدواژه های علامت گذاری شده و شماره صفحه های آنها را

استخراج کندویک برنامه یا ابزار ساده دیگر نمایه را ترتیب بندی و مرتب کند. روشهای نمایه سازی با کمک رایانه دارای مزایای سرعت ورود اصطلاحات نمایه ای، ظرفیت ذخیره سازی بالا، سهولت اصلاح، حروفچینی و تایپ، ویرایش و چاپ نمایه هستند و از لحاظ اقتصادی مقرون به صرفه می باشند. در نمایه سازی به کمک رایانه، نمایه ساز از رایانه برای انجام امور نمایه سازی استفاده می کند. اولین نمایه های رایانه ای مانند کوئیک و کواک از نوع نمایه سازی به کمک رایانه بوده اند.

نمایه سازی کوئیک: kwic

Keyword In Context:

(پراشر 1989):

این نظام نمایه سازی بر این اصل استوار است که عنوان یک مدرک، محتوای آن را نشان می دهد. در این نظام نمایه سازی عنوان و واژه های آن نشاندهنده موضوع مدرک می باشد.

(نایت 1970):

این نوع نمایه، عنوان های مدارک را بر اساس کلیدواژه های موجود در آن مرتب کرده و هر کلیدواژه در نظامی الفبایی، نقش مهمی به عنوان مدخل دارد.

(لنکستر 1991):

نمایه کوئیک نمایه ای گردشی است که اغلب از عنوان های انتشارات به وجود می آید در نظر گرفته می شود؛ (Access Point) بدین مفهوم که هر کلیدواژه موجود در عنوان به منزله نقطه بازیابی مثال:

مبانی نظری و توصیه های عملی تحلیل موضوعی و نمایه سازی /
توصیه های عملی تحلیل موضوعی و نمایه سازی / ، مبانی نظری و
تحلیل موضوعی و نمایه سازی / ، مبانی نظری و توصیه های عملی
نمایه سازی / ، مبانی نظری و توصیه های عملی تحلیل موضوعی و

Keyword Out of Context: نمایه کواک kwoc

نمایه کواک همان نمایه برون بافتی است.

در این نمایه هر کلیدواژه به ترتیب از عنوان خود خارج شده و مقدم بر سایر اجزای عنوان قرار می‌گیرد. سپس عنوان مدرک به ترتیب طبیعی خود و به طور کامل در زیر این واژه یا به دنبال آن می‌آید. به این ترتیب برای هر واژه مهم یک مدخل ساخته می‌شود.

(لنکستر 1991):

همان نمایه کوئیک است با این تفاوت که کلیدواژه‌های آن در خارج از متن می‌آید.

(نیزو 1997):

نمایه‌ای که در آن هر واژه مهم در زنجیره‌ای از متن، به عنوان اصطلاح هدایت‌کننده یا نقطه دسترسی و به دنبال زنجیره کامل آمده است.

(راولی نقل در نیزو، 1982):

- 1- در نمایه کواک همه واژه‌ها که به صورت مدخل ظاهر می‌شوند، از عناوین مدرک استخراج می‌گردند
- 2- در این نمایه بعضی مدخل‌ها معمولاً اصطلاحات تک واژه‌ای هستند

مثال:

“مبانی نظری و توصیه‌های عملی تحلیل موضوعی و نمایه‌سازی”

تحلیل موضوعی

مبانی نظری و توصیه‌های عملی تحلیل موضوعی و نمایه‌سازی

توصیه‌های عملی

مبانی نظری و توصیه‌های عملی تحلیل موضوعی و نمایه‌سازی

مبانی نظری

مبانی نظری و توصیه‌های عملی تحلیل موضوعی و نمایه‌سازی

نمایه‌سازی

مبانی نظری و توصیه‌های عملی تحلیل موضوعی و نمایه‌سازی

نمایه‌سازی رایانه‌ای:

یکی از انواع پرکاربرد نمایه‌سازی که سرعت و یکدستی را به دنبال دارد، نمایه‌سازی ماشینی است.

نمایه سازی ماشینی : عبارتست از انتخاب واژه های کلیدی يك اثر بوسیله روش های ماشینی استفاده از کامپیوتر برای بیرون آوردن و نشان دادن واژه های نمایه بدون دخالت انسان در حالی که يك بار برنامه، و سیاست کار آن، به ماشین داده شده است.

نمایه سازی رایانه ای ، انجام کلیه مراحل نمایه سازی اعم از :

✓ -انتخاب و استخراج اصطلاحات نمایه ای از متن ،

✓ -مدخل آرایه و

✓ -ارائه جاینامه های هر مدخل و

✓ -چاپ نمایه توسط رایانه و بدون دخالت انسان را نمایه سازی رایانه ای گویند. نمایه سازی کامپیوتری، رایانه ای ، ماشینی و خودکار همگی با هم مترادفند و بجای یکدیگر به کار می روند.

نکات :

وقتی شناسه ها توسط نمایه ساز انتخاب می شوند نمایه سازی دستی گفته می شود ، اگر این کار را کامپیوتر انجام دهد، نمایه سازی خودکار یا ماشینی نامیده می شود. جهت استفاده از روش نمایه سازی ماشینی داده ها باید به صورت ماشین خوان درآیند. در این نوع نمایه سازی همه امور از انتخاب کلید واژه، شماره گذاری ، ترتیب بندی و غیره توسط کامپیوتر انجام می گیرد.

نمایه سازی که با اصول تحلیل و طراحی سیستم و برنامه سازی کاربردی آشنایی دارند، می توانند در زمینه

کامپیوتری سازی نمایه سازی موثر باشند

روال کار :

کامپیوتر مفاهیم مهم را که بارها تکرار شده اند و جزو کلمات غیر موضوعی زبان نیستند، به علاوه اسمهای اشخاص، مکانها و غیره را به عنوان کلید واژه در نظر می گیرد. این کار توسط يك نرم افزار به نام نرم افزار بسامدی استخراج می کند. مثلاً يك مفهوم که در ذیل يك بخش از متن بارها تکرار شده باشند، محتملاً کلید واژه است، مگر آن که جزو واژگان غیر موضوعی باشد. واژگان غیر موضوعی مانند(است، که ، را و...) توسط يك فهرست به نام فهرست ایستا مشخص و نادیده گرفته می شود. قسمت دیگر نرم افزار حاوی اسمهای خاص است. نرم افزار با استفاده از این دادگان ، اسمهای خاص متن را تشخیص می دهد و به عنوان کلید واژه در نظر می گیرد.

فرایند انتخاب کلید واژه درنمایه سازی خودکار :

1. شناسایی واژه های انفرادی از متن که تحلیل واژگان نامیده میشود.
2. برداشتن واژه های با بسامد تکرار بالا که در ارائه محتوای متن بی تأثیرند، با استفاده از فهرست واژه های غیرمجاز.
3. تبدیل واژه های باقی مانده به شکل ریشه آنها؛ یعنی حذف پسوندها یا پیشوندها تا هر کلمه تا حد ریشه اش کوتاه مبدع آن است. porter . شود
4. محاسبه رایانه ای بسامد ریشه هایی که در متن تحلیل شده اند، به منظور تعیین تابع ارزشگذاری هر ریشه.
5. ریشه هایی که ارزشگذاری بزرگتری دارند، برای متنی که در آن ظاهر شده، به عنوان کلیدواژه تعیین می شود.

شرکت زافتکس نوعی برنامه کامپیوتری با عنوان IDX طراحی کرده که به ارائه خدمات نمایه سازی ماشینی می پردازد. این نرم افزار از روش استفاده از واژه نامه بهره می برد. این نرم افزار امکان تبدیل واژگان به ریشه آنها را جهت بازیابی بعدی فراهم می کند علامتگذاری و محدودکردن واژه های ناخواسته را انجام م دهد. شکستن واژه های مرکب و ترجمه و انجام عمل ارجاع و مترادف سازی و ساخت عبارات را نیز انجام می دهد ماشین امکان تشخیص را تنها از طریق تطبیق واژه های استخراج شده از متن یا منتسب شده به متن با فهرستی که واژه های غیرمجاز نامیده می شود، به دست می آورد.

در اختیار داشتن فهرستی از این واژه ها و ارائه آنها به برنامه رایانه ای برای ممانعت از ورود آنها به فهرست واژه های مفهومی

مطلوب برای نمایه شدن، یکی از اقدامهای سودمند در نمایه سازی خودکار مبتنی بر کلیدواژه هاست.

روش های نمایه سازی ماشینی

الف. روش زبان شناختی: این روشها می کوشند با کمک تحلیل های شکل شناسی و ساختار نحوی مدرک توصیفگرها را استخراج نمایند.

- ✓ تجزیه و تحلیل ریخت شناسی که بر مبنای ریشه لغات عمل می کند.
- ✓ کلمات بدون بار معنایی موجود در سیاهه بازدارنده را حذف می کند.
- ✓ شکل های دستوری صرف کلمه را به یک شکل می آورد.
- ✓ ضمایر را بر اساس اسامی مربوط به آنها مرتب می کند.
- ✓ تجزیه و تحلیل نحوی که در سطح جملات امکانپذیر است و بر مبنای علامتهای نحو لغات انجام می شود.
- ✓ تجزیه و تحلیل معنا شناختی که در سطح مدارک مشترك در يك پایگاه صورت می گیرد . ارتباطات معنایی موجود

✓ در يك مدرک شناسايي مي شوند تا متون مشترك بتوانند به صورت واحدهاي هم معني تجزيه شوند.

ب. روش هاي آماری: مشخص مي کند که معني هر مفهوم منفرد در مدرک با حضور آن در جاگاههاي مختلف مدرک ارتباط تنگاتنگ دارد. بنابراین لغات درون متن شمارش مي شوند و ارتباط آنها ارزش گذاري مي شود. هدف آماری از اطلاعات آن است که لغات داراي بارمعنایي در مدرک به عنوان توصيفگر انتخاب شوند. روش هاي آماری عملاً براي بالابردن جامعيت به کار گرفته مي شوند. درحاليکه روشهاي زبان شناختي در جهت بهبود مانعيت کاردارند.

ج. روش هاي مبتني بر احتمالات: در اين روشها تئوري احتمالات براي مدل سازي رياضي مراحل بازيابي به کار گرفته مي شود. درحاليکه در توزيع آماری اصطلاحات يك مدرک مورد استفاده قرار مي گيرد. اين روش با عمليات رياضي مفروضات ساده و مطمئني را ارائه مي دهد. فرض بر آن است که مدارک بر اساس ميزان ربط در هنگام بازيابي مورد ارزشي قرار مي گيرند.

فهرست ايستا , يا واژگان غير مجاز

تحليل کلمات يك متن نشان مي دهد گروهي از کلمات بي اهميت وجود دارد که به فراواني در متن ظاهر مي شود (مانند

يك، به، نه، براي، با، چه کسي، چه موقع، است، آن). گروهي نيز وجود دارد که بندرت در متن مي آيند و ممکن است نشان دهنده محتوای اطلاعاتي متن نباشند

اين دسته از واژه ها به تنهائي بارمعنایي ندارند , بود يا نبود آنها نه تنها در پرسش کاربر تأثيري ندارد بلکه در ميزان ربط يا عدم ربط مدارک بازيابي شده نيز تأثيري ندارد. اين واژه ها با عنوان واژههاي غيرمجاز براي ورود به نمايه معرفي مي شوند.

مزيابي تهيه ليست واژه هاي غير مجاز

در صورتي که واژه هاي غيرمجاز قبل از فرايند نمايه سازي مدارک مشخص و فهرست آنها براي کنترل به رايانه داده شود، علاوه بر صرفه جويي در زمان و حجم بايگانيهاي نمايه، به ميزان زيادي از بازيابي مدارک نامرتب و ريزش کاذب در جستجو جلوگیری خواهد شد.

چند نمونه از نرم افزارهاي معروف نمايه سازي که در سيستمهاي کامپيوتر هاي شخصي کار مي کنند به

شرح زیر است:

INDEXING RESEARCH

محصول شرکت CINDEX

INDEX AID

محصول شرکت Santa Barbara Software Products

INQUIRY

محصول شرکت Indexer assistant

Indexit

محصول شرکت Norman Swartz

Bayside Indexing Service h

MARCSEX محصول شرکت

Newberry Library

NLCINDEX محصول شرکت

-Watch City Software

Windex محصول شرکت

طرح های موفق که در اجرای نمایه سازی ماشینی ارائه گردیده :

1. طرح AIMS از کتابخانه ملی پزشکی آمریکا
2. طرح AIR/X در دانشگاه فنی دارمشتات با موضوع فیزیک برای پایگاه اطلاعاتی فیزیک
3. طرح LISA و طرح COPSY توسط زیمنس اجرا شد .
4. 4PASSAT توسط زیمنس، سیستم نمایه سازی برای زبان آلمانی
5. Saphir در دانشکده پزشکی هاروارد، سیستمی برای یک سیستم بازیابی هوشمند اطلاعات در مدارک بیولوژی و پزشکی بر پایه روش های آماری و زبان شناسی
6. SPECIALIST توسط کتابخانه ملی پزشکی در آمریکا

عوامل موثر در نمایه سازی خودکار

1 - محدوده رکورد

اولین تصمیم گیری مهم برای تهیه هر نوع نمایه ، گزینش حد و حدود رکوردی است که واحد قابل جستجو را تعریف می کند. این تصمیم گیری در بازیابی کارآمد نقشی حیاتی دارد.

2 - محدوده اصطلاحات

تعیین حد و حدود یک واژه از دیگر مسائلی است که در نمایه سازی خودکار باید به آن توجه شود. در نظام های نمایه سازی دستی ، گزینش کلمات برای نمایه به سهولت انجام می شود. اما در نمایه سازی خودکار از آنجا که ماشین از هوشمندی لازم برای انتخاب کلمات برخوردار نیست بنابراین باید حدود کلمه را تعریف کرد. معمول حدود کلمات نمایه را با استفاده از علائم نقطه گذاری تعریف می کنند. به طور معمول ، فاصله بین کلمات و علائم دستوری و نقطه گذاری به عنوان مرز کلمات در نظر گرفته می شود. روش های تعیین حد و حدود کلمات در نمایه سازی خودکار ، بر اساس نوع برنامه و میزان پیشرفتگی آنها متفاوت است

کستر نمایه سازی خودکار را به دو دسته استخراجی و تخصیصی تقسیم می کند:

نمایه سازی استخراجی

ساده ترین روش نمایه سازی در پایگاه های اطلاعاتی ، روش نمایه سازی استخراجی است که در آن واژه ها برای قرار گرفتن در نمایه ، توسط رایانه از متن استخراج می شوند. در این روش عموماً بسامد تکرار واژه در هر رکورد یا مقاله تعیین شده و کلماتی که بسامد تکرار آنها زیاد است در متن نمایه قرار می گیرند.

نمایه سازی تخصیصی

انتسابی فرآیند پیچیده ای است که با استفاده از تحلیل های آماری ، کلمات و اصطلاحات به مدرک منتسب می شوند رایانه برای نمایه سازی از اصطلاحنامه یا کنترل واژگان بهره می گیرد.

انواع نمایه استخراجی

نمایه سازی با استفاده از فهرست کلمات ممنوعه : در این روش در هنگام نمایه سازی ، رایانه تمام کلمات متن را استخراج می کند، سپس کلمات ممنوعه را حذف و بقیه کلمات را در يك نظام الفبایی مرتب می کند. نمایه سازی بسامدی : در این روش، بسامد تکرار کلمات در هر رکورد مورد بررسی قرار می گیرند و براساس بسامد تکرار در فهرست کلمات نمایه قرار می گیرند.

نمایه سازی ریشه یابی : در بعضی از سیستم های نمایه سازی استخراجی ، از پسوند یا ریشه کلمات استفاده می شود. در این روش ریشه یا پسوند کلمات جایگزین مجموعه ای از کلمات هم ریشه یا پسوند مشترک می شود. الگوریتم های ریشه یابی مختلفی چون الگوریتم های موضوعی خاص مانند الگوریتم های پزشکی وجود دارند الگوریتم Lovins. پسوند 260

نمایه سازی بر اساس وزن دهی : در این روش ، کلمات بر اساس محل قرار گرفتن خود در متن (مثل عنوان ، چکیده و) ... امتیازدهی می شوند. حضور کلمات در بخش های مختلف رکورد، امتیازات متفاوتی دارد. در نظام های رایانه ای بیشتر از روش های نمایه سازی استخراجی استفاده می شود. یکی از عمده ترین مشکلات این نمایه ها به ویژه هنگام استفاده در پایگاه های اطلاعاتی ، عدم بازیابی اطلاعات به دلیل نبودن آن کلمه درخواستی در نمایه پایگاه اطلاعاتی است. دلیل این امر آن است که بهره گیران ، یا همه کلمات مترادف با اصطلاح موجود در درخواست را وارد نکرده اند و یا از مترادفات آن بی خبرند. بنابراین ، بسیاری از مدارک مرتبط از دست می روند. برای رفع این معضل طراحان پایگاه های اطلاعاتی توانایی های نمایه ای را با توانایی نرم افزاری در هم می آمیزند: یکی از روش ها، امکان نمایش نمایه و انتخاب واژه درخواستی توسط خود بهره گیر است.

روش دیگر، استفاده از نظام بازخورد مرتبط است . این روش به بهره گیران اجازه می دهد تا مدارک مرتبط را برگزینند. سپس از سیستم می خواهند تا با توجه به این مدارک ، مدارک مرتبط بیشتری را بازیابی نمایند. امروزه این روش در اینترنت و پایگاه های اطلاعاتی تمام متن کاربرد فراوانی دارد نمایه سازی تخصیصی

درواقع در نمایه سازی تخصیصی برای هر واژه "پرونده ای" از کلمات و عبارات مرتبطی که به نظر می رسد تهیه می شود. بنابراین می توان از برنامه ای رایانه ای برای انطباق عبارت های مهم در مدرک با این مجموعه پرونده ها استفاده کرد و در صورت انطباق واژه موجود در مدرک با واژه های موجود در پرونده های کلمات ، اصطلاح نمایه ای را انتخاب نمود موسسه Biosis . (15000 اصطلاح زیست شناسی)

- هر چه نمایه سازی تخصیصی تر باشد صرفه اقتصادی آن در يك سیستم خبره بیشتر می شود .

• سیستمی که از يك نمایه سازی تخصصی استفاده می کند می تواند بعنوان يك سیستم کمی خبره از آن نام برد. نتیجه تحقیق

پژوهشی را با موضوع نمایه سازی ماشینی متون فارسی براساس قانون زیف انجام دادند. نتایج نشان داد: « داورپناه و بلندیان »

✓ توزیع فراوانی واژگان در متون فارسی دارای الگوی پیش بینی پذیر است.

✓ کاربرد واژه‌های با بسامد بالا و بسامد پایین در مقاله های فارسی، از قانون زیف پیروی میکند.

✓ بسامد واژگانی می تواند به عنوان معیاری برای نمایه سازی ماشینی متون فارسی در نظر گرفته شود.

همخوانی کامل بین بسامد واژگانی و کلیدواژه های موضوعی در شیوة تفکیک صرفاً ماشینی بدون دخالت عامل انسانی به طور 21 % است. در شیوة تفکیک ماشینی با دخالت عامل انسانی، میزان همخوانی / متوسط در کل مقاله های مورد بررسی به میزان 50 به 52 % میرسد.

وضعیت همخوانی کامل بسامد واژگانی با کلیدواژه‌های عنوانی در شیوة صرفاً ماشینی بدون دخالت عامل انسانی، به طور متوسط 9 % است که در شیوة ماشینی با دخالت عامل انسانی این میزان بیشتر از 5 برابر شده و به 20 / 9 % در کل، مقاله های مورد بررسی 14/54 میرسد .

قانون زیف

جورج کینگزلی زیف استاد زبان شناسی دانشگاه هاروارد، در سال 1949 با آزمایش کلمات کتاب اولیس جیمز جویس به نتایجی در مورد کلمات و میزان تکرار آنها در متن رسید. نتایج او به این صورت بود که: هر کلمه با فراوانی (بسامد) \circ اگر تمام کلمات يك کتاب را بشماریم و از زیاد به کم مرتب کنیم به این نتیجه می رسیم که رتبه همان کلمه در متن رابطه معکوس \circ همان کلمه نسبت عکس دارد، یعنی تعداد بارهایی که هر کلمه در متن ظاهر می شود با رتبه دارد. این نسبت در کلمات کل متن برقرار است. که به قانون زیف معروف شده است. قرار دارد و 3 برابر \circ 1 قرار دارد دوبرابر بیشتر از کلمه ای در متن ظاهر می شود که در رتبه \circ بر طبق زیف کلمه ای که در رتبه 3 قرار دارد و همینطور تا آخر \circ بیشتر از کلمه ای ظاهر م ی شود که در رتبه او این قضیه را با اصل کمترین کوشش توجیه کرد. انسانها بر اساس این اصل تمایل دارند کارهای خود را به گونه ای ساده تر انجام دهند و در نوشتن متنی سعی دارند بیشتر از کلمات تکراری استفاده کنند. و به همچنین در هنگام صحبت کردن و سخنرانی سعی دارند کلمات کمتری را بیشتر تکرار کنند. آن در کل هر متن \circ حاصل ضرب فراوانی (بسامد) واژه در رتبه $k = r * f$: برقرار است که r و رتبه f این رابطه بین فراوانی عددی (تقریباً) ثابت است.

رابطه لگاریتمی آن شناخته شده تر است و کاربرد بیشتر دارد: $\log r + \log f = \log c$
 این رابطه به جز کلمات در بسیاری از دیگر محیطها از جمله جمعیت شهرها، میزان بازدید از صفحات اینترنت، شرکت ها و کارکنان آن و نیز در نمایه سازی خودکار و... استفاده می شود. از نظر محققان بسیار عجیب است که چطور و چرا همچنین ساده ای در بسیاری محیطهای پیچیده اتفاق می افتد. رابطه اما به هر حال قانون زیف بسیار ساده است چراکه خود زیف زبان شناس بود و به مسائل ریاضی چندان علاقه نداشت. و این فرمول نتوانست محیطهای خیلی پر تکرار را به درستی نشان دهد. بعد از زیف سه عدد ثابت به این فرمول اضافه شد و کمی محاسبه آن را انعطاف پذیرتر کرد. که به زیف مندلیبرت معروف است.

$$f = (r+m)^B$$

نمایه سازی وب

وب به عنوان یکی از جذاب ترین بخش های اینترنت کاربردهای فراوانی دارد و مجموعه ای است از صفحات به هم پیوسته که حاوی اطلاعات مفیدی در زمینه های موضوعی متفاوت است. همه دلایل مربوط به چرایی سازماندهی اطلاعات در محیط چاپی با شدت بیشتری در محیط الکترونیک صادق است.

در حال حاضر موتورهای کاوش، وب راتحت ضابطه درآورنده اندو با نمایه سازی صفحات، پاسخی برای پرس وجوی کاربران فراهم می آورند.

تحولات اینترنت در حوزه نمایه سازی :

اینترنت بر سرعت، دقت، هوشمندی، قدرت مشارکت، کاربرپسندی، جهانی بودن، جامعیت و چندزبانی نمایه ها افزوده است.

- حجم مجموعه

حجم اطلاعات موجود بر روی اینترنت و به طور خاص بر روی وب هر روزه در حال افزایش است. بدیهی است هر چه حجم

مجموعه افزایش یابد نمایه سازی آن دشوارتر شده و نیازمند استفاده از روش های کامل تر و پیشرفته تری است.

- تفاوت های دایره لغات

معمولاً در یک مجموعه یا پایگاه اطلاعاتی یکدستی - چه از نظر ساختار و چه از نظر موضوعی - وجود دارد. وب برخلاف پایگاههای اطلاعاتی، تقریباً از هیچ تجانس موضوعی و زبانی برخوردار نیست. علاوه بر این،

بهره گیران از اینترنت نیز بسیار متنوع اند. در يك تحقیق مشخص شد که احتمال استفاده یکسان دو فرد از اصطلاحی واحد در اینترنت در حدود 20 درصد است .

- راهبردهای کاوش

کاهش مداخله کاوشگران متخصص در بازیابی اطلاعات، تنوع بهره گیران ، فقدان آگاهی درباره شیوه و راهبرد کاوش در اینترنت ، تولید نمایه منسجم را دشوار ساخته است .

- رسانه های جدید ارتباطی در اینترنت

امکانات تازه ای که در وب پدید آمده مانند امکانات فرامتنی ، لینکها ، خدمات چندرسانه ای و چندزبانه بودن اطلاعات اینترنت ، بیش از گذشته ایجاد نمایه را با دشواری مواجه ساخته است . نکته مهم برای نمایه سازان این است که نحوه و ساختار ذخیره و بازیابی اطلاعات چگونه است و از چه الگوریتمی برای ذخیره و بازیابی اطلاعات در آنها استفاده شده است . تفاوتی موجود در قابلیت های جستجو و بازیابی اطلاعات و نحوه نمایش و رتبه بندی اطلاعات در هر پایگاه اطلاعاتی نشان دهنده تفاوت در به کارگیری الگوریتمهای مختلف است . سعی شده از روابط الگوریتمی جهت جایگزین نمودن با پردازش فکری انسان برای نمایه سازی ماشینی استفاده شود . یعنی کلید واژه ها بر اساس الگوریتم خاصی در نمایه سازی ماشینی انتخاب شوند . اکثر موتورهای کاوش و پایگاه های تجاری به دلیل ماهیت تجاری بودن ، الگوریتمهای نمایه سازی خود را به راحتی در اختیار کاربر قرار نمی دهند. هر يك با الگوریتمها و سیاستهای متفاوتی به مقوله نمایه سازی می نگرند.

تفاوت در نتایج بازیابی شده در موتورهای کاوش مختلف نشان بارزی از وجود تفاوت در الگوریتمهای نمایه سازی و یا پایگاه های داده آنهاست. حجم نمایه در بین موتورهای کاوش هنوز هم یکی از نکات مهم و اصلی رقابت در بین تولیدکنندگان موتورهای کاوش است گرچه امروزه تهیه بهترین نمایه و نه بزرگترین آن مورد توجه قرار می گیرد. نمایه های وب برای نشان دادن لیستهای منابع خود از پیوندهای فرا متن استفاده می کنند، آنها این امکان را دارند تا از طریق لینک ها صدها و بلکه هزاران منبع را در بر بگیرند. در محیط وب به دلیل حجم وسیع اطلاعات منتشر شده عملاً نمی توان به شیوه های دستی نمایه سازی متوسل شد. بنابراین نمایه سازی در وب به صورت خودکار و توسط موتورهای کاوش انجام می شود. موتورهای کاوش پایگاههای اطلاعاتی قابل جستجویی هستند که از طریق برنامه های کامپیوتری به شناسایی و نمایه سازی خودکار صفحات وب می پردازند. موتورهای کاوش برنامه های خودکاري هستند که هیچ گونه وابستگی به نیروی انسانی ندارند. نوع جمع آوری اطلاعات در اینترنت به دو گروه اصلی تقسیم میشود:

موتورهای جستجوگر یا Search Engine

فهرست های وب یا Web Directory

هر دو مورد اطلاعات را در اختیار کاربران قرار میدهند اما تفاوت اصلی آنها در روش جمع آوری اطلاعات است. در موتورهای توسط Web Directory جستجوگر اطلاعات توسط نرم افزار جمع آوری و طبقه بندی میشوند اما در فهرست های وب یا عوامل انسانی انجام میشود اما با توجه به رشد بسیار سریع اینترنت عملاً وب دایرکتوری ها کاربران چندانی ندارند

Yahoo directory

Dmoz.org

دلایل استفاده از راهنماهای موضوعی

- ✓ رده بنده ی دستی تضمین کننده ربط مدارك در موضوع می باشد
- ✓ کانون تمرکز ابتدایی فرایند جستجو انجام شده است
- ✓ با مرور حوزه موضوعی به جنبه های تخصصی تر به صورت خود کار دست می یابیم
- ✓ مدارك با موضوعات مشابه با هم گروه بندی شده و ارتباط مدارك انجام می گیرد
- معایب راهنماهای موضوعی
- ✓ مشکل در روز آمد شدن اطلاعات
- ✓ يك مدارك ممکن است زیر يك حوزه موضوعی قرار گیرد اما جستجوگر در حوزه موضوعی دیگر به دنبال آن باشد

- ✓ شناخت طبقه آغازین مرور ممکن است مشکل باشد

فناوری موتورهای جستجو از دو فرایند مجزا اما یکپارچه تشکیل شده است

- 1 – ایجاد نمایه ای از مدارك وب : وقتی دکمه جستجو فعال می شود ، این نمایه جستجو می شود (فایل مقلوب)
- 2 – مجموعه ای از مدارك در قالب سیاهه ای که بر مبنای میزان ربط با پرسش جستجو رتبه بندی شده نشان داده می شود

روش کار موتورهای جستجوگر

Spider , crawler, در موتورهای جستجوگر کار جمع آوری و طبقه بندی اطلاعات بر عهده نرم افزار است. این نرم افزارها نامیده میشوند. اسپایدرها همیشه مشغول کار هستند و اطلاعات را از گوشه و کنار جمع worm wanderer, Robot, آوری میکنند حال تفاوتی ندارد که این اطلاعات صفحات وب جدید باشند و یا اطلاعاتی

باشند که قبلا وجود داشتند اما حالا بروزرسانی شده اند. هدف تمامی موتورهای جستجوگر وب در واقع واحد است. همه آنها سعی در جمع آوری اطلاعات مورد نیاز برای کاربران خود را با بالاترین میزان دقت دارند. این مسئله که يك وب سایت جستجوگر چگونه اطلاعات را بهتر طبقه بندی کند بستگی مستقیم به نوع موتور جستجوی سایت و الگوریتم آن دارد همه جستجوها حتی ساده ترین آن ها ، نمایه موتور جستجو یا پایگاه داده در يك کامپیوتر دور را جستجو می کند
چه اقلامی نمایه سازی می شوند

- ✓ بزرگترین موتورهای جستجو هم قادر به نمایه سازی کامل وب نیستند.
- ✓ برخی از این موتورها نظام نمایه سازی تمام متن دارند و هر واژه موجود در متن به جز واژه های فاقد بار اطلاعاتی مانند حروف
- ✓ اضافه، ربط و تعریف را نمایه می کند .
- ✓ برخی دیگر، سرعنوانها، عناوین فرعی و فرابوندها را همراه 20 خط ابتدای متن و 100 کلمه ای که از بسامد بالایی برخوردار
- ✓ است، نمایه می کنند.

- ✓ برخی توصیف متنی از يك شکل را نمایه می کنند
- ✓ برخی از موتورهای جستجو در ازای هزینه ای ، صفحات را بی درنگ نمایه سازی می کنند
- ✓ برخی از اسپایدر ها در هنگام گردآوری صفحات پیوندهای فرامتن را دنبال می کنند که زمان نمایه سازی را افزایش می دهد
باز نمود نمایه در وب
نمایه وب به سه صورت نمایش داده می شود:

1. نمایش تیترا حرف

2. نمایش فرم

3. نمایش سنتی

الگوریتم موتورهای جستجوگر چیست؟

الگوریتم در واقع مجموعه ای از دستورالعمل های گوناگون است که ترتیب قرارگیری سایت ها را در موتورهای جستجوگر تعیین می کند. برای این که مشخص شود که کدام سایت ها باید در لیست نمایش نتایج جستجو در ابتدا نمایش داده شوند ، موتور جستجو ب اساس پارامترها عمل میکنند. پس در واقع دو مورد

الگوریتم و پارامترها در نوع نمایش نتایج جستجو موثر هستند و اما معماری هر دو کاملاً محرمانه است. موتورهای جستجو پیوسته الگوریتم های نمایه سازی خود را بهبود می بخشند. الگوریتم وزن دهی:

- ✓ پاره های اطلاعاتی را جدا کرده سپس به آن ها بر اساس معیارها مقدار و وزن داده و در نمایه قرار می دهند.
- ✓ وزن واژگان عامل مهمی در تعیین ربط مدارک یافت شده و نتایج جستجو به شمار می رود (رتبه بندی مدارک)
- ✓ در بخش FAQ یا راهنمای جستجو موتورهای جستجو، اطلاعاتی در خصوص معیارهای رتبه بندی و ربط دارند

معیارهای ربط یا وزن دهی

تحلیل پیوند clustering

خوشه سازی

ربط آماری

معیارهای ربط آماری

- ✓ بسامد تکرار واژه ها
- ✓ مکان ظهور واژه مثلاً " وجود واژه در برجسبهای عنوان
- ✓ طول مدرک (واژه در مدرک کوتاه وزن بیشتر در مقایسه با مدرک طولانی دارد)
- معیار تحلیل پیوندی شهرت (نخستین روش رتبه بندی در موتورهای جستجو)
- pagerank شمارش تعداد پیوند ها به صفحه از نمایه سازی استنادی گرفته شده است (هرچه به اثری بیشتر استناد گردد، اهمیت و اعتبار بیشتری دارد)

به کارگیری این معیار در گوگل مبنای I'm feeling lucky

معیار خوشه سازی

موتور جستجو بر اساس يك الگوریتم، اصطلاحات مشترك موجود در مدارک را در قالب خوشه های موضوعی مرتب می نماید.

(related terms)

Ask, clusty, google

معیار پیوندهای مرتبط

مشابه استفاده از " نیز نگاه کنید " به معنی مدارک بیشتری درست همانند این مدرک بدهید.

در گوگل Similar pages

قابلیتهای جستجو در موتورهای کاوش

جستجوی ساده جستجوی پیشرفته

عملگرهای بولی مجاورت

کوتاه سازی جستجو در فیلد

انواع موتور کاوش

موتور کاوش عمومی

yahoo, google ✓

موتور کاوش تخصصی: الگوریتم آنها مدارک مربوط به حوزه موضوعی را هدف قرار داده و توسعه می دهند .

Search-engine-index.co.uk ✓

citeseer.ist.psu.edu (کتابخانه دیجیتال حاوی مقالات داوری شده) ✓

بخشی برای جستجوی وب پنهان (مقالات مجلات) دارد scirus.com ✓

Techxtra.ac.uk ✓

ابر موتور کاوش: پرسش را همزمان به چند موتور کاوش و راهنما می دهند

Dogpile , mamma, metacrawler ✓

اینترنت همانگونه که ابزاری قدرتمند برای دسترسی به اطلاعات است به همان نسبت نیز ابزاری است که

اطلاعات گمراه کننده را اشاعه می دهد.

تهیه کنندگان hotbot

”ایجاد صفحاتی که مغرضانه مرورگرهای کاوش را گول زده و باعث می شوند تا صفحات نامرتب با کاوش را

بازیابی کنند یا امتیاز بالایی را به مدارک آنها اختصاص دهند بسیار رواج یافته است . تکرار هزاران باره کلمات

در بخش کلیدواژه ها یا اظهارنظرها , یا گنجاندن حجم زیادی از کلمات نامرئی در فونتهای ریز یا به رنگی

مشابه با رنگ پس زمینه صفحه مثال هایی از حقه زدنهایی است که عمومیت دارند . چنین حقه هایی باعث می

شود تا هات بات به دو روش به این صفحات امتیازات کمتری را اختصاص دهد : اول حجم این مدارک طولانی

تر است , دوم اگر هات بات يك روش حقه زدن متداول را تشخیص دهد به شدت امتیاز آن صفحه را کم خواهد

کرد ”

وب پنهان : بخشی از وب که توسط نرم افزار خزنده موتورهای جستجوی عمومی نمایه نمی شوند یا نمی توانند

نمایه شوند , شامل :

- ✓ مبتنی بر پایگاه اطلاعاتی html صفحات
 - ✓ کتابخانه های دیجیتال
 - ✓ pdf مانند فایل html قالبهای غیر
 - ✓ اطلاعات ناپایدار مانند خبرهای جاری، آگهی ها
 - ✓ آثار غیر تجاری مانند گزارشهای دولتی، خبرنامه ها
- موتورهای جستجو با افزایش قابلیت‌های جستجوی خود با چالش وب پنهان روبرو شده اند. از جمله گوگل اسکولار رویکردهای نمایه سازی خودکار دروب:
- ✓ محتوا محوری
 - ✓ معنا محوری

محتوا محوری: اغلب موتورهای کاوش حاضر از روش نمایه سازی بر مبنای کلیدواژه های متن استفاده می کنند. در این شکل فرایند نمایه سازی سه مرحله خواهد داشت: شکستن کلمات، تعدیل و حذف کلمات غیرموضوعی، استفاده از الگوریتم ریشه ساز جهت تولید ریشه های مفاهیم. در مرحله شکستن کلمات، داده هایی که به صورت رشته ای از کاراکترها هستند مورد بررسی قرار گرفته و حدود کلمات و فاصله میان آنها مشخص می گردد.

در مرحله تعدیل کلمات مزاحم، بزرگ نویسی، نقطه گذاری و مواردی از این دست مدیریت می شود. کلمات مزاحم در نمایه سازی کلماتی هستند که بار معنایی خاصی ندارند و تنها برای ایجاد پیوستگی و ارتباط در جمله ها به کار می روند. بعد از این مرحله از الگوریتمی جهت تولید ریشه ها و مفاهیم استفاده می شود. معنا محوری:

بعضی موتورهای کاوش رویکردی مکانیکی دارند و به مفاهیم، الگوها و کلیدهایی که به فهم مفاهیم می انجامد توجهی ندارند. در این ابزارها، جست و جوی واقعی صرفاً بر مبنای کلیدواژه هاست. استفاده از فهرست مترادفها و بهره گیری از جست و جوی فازی از راهبردهای مطرح شده جهت رفع مشکلات جست و جوی کلید واژه ای است. روشی که در اینجا مطرح است بهره گیری از نمایه سازی معنایی پنهان جهت بهبود مانعیت، جامعیت و رتبه بندی نتایج کاوش

است. نمایه سازی معنایی پنهان به کاربران این اجازه را میدهد که جست و جوی خود را به مفاهیم و نه فقط کلید واژه ها محدود کنند. در زمینه نمایه سازی خودکار دروب حرکت‌های جاری به سمت بهره گیری از داده های ساختار یافته و تحقق وب معنایی است، اما این حرکتها به طور کامل به انجام نرسیده است و هنوز مشکلات حل نشده فراوانی در این مسیر وجود دارد.

وب معنایی (Semantic web)

وب معنایی نسبت به وب يك انقلاب محسوب مي شود كه در آن اطلاعات قابل خواندن و تجزيه و تحليل توسط ماشين است در حالي كه صفحات وب كنوني را فقط انسان مي تواند بخواند وب معنایی اين اجازه را به مرورگرها و ديگر نرم افزارها مي دهد تا اطلاعات را خوانده به راحتی تجزيه و تحليل كنند فضايي جهاني از جنس محاسبات هوشمند ماشيني را مي توان تصور كرد كه در آن تمامي پايگاه هاي دانش (Knowledge bases) به صورتي معني گرا و با توانايي درك مفهومي همدیگر در کنار هم قرار خواهند گرفت. آینده وب نه فقط توسط انسانها قابل فهم است بلکه توسط ماشين ها نيز قابل درك و پردازش است.

در زیر سه تعريف مختلف از وب معنایی ارائه شده است :

- ✓ پروژه اي با هدف ايجاد رسانه اي جهاني براي رد و بدل كردن اطلاعات بصورتي كه براي كامپيوتر قابل فهم و پردازش باشد .
 - ✓ وب معنایی، شبکه اي از اطلاعات در مقياس جهاني است به نحوي است كه پردازش آنها توسط ماشين ها به سادگي امكان پذير است .
 - ✓ وب معنایی شامل داده هاي هوشمند وب است كه توسط ماشين ها قابل پردازش است از دیدگاه مدیریت نظام اطلاعاتي چهار ویژگی اصلي براي اطلاعات بازیابی شده قابل لحاظ است :
 - ✓ دقت، پیوستگی زمانی، بهنگام بودن و مرتبط بودن.
 - ✓ لذا نمایه سازي باید بر مبنای چهار ویژگی مذکور انجام پذیرد. بنابراین جهت رسیدن به نمایه سازي جامع و مانع وب به نمایه سازي دقيق موتورهاي جستجو نیاز داریم.
- در نمایه سازي همواره اين سؤال مطرح است كه يك نمایه خوب چگونه نمایه اي است؟
- نمایه مؤثر كاربر را به اطلاعات دقيق، بدون دشواري، خطا و موارد نامرتبط هدايت مي كند و به ندرت منجر به بازیابی اطلاعات
- سطحي مي شود. نمایه از ابزارهاي مهم تحليل مدرك است بنابراین كنترل كيفيت اين ابزار در نظام اطلاع رسانی مهم است هدف هر نمایه، بازیابی ركورد ها يا مداركي است كه به وسیله فرایند نمایه سازي، ذخیره و سازمان دهی شده اند. هدف از ارزیابی نمایه، تعیین میزان اثربخشي و كارایی آن است
- كلوند Cleveland معيارهاي ارزیابی در سه سطح ارزیابی کرده است:

1. فني

2. معنایی

3. میزان کارایی

در سطح فنی، نمایه باید دارای زبان مناسب و شکل قابل درک باشد و به سادگی بتوان از آن استفاده کرد
در سطح معنایی، واژه ها باید معانی را بدون ابهام منتقل کنند
سطح سوم، نمایه باید اطلاعات مرتبط را به درستی شناسایی کند و در بازیابی اطلاعات مؤثر باشد

معیارهای ارزیابی نمایه از دید لنکستر

1. دامنه آن

2. توانایی بازیابی گزینه های مورد نیاز (جامعیت)

3. توانایی پیشگیری از بازیابی گزینه های ناخواسته (مانعیت)

4. مدت زمان پاسخگویی نظام: از یک نظام نمایه سازی خوب، بازیابی سریع مدارک مرتبط انتظار می رود. اما باید به خاطر

داشت که زمان پاسخگویی به نوع پرسش جست و جو بستگی دارد.

5. میزان تلاشهای مورد نیاز کاربران: اگر دسترسی به مدارک آسان باشد، یعنی نظام نمایه سازی اثربخش است.

6. هزینه: در نظام نمایه سازی، هزینه در مقابل سودمندی های نظام، همیشه قابل قبول است. اما نظام نمایه سازی مؤثر حداکثرمزایا را با حداقل هزینه ها فراهم می کند.

کیفیت نظام نمایه سازی به عواملی مانند :

✓ دامنه پوشش، جامعیت، مانعیت، تازگی و صحت داده هابستگی دارد.

پراشر عوامل زیر را بر کارایی نمایه سازی مؤثر می داند:

✓ عوامل مخل

✓ ضریب ریزش

✓ ضریب تازگی و نوظهوری

✓ کل نگری و جزءنگری.

ضریب عوامل مخل، مکمل ضریب مانعیت است که تعداد مدارک نامرتب در کل مدارک بازیابی شده را نشان می دهد و از لحاظ ریاضی به صورت زیر نشان داده می شود:

$$100\% \times \text{تعداد کل مدارک نامرتب بازیابی شده} = \text{ضریب عوامل مخل}$$

تعداد کل مدارک بازیابی شده

هر چه ضریب عوامل مخل کم تر باشد، کارایی نظام نمایه سازی بیشتر است.

ضریب ریزش نشان می دهد که چه تعداد مدارک نامرتب از کل مدارک نامرتب پایگاه بازیابی شده اند و از لحاظ ریاضی فرمول آن به صورت زیر است:

$$100\% \times \text{تعداد کل مدارک نامرتب بازیابی شده} = \text{ضریب ریزش}$$

تعداد کل مدارک نامرتب موجود در پایگاه

ضریب تازگی و نوظهوری بخشی از مدارک جدید است که برای اولین بار مورد توجه کاوشگر اطلاعات قرار گرفته است. در میان کل مدارک مرتبط، درصد کمی از مدارک جدید وجود دارند که اطلاعات جدیدی ارائه کنند. کارایی بازیابی در نظام نمایه سازی با اصلاح زبان نمایه سازی و زبان پرس و جو افزایش می یابد. جزءنگری و کل نگری فنونی هستند که برای چنین اصلاحی مورد استفاده قرار می گیرند. سطح بالاتر کل نگری، جامعیت را افزایش و

مانعیت را کاهش می دهد. اما تعدیل جامعیت و مانعیت برای حفظ کارایی بهینه نمایه سازی الزامی است

ضعف در نمایه سازی

ضعف در تحلیل مفهومی :

1 - ضعف در تشخیص موضوعی که مورد علاقه جامعه بهره گیر است

2 – تفسیر نادرست از جنبه هایی که مدرک واقعا در مورد آن بحث کرده است .

ضعف در ترجمه :

1 – کوتاهی در استفاده از اصطلاح اخص تر برای بعضی از موضوعات

2 – استفاده از اصطلاحی نامناسب با محتوای موضوعی

استانداردهای حوزه نمایه سازی

ANSI /Z. 39.4 معیارهای اساسی برای نمایه ها

ANSI/NISO Z.39.50 بازیابی اطلاعات: کاربرد تعریف خدمات و پروتکل

BS .3700 آماده سازی نمایه برای کتابها، نشریات ادواری و دیگر مدارک

BS 6529 1984 بررسی مدارک، تعیین موضوعهای آن ها و انتخاب اصطلاحات نمایه ای

ISO 5963 تحلیل مدارک، تعیین موضوعها، انتخاب اصطلاحات نمایه ای

ISO 999 دستور عملهایی برای محتوا، ساختار و نمایش نمایه ها

1. دامنه زیر پوشش نمایه کامل باشد؛

2. در انتخاب اصطلاح، انسجام داشته باشد؛

3. اصطلاحات انتخاب شده با سطح کاربران متناسب باشد؛

4. ارجاعات به میزان کافی وجود داشته باشند؛

5. زنجیره های بسیار طولانی از سرعنوانهای فرعی به هم پیوسته، وجود نداشته باشد؛

[/sra.blogsky.com9h](http://www.sra.blogsky.com9h)

Page | ۴۹

6. سرعنوانهای فرعی به درستی سرعنوانهای اصلی را نشان دهند؛

7. مکان نمایی نادرستی وجود نداشته باشد؛

8. زنجیره بسیار طولانی از مکان نماها وجود نداشته باشد؛

9. نظم الفبایی در متن، یکپارچه و صحیح باشد؛

10. اشتباه املائی وجود نداشته باشد؛

11. ارجاع نادرست و دوطرفه وجود نداشته باشد.

اصطلاحنامه

گنجواژه یا اصطلاحنامه، مجموعه اصطلاحات يك رشته است که میان آنها روابط معنایی، رده ای، و سلسله مراتبی برقرار شده و توانایی آن را دارد که موضوع آن رشته را با همه جنب‌های اصلی و فرعی و وابسته، ب‌ه‌گونه‌ای نظام‌یافته و ب‌هم‌منظور ذخیره و بازیابی اطلاعات ارائه دهد. اصطلاحنامه معادل فارسی واژه انگلیسی Thesaurus است

اصطلاحنامه: از نظر وظیفه و کارکرد، ابزار کنترل واژه‌ها به منظور برگرداندن زبان طبیعی مدارک به زبان مفید است. از نظر ساختار، واژگان کنترل شده و پویای زمینه‌ای خاص از دانش بشری است که برای ذخیره و بازیابی اطلاعات آن حوزه به کار می‌رود.

اهداف اصطلاحنامه

اصطلاحنامه دارای هدف‌های اساسی زیر است:

نمایاندن ساختار زمینه معینی از دانش چنان که هم‌نمایه ساز و هم‌جست‌وجوگر بتوانند از گستره آن زمینه و ارتباط میان مفاهیم آن با اندیشه‌های مرتبط آگاهی یابند

- ✓ ارائه اصطلاحات استاندارد در زمینه‌ای معین
- ✓ برقراری نظام ارجاعات میان اصطلاحات و رده‌بندی اصطلاحات به صورت سلسله‌مراتبی
- ✓ تأکید بر توجه به نیازهای اطلاعاتی استفاده‌کنندگان
- ✓ تعیین اصطلاحات مجاز و مشخص کردن حدود معانی اصطلاحات به منظور ایجاد هماهنگی در نمایه‌سازی

روابط میان اصطلاحات

باید توجه داشت که ویژگی ذاتی اصطلاحنامه، توانایی تعیین و نمایش روابط معنایی میان واژه‌هاست و يك رابطه يك‌سویه نداریم و رابطه همه‌جانبه وجود دارد. این روابط ممکن است یکی از این سه نوع باشد:

الف) رابطه هم‌ارزی Equivalence Relation:

میان اصطلاح پذیرفته شده و اصطلاح پذیرفت‌هنشده برقرار می‌شود

مثال: اصطلاح گیاهان به جای نباتات

ب) رابطه سلسله‌مراتبی Heirachial Relation

بیان‌کننده رابطه اعم و اخص میان مفاهیم است که در واقع، اصطلاحنامه‌ها را از واژه‌نامه‌های متداول متمایز می‌کند.

- گیاهان ا.خ درختان
- صنوبر ا.ع درختان

رابطه همبسته یا همایند: Associative Relation

رابطه میان دو اصطلاح که به دلیل وابستگی معنایی، وجود یکی دیگری را نیز به ذهن متبادر م‌یکند. مثال: دو اصطلاح اسب و سوارکاری

- ✓ پشتیبانی اصطلاحنامه‌ها بسیار گران است و به دانش خاص نیاز دارد
- ✓ اصطلاحنامه‌ها نقش مهمی در نظامهای ذخیره و بازیابی اطلاعات دارند.
- ✓ ظهور وب، همراه با توسعه و پیشرفتهای اخیر در کاربرد اصطلاحنامه‌ها به عنوان ابزارهای بازیابی اطلاعات، باعث تولد نسل جدیدی از اصطلاحنامه‌ها شده است.
- ✓ مشکلات روشهای آماری و زبانشناسی رایانه‌ای باعث شده است که این پروژه‌ها فکر استفاده از ابزارهای دیگری را در سر بپرورانند.
- ✓ اصطلاحنامه‌های وب محور، راه خود را به محیطهای بازیابی و سازماندهی اطلاعات وب محور باز نموده و در تهیه ابر داده‌ها، نمایه‌سازی صفحات، سایت‌های وب، پایگاههای داده و موتورهای جستجو استفاده می‌شوند.
- ✓ ساختارهای معنایی موجود در اصطلاحنامه‌ها می‌توانند هم در سازماندهی و هم در بازیابی اطلاعات وب و منابع دانش نقش داشته باشند.
- ✓ اصطلاحنامه‌های پیوسته یا به کمک رایانه تدوین و بازنمایی می‌شوند یا به صورت رایانه‌ای تولید می‌گردند. در این دو رویکرد کلی تدوین، حرکت از ذهنیت به عینیت است، بدین ترتیب که روابط اصطلاحنامه‌ای ممکن است توسط متخصصان تنظیم شوند و از رایانه برای بازنمایی آن‌ها استفاده گردد که در این صورت، رسوخ ذهنیت در تدوین اصطلاحنامه، بیشتر است. ممکن است روابط به وسیله استدلال‌گرهای خاصی که مبتنی بر قواعد ریاضی و جبر هستند تولید شوند که در این صورت، رسوخ ذهنیت در تنظیم روابط بسیار کاهش می‌یابد، ولی انتخاب اصطلاحات همچنان به وسیله انسان صورت می‌گیرد. در روش‌های جدید تولید خودکار اصطلاحنامه از آمارهای متنی و الگوریتم‌های پیچیده برای انتخاب اصطلاحات و تشخیص روابط اصطلاحنامه‌ای و خوشه‌بندی لغات استفاده می‌شود. رسوخ ذهنیت محدود به تدوین الگوریتم‌های تولید اصطلاحنامه است و پس از این مرحله، اصطلاحنامه به صورت خودکار و بدون دخالت انسان تولید می‌شود. هر دوی این دورویکرد کلی می‌توانند معایب و مزایایی داشته باشند که با توجه به حوزه‌ها و موقعیت‌های مختلف، باید رویکرد مناسب برای تدوین اصطلاحنامه انتخاب شود. اصطلاحنامه‌هایی که به صورت خودکار تولید می‌شوند ممکن است با نیازهای کاربران و اصطلاحات مورد نظر آن‌ها برای انجام جست‌وجو و بازیابی فاصله داشته باشند و نیازهای خاص آنان لحاظ نشده باشد. همچنین برای حوزه‌های فعالیت کوچک، مقرون به

صرفه نیستند . در عوض سرعت تولید، یکدستی روابط اصطلاحنامه ای، قابلیت روزآمدسازی، ویرایش سریع و کارآمد، و بازنمایی مناسب این نوع اصطلاحنامه ها در وب، از نقاط قوت آن ها محسوب می شوند.

آنتولوژی

✓ به معنای هستی شناسی است. (onto+logy) آنتولوژی برگرفته از ترکیب یونانی

✓ هستی شناسی : فهم چگونگی و علت امور در دنیای بیرونی

✓ همه علوم بالاخص علوم پایه به دنبال شناخت هستی اند

✓ اطلاعات و منابع موجود در وب بصورت فزاینده ایی رو به رشد هستند و استفاده کنندگان وب نیازمند یک درک مشترک از

✓ آنها دارند.

✓ آنتولوژی نقش اصلی را در مبادله اطلاعات و توسعه وب لغوی بسمت وب معنایی دارد.

✓ آنتولوژی یک مدل مفهومی است که موجودیتهای واقعی در یک دامنه خاص و روابط بین آنها را به صورت

صریح و رسمی

✓ مدلسازی می کند.

در سال 1980 ، مجمع هوش مصنوعی از لغت آنتولوژی برای دو منظور استفاده کرد:

✓ نظریه ای در مورد جهان مدل شده

✓ مؤلفه ای از سیستم های دانش.

آنتولوژی در هوش مصنوعی و همچنین علوم کامپیوتر به مجموعه ای از لغات و فرضیات (عموماً در منطق

مرتبه ی اول) گفته می شود که با توجه به معنی آن لغات ایجاد شده اند و به منظور توصیف یک واقعیت خاص

طراحی شده اند. امروزه آنتولوژی در هوش مصنوعی، مهندسی نرم افزار، مهندسی سیستمها و معماری

اطلاعات کاربرد دارد.

وب معنایی یا وب 3 که امروزه از مباحث تازه و جذاب در حوزه کامپیوتر و ارتباطات است، موجب می شود که

جستجو در اینترنت برای کاربران سریعتر و راحتتر گردد. در این خصوص عامل اصلی ظهور و موفقیت و

بمعنایی و آنچه که در پس این فناوری قرار دارد، آنتولوژی می باشد. در دهه اخیر این واژه و کاربرد آن در

حوزه فناوری اطلاعات و ارتباطات وارد شده و نقش مهمی در تبدیل وب به وب معنایی دارد .

تفاوت آنتولوژی در فلسفه و کامپیوتر

✓ در فلسفه، از نظم و ترتیب میان مفاهیم به آنتولوژی میرسیم اما در علوم کامپیوتر، آنتولوژی را از روی ترتیبی

که خود برای

مفاهیم در نظر میگیریم، استخراج میکنیم.

✓ نگاه آنتولوژی در فلسفه نگاهی جامع و جها نشمول است، درحالیکه آنتولوژی در کامپیوتر دارای دامنه بسیار کوچکتري است.

آنتولوژی در وب

فرض کنید که می خواهید در مورد موضوعی با کسی صحبت نمایید. برای اینکه طرف مقابل، حرف شما را کامل و درست متوجه شود، احتیاج است که حوزه بحث کاملاً مشخص باشد. واژه‌های مختلف در حوزه‌های گوناگون، معانی یا تعبیرهای متفاوتی دارند و حتی گاهی ممکن است که با وجود مشخص بودن حوزه بحث، یک کلمه خاص در ذهن افراد مختلف دارای تفاوت‌های اندکی باشد. واضح است که تنها راه‌هایی از چنین وضعیتی، یک مجموعه واژگان مشترک بین افراد است. آنتولوژی در وب معنایی دقیقاً چنین کاربردی را دارد. در هر آنتولوژی، تمام موجودیتهای یک حوزه به صورت کامل و با ذکر تمام ویژگیها فهرست می‌شوند. بعد از اینکه موجودیتهای حوزه را شناسایی کردید، باید ارتباطات بین آنها را نیز بیان نمایید.

مرحله بعدی این است که تمام اطلاعات فوق را با یک فرمت خاص درون مستندات اینترنتی قرار داده و اطلاعات موجود در آن مستند را به آنتولوژی اتصال دهید..

یک آنتولوژی، لغات و مفاهیمی را که در تعریف و نمایش محدود ه ای از دانش به کار میروند، تعیین کرده و بنابراین معانی را استاندارد میکند.

آنتولوژی توسط مردم، پایگاه‌های داده و برنامه‌های کاربردی که نیاز به اشتراک‌گذاری اطلاعات یک دامنه خاص را دارند، به کار برده می‌شود. آنتولوژی در وب معنایی واژه‌ها و ارتباطات بین آنها را در دامنه‌ای که استفاده می‌گردند، نشان می‌دهد.

عناصر اصلی تشکیل دهنده آنتولوژی عبارتند از:

1 - مفاهیم

2 - ارتباط بین مفاهیم

3 - خصوصیات آنها

به عبارت دیگر آنتولوژی ارتباط بین مفاهیم در اسناد وب و دنیای واقعی را مشخص میکند که با این کار اسناد مربوطه توسط ماشینها قابل پردازش و فهم می‌شوند و اشتراک‌گذاری بین عاملها را تسهیل مینماید. در واقع میتوان گفت:

Vocabulary+Structure=Taxonomy ✓

Taxonomy+Relationship و Constraints&Rules=Ontology ✓

✓ $\text{Ontology} + \text{Instance} = \text{Knowledge}$

✓ واژگان + ساختار = طبقه بندی

✓ طبقه بندی + ارتباط و محدودیت ها و قوانین = هستی شناسی

✓ هستی شناسی + نمونه = دانش

5 مرحله در طراحی هستی شناسی ها :

1 - تعیین هدف و دامنه هستی شناسی؛

2 - طراحی هستی شناسی در يك فرآیند سه مرحله ای شامل:

✓ گردآوری هستی شناسی (تعیین و تعریف مفاهیم و روابط اصلی)؛

✓ کدگذاری هستی شناسی (به کار گرفتن واژ ههای اصلی برای هستی شناسی (رده، موجودیت، رابطه)؛ انتخاب

يك زبان

بازنمون؛ نوشتن کد)؛

را به عنوان زبان نشان هگذاری معنایی به منظور انتشار OWL کنسرسیوم وب جهانی در نوامبر 2002 ،

زبان با تعریف کلا سها، نمونه ها و روابط به طور OWL . است RDF هستی شناسی های وب پیشنهاد کرد.

این زبان بر مبنای امکانات RDFS و XML،RDF نسبت به OWL . واضح و رسمی در توسعه و ساخت

هستی شناسی ها به کار می رود بیشتری برای بیان مفاهیم و معانی دارد و به دلیل قابلیت نمایش محتوای

میانکش پذیر رایانه ها در وب، برتر از سایر زبان ها است که هر يك ویژگی های خاص خود را OWL Full

و OWL DL،OWL Lite است. این زبان دارای سه زبان فرعی دارند و برای گروه خاصی از کاربران

طراحی شده اند.

✓ یکپارچه سازی هستی شناسی های موجود؛

3 - ارزیابی هستی شناسی؛

4 - مستندسازی؛

5 - ارائه راهنماها و دستورالعمل هایی برای هر يك از مراحل قبل

آنتولوژی به عنوان ابزاری قدرتمند برای نمایش و بیان دانش مربوط به يك حوزه، در قالبی رسمی و قابل

پردازش توسط ماشین مطرح است.

به کمک آن میتوان ارتباط بین سیستمهای ناهمگون را برقرار کرد و تعامل و ارتباط متقابل بین برنامه ها، ماشینها

و سیستمهای ناهمگون را بهبود بخشید.

موتور جستجوی وب معنایی، برای جستجوی هستی‌شناسی‌ها، مدارک، واژه‌ها و داده‌های منتشر شده در وب. این موتور جستجو دارای سیستم جستجوی مدارک RDF, Html میباشد

نمایه‌سازی معنایی پنهان

چکیده

این مقاله به معرفی و توصیف روش نمایه‌سازی معنایی پنهان (ال. اس. آی) می‌پردازد که یکی از روش‌های نوین نمایه‌سازی خودکار است. ابتدا در مقدمه‌ای کوتاه به نمایه‌سازی و چالش‌های آن اشاره می‌شود سپس در بخش دودل‌های فضایی برداری که نمایه‌سازی معنایی پنهان یکی از گسترش‌های آن است، توصیف می‌شود. در بخش بعدی ضمن تشریح مفهوم نمایه‌سازی معنایی پنهان، کاربردها و موارد استفاده از آن بیان می‌گردد و سپس مبانی ریاضی آن که روش آماری تجزیه مقادیر منفرد است با مثالی تشریح می‌شود. در بخش بعدی فرآیندکار نمایه‌سازی معنایی پنهان همراه با مثال توضیح داده می‌شود و در نهایت طرح‌ها و برنامه‌هایی که هم‌اکنون در این زمینه اجرامی شوند معرفی و به بعضی پیشرفت‌های فناورانه موثر در بهبود عملکرد نمایه‌سازی معنایی پنهان اشاره می‌شود.

کلید واژه‌ها:

نمایه‌سازی معنایی، نمایه‌سازی پنهان، بازیابی اطلاعات، تجزیه مقادیر منفرد

امروزه نمایه‌سازی و به‌طور کلی حوزه بازیابی اطلاعات، به سبب تغییرات قابل توجه در محیط پیرامون خود متحول شده‌اند. می‌توان بعضی از این تغییرات را این‌گونه برشمرد گسترش قابل توجه دامنه‌های بازیابی اطلاعات با ظهور چند رسانه‌ای‌ها، اینترنت، و اطلاعات جهانی ظهور پایگاه‌ها به داده و بعضی نظریه‌های جدید هوش مصنوعی، زبان‌شناسی و ریاضیات وویکرد تحلیلی مباحث ریاضی و پیچیده به کار رفته در بازیابی اطلاعات طراحی نظام‌های بازیابی اطلاعات همانند نظام‌های پایگاه داده‌ای رابطه‌ای، و بروز امکانات جدید ناشی از فناوری‌های نوینی چون وب مانند پالایشی اشتراکی (که با عنوان نظام‌های توصیف هاف یا شخصی شناخته می‌شود). علاوه بر چالش‌های فوق باید به این نکته اصلی توجه داشت که در بازیابی اطلاعات مفهوم نسبتاً مبهم و گنگی به نام ربط وجود دارد که به روش‌های پیچیده تشخیص نیت کاربر و ماهیت مدرکی بستگی دارد. طبق گفته پاپادیمتریو و همکارانش نظریه‌های اندکی در این زمینه ارائه شده است برای درک بیشتر این

چالش بررسی متون کلاسیکی چون آثار ریجسبرگن و سا لتون مفید خواهد بود باید به یاد داشت که چالش های جدید، علاوه بر چالش های پیشین چون ذهنی بودن فرایند نمایه سازی است که ناشی از نظری بودن این عمل است. اصلی ترین روش تعیین ذهنی بودن نمایه سازی، بررسی یکدستی آن هنگام تحلیلکار چند نمایه ساز از یک مدرک یا حتی یک نمایه ساز در زمان های مختلف است. همه این عوامل نشان دهنده پیچیدگی قابل توجه فرایند نمایه سازی است که در مقاله اندر سون و پرز-کاربالو به طور مفصل تشریح شده است.

وجود این چالش ها موجب شد که پژوهشگران از دهه 1950 تاکنون به دنبال روش های نوین و کارآمد برای رویا رویی با چالش ها باشند و از دهه 1970 تاکنون برای رفع این مسئله اهتمام بسیار داشته اند. همان طور که پالگدرین و کیل ه لیو در پژوهش خود نشان دادند که بیش از 800 اثر پژوهشی طی سال های 1956 تا 2000 در مورد نمایه سازی خودکار، نیمه خودکار و رایانه ای نگاشته شده اند

با توجه به نکات فوق، نوعی توافق عمومی وجود دارد مبنی بر اینکه بازیابی سنتی اطلاعات نمی تواند جوابگوی این چالش ها باشد. نتایج بازیابی نظام های سنتی بازیابی اطلاعات به دو دلیل عمده ناکارا و غیر دقیق بود. اول اینکه در این نوع نظام ها مفهوم واحدی را می توان به روش های مختلف توصیف کرد. عبارت های پرسش کاربر ممکن است در مدرک مربوط وجود نداشته باشد.

دوم اینکه بیشتر کلمات بیشتر از یک معنا دارند در نتیجه مطابقت واژگانی عبارت های پرسشی کاربر ممکن است به بازیابی مدرکی نامربوط بینجامد دانشمندان اصلی ترین علل محدودیت مدل ها و نظام های بازیابی اطلاعات را ابهام و گویانبودن واژه ها، ناکارآمدی بازنمون مدارک مجموعه، و در سویی دیگر اغتشاش و عدم صراحت پرسش های کاربر می دانند. راه حل های گوناگونی برای این مسائل پیشنهاد شده است که بازیابی اطلاعات براساس هستی شناس یا استفاده از بازنمون معنایی مدارک و پرسش ها از مهم ترین راه حل ها هستند در این مقاله سعی می شود نمایه سازی معنایی پنهان که نوعی بازنمون معنایی مدارک و نوع گسترش یافته مدل های بازیابی برداری است و از مباحث جبر خطی و ماتریسهای ریاضی و فن تجزیه مقادیر منفرد استفاده می کند، معرفی و ابعاد مختلف استفاده از آن بیان شود

مدل فضا برداری

از آنجا که نمایه سازی معنایی پنهان یکی از راهکارهایی بود که برای رفع مشکلات مدل فضای برداری به وجود آمد، در این بخش سعی می شود خلاصه ای از این مدل و مشکلات آن بیان شود.

مدل فضایی برداری یکی از چند روش تشخیص تشابه میان دو مدارک است که در سال 1975 به وسیله سا لتون گسترش یافت و چارچویی تأثیرگذار و قدرتمند برای ذخیره، تحلیل و ساختن مدارکی است. این مدل از ابتدا به منظور بازیابی اطلاعات گسترش یافت و در نمایه سازی مدارکی مبتنی بر فراوانی عبارت ها، به طور گسترده ای به کار رفت. سه مرحله این مدل عبارتند از نمایه سازی مدارک، وزن دهی عبارت، محاسبه ضریب تشابه:

1. *نمایه سازی مدارک:* هر مدارک یا (پرسش) به صورت یک بردار در فضایی با ابعاد بزرگ ترترسیم می شود. تعداد عبارت های منحصر به فرد مجموعه مدارک محاسبه می شود. واژگان غیرمهم در بردار مدارک حذف می شوند. در سیاهه واژگان غیرمجاور که واژگان عمومی را در خود جای می دهد، برای حذف واژه های پرکاربرد استفاده می شود (که در کل 40 تا 50 درصد کل واژگان مدارک حذف می شوند).

2. *وزن دهی عبارت:* وزن دهی برای نشان دادن میزان اهمیت عبارت ها در باز نمون مدارک انجام می شود. پیمایش فرض اغلب روش های وزن دهی چون فراوانی معکوس مدارک این نکته است که اهمیت یک عبارت با افزایش میزان رخداد آن عبارت متناسب است. برای پیشگیری از بازیابی مدارک طولانی تر از شیوه هنجار سازی استفاده می شود (زیرا با توجه به نکته فوق، احتمال بازیابی مدارک طولانی بیشتر از احتمال بازیابی مدارک کوتاه تر خواهد بود).

3. *محاسبه ضریب تشابه:* تشابه میان دو مدارک (یا میان یک پرسش و یک مدارک) با فاصله میان بردارها در فضایی با ابعاد بزرگ تر تعیین می شود. بدین صورت که همپوشانی واژه، نشان دهنده تشابه خواهد بود. شناخته شده ترین مقیاس تشابه، ضریب کسینوس است و تشابه میان دو مدارک را با کسینوس زاویه میان دو بردار نشان می دهند

در این مدل هنگام ورود یک پرسش توسط کاربر، آن پرسش مانند دیگر مدارک فرض می شود و به صورت یک بردار نشان داده می شود. نظام در فرایند بازیابی، موقعیت پرسش را نسبت به محل هر مدارک در فضای برداری مقایسه می کند و مدارک را با توجه به میزان تشابه با پرسش کاربر رتبه بندی می کند. پس به این ترتیب یا تعداد مدارکی را که بیشترین تشابه را با پرسش دارند بازیابی می کند یا همه مدارکی را که تا حدی از آستانه تشابه (که قبلاً معین شده است) بیشتر باشد، بازیابی می کند. اصلی ترین مزیت مدل فضای برداری، رتبه بندی کارآمد و دقیق مدارک مرتبط با توجه به میزان تشابه آنها با پرسش است، در صورتی که در روش های مبتنی بر مطابقت واژگانی یا هیچ طرحی برای رتبه بندی وجود ندارد یا اگر هم وجود داشته باشد فاقد امکان رتبه بندی موارد پیچیده هستند. چنانکه ممکن است به واژه ای که در اول عبارت پرسش ظاهر شده است رتبه ای بالاتر اختصاص داده شود.

کاربران در سیاهه مدارک بازیابی شده، همیشه مدارک با بیشترین میزان تشابه را مطالعه می کنند به این دلیل که از نظر معنایی بیشتر به پرسش مربوط هستند. با وجود این سنجش میزان تشابه میان بردارها، به دستورالعملی برای چگونگی وزن دهی هر عبارت نمایه نیاز دارد.

روش های مختلفی برای این امر وجود دارد که توابع دودویی، فراوانی عبارت و لگاریتمی بیشتر از بقیه شناخته شده اند

البته باید توجه داشت که استفاده از این مدل مشکلاتی را نیز دارد. اصلی ترین مشکل مدل های فضای برداری، عدم مطابقت واژگان است. مدل های فضای برداری، با عبارت های نامشابه مانند اقلام نامرتب برخورد می کنند. برای مثال رایانه و لپ تاپ اگرچه عبارت های مرتبطی هستند اما مدل های فضای برداری از کشف چنین رابطه ای عاجز هستند. پس اگر میان پرسش و مدارک در مجموعه متون هیچ کلمه

مشترکی وجود نداشته باشد حتی اگر بعضی مدارک با پرسش مرتبط باشند، مقدار تشابه عملاً منجر خواهد بود و در نتیجه هیچ مدرکی بازیابی نخواهد شد. دومین اشکال این مدل ها در مجموعه مدارک بزرگ پیس می آید که در صورت تشکیل ماتریس عبارت مدرک، ماتریس حاصل بزرگ و پراکنده خواهد بود که فضای ذخیره زیادی را می طلبد و مدت زمان پردازش و محاسبه آن نیز طولانی می شود. برای رفع این مشکلات از یکی از شاخه های این مدل با عنوان نمایه سازی معنایی پنهان استفاده می شود که در سال های اخیر رواج یافته است

به طور کلی می توان گفت که بازیابی اطلاعات سنتی با دو مشکل قدیمی متران ها (مانند مدارک حذف شده مربوط به اتومبیل به هنگام پرسش درباره ماشین) و چند معنایی ها (مثل بازیابی مدارک درباره اینترنت به هنگام پرسش درباره موج سواری) روبرو است. در این مدل هیچ تفاوتی میان crane به معنای پرند ماهی خوار crane به معنای جرثقیل نیست. برای مقابله با این دو مشکل و سایر مشکلات بازیابی سنتی اطلاعات، سعی می شود مدارک (و پرسش ها) نه با استفاده از خود عبارت ها (همان طور که در روش های برداری معمول است)، بلکه با استفاده از مفاهیم ضمنی (پنهان) آن عبارت ها، بازنمون شوند. این ساختار پنهان، نقشه ای ثابت میان عبارت ها و مفاهیم نیست بلکه به کل مدارک مجموعه و رابطه عبارت باکل آن بستگی دارد

نمایه سازی معنایی پنهان

نمایه سازی معنایی پنهان یکی از فنون نمایه سازی مفهومی است که برای غلبه بر مشکلات ناشی از عدم مطابقت واژگان به وجود آمده است. همان طور که گفته شده در نظام های سنتی بازیابی، اطلاعات از مطابقت واژه به واژه عبارت های مدارک با پرسش کاربر بازیابی می شود، ولی شواهدناکارآمدی این روش را نشان داده اند. همچنین از آنجا که معمولاً روش های متعددی برای بیان یک مفهوم وجود دارد (ترادف) احتمال دارد که میان عبارت ها و واژگان پرسش کاربر با بازنمون واژگانی مدرک هیچ اشتراکی وجود نداشته باشد. علاوه بر آن بیشتر واژه ها چندین معنا دارند (چندمعنایی) به طوری که نتیجه بازیابی با عبارت های پرسش کاربر، مدرکی نامربوطی را دربرخواهد گرفت. رزاریو معتقد است روشی که کاربر را قادر به بازیابی اطلاعات براساس مفهوم یا معنای یک مدرک خواهد نمود و مشکلات ذکر شده را نیز برطرف می کند، نمایه سازی معنایی پنهان است

تعاریف مختلفی از نمایه سازی معنایی پنهان ارائه شده است. رزاریو آن را این گونه تعریف می کند: "نمایه سازی معنایی پنهان، فنی است که پرسش ها و مدرکی را به فضایی با ابعاد معنایی پنهان وارد می سازد. پایادیمیتریو و همکارانش آن را فنی برای بازیابی اطلاعات بر اساس تحلیل طیفی ماتریس عبارت مدرک تعریف می کنند که پیشرفت های تجربی پیشی از آن فاقد هرگونه پیشی بینی و توصیف محکم بود. به زعم آنها نمایه سازی معنایی پنهان نوعی روش بازیابی اطلاعات است که تلاش دارد تا ساختار معنایی پنهان مجموعه مدرکی را با استفاده از فنون برگرفته از جبر خطی کشف کند

نمایه سازی معنایی پنهان اولین بار توسط گروهی از پژوهشگران (دیر وستر و همکارانش) در بلکور در بازیابی اطلاعات به کار رفت (18) و نمایه سازی معنایی پنهان نامیده شد. پژوهشگران به این دلیل واژه "پنهان" را به آن می افزایند که عبارت های جدید که بازنمون اطلاعات معنایی هستند، مستقیماً از مدرکی یافت نمی شوند بلکه حاصل بررسی کل مجموعه مدرکی و استفاده از روش ریاضی خاصی با عنوان تجزیه مقادیر منفرد (اس. وی. دی) است. به عقیده چنین ساختار معنایی پنهان مدرک با توجه به الگوی استفاده از واژگان (با توجه به امکان انتخاب چندین واژه) شکل می گیرد. مدل نمایه سازی معنایی پنهان از فنون آماری خاصی برای نشان دادن ساختار پنهان معنایی و زدودن زواید ناشی از امکان انتخاب چند واژه به جای یک مفهوم که دیر وستر و همکارانش به طور مبسوط آن را توصیف کرده اند،

استفاده می کند و باکشف الگوی استفاده از واژه ها، این ساختار را آشکار می کند و باعث حذف نوفه (پارازیت) می شود (13).

پیشی فرض نمایه سازی معنایی پنهان این است که معمولاً کل محتوای معنایی یک سند چون پاراگراف، چکیده یا کل مدرک با مجموع معانی واژه های آن به طور تقریبی برابر است. یعنی بدین صورت که: معنای واژه اول - معنای واژه دوم - معنای واژه سوم... معنای واژه «ام-معنای پاراگراف همچنین می توان با احتساب هر متن به صورت یک معادله خطی و کل مجموعه مدارک به صورت نظامی از معادلات هم زمان، معنای پایدار بازنمون های واژه ها را از کل یک مجموعه مدرک بزرگ به دست آورد (

نمایه سازی معنایی پنهان با استفاده از بستر متن، مترادف (یعنی امکان انتخاب چند واژه به جای یک مفهوم) و چند معنایی ها (یعنی وجود چند معنا برای عبارتی واحد) را کنترل و بدین ترتیب از ریزش کاذب پیشگیری می کند. یک عبارت پنهان ممکن است مرتبط به یک مفهوم نمایان (مانند مفهوم تعامل انسان و رایانه) باشد که با چند کلیدواژه توصیف می شود و ترکیبی از چند واژه است بنا به ادعای پژوهشگرانی چون دیروستر، دو میس، لندئر، فارناس و هارشمن، نمایه سازی معنایی پنهان شناخته شده ترین الگوریتم بازیابی اطلاعات است و برای مقاصد گوناگونی چون جستجو و بازیابی رده بندی و پالایش به کار می رود. نمایه سازی معنایی پنهان یکی از روش های فضایی برداری برای مدل سازی مدارک است و طبق نظر دیروستر و دیگران دو میس و کنتستاتیس و پتنگر معنای "پنهان" مجموعه مدارک را آشکار می سازد

پژوهش های دومیس نشان دادند که نمایه سازی معنایی پنهان با ترجیح مفهوم معنایی مدرک بر واژگان آن، مدل های فضایی برداری را بهبود می بخشد

در نمایه سازی معنایی پنهان برخلاف مدل های برداری مدارک که عبارت ها و واژه ها را مستقل فرض می کنند، سطوح مختلفی از همبستگی، وابستگی یا پیوستگی برای آنها در نظر گرفته می شود و این پیوستگی های بین عبارت ها با تشکیل مجموعه جدیدی از عبارت ها با استفاده از روش آماری تجزیه مقادیر منفرد مشخص می شوند همچنین در فضای معنایی پنهان، یک پرسش و یک مدرک، حتی در صورت نداشتن عبارت مشتری، می تواننا تشابه کسینوس زیادی داشته باشند زیرا عبارت های آنها از نظر معنایی، مشابهت "مفهومی دارند در صورتی که در مدل های برداری چنین چیزی امکان پذیر نیست

مهمترین نکته قوت نمایه سازی معنایی پنهان کار آمدی بازیابی مبتنی بر پرسش کاربر است که از طریق محاسبه ماتریس حاصل می شود. همان طور که با استفاده از این روش مدرک مربوط، حتی در صورت عدم مطابقت واژگان محتوای دو مدرک با یکدیگر، بازیابی می شوند

به طور کلی می توان گفت که نمایه سازی معنایی پنهان یکی از فنون روبه رشد نمایه سازی مجموعه مدرک بزرگ است که سعی دارد از یادگیری آماری ماشین در تحلیل متون استفاده کند (33).

محققان مزایای مختلفی برای نمایه سازی معنایی پنهان برمی شمردند که در ادامه به چند مورد اشاره می شود: چنگ، یکی از اصلی ترین مزایای نمایه سازی معنایی پنهان را استفاده از مفاهیم معنایی به جای تک تک واژه

ها در نمایه سازی می داند که بدین ترتیب مدرکی مربوط حتی در صورت نبوت واژه مشتری با عبارت پرسش بازیابی می شوند دیروستر نیز مزیت اصلی استفاده از بازنمون نمایه سازی معنایی پنهان را کار آمدی آن در رویا رویی بامدارکی می داند که حاوی مترادف ها، چند معنایی ها، و عبارت های وابسته بههمدیگر است به نظر هازبندز، سیمون و دینگ نمایه سازی معنایی پنهان چون یکی از روش های گسترش پرسش است می توانا جامعیت را بهبود بخشد.

موارد استفاده کاربرد های نمایه سازی معنایی پنهان پژوهشگران استفاده از نمایه سازی معنایی پنهان را در حوزه های مختلف ارزیابی کرده اند و آن را یکی از پرکار بردترین روش های نمایه سازی دانسته اند.

موارد استفاده از نمایه سازی معنایی پنهان عبارتند از:

1. *بازیابی اطلاعات*: همان طور که گفته شد نمایه سازی معنایی پنهان اولین بار توسط گروهی از پژوهشگران در بلکور در بازیابی اطلاعات به کار رفت و به همین نام شناخته شد. این روش بهتر از روش های برداری استاندارد عمل می کند و حتی در صورت نبوت اشتراك واژگانی میان سوال و مدرک، کارایی خود را حفظ می کند. استفاده از نمایه سازی پنهان معنایی در بازیابی اطلاعات در مقالات متعددی توسط محققانی چون بری، دومیس و ابریان و دومیس هال و هازبندز، سیمون و دینگ تاکید شده است محققان دیگری چون اندو و لی، بارتل، کنترل و بلو دومیس و ژا، مار کوس و سیمون نیز در پژوهش های خود نشان دادند که با استفاده از نمایه سازی معنایی پنهان جامعیت و مانعیت نظام بازیابی به طور قابل توجهی افزایش می یابد بشیری نیز از این روش در بازیابی متون فارسی استفاده کرده است
2. *بازخورد ربط*: بیشتر آزمون هایی که از نمایه سازی معنایی پنهان برای بازخورد ربط استفاده می کنند از روشی استفاده می کنند که در آن، حاصلجمع برداری مدارکی که توسط کاربران مربوط شنا سایر شده اند، جایگزین پرسش اولیه می شوند. پژوهش ها نشان داده اند که جایگزین سازی اولین مدرک مربوط با پرسش اولیه، عملکرد را در حد 33 درصد و جایگزین سازی سه مدرک مربوط به آن عملکرد را در حد 67 درصد بهبود می بخشد. نمایه سازی معنایی پنهان این قابلیت را دارد که با استفاده از فنون گسترش پرسش حتی بدون استفاده از بازخورد ربط اقدام با شنا سایی مدرکی مربوط کند ولی با استفاده از اطلاعات حاصل از بازخورد ربط، بازدهی عملکرد را به طور قابل توجهی افزایش می دهد.

3. *پالایش اطلاعات*: استفاده از نمایه سازی معنایی پنهان در پالایش اطلاعات بسیار آسان است. ابتدا نمونه ای از مدرک با استفاده از ابزارهای استاندارد نمایه سازی معنایی پنهان و تجزیه مقادیر منفرد تحلیل می شوند و بدین ترتیب علایق کاربر به صورت برداری با ابعاد کم نشان داده می شود و سپس مدرکی موجود

در مجموعه با آن بردار مطابقت داده شده و در صورت تشابه مدرک با آن برودر، به کاربر توصیه می شود. البته می توان با گذشت زمان با استفاده از روش های یادگیری نظیر بازخورد ربط برای بهبود بازنمون بردارهای علایق استفاده کرد (زیرا این علایق در طول زمان تغییر می یابند)

4. *همایش ارزیابی*: ارزیابی متنی اخیراً از نمایه سازی معنایی پنهان در زمینه های پالایش و ارزیابی اطلاعات همایش ارزیابی ارزیابی متنی استفاده می شود. پرسش های این برنامه بسیار طولانی و حاوی تومیغات مفصلی هستند که گاه طول آنها به 50 کلمه نیز می رسد. پرسش های این برنامه، دارای ابزارهای قوی تری نسبت به محاسن نمایه سازی معنایی پنهان یا سایر روش هایی که سعی در غنی سازی پرسش های کاربران دارند، است.

اصلي ترین چالش این مجموعه، گسترش ابزارهای نمایه سازی معنایی پنهان برای کنترل مجموعه بود که این نتایج، کاملاً دلگرم کننده بودند. از آنجا که در زمان همایش های ارزیابی ارزیابی متنی محاسبه کل مدرکی مجموعه منطقی نبود، از نمونه ای در حدود 70 هزار مدرک و 90 هزار عبارت استفاده شد. چنین ماتریسهای عبارت مدرکی بسیار پراکنده هستند و فقط 0.001 تا 0.002 درصد ورودی ها را دربر می گیرند، برای مثال محاسبه دویستمین مقدار فردی بزرگ تر، به حدود 18 ساعت زمان کار واحد پردازش مرکزی در یک پایگاه سانس پارک با ده ایستگاه کاری نیاز داشت. با وجود دشواری بسیار مقایسه تفصیلی نظام ها (به علت تفاوت های قابل توجه پیش پردازش، بازنمون و مطابقت)، عملکرد نمایه سازی معنایی پنهان بسیار خوب گزارش شد. استفاده از اطلاعات درباره مدرکی شناخته شده مربوط برای ایجاد برای هر پرسش، در پالایش بسیار سودمند بود. مزیت بازیافت 31 درصدی تا حدی کمتر از میزان مشاهده شده در سایر آزمون های پالایش بود که به پرسش های اولیه این همایش نسبت داده می شود. در ارزیابی نیز با استفاده از نمایه سازی معنایی پنهان در مقایسه با روش های برداری کلید واژه ای 16 درصد بهبود مشاهده شد

5. *بازیابی در متون چند زبانه*: چون نمایه سازی معنایی پنهان از دستور زبان انگلیسی استفاده نمی کند می توان از آن در هر زبانی استفاده کرد. به علاوه می توان از آن در بازیابی از متونی که چند زبانه هستند استفاده نمود و پرسش های کاربران (در هر زبان موجود در متون) با مطابقت زبان پرسش با مدرکی هم زبان صورت می گیرد. لند وئر و لیتمن روشی را برای ایجاد فضایی مشترک که واژه های هر زبان در آن ارائه شوند با استفاده از نمایه سازی معنایی پنهان توصیف می کنند. در این نوع موارد ماتریس اولیه عبارت مدرک، با استفاده از مجموعه ای از چکیده های چندین (که در پژوهش آنها فرانسویو انگلیسی بود) تشکیل می شود و به هر چکیده به صورت ترکیبی از نسخه های انگلیسی و فرانسوی، نگاه می شود. تجزیه مقادیر منفرد کاهش یافته، ماتریس عبارت چکیده ترکیبی محاسبه می شود. فضای حاصل از چکیده های ترکیبی انگلیسی و فرانسوی تشکیل می شود. بدین ترتیب در این فضای کاهش یافته، واژه های انگلیسی و فرانسوی که در چکیده های ترکیبی مشابه ظاهر می شوند، در کنار همدیگر قرار خواهند گرفت. بعد از این تحلیل، چکیده های تک زبانی وارد می شوند چکیده فرانسوی به سادگی در بردار حاصل جمع واژه های اجزای اصلی که قبلاً در فضای نمایه سازی معنایی پنهان وارد شده بودند قرار می گیرد. پرسش های انگلیسی یا فرانسوی با چکیده های انگلیسی یا فرانسوی مطابقت داده می شوند. در فضای نمایه سازی معنایی پنهان نیازی به ترجمه نیست. تجربه ها نشان داده اند که فضای چند زبانی کاملاً خودکار، حتی بهتر از فضایی با زبان واحد عمل می کند و نتایج بهتری عاید کاربران می کند. عملکرد ارزیابی مدارک فرانسوی زبان با پرسش های به زبان انگلیسی (و برعکس) برابر با عملکردی است که ابتدا پرسش ها به فرانسوی ترجمه و سپس در پایگاه داده منحصر فرانسوی جستجو شوند. این روش نتایج خوبی را در بازیابی چکیده های انگلیسی و اندیشه نگاشت های ژاپنی کنجی و ترجمه های چند زبانی (انگلیسی و یونانی) ارائه داد.

6. *مطابقت افراد به جای مدارک*: از دیگر کاربردهای نمایه سازی معنایی پنهان یافتن افراد خبره در یک حوزه

خاص با استفاده از مقالات و آثار آنهاست. در این کاربرد، اشخاص براساس مقالاتی که نوشته اند معرفی می شوند. برای مثال در پژوهش فرناس و دیگران که بلکور ادوایز نام دارد نظامی برای یافتن خبرگان محلی باتوجه به پرسش های کاربران طراحی شد. یک پرسش با جدیدترین مدارك مطابقت داده می شد و توصیف و مشخصات مربوط به پدیدآوران به عنوان مربوطترین مطلب ارائه می شد. در موردی دیگر، از نمایه سازی معنایی پنهان برای تخصیص مقالات به افراد به منظور ارزیابی داوری آنها استفاده شد. برای انجام اینکار ابتدا صدها منتقد بر اساس متون تألیفیشان توصیف شدند که اینکار بنیان تحلیل نمایه سازی معنایی پنهان شد. صدها مقاله نیز براساس چکیده آنها نشان داده شدند و با نزدیکترین و مرتبطترین داوران مطابقت داده شدند. استفاده از امکان تشابه یابی این برنامه برای تخصیص مقالات همایش 'تعامل انسان رایانه' به داوران انجام شد. تحلیل های بعدی نشان داد که اینکار کاملاً خودکار، به خوبی کاری بود که توسط افراد خبره انجام شده بود.

ابرداده (فراداده)

Metadata

ابرداده

برای ابرداده تعاریف بسیاری آمده که اکثر آنها بر کاربرد آن تکیه دارند و به عبارت دیگر، تعاریفی عملگرا به شمار می آیند.

به طور خلاصه ابرداده را «داده ای برای داده» یا «داده ای درباره داده» تعریف می کنند.

فراداده، مکرراً "داده درباره داده" یا "اطلاعات درباره اطلاعات" تعریف شده است. به عبارت دیگر، فراداده، داده ای است که منابع اطلاعاتی را توصیف می کند.

این تعریف گسترده، سطوح گوناگونی از توصیف (ساده تا پیچیده) را در بر می گیرد: یادداشت توصیفی کوتاهی در مورد یک کتاب، توصیف غیر رسمی را که موتورهای کاوش در مورد پیشینه های باز یابی شده ارائه می دهند، یک پیشینه کتابشناختی چاپی (فهرستبرگه) یا الکترونیکی (پیشینه مارک)، چکیده ها و اصطلاحات نمایه ای، و حتی یک استناد کتابشناختی.

ابرداده

بطور کلی می توان سه نوع ابرداده برای توصیف یک شیء دیجیتال ایجاد کرد: ابرداده **توصیفی**،

ساختاری، و **اجرایی**

● ابرداده توصیفی به صفاتی از عناصر اطلاعاتی مانند عناصر کتابشناختی نظیر

عنوان، نویسنده، ناشر و غیره، اشاره دارد.

● ابرداده ساختاری، ساختار و رابطه مجموعه ای از عناصر دیجیتال را توصیف می

- کند. به عنوان مثال، برای تولید ابرداده ساختاری يك کتاب، باید شیوه آرایش صفحات، فهرست مندرجات و ارتباط بین بخش ها و فصل ها را ثبت کرد.
- ابرداده اجرایی، همه اطلاعاتی است که در تمام دوره حیات یک شیئی دیجیتال، برای مدیریت آن مورد نیاز است و همه اطلاعات مورد نیاز برای حفاظت از آن را نیز در بر خواهد داشت، مانند تاریخ ثبت، ساختار فایل، حق مؤلف و غیره.

ابرداده

○ کاربردهای مهم ابرداده:

- سازماندهی منابع اطلاعاتی؛
- توصیف منابع اطلاعاتی اعم از متن، تصاویر، فایل‌های صوتی و نظایر آنها؛
- تحلیل محتوا و نمایه سازی؛
- تطبیق، اشتراک و یکپارچه سازی منابع اطلاعاتی ناهمگن؛
- زمینه سازی برای استفاده مجدد از اطلاعات توزیع یافته در سایر کتابخانه های دیجیتال و محیط وب؛
- ایجاد امکان دسترسی به اطلاعات دقیق و مرتبط توسط کاربران؛
- مدیریت دقیق تر بر حجم گسترده ای از اطلاعات در کتابخانه های دیجیتال

ابرداده

- یکی از تعاریف جامعی که از ابرداده ارائه شده، تعریفی است که «هینز» پس از ارائه تعاریف مختلف از آن، آورده است:
- «ابرداده داده‌ای است که محتوا، شکل یا خصوصیات یک رکورد داده‌ای یا یک منبع اطلاعاتی را توصیف می‌کند. ابرداده را می‌توان در توصیف منابع کاملاً ساختاریافته، یا اطلاعات ساختار نیافته از قبیل مدارک متنی به کار گرفت. همچنین می‌توان برای توصیف منابع الکترونیکی، داده‌های رقمی (شامل تصاویر رقمی) و مدارک چاپی از قبیل کتاب‌ها، مجلات و گزارش‌ها، مورد استفاده قرار داد. می‌توان آن را در درون يك منبع اطلاعاتی (مانند منابع وب) جای داد یا به طور جداگانه در

یک پایگاه اطلاعاتی نگهداری کرد» (Haynes 2004).

ابرداده

- می توان فراداده را، "داده های ساختارمند درباره دیگر داده ها" دانست. بنابراین، فراداده عبارت است از مجموعه ساختارمندی از عناصر که منابع اطلاعاتی را به منظور شناسایی، کشف، و استفاده، توصیف می نماید.
- اگر چه این اصطلاح، اصطلاحی نسبتاً جدید در حوزه کتابداری و اطلاع رسانی است اما، مفهوم توصیف داده ها یا منابع اطلاعاتی از زمان بوجود آمدن طرح های اولیه سازماندهی اطلاعات مشاهده می شود.
- همه ابزارهای کتابخانه ای که برای توصیف محتوای داده ها استفاده می شوند را می توان فراداده نامید.
- برای مثال، قواعد فهرست نویسی انگلو – امریکن ابزاری است که برای توصیف کتابشناختی و کمک به دسترسی سازمان یافته به داده ها و منابع الکترونیکی تدوین شده است .

- Data about data

- Information about information

- Structured data about data

تاریخچه فراداده

اصطلاح فراداده نخستین بار توسط جک ای. مایرز و اس. کی. کول جیان در سال 1969 به عنوان نام شرکتی تجاری (آمریکایی) برای ایجاد و توسعه محصولات مربوط به "فرا الگوها" به کار رفت و در یک جزوه معرفی محصول در سال 1972 چاپ شد. چند سال بعد، یعنی در

سال 1986، علامت تجاري MEADATA® براي شرکت فراداده ثبت گردید.

اصطلاح فراداده بعدها بوسیله متخصصان علوم رایانه، آمار، پایگاه های اطلاعاتی، و جامعه کتابداری و اطلاع رسانی به شکل های "فراداده"، "فرا داده"، و "فرا - داده" استفاده شد. اما کاربرد مکرر این اصطلاح بیشتر به دهه 1990 و شکل گیری شبکه جهانی وب در سال 1993 مربوط می شود.

انواع قالب های مارکی که در دهه های 1960، 1970، و 1980 شکل گرفته را می توان نخستین قالب های فراداده ای دانست. تاریخچه شکل گیری سایر قالب های فراداده ای به دهه 1990، به ویژه 1993-1996 باز می گردد.

از شناخته شده ترین قالب های فراداده ای در دهه 1990 می توان به طرح کدگذاری متن، خدمات مکان یابی اطلاعات دولتی، تبادل رایانه ای اطلاعات موزه ها، چهارچوب توصیف منبع، توصیف کدگذاری شده آرشیوی، و طرح فراداده ای هسته دوبلین اشاره کرد.

- با توجه به ویژگیهای منابع الکترونیکی و محیط هایی که این منابع ایجاد کرده اند و یا در آنها بکار می روند، نیز تنوع منابع و محصولات اطلاعاتی، ضعف ابزارهای جستجوی اطلاعاتی (مانند موتورهای کاوش)، تقاضای های فزاینده استفاده کنندگان، تولید کنندگان و عرضه کنندگان منابع الکترونیکی و شبکه ای، نیاز به ابزارهایی جدید برای سازماندهی و یا مطابقت ابزارهای سنتی با این تحولات کاملاً بدیهی می نماید.
- این همان جایی است که بحث فراداده مطرح می شود. بر همین اساس، فراداده ها می توانند با افزوده شدن (جاسازی) و پیوند یافتن با منابع، باعث توصیف، شناسایی، مکان یابی، و مدیریت منابع گردند.

کارکرد فراداده

در بسیاری از بحث های مربوط به فراداده بویژه در حوزه مدیریت اطلاعات (از دیدگاه علم کتابداری و اطلاع رسانی)، تمایل به گروه بندی عناصر فراداده ای بر اساس کارکردهایی که

این عناصر پشتیبانی می کنند، وجود دارد.
 نتایج این بحث ها شناسایی انواع متفاوت فراداده یا گروه های فراداده ای است. در ذیل جدولی از مشهورترین انواع فراداده به همراه کارکردهای آنها ارائه شده است:

مزایا و کاربردهای فراداده

- اساس کاربرد فراداده، کاوش، بازیابی، دسترسی، کشف، مستند سازی، ارزیابی و انتخاب منابع الکترونیکی بویژه شبکه ای است که باعث افزایش دقت بازیابی می گردد. به تعبیر دقیق تر می توان مهمترین کاربردهای فراداده را این چنین توصیف کرد:

دو عبارتی که معمولاً با فراداده همراه هستند

● زبان نشانه گذاری گسترش پذیر

● eXtensible Markup Language (XML)

● قالب توصیف منبع

● Resource Description Framework (RDF)

- زمانیکه برخی افراد درباره XML و RDF صحبت می کنند گویی که این دو آغاز خودشان قالب ابر داده ای هستند، اما این آشوبی میان قالب و محتوا است.
- در عمل XML و RDF قالبهایی داده ای عمومی هستند که می توانند در موارد گوناگون استفاده شوند. خصوصاً از XML اغلب به عنوان قالب مدرک استفاده می شود و قالب کلی تری است که از آن HTML حاصل می شود.

زبان نشانه گذاری گسترش پذیر XML

- هرگاه شما فیلدهای رکورد مارک را با تگ هایی مانند استفاده از "245" به معنی "عنوان" در نظر بگیرید:

● a Hamlet, Prince of Denmark\$ 245

● در این صورت XML تنها روش دیگری برای تگ زدن یک قطعه داده است، اگرچه آن شامل گذاشتن یک تگ شروع و یک تگ پایانی (با یک "/" قبل از نام تگ) در اطراف هر عنصر داده ای است:

● `<title>Hamlet, Prince of Denmark</title>`

● تگ ها می تواند هر چیزی که شما تمایل دارید باشد، به شرط اینکه شما آنها را در یک قالب داده ای تعریف ساختار، از قبل تعیین کرده باشید. بنابراین اگر شما ترجیح می دهید، تعریفان می تواند دارای هر کدام از تگ های زیر برای "عنوان" باشد.

● `Hamlet, Prince of Denmark</245><245>`

`<ti>Hamlet, Prince of Denmark</ti>`

● در اصل XML، مشابه تگ ها و فیلدهای فرعی مارک، سلسله مراتبی است. مزایای آن بیش از مارک 21 است که می تواند به اندازه ای که سطوح سلسله مراتبی لازم است را دارا باشد. برخلاف مارک 21 که می تواند دارای دو سطح تگ و فیلد فرعی باشد.

قالب توصیف منبع RDF

- RDF یک مرحله یا دو مرحله بعد از XML است.
 - بر روابط میان عناصر داده ای تأکید دارد.
 - همانطور که از نام RDF برمی آید ایجاد یک مکانیزمی برای توصیف منابع و اسناد بر روی اینترنت است و رابطه کلیدی آن توصیف می باشد.
 - RDF عنصر ضروری و مهمی از وب معنایی است که توسط ائتلاف شبکه جهانگستر وب برای افزودن عنصر معنایی و در چارچوب توصیف منابع برای اشتراک داده ها روی اینترنت ایجاد شده است.
 - RDF پیچیده تر و نسبت به XML کمتر استفاده شده است و هنوز مشخص نیست، آیا آن به عنوان یک زبان عمومی در توصیف جهانی وب موفق است.
- برخی از استانداردهای فراداده مهم
- مارک

- هسته دوبلین
- مودس
- مدس
- متز

مارک MARC

قالب مارک به منظور ماشین خوان و ماشین فهم نمودن پیشینه های کتابشناختی طراحی شده است. قالب مارک در واقع قالبی است برای ثبت داده های فهرست نویسی و استفاده بعدی از آنها برای تبادل، جستجو، بازیابی، نمایش، اطلاعات کتابشناختی، و یا چاپ پیشینه های کتابشناختی. مارک دارای قالب های مختلف داده ای است که هر یک به منظور استفاده خاصی ایجاد شده اند

هسته دوبلین (Dublin core)

- استاندارد فراداده ای هسته دوبلین، طرحی بین المللی و میان رشته ای است که مجموعه عناصری ساده و کارآمد برای توصیف طیف گسترده ای از منابع اطلاعاتی شبکه ای ارائه می دهد.
- پدید آمدن این قالب فراداده ای مانند دیگر قالب های فراداده ای برخاسته از اهداف و نیاز هایی بوده که توسعه دهندگان و پدیدآورندگان آن در نظر داشته اند.
- این طرح استاندارد دو سطح دارد: ساده و ویژه (مقید به توضیحگرها).
- سطح ساده شامل پانزده عنصر، و سطح ویژه شامل هفت عنصر (مخاطب، منشا (ریشه)، نگهدارنده حقوق) بیشتر یعنی بیست و دو عنصر می باشد.
- هسته دوبلین، همچنین دارای گروهی از توضیحگرها (پالایش های عنصر و طرح های کدگذاری عنصر) است که عناصر را از لحاظ معناشناختی به منظور فرایند کشف منبع پالایش می کنند.
- مباحث مربوط به معناشناختی عناصر هسته دوبلین توسط گروهی از متخصصان رشته های کتابداری و اطلاع رسانی، علم رایانه، کدگذاری متن، جامعه موزه ها، و دیگر حوزه های مرتبط پرداخته شده است.

4-5. مجموعه عناصر فراداده‌ای هسسه دوبین
4-5-1. عناصر سطح ساده

نام عنصر	شرح عنصر
عنوان	Title
عنوان	Title
پدیدآورنده	Creator
موضوع	Subject
توصیف	Description
ناشر	Publisher
همکار	Contributor
تاریخ	Date
نوع منبع	Type

فالب	Format
شناسه‌گر	Identifier
زبان	Language
ارتباط	Relation
منبع اصلی	Source
پوشش	Coverage
حقوق	Rights

عناصر سطح ویژه

طبقه ای از یک موجودیت که منبع برای او تهیه شده یا مفید است	Audience	مخاطب
شرح هر تغییری در مالکیت و حفاظت منبع پس از ایجاد آن که برای صحت و اعتبار، یکپارچگی، و تفسیر منبع مهم قلمداد می شود	Provenance	منشا (ریشه)
شخص یا سازمانی که حقوق منبع را مالکیت یا مدیریت می کند	Rightsholder	مالک حقوق
فرآیندی که برای تولید دانش صورت می گیرد، با گرایش ها و مهارت هایی که منبع به منظور پشتیبانی آنها اختصاص یافته است؛ شامل: روش های ارائه مطالب آموزشی یا هدایت فعالیت های مرتبط (مانند الگوهای تعامل یادگیرنده- یادگیرنده؛ یادگیرنده - آموزنده و ...)	Instructional Method	شیوه آموزشی
روش افزوده شدن آثار به یک مجموعه. بهترین روش توصیه شده، استفاده از واژگان های کنترل شده است.	Accrual method	شیوه گسترش
بسامد افزوده شدن آثار به یک مجموعه. بهترین روش توصیه شده، استفاده از واژگان های کنترل شده است.	Accrual periodicity	تناوب گسترش
خط مشی افزوده شدن آثار به یک مجموعه (خط مشی مجموعه سازی). بهترین روش توصیه شده، استفاده از واژگان های کنترل شده است.	Accrual Policy	خط مشی گسترش

یک پیشینه فراداده ای در قالب استاندارد هسته دوبلین و در بستر نحوی "زبان نشانه گذاری فرامتن (HTML)"

<link rel="schema.DC" href="http://purl.org/dc/elements/1.1/" />

<link rel="schema.DCTERMS"

href="http://purl.org/dc/terms/" />

<meta name="DC.title" lang="فارسی" content="پژوهشگاه علوم و

فرهنگ اسلامی" />

<meta name="DC.creator" content="پژوهشگاه علوم و فرهنگ

اسلامی" />

<meta name="DC.subject" lang="فارسی" content="اخبار علمی،

دفتر تبلیغات اسلامی حوزه علمیه قم - فعالیتهای پژوهشی،وب سایتها" />

<meta name="DC.description" lang="فارسی" content="در سال

1380 همزمان با برنامه‌ریزی جامع دفتر تبلیغات اسلامی و تدوین بیانیه مأموریت آن، معاونت پژوهشی نیز به بازنگری و اصلاح ساختار، اهداف و خطمشی‌های خود پرداخت و سرانجام، چشم‌انداز و سیاست‌های کلان پژوهشی و همچنین اهداف و برنامه‌های پژوهشکده‌ها و گروه‌های پژوهشی خود را تنظیم و تدوین کرد. در ساختار جدید دفتر مقرر گردید که کلیه مراکز پژوهشی در قالب یک پژوهشگاه جامع سازمان‌دهی شود. پس از چند سال پی‌گیری، سرانجام با تصویب چهار پژوهشکده در شورای عالی گسترش، مجوز قطعی «پژوهشگاه علوم و فرهنگ اسلامی» از وزارت علوم، تحقیقات و فناوری گرفته شد. هم‌اینک پی‌گیری اخذ مجوز رسمی برای «مرکز فرهنگ و معارف قرآن» و «پژوهشکده علوم اجتماعی و اقتصادی» در دست اقدام است />

<meta name="DC.publisher" content="علمیه قم" />

<meta name="DC.contributor" content="اسلامی مرکز اطلاعات و مدارك" />

<meta name="DC.date" scheme="DCTERMS.W3CDTF" content="2007-04-25" />

<meta name="DC.type" scheme="DCTERMS.DCMiType" content="Text" />

<meta name="DC.format" content="text/html" />

<meta name="DC.format" content="54591 bytes" />

<meta name="DC.identifier" scheme="DCTERMS.URI" content="http://www.isca.ac.ir" />

<meta name="DC.language" scheme="DCTERMS.URI" content="فارسی" />

<meta name="DC.relation" scheme="DCTERMS.URI" content="http://www.islamicdoc.org" مدارك اسلامی مرکز اطلاعات و مدارك" />

<meta name="DC.rights" scheme="DCTERMS.URI"

content="میکلیه حقوق وب سایت متعلق به دفتر تبلیغات اسلامی حوزه علمیه قم

می باشد /> "

مارک در زبان نشانه گذاری گسترش پذیر (MARCXML)

"دفتر استانداردهای مارک و توسعه شبکه" کتابخانه کنگره، در حال گسترش چارچوبی برای پیاده سازی داده-های مارک در محیط "زبان نشانه گذاری گسترش-پذیر-XML" است. انعطاف-پذیری و گسترش پذیری این چارچوب، امکان پاسخگویی به نیازهای مختلف و خاص را برای کاربران فراهم می-کند. این چارچوب خود شامل طرح-ها، نرم-افزارها، و الگوهای از پیش تعریف شده-ای می-شود. در ذیل به مهمترین آنها اشاره می گردد:

توصیف آرشیوی رمزگذاری شده EAD

استاندارد ابرداده ای که توسط "دفتر توسعه شبکه و استانداردهای مارک کتابخانه کنگره" برای جامعه آرشیویست های آمریکا تهیه شده است. باعث می شود که اطلاعات آرشیوی به صورت ماشینی ذخیره، پردازش، بازیابی و قابل انتقال شوند.

مبنتی بر زبان ایکس ام ال است.

www.loc.gov/ead

طرح فراداده ای توصیف شی (MODS)

- این طرح برای مجموعه عناصر کتابشناختی که با اهداف گوناگون، به خصوص کاربردهای کتابخانه ای استفاده می شوند ، تهیه شده است.
- به عنوان یک طرح "زبان نشانه گذاری گسترش-پذیر"، امکان انتقال داده های کتابشناختی گزیده از پیشینه های مارک موجود، نیز توانایی ایجاد پیشینه-های توصیف برای منابع جدید را فراهم می آورد.
- این طرح بخشی از مناطق (فیلدهای) مارک 21 را شامل شده و از برچسب های زبان-

پایه بیش از برچسب های عددی استفاده می کند.

طرح فراداده ای توصیف شی (MODS)

طرح مذکور با سایر طرح ها و ابزارهای "مارک در قالب زبان نشانه گذاری گسترش پذیر (MARCXML)" که از طریق وب سایت مارک دسترس پذیرند، سازگار است و به طور ویژه از فهرست نویسی منابع الکترونیکی پشتیبانی می نماید.

این طرح بوسیله "دفتر استانداردهای مارک و توسعه شبکه" و با همکاری گروهی از کاربران علاقه-مند به مارک در سال 2002 گسترش یافت.

قالب مارک استاندارد پیچیده برای رمزگذاری اطلاعات کتابشناختی است.

در محیط شبکه ای در شرایطی که ابرداده توصیفی می تواند بر روی نظامها ارسال شود و می تواند درون یا همراه با انواع دیگر ابرداده باشد، استفاده از رکوردهای مارک با این هدف، ایده آل به نظر خواهد رسید.

با این وجود، گنجاندن رکوردهای مارک درون ابرداده، استفاده از ساختار داده ای XML را ضروری ساخته است، درحالیکه مارک یک رکورد XML نیست.

کتابخانه کنگره روشی را برای تبدیل رکورد مارک به XML ایجاد کرده است، اما آن با اقبال زیادی روبرو نشد و احتمالاً به این دلیل که رکورد مارک بزرگتر و جزئی تر از اکثر نظامهای مورد نیاز است، و استفاده اش از تگ های عددی و کدهای فیلد فرعی موجب شده درک آن را بدون آموزش قابل توجهی مشکل سازد.

چیزی که مورد نیاز بود نسخه ای خوشایندتر و باکیفیت تر از مارک بود تا بتواند عناصر کلیدی داده ای از رکورد مارک را قبول نماید و آنها را به یک قالب آسان و قابل درک XML ارسال نماید. بنابراین استاندارد توصیف ابرداده ای شی گرا (مودس) متولد شد.

مودس از تگ های قابل درک بشری بجای تگ های سه رقمی و کدهای فیلد فرعی مارک (یعنی "عنوان" بجای "245") استفاده می کند.

بیشتر عناصر داده ای فیلد ثابت، به استثنای کدهای قالب فیزیکی (از 007) و بیشتر کدهای نوع (از 008) را نادیده می گیرد.

برخی قابلیت ها و نوآوری ها را معرفی می کند.

مودس یک ساختار نام گذاری شده "نام" را تعیین می کند که فیلدها و فیلدهای فرعی برای اسامی شخصی، سازمانی و برای کنفرانسها را ارائه می نماید. این ساختار می تواند در هر جایی که نامها ظاهر خواهند شد، یا به شکل سرشناسه، شناسه های افزوده یا موضوعات استفاده شوند.

نمونه ای از طرح فراداده ای توصیف شی (MODS)

```
name type="personal">
```

```
<namePart>Shakespeare, William</namePart>
```

```
<namePart type="date">1564-1616</namePart>
```

```
</name>
```

می تواند به عنوان یک فیلد نویسنده استفاده شود یا آن می تواند قسمتی از یک سرعنوان موضوعی باشد

```
<subject authority="lcsch">
```

```
<name type="personal">
```

```
<namePart>Shakespeare, William</namePart>
```

```
<namePart type="date">1564-1616</namePart>
```

```
</name>
```

```
<topic>Bibliography</topic>
```

```
<topic>Periodicals</topic>
```

```
</subject>
```

طرح فراداده-ای توصیف مستند (MADS)

طرح فراداده-ای توصیف مستند" برای مجموعه عناصر مستندی که در تهیه پیشینه-های فراداده-ای به منظور توصیف افراد حقیقی و حقوقی، رویدادهای مهم، و شناسه-های موضوعی و جغرافیایی، و... بکار می-روند، و با هدف پیاده سازی قالب مستندات مارک 21 در بستر نحوی "زبان نشانه گذاری گسترش-پذیر" و استفاده از مزایای این محیط توسعه یافته است. این طرح نیز محصول "دفتر استانداردهای مارک و توسعه شبکه" و برخی متخصصان علاقه-مند به قالب مارک می-باشد.

استاندارد انتقال و کدگذاری فراداده-ها (METS)

پروژه MOA2، در زمینه تهیه یک قالب کدگذاری برای فراداده-های توصیفی، مدیریتی، و ساختاری برای منابع متنی و تصویری به منظور پاسخگویی به نیازهای سازماندهی کتابخانه-های دیجیتال فعالیت می-کند. در همین راستا، "استاندارد انتقال و کدگذاری فراداده-ها" که محصول "فدراسیون کتابخانه-های دیجیتال" است برای تحقق اهداف پروژه MOA2، با ارائه یک قالب مدرک در بستر نحوی "زبان نشانه گذاری گسترش-پذیر" به منظور ایجاد فراداده-های ضروری برای مدیریت اشیای محتوایی کتابخانه دیجیتال یک سازمان و تبادل این اشیای محتوایی میان چند سازمان و کاربران آنها، تلاش می-کند. بسته به نوع استفاده، این استاندارد می-تواند در بسته-های گوناگونی بکار گرفته شود.

مدیریت کتابخانه های متشکل از شیء های دیجیتالی نیازمند مدیریت ابرداده های آن شیء هاست. ابرداده های ضروری برای مدیریت موفق و استفاده از شیء های دیجیتالی، گسترده تر و متفاوت تر از ابرداده های مورد استفاده در مدیریت مجموعه آثار چاپی و سایر منابع فیزیکی می باشد. در فهرستهای کتابخانه ها ابرداده توصیفی هر کتاب ثبت می شود. اگر کتابخانه در ثبت دقیق ابرداده ساختاری سازمان منتشردهنده کتاب ناموفق باشد، به هم پیوستگی صفحه های کتاب مذکور از بین خواهد رفت. همچنین اگر کتابخانه در بیان نحوه چاپ کتاب و استفاده از ابزارهای ویژه دقت نکرده باشد، در بازیابی کتاب توسط محققان آثار منفی در پی خواهد داشت

اما این مسائل درخصوص نسخه دیجیتالی همان کتاب صحیح نیست. بدون ابرداده ساختاری، تصویر صفحه یا فایل متنی به وجود آورنده اثر دیجیتالی کمتر مورد استفاده قرار خواهد گرفت. برای اهداف مدیریت داخلی، کتابخانه باید ابرداده فنی مناسب جهت روزآمدسازی و انتقال متناوب داده و اطمینان از ماندگاری منابع باارزش دسترسی داشته باشد.

این قالب استاندارد با فراهم کردن یک قالب کدگذاری برای ابرداده توصیفی، مدیریتی و ساختاری جهت آثار متنی و تصویری به حل این مسئله کمک کند. نتیجه و حاصل از اینکار، مدارک با قالب XML می باشد که از این طریق کدگذاری ابرداده های ضروری جهت مدیریت شیء های دیجیتالی در مخزن و نیز تبادل چنین شیء هایی بین مخزنها (یا بین مخزنها و کاربران آنها) فراهم می شود(6).

با این اوصاف، مدرک ابرداده ای وجود دارد که هدف "توصیف" به مفهوم فهرست نویسی نیست. یک مثال از قالب ابرداده ای که در حال استفاده توسط کتابخانه ها و آرشیوهای دیجیتالی است، استاندارد رمزگذاری و انتقال ابرداده (متس) می باشد. از متس به عنوان "یک روکش" برای حفظ فایلها از

یکدیگر استفاده می شود تا با همدیگر مخلوط نشوند. یک شیء دیجیتالی ایجاد شود. برخلاف جلدکتاب، مدارک دیجیتالی اغلب از تعدادی فایل‌های جداگانه ارائه شده در صفحات یا واحدهای دیگر ساخته می شود. و بر خلاف کتاب فیزیکی، جلد یا عنوان صفحه قابل رویت نیست، نه می توان از انگشت شست برای یافتن صفحه ای خاص در کتاب استفاده نمود. متس را به عنوان صحافی، جلد و ردیابی یک گروه از فایل‌های دیجیتالی تصور کنید (3). همچنین متس شامل اطلاعات فنی است که نیاز به کنترل و شناخت این فایلها خواهد داشت، مانند قالب‌های فایل، فناوری استفاده شده در پوشش، و تغییر شکل دیجیتالی و فشردگی که بر روی فایلها استفاده شده است. چیزی که متس تعیین نمی کند ابر داده توصیفی است. در عوض امکان ایجاد رکوردهای متس برای درج در هر ابر داده توصیفی وجود دارد. این ویژگی مهم، دنیای ابر داده را توصیف می کند، که ما آن را در نمونه کار کریپتوکامانز مشاهده نمودیم: ابر داده می تواند دوباره استفاده شود بجای اینکه دوباره اختراع شود. معمولاً ابر داده توصیفی در دابلین کور یا در مودس شامل رکوردهای متس هستند

فراداده های مبتنی بر XML برای تصاویر (MIX)

دفتر استانداردهای مارک و توسعه شبکه، با همکاری کمیته استانداردهای فراداده ای فنی برای تصاویر ثابت وابسته به "سازمان استانداردهای اطلاعات ملی (NISO)"، و دیگر متخصصان علاقه مند، در حال گسترش یکه فرانما XML برای مجموعه ای از عناصر داده ای فنی مورد نیاز برای مدیریت مجموعه تصاویر دیجیتالی هستند. این فرانما، قالبی برای مبادله و/یا ذخیره سازی داده های معین در "فرهنگ داده ها- فراداده های فنی برای تصاویر ثابت دیجیتالی (ANSI/NISO Z39.87-2006)" فراهم می نمایند. فرانمای یاد شده، در حال حاضر به "فراداده مبتنی بر XML برای تصاویر" مربوط به "سازمان استانداردهای اطلاعات ملی (NISO)" ارجاع می دهد. MIX با استفاده از زبان فرانمای XML مربوط به کنسرسیوم وب بیان می گردد. و بوسیله دفتر استانداردهای مارک و توسعه شبکه به سفارش "سازمان استانداردهای اطلاعات ملی (NISO)"، با دروندادی از کاربران نگهداری می شود.

ابر داده برای آرشیو الکترونیکی

مدیریت رکوردها از طریق تنظیم منابع و فراهم کردن بررسی خلاصه مذاکرات، نقش مهمی در مستندسازی تصمیمات دولتی ایفا می کند. مدیریت رکوردها شامل مدیریت فایل های فیزیکی و

رکوردهای الکترونیکی می‌باشد؛ هرچند که آن‌ها محل رسانه‌ای متفاوتی داشته باشند، اما اکثر فرایندهای مشابه بر روی هر نوع رکوردی اعمال می‌شود. مدیریت خوب رکوردها به ایجاد رکوردهای درست و صحیح بستگی دارد. این قسمت بر کاربرد ابر داده در مدیریت چرخه عمر رکوردها تأکید می‌کند.

زمانی که رکورد ایجاد یا دریافت می‌شود، به ثبت می‌رسد یا با یک نام معین ذخیره می‌شود. در طول فرایند فراهم‌آوری، ابر داده‌های مربوط به رکورد که شامل تاریخ ایجاد، مالک، تاریخ بازبینی و رده امنیتی آن است، ایجاد می‌گردند.

تعدادی از عناصر داده‌ای که اختصاصاً برای مدیریت رکوردها طراحی شده‌اند، در پروفایل‌های کاربردی ابر داده بر اساس «دوبلین کور» [7] تعریف شده‌اند. مثلاً «آرشیو ملی انگلستان» به عنوان یک سازمان موفق در اداره رکوردهای عمومی و دست‌نوشته‌های تاریخی، گروهی از عناصر داده‌ای را اختصاصاً برای مدیریت رکوردها در دولت مرکزی تعریف کرده است. اما این عناصر را می‌توان در یک پروفایل کاربردی عمومی‌تر نیز به کار بست (Public Records Office)

(2002). این عناصر داده‌ای عبارت‌اند از: Subject

Title.

Identifier.

Record Type.

Addressee.

Date.

Rights.

Location.

Language.

Description.

Mandate.

Preservation.

Aggregation.

Relation.

Creator.

Digital Signature .

موتورهای جستجو

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. مهمترین مرجع الکترونیک امروزی هستند

برخی از جویشرها برای تنها یک وب‌گاه (پایگاه وب) اینترنت به کار برده می‌شوند و در اصل جویشرهای 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 مورد نظر آن‌ها را در خود دارد پی ببرند.

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 تغییراتی 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.

جویشرهای پیمایشی

وجه مشخصه این گروه از جویشرها وجود نرم‌افزار 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 دست‌نویس شده یا Human-Powered (Directories) مانند فهرست بازی Open (Directory) 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.

فهرست‌های دست‌نویس شده

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 شانس بیشتری نسبت به یک سایت با پایگاه داده ضعیف دارد.

البته در مورد جویسگرهای 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 و بگانه خود در فهرست بندی هستند، می‌توانند با پرداخت پول به این جویسگرها به هدف خویش برسند.

جویسگرهای ترکیبی با نتایج مختلف

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

ابرجویسگرها

این گونه جدید از جویسگرها که قدمت چندانی نیز ندارند، بصورت هم‌زمان از چندین جویسگر برای کاوش در شبکه برای کلید واژه مورد نظر استفاده می‌کنند. بدین معنی که این جویسگر عبارت مورد نظر شما را در چندین جویسگر دیگر جستجو کرده و نتایج آنها را با هم ترکیب کرده و یک نتیجه کلی به شما ارائه می‌دهد. به‌عنوان مثال جویسگر داگ پایل از نتایج حاصل از موتورهای [Google - Yahoo - MSN - ASK](#) استفاده کرده و نتیجه حاصله را به شما ارائه می‌دهد.

نوجویسگرها

این گونه از جویسگرها، نسل جدید و متفاوتی از جویسگرهای گذشته هستند. امکان ثبت جستجو و مدل‌سازی فعالیت‌های کاربر و ارائه نتایج جدید به‌کاربر، به‌صورت متفاوت و تفکیک شده، از امکانات نوجویسگرها است

بررسی یک جویشرگر پیمایشی

جویشرگرهای 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 جویشرگر 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 جویشرگر 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 جویشرگر، سومین عنصر یک جویشرگر است

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 تجارب 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 نظرتان 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

عوامل خارج از صفحه

سرفصل‌های بهینه‌سازی

مکان و بسامد

یکی از قوانین اصلی در الگوریتم‌های رتبه‌بندی موقعیت و بسامد (تعداد تکرار) واژه‌هایی است که

در صفحه مورد استفاده قرار گرفته‌اند

بطور خلاصه روش مکان-بسامد (Location/Frequency Methode) نامیده می‌شود.

کتابدار مذکور را به خاطر می‌آورد؟ لازم است که او کتاب‌های در رابطه با واژه «سفر» را

طبق درخواست شما بیابد. او در حله اول احساس می‌کند که شما به دنبال کتاب‌هایی هستید که

در نامشان کلمه «سفر» را شامل شوند. جویسگرها هم دقیقاً همان کار را انجام می‌دهند. آنها هم صفحاتی را برایتان فهرست می‌کنند که در پرچسب عنوان (Title) موجود در کد زبان نشانه‌گذاری آپرمتی (زنگام) (HTML حاوی واژه «سفر» باشند.

جویسگرها 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 جویسگرها از طریق آن صفحات مربوط را شناسایی می‌نمایند.

جویسگرها 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 جویسگرها درباره روش‌هایی از قبیل مکان-تکرار فاش نمی‌شود و هر جویسگری 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 جویسگرهای معتبر و بزرگ همچنان محرمانه نگهداری می‌شوند. برخی جویسگرها نسبت به برخی دیگر صفحات بیشتری را فهرست کرده‌اند. نتیجه این خواهد شد که هیچ جویسگری نتیجه جستجوی مشترکی با موتور دیگر نخواهد داشت و شما نتایج متفاوتی را از آن‌ها دریافت می‌کنید.

جویسگرها 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 جویشگرها داشته باشند.

فرستادن هرزنامه (Spamming) روشی است که برخی از صفحات برای احراز رتبه بالاتر در جویشگرها در پیش می‌گیرند

به این صورت است که با تکرار بیش از حد واژه‌ها و یا بزرگ نوشتن یا بسیار ریز نوشتن متنها بطور عمدی کوشش در بر هم زدن تعادل و در نتیجه فریب جویشگرها 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 جویشگرها راه‌های متنوعی برای جلوگیری از فرستادن هرزنامه دارند و در این راه از گزارش‌های کاربران خود نیز بهره می‌برند.

امروزه بهینه‌سازی سایت‌های اینترنت برای جویشگرها یکی از مهم‌ترین روش‌های جلب بازدیدکننده به سایت است.

عوامل خارج از صفحه

جویشگرهای 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 جویشگرهای 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.

عوامل خارج از صفحه

به وسیله تجزیه صفحات، جویشگرها 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 جویشگرها جلوگیری می‌نماید.

علاوه بر آن جویشگرها 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.

تکنیک‌های جستجو برای جستجوی بهینه

یکی از کاراترین و مقتدرترین روش‌های جستجوی اطلاعات در دنیای وب استفاده از واژه‌هایی است که اصطلاحاً کلمات کلیدی نامیده می‌شوند.

اغلب کاربران حرفه‌ای و جستجوگران ورزیده دنیای اینترنت می‌توانند با طرح بهترین کلمات کلیدی و بکار بستن قوانین ترکیب آن‌ها با هم برای نیازهای اطلاعاتی خود پاسخی در خور بیابند. در این روش توصیه‌های زیر برای انتخاب کلمات کلیدی و نیز جستجوی دقیق و مفید پیشنهاد می‌شود که بشرح ذیل است :

تکنیک‌های جستجو برای جستجوی بهینه

- حتی‌المقدور سعی شود کلمات کلیدی از میان اصطلاحات منحصر به فرد و اسامی خاص انتخاب شود .
- حتی‌المقدور از آوردن کلمات عمومی که عناوین بسیاری را در زیر مجموعه خود شامل می‌شوند، جداً خودداری کنید .
- همیشه اسم شخص یا نام شی یا هر چیز دیگری را که مد نظر دارید به‌طور کامل وارد کنید .
- دقت کنید که اگر موتور جستجو میان حروف بزرگ و کوچک تفاوتی می‌گذارد، این مسأله را در طرح کلمات کلیدی خود مدنظر داشته باشید .
- در نظر داشته باشید اگر نتیجه جستجو صفر بود به احتمال زیاد می‌تواند از يك اشتباه تایپی باشد .
- اگر املاي صحیح و کامل کلمه‌ای را نمی‌دانید از کارکتر جانشین که اغلب * و یا ؟ است

استفاده کنید .

تکنیک های جستجو برای جستجوی بهینه

7. اگر يك کلمه کلیدی را برای طرح دقیق و تمام و کمال يك مورد جستجو کفایت نمی‌کند، از تکنیک‌های جستجوی عبارتی، استفاده از اپراتورهای جبر بولین (AND, OR, NOT) استفاده کنید. جستجوی عبارتی یکی از مهم‌ترین و قدرتمندترین امکانات جستجو در اغلب موتورهای جستجو می‌باشد و می‌توان يك عبارت یا جمله مشخص را به همان ترتیبی که کلمات وارد شده‌اند مورد جستجو قرار داد. برای این روش جستجو عبارت مورد نظر را داخل گیومه "" بگذارید .

8. استفاده از عملگر AND : AND به مفهوم "و" برای محدود کردن دامنه جستجو از طریق ترکیب کلید واژه‌های مختلف به کار می‌رود و برای ترکیب کلیدهای جستجو زمانی که برای شما مهم است که دو یا چند کلمه کلیدی حتماً وجود داشته باشد و علامت آن در پایگاه‌های مختلف به صورت استفاده از عبارت AND ، استفاده از + ، انتخاب عبارت ALL THE WORD از منو، انتخاب عبارت (MATCH ON ALL WORDS AND) به‌وسیله کلیک کردن بر روی دکمه‌های رادیویی است .

تکنیک های جستجو برای جستجوی بهینه

9. استفاده از عملگر OR: OR اپراتور OR به مفهوم "یا" و برخلاف عملگر AND باعث گسترش دامنه جستجو و بازیابی اطلاعات بیشتر شده برای ترکیب کلید واژه‌های جستجو زمانی که انتظار دارید تنها يك، دو یا چند کلمه کلیدی حضور داشته باشند و علامت آن استفاده از عبارت OR ، نحوه‌ی اجرای ساده و معمولی آن، انتخاب عبارت ANY OF THE WORDS از منو، انتخاب عبارت (MATCH ON ANY WORDS OR) با کلیک بر روی دکمه‌های رادیویی می‌باشد. یکی از کاربردهای مهم این عملگر پوشش مفاهیم یا اصطلاحات مترادف، مرتبط یا با املاهای متفاوت است .

10. استفاده از عملگر NOT: NOT اپراتور NOT به مفهوم "نه" و یا به جز که در این صورت تمامی جواب‌های بازگشتی که حاوی عبارت یا کلمه کلیدی هستند حذف خواهند گردید و برای اجرای آن تنها کفایت که NOT را قبل از عبارت یا کلمه کلیدی مورد نظران با يك فاصله بیاورید .

تکنیک های جستجو برای جستجوی بهینه

- 10 استفاده از عملگر نزدیک‌یابی: در بسیاری از موارد استفاده از عملگر AND باعث بازیابی اطلاعاتی می‌شود که برای ما مفید نیست. به این دلیل که این عملگر کلید واژه‌ها را در هر کجای متن که باشند بازیابی می‌کند. در این موارد استفاده از تکنیک نزدیک‌یابی می‌تواند از ریزش کاذب اطلاعات و یا بازیابی اطلاعات غیرمرتبط جلوگیری نماید. همه موتورهای جستجو قابلیت استفاده از این تکنیک را ندارند ولی به عنوان مثال در موتور جستجوی آلتاویستا می‌توان با استفاده از عملگر NEAR از این تکنیک استفاده نمود.
12. جستجوی ترکیبی با استفاده از پرانتز: این تکنیک یکی از مهمترین تکنیک‌های جستجو می‌باشد که به وسیله آن می‌توان تا حدود زیادی از بازیابی موارد غیرمرتبط در محیط وب جلوگیری کرد. در این روش می‌توان از همه عملگرهای جستجو که در بالا گفته شده یکجا استفاده کرد و آن‌ها را با هم‌دیگر ترکیب نمود.
- تکنیک‌های جستجو برای جستجوی بهینه
13. جستجوی کلیدواژه در عنوان صفحات وب: این تکنیک با این پیش فرض که عنوان یک صفحه وب تا حدود زیادی نمایانگر محتوای اطلاعات موجود در آن است به جستجوی واژه‌های کلیدی در عنوان سایت‌ها می‌پردازد. علامت آن در موتورهای جستجو متفاوت است ولی اغلب موتورهای جستجو از طریق فهرست انتخابی و یا گزینه‌های دیگر این امکان را فراهم می‌آورند.
14. جستجوی حوزه سایت‌ها: با توجه به این که به صورت قراردادی هر کشوری حوزه خاصی در محیط وب دارد، قابلیت جستجوی حوزه سایت‌ها به ما این امکان را می‌دهد که فرایند جستجو را به حوزه خاصی نظیر سایت‌های وب ایران (IR) و یا سایت‌های وب سازمان‌های غیر انتفاعی (ORG) محدود کنیم. دستورات استفاده از این تکنیک در موتورهای جستجو مختلف می‌باشد.
- تکنیک‌های جستجو برای جستجوی بهینه
15. محدود کردن جستجو به زبان‌های مختلف باعث می‌شود نتایج جستجو به زبان‌های دیگر آورده نشود و انتخاب مطلب مورد نظر آسان‌تر است.
16. محدود کردن جستجو به تاریخ انتشار منابع در وب: تاریخ انتشار یا به اصطلاح روزآمدی مطلب به خصوص در منابع علمی اصل مهمی است و این‌گونه محدودیت باعث می‌شود بنا به نیاز کاربر جدیدترین و یا قدیمی‌ترین منبع بازیابی بشود.

17. جستجوی رسانه‌های مختلف؛ موسیقی، عکس، ویدئو: زمانی که فقط نوع خاصی از رسانه مورد نیاز است به عنوان مثال زمانی که به عکس يك شخصیت نیاز داریم، جستجو در میان عکس‌ها باعث می‌شود نتیجه جستجو شامل اطلاعات دیگری در مورد آن شخصیت نباشد.

تکنیک‌های جستجو برای جستجوی بهینه

18. جستجوی صفحات با فرمت‌های مختلف: PDF, WORD, MP3, MPEG,: زمانی که فرمت خاصی مورد نظر است می‌توان از این تکنیک استفاده کرد. به عنوان مثال اگر مایل باشیم منبع بازیابی شده در فرمت PDF باشد، این تکنیک می‌تواند مفید باشد.

19. آگاهی از پیش‌فرض‌های جستجو در موتور جستجو: با توجه به این که هر موتور جستجو برای ترکیب واژه‌ها يك پیش‌فرض دارد و اگر از هیچ‌گونه عملگری استفاده نشود، کلید واژه‌ها را به صورت پیش‌فرض با یکی از عملگرهای جبر بولی ترکیب می‌کند؛ آگاهی از این پیش‌فرض موتورهای جستجوی مختلف مهارت ما را در جستجو بالا می‌برد.

20. وب نامریی: وب نامریی به دو دلیل کمی و کیفی اهمیت دارد کمی از این نظر که موتورهای جستجو فقط قادر هستند حدود 16 درصد از اطلاعات موجود در اینترنت را بازیابی کنند و اندازه وب نامریی تقریباً 500 برابر وب مرئی است و کیفی از این نظر که منابع اطلاعاتی موجود در وب عمیق معمولاً ارزشمند و مفید هستند و در بسیاری از موارد پاسخگویی نیاز کاربران می‌باشند.

آشنایی با ابزارهایی که برای شناسایی منابع وب نامریی به وجود آمده‌اند و کاربران را به سایت‌های مناسب راهنمایی می‌کنند، باعث دسترسی به این بخش عظیم از اطلاعات مفید و ارزشمند می‌شود. مثل سایت INVISIBLEWEB که فهرستی از منابع نامریی را و سایت COMPLETEPLASET که فهرستی از تقریباً 40000 پایگاه اطلاعاتی وب نامریی را ارائه می‌دهد.

جایگاه کتابخانه‌ها و مراکز اطلاع‌رسانی در جامعه اطلاعاتی

کتابداری و اطلاع‌رسانی جزو مشاغل اطلاعاتی است

کتابخانه‌ها به عنوان نهادهایی که وظایف گردآوری، سازماندهی و اشاعه اطلاعات را بر عهده دارند از عناصر اصلی **جامعه اطلاعاتی** و **زیرساخت‌های اطلاعاتی** محسوب می‌شوند.

جوامع توسعه یافته، کتابخانه های توسعه یافته تری نیز نسبت به جوامع در حال توسعه دارند.

راهنماهای موضوعی

کتابداران، اطلاع رسانان، آموزشگران و پژوهشگران کتابداری و اطلاع رسانی، صرفنظر از این که در چه حوزه ای فعالیت می کنند، می توانند با بهره گیری از این رویکردها، خدمات خود را نظام مند کرده و به فعالیت‌های خود جهت دهند.

راهنماهای موضوعی...

راهنماهای متون با عناوین مختلفی نظیر راهنماهای منابع اطلاعاتی و راهنماهای منابع مرجع و دروازه های موضوعی شناخته می شوند و منابع اطلاعاتی را در یک یا چند زمینه موضوعی، معرفی و ساختار آنها را توصیف می کنند.

منابعی نظیر *Walford's Guide to Reference* و *Guide to Reference Book*

Materials که کتابداران با آنها آشنا هستند، نمونه هایی از راهنماهای ارزشمندی هستند که تمام زمینه های موضوعی را پوشش می دهند.

راهنماهای موضوعی...

می توان این راهنماها را با استفاده از دروازه های موضوعی (*gateways*) مبتنی بر وب نیز تهیه کرد چون بیشتر پژوهشگران تمایل دارند منابع مرتبط با موضوعات تخصصی مورد علاقه شان را بدون صرف زمان زیادی، از طریق اینترنت بدست آورند و کتابخانه ها و کتابداران همواره درگیر توسعه و طراحی دروازه های موضوعی در وب بوده اند.

مطالعات کاربران

دو نوع اصلی از پژوهش های مربوط به مطالعات کاربران، **نیاز اطلاعاتی** و **رفتار اطلاع یابی** است.

جستجوی هدفمند اطلاعات بعنوان پیامد یک نیاز اطلاعاتی در راستای برآوردن یک یا چند هدف. جستجوی اطلاعات از سوی فرد ممکن است از طریق تعامل وی با نظام های اطلاعاتی دستی (نظیر یک روزنامه یا یک کتابخانه) یا از طریق تعامل با نظامهای رایانه محور (نظیر وب جهان گستر) صورت پذیرد

مطالعات کاربران

اهمیت دیگری که بررسی نیاز اطلاعاتی کاربران دارد در طراحی نظام های بازیابی اطلاعاتی است. این نظام ها را، بدون یافتن درکی روشن از آنچه استفاده کنندگان نیاز دارند یا می خواهند بدانند، چگونگی جستجوی اطلاعات توسط آنها و چگونگی ارزیابی اطلاعاتی که دریافت می کنند

مطالعات کتابسنجی

- مطالعات کتابسنجی
- اطلاع سنجی
- علم سنجی
- وب سنجی

افزایش حجم سریع اطلاعات، انفجار اطلاعات، آلودگی اطلاعات، اضافه با اطلاعاتی و ... سردرگمی متخصصان در تشخیص اطلاعات سره از ناسره

وب سنجی

- تحلیل محتوای صفحات وب
- تحلیل فناوریانه وب (عملکرد موتورهای کاوش)
- تحلیل استنادی وب
- بررسی ضریب تاثیرگذاری وب (web impact factor)

و ...

مهارتهای اطلاع یابی در محیط وب
با تأکید بر پایگاه های اطلاعاتی

مقدمه

افزایش روزافزون اطلاعات
انتشار منابع اطلاعاتی متعدد

جامعه علمی را با مشکلات و چالش های جدی روبرو ساخته است.

در میان این حجم انبوه اطلاعات (چاپی، غیر چاپی، و الکترونیکی) پیدا کردن اطلاعات کاملاً مرتبط و معتبر کار آسانی نیست و نیازمند مهارتهای ویژه ای است.

مهارت های اطلاعاتی پیش نیاز تحقیق است.

تعریف وب مرئی و نامرئی

تعریف فنی

وب مرئی

صفحات ثابت با دسترسی آزاد

وب نامرئی

تشکیل شده از صفحات وب متحرک و یا قابل دسترس به صورت محدود (مثلاً با استفاده از رمز

ورود ..)

تقسیم بندی وب مرئی و نامرئی

مهارت های اطلاع یابی

- مهارت در بازیابی اطلاعات؛
- مهارت در ارزیابی اطلاعات؛
- مهارت در سازماندهی اطلاعات؛
- و مهارت در تبادل ارتباط.

ابزارهای جستجو

روشهای جستجوی اطلاعات در محیط وب

- ساده

- پیشرفته
- کنترل واژگان
- اصطلاحنامه

فرایند جستجو

ارزیابی پایگاه ها

- جامعیت
- مانعیت

چگونه می توان مدارك بیشتر و مرتبط تر را بازیابی کرد؟

بخش های مختلف وب نامرئی

وب مات یا تاریک

۱ بخشی از فضای وب نامرئی به وب مات موسوم گردیده که می توانسته مورد استفاده کاربران قرار گیرد، اما به دلایل زیر این اطلاعات در خارج از دسترس کاربران قرار گرفته و موتورهای کاوش نمی توانند آن ها را بازیابی کنند :

۱. از آنجا که اولاً محیط وب دائماً در تغییر است و هر روز منابع و اطلاعات جدید به آن افزوده می گردد و ثانیاً صفحاتی در وب وجود دارند که هیچ پیوندی بین آن ها با منابع دیگر برقرار نشده، خزنده های موتورهای جستجو قادر به یافتن این صفحات و همگام نمودن خود با این حجم عظیم اطلاعات نیستند .

وب مات یا تاریک

۲. به دلیل محدودیت توانایی، نرم افزارهای خزنده فرصت کافی برای روزآمدسازی صفحات جدید وب را ندارند. موتورهای کاوش نیز امکان روزآمدسازی حجم عظیمی از اطلاعات و منابع جدید را ندارند و به همین دلیل بسیاری از این اطلاعات از حوزه موتورهای کاوش دور می مانند .

۳. محدودیت توان مالی بسیاری از موتورهای کاوش سبب گردیده که موتورهای کاوش نتوانند تمام صفحات وب سایت ها را نمایه سازی کنند، چرا که برای آن ها هزینه های زیادی دارد و

بنابراین موتورهای کاوش بنا بر سیاست های خودشان، تنها بخشی از وب سایت ها یا لایه های بیرونی آن ها را نمایه سازی می کنند. بنابراین همیشه بخش عظیم لایه های درونی وب سایت ها پنهان می مانند .

۲. وب عمیق

به مجموعه ای از اطلاعات الکترونیکی پیوسته اطلاق می شود که بسیاری از پایگاه های اطلاع رسانی، آن ها را از طریق شبکه جهانگستر وب در دسترس عموم قرار داده اند. برخی این اطلاعات را به رایگان، و برخی دیگر را با دریافت هزینه در دسترس عموم قرار می دهند.

مندرجات این پایگاه ها معمولاً خارج از حوزه جستجوی موتورهای کاوش قرار دارند هر یک از این پایگاه ها صفحه جستجوی مبتنی بر وب دارند. که امکان جستجو در آن ها برای کاربران را فراهم می کند، اما خزنده های موتورهای جستجو توان ورود به آن ها را ندارند و در نتیجه حجم انبوهی از اطلاعات، نمایه نشده باقی می ماند.

۲. وب عمیق

به عنوان نمونه اگر يك متخصص موضوعي (مثلاً يك دانشجوي رشته پزشکی) بخواهد خود را به موتورهای کاوش معمولی محدود کند و نتواند به پایگاه های اطلاعاتی تخصصی مراجعه نماید یا از وجود آن ها آگاه نباشد، از دسترسی به حجم انبوهی از اطلاعات محروم خواهد ماند. بنابراین کاربر باید در این موارد از طریق موتورهای جستجو، پایگاه های مرتبط با موضوع خود را شناسایی کند و سپس، جداگانه به جستجو در آن ها بپردازد تا از دسترسی به وب عمیق باز نماند .

۳. وب خصوصی و وب ملکی

بخشی دیگر از وب نامرئی وجود دارد که چون اطلاعات موجود در آن جزو دارایی های شخصی یا خصوصی سازمان ها یا افراد می باشد، از حوزه دسترسی موتورهای جستجو پنهان است. مثلاً در برخی از سازمان ها و مؤسسات خصوصی یا دولتی، به دلایل امنیتی از اطلاعات مربوط به مسائل کاری و سازمانی و پرسنلی خود حفاظت می کنند اجازه دسترسی به آن ها را به دیگران نمی دهند و فقط کسانی که دارای اسم کاربر و گذرواژه هستند می توانند از آن ها استفاده کنند؛ این

بخش، وب خصوصی محسوب می‌گردد (منصوریان، ۱۳۸۲). بخش دیگر، منابع اطلاعاتی از قبیل نشریات الکترونیکی مبتنی بر وب می‌باشند که دسترسی به آن‌ها از طریق پرداخت حق اشتراک و خرید محصولات اطلاعاتی شرکت‌های مختلف صورت می‌گیرد و «وب ملکی» نامیده می‌شود.

۴. اینترنت واقعاً پنهان

بخش دیگری از وب پنهان وجود دارد که بنا به مسائل فنی و ناکارآمدی ابزارهای جستجو، از دسترسی کاربران دورمانده است. بسیاری از موتورهای جستجو قادر به بازیابی اطلاعات متنی اچ تی ام ال هستند، ولی توانایی بازیابی فایل‌های پی‌دی‌اف را ندارند، یا به دلیل کمبود منابع مالی و فنی از جستجوی فایل‌های غیرمتنی صرف نظر کرده‌اند. بنابراین منابع اطلاعاتی متنوعی نیز در وب وجود دارند که تنها به دلیل محدودیت‌های فناوری‌های موتورهای جستجو، از حوزه کاوش آن‌ها و در نتیجه از دسترس کاربران دور مانده‌اند.

منابع موجود در وب نامرئی

بخش بزرگی از وب وجود دارد که عنکبوت‌های موتورهای جستجو آن‌ها را نمایه نمی‌کنند یا نمی‌توانند نمایه کنند و عبارت‌اند از سایت‌های دارای رمز عبور، اسناد موجود در پشت سامانه‌های حفاظتی، فایل‌های پی‌دی‌اف از متون آرشیو شده، و ابزارهای تعاملی نظیر ماشین حساب‌ها و برخی واژه‌نامه‌ها و همچنین محتویات بعضی از پایگاه‌های اطلاعاتی، منابع محافظت شده از طریق اسم کاربر و گذرواژه، منابع و صفحات وب بدون پیوند، و صفحات افزون بر حداکثر تعداد صفحات قابل مرور در نتایج بازیابی.

اهمیت وب نامرئی

به دو دلیل می‌توان گفت که وب نامرئی اهمیت دارد. نخست از نظر کمی، باید گفت که حجم اطلاعات موجود در این بخش خیلی بیشتر از سطح آشکار است [۱۷]. موارد زیر، اهمیت وب نامرئی را از نظر کمی نشان می‌دهند

۱. بهترین موتورهای کاوش فقط قادر هستند که حدود ۱۶ درصد از اطلاعات موجود در وب را بازیابی کنند و بنابراین ۸۴ درصد آن‌ها جزو وب نامرئی به حساب می‌آیند.

۲. اندازه وب نامرئی تقریباً ۵۰۰ برابر وب مرئی است: وب نامرئی ۵۵۰ میلیون سند، و وب مرئی تقریباً ۱ میلیون سند را دارا می باشد .

دوم اینکه از نظر کیفی، اطلاعات بخش های مختلف این مجموعه بویژه منابع اطلاعاتی موجود در وب عمیق، معمولاً منابع ارزشمند و مفید هستند و در بسیاری از موارد پاسخگویی نیاز کاربران می باشند. تقریباً بیش از نیمی از وب نامرئی را پایگاه های اطلاعاتی موضوعی تشکیل می دهند .

دلایل عدم بازیابی و نمایه سازی وب نامرئی توسط موتورهای کاوش

۱. دلایل فنی: بسیاری از موتورهای کاوش به دلیل محدودیت های نرم افزاری توانایی روزآمدسازی اطلاعات جدید وب را ندارند. باید یادآور شد که هنوز هیچ موتور کاوشی ادعا نکرده است که قادر به گسترش حوزه کاوش خود به تمام محیط وب می باشد و همیشه این موتورها يك گام از سرعت روزافزون اطلاعات عقب تر هستند .

۲. دلایل بودجه ای: همانطور که قبلاً اشاره شد فرآیند نمایه سازی تمام صفحات وب، هزینه بر خواهد بود و موتورهای کاوش نیز بنا به محدودیت بودجه ناگزیرند فقط بخشی از وب سایت ها را نمایه سازی کنند .

۳. دلایل اجتماعی و حقوقی: از آنجا که اطلاعات موجود در وب در دسترس عموم قرار می گیرد، بسیاری از افراد و سازمان ها به دلیل صرف بودجه های کلان در راه اندازی سایت ها و پایگاه های اطلاعاتی خود، حاضر نیستند این اطلاعات را به صورت رایگان در اختیار همه بگذارند. البته این از لحاظ اجتماعی و حقوقی حق مسلم آن ها است .

شیوه های اطلاع یابی در وب نامرئی

در حال حاضر ابزارهایی به وجود آمده اند که منابع وب نامرئی را شناسایی، و کاربران را به سایت های مناسب راهنمایی می کنند. این رویکرد توسط بزرگراه های اطلاعاتی و کتابخانه های مجازی پذیرفته شده است؛ بطوری که فقط توصیفی از پایگاه های اطلاعاتی و مجلات نامرئی را ارائه می کنند؛ مثل سایت **Invisibleweb** که فهرستی از منابع نامرئی را، و سایت **Completeplanet** که فهرستی از تقریباً ۴۰۰۰۰۰ پایگاه اطلاعاتی وب نامرئی را ارائه می دهند. برخی دیگر از ابزارهای اطلاع یابی نیز که تاکنون معرفی شده اند و ما می توانیم با استفاده از آن ها به این اطلاعات دسترسی پیدا کنیم به شرح زیر است :

دروازه های اطلاعاتی

- دروازه های اطلاعاتی مجموعه ای از پایگاه ها و سایت ها هستند که به وسیله متخصصان اطلاعاتی و معمولاً کتابداران گردآوری، بررسی و براساس موضوع مرتب شده اند و معمولاً به کاربران نیز توصیه می شوند. نمونه ای از این دروازه ها عبارت اند از :

- Academic Information
- Digital Librarian
- Gaary price direct search
- Infomine
- Internet public Library

پرتال ها یا پایگاه های اطلاعاتی خاص موضوعی

مجموعه ای از پایگاه های اطلاعاتی خاص موضوعی هستند که به يك موضوع خاص اختصاص دارند و به وسیله دانشمندان، محققان، متخصصان، مؤسسات دولتی، شرکت های بازرگانی و کارشناسان موضوعی، افراد بسیار علاقه مند یا دارای دانش حرفه ای و اطلاعات وسیع در حوزه خاص ایجاد می شوند (Oxford uni. Libraries) ۲۰۰۰. (از ورتال ها در هنگام جستجو برای موضوعات خاص مانند پیوندهای خبری، فایل های چنדרسانه ای، آرشیوها، فهرست های پستی اشخاص، شغل یاب ها و هزاران پایگاه اطلاعاتی که به موضوعات خاص اختصاص دارند استفاده می شود. از مزیت های عمده استفاده از دروازه های اطلاعاتی این است که برای ایجاد آن ها افرادی با دانش موضوعی خاص، در اینترنت جستجو کرده اند و به پالایش اطلاعات مفید از غیرمفید پرداخته اند .

ویژگی پایگاه های اطلاعاتی :

۱. جستجوپذیر و مرورپذیر، حاوی توصیف منابع اینترنتی در يك زمینه موضوعی خاص؛
۲. مشتمل بر معیارهای شفاف و تعریف شده برای ارزیابی کیفیت منابع اطلاعاتی، به جای انتخاب بدون ارزیابی؛

۳. دخالت کتابداران و متخصصان موضوعی در ایجاد آن؛
۴. تولید دستی رکوردهای آن، برای توصیف بامعنا و اطلاع بخش از منابع اطلاعاتی؛
۵. فهرست نویسی و رده بندی منابع اطلاعاتی با استفاده از شیوه های سنتی کتابخانه ای به منظور بازیابی موثر .

ابرموتورهای کاوش :

گسترش پذیری حوزه های جستجو نیز یکی از شیوه های دسترسی به وب نامرئی شمرده می شود که نمونه آن، استفاده از ابرموتورهای کاوش است. این ابرموتورها خود، موتورهای جستجوی واقعی نیستند. بلکه به کاربران این امکان را می دهند که کلیدواژه های خود را همزمان توسط چند موتور، مورد کاوش قرار دهند و نتایج جستجوی همه آن ها را با هم ببینند .

عوامل هوشمند :

این عوامل هوشمند از ابزارهای بازیابی در اینترنت هستند که برای اجرای کارهای بخصوص به کارگرفته می شوند. این عوامل توانایی جستجو و مقایسه و انتخاب منابع اطلاعاتی بر اساس نیاز مطرح شده توسط کاربر را دارند .

توصیه ها و ابزارهایی برای جستجوی وب نامرئی

- متخصصان جستجوی وب میگویند موتورهای جستجو مثل گوگل و یاهو فقط 1 درصد (surface Web) از اطلاعات وب را برای جستجو اطلاعات اینترنت را مورد جستجو قرار میدهند مابقی اطلاعات اینترنت را وب نامرئی یا Deepnet ،invisible Web ،dark Web ،hidden Web می نامند. خوشبختانه برای جستجو قسمتهایی از وب پنهان ابزارها و موتور جستجوهای متعددی ایجاد شده و در یافتن اطلاعات دلخواه کمک خوبی هستند. این موتور جستجوها از الگوریتم های پیشرفته رنگینگ و language-analysis استفاده می کنند که اطلاعات مفید و مرتبط را چکیده می کنند.

موتور جستجوهای وب نامرئی:

DeepDyve: یکی از جدیدترین موتور جستجوها که بر وب نامرئی تمرکز کرده است

CluserLook Search: جستجوگری برای اطلاعات پزشکی ، سلامت ، داروها و...
CompletePlanet: بیش از هفت هزار دیتابیس و موتور جستجو در این پایگاه در دسترس است و یکی از راه های عالی برای جستجوی وب نامرئی می باشد
Daylife: جستجوی اخبار به همراه تصاویر ، مقالات و...
spock: جستجوی افراد در وب ؛ پیدا کردن پروفایلها و تصاویر مخفی دوستانان
The WWW Virtual Library: یکی از قدیمی ترین پایگاه اطلاعاتی وب که با کلمات کلیدی یا دسته بندی ها می توانید به اطلاعاتش دست یابید

موتور جستجوهای وب نامرئی:

pipl: طراحی شده برای جستجو وب نامرئی ؛ نتایجی که این وب سایت در دسترس قرار می دهد از نظر بعضی از کاربران خطرناک است (حریم شخصی)
CustomSearchEngine: آیسیتی از جستجوگرهای سفارشی سازی شده گوگل
SurfWax: موتور جستجوی ساده و پر محتوا
Freebase: مجموعه ای از میلیونها بانک اطلاعاتی در موضوعات متنوع
RefSeek: موتور جستجویی برای دانش جویان و محققان ؛ بیشتر از یک میلیارد سند که شامل صفحات وب ، مجلات ، روزنامه ها ، کتاب ها و دایره المعارف را جستجو می کند

چه کنیم که اطلاعاتمان جزء وب نامرئی محسوب نشود؟

1-بازیابی صفحات وب عمیق

یک صفحه P پویا گفته می شود اگر بعضی یا تمام محتوای آن در زمان اجرا (زمان بعد از دریافت درخواست صفحه در خدمتگذار) توسط برنامه ای بر روی خدمتگذار یا مشتری تولید شود.

تشخیص فرمها

تشخیص اسکرپیتها

2-تشخیص فیلدهای فرم

کار بسیار دشواری است!

● اکثر تکنیکها براساس روشهای مکاشفه ای است.

خوشبختانه، در اکثر فرمها از عناصر یکسانی استفاده شده است.

3- تکمیل خودکار فیلدها

انتصاب مقدار مناسب به فیلهای استخراج شده برای کشف محتویات داخل داده‌پایگاه
تکنیک اول:

استفاده از پرس‌وجوهای از پیش تعیین شده

تکنیک دوم:

استفاده از مقادیر موجود در فیلدها بصورت جایگشتی

تکنیک سوم:

آموزش خزشگر با استفاده از پالایش صفحه‌ی جستجو

3- تکمیل خودکار فیلدها - ادامه

تکنیک چهارم:

استفاده از تکنیک‌های یادگیری ماشین

تکنیک پنجم:

تبدیل مسأله‌ی پیدا کردن بهترین پرس‌وجو برای یک داده‌پایگاه به مسأله‌ی پوشش مجموعه

در گراف‌ها

تبدیل به مسأله‌ی مجموعه‌ی غالب وزن‌دار کمینه

تکنیک ششم:

کار با اسکرپیت‌های سمت مشتری

4- آنالیز نتایج دریافتی از داده‌پایگاه‌ها

دلایل:

بدست آوردن کلمات کلیدی جدید

تخمین تعداد مستندات داده‌پایگاه در یک زمینه‌ی خاص

5- دسته‌بندی یا خوشه‌بندی داده‌پایگاه‌ها

دسته‌بندی یا خوشه‌بندی؟

روش‌های مختلف دسته‌بندی:

دسته‌بندی مبتنی بر پرس‌وجو

دسته‌بندی مبتنی بر خزش

دسته‌بندی با استفاده از توصیفات کلاس سرویس

