



دانشگاه صنعتی امیرکبیر

دانشکده مهندسی کامپیوتر و فناوری اطلاعات

پایان نامه کارشناسی ارشد رشته مهندسی کامپیوتر گرایش نرم افزار

بازیابی معنایی اطلاعات با استفاده از بسط مفاهیم حاصل از جستجوی مبتنی بر

کلید واژه

نگارش

وحید جلالی

استاد

دکتر محمدرضا مطش بروجردی



شماره مدرک:

فرم اطلاعات
پایان نامه
کارشناسی- ارشد
و دکترا

مشخصات دانشجو		نام خانوادگی : جلالی	نام : وحید	شماره دانشجویی: 85131049
عنوان		بازیابی معنایی اطلاعات با استفاده از بسط مفاهیم حاصل از جستجوی مبتنی بر کلید واژه		
Title	Semantic information retrieval using result concepts of a keyword based query			
استاد راهنما	نام خانوادگی : مطش بروجردی	درجه و رتبه	استاد راهنما	نام خانوادگی : درجه و رتبه
	نام : محمد رضا	استاد ديار		نام :
استاد مشاور	نام خانوادگی	درجه و رتبه	استاد مشاور	نام خانوادگی: درجه و رتبه
	نام:			نام:
دانشنامه	<input type="radio"/> ارشد • <input checked="" type="radio"/> دکترا <input type="radio"/> کارشناسی		سال تحصیلی : 87	
نوع پروژه	<input checked="" type="radio"/> کاربردی • <input type="radio"/> نظری • <input type="radio"/> توسعه ای • <input type="radio"/> بنیادی			
مشخصات ظاهری	تعداد صفحات 107	تصویر • جدول • نمودار • نقشه • واژه نامه •	تعداد مراجع 44	ضمائم • تعداد صفحات
زبان متن	<input checked="" type="radio"/> فارسی • <input type="radio"/> انگلیسی		چکیده	<input checked="" type="radio"/> فارسی • <input type="radio"/> انگلیسی
یادداشت				
توصیفگر				
کلید واژه فارسی	بازخورد ارتباطی مفهومی، بازیابی اطلاعات مبتنی بر هستان شناسی، گسترش جستار، معیارهای شباهت معنایی.			
Key word of English	<i>Concept-based Pseudo Relevance Feedback, Ontology-based Information Retrieval, Query Expansion, Semantic Similarity Measures.</i>			

تقدیم به

دوستانی که شناختم.

با تشکر از زحمات و حمایت‌های بی‌دریغ دکتر بروجردی عزیز

چکیده

اگرچه بازیابی اطلاعات به طور سنتی مستندات و جستارها را از نقطه نظر کلیدواژه‌های به کار رفته در آنها مورد بررسی قرار می‌دهد، در نظر گرفتن روابط معنایی آنها در کنار شباهت‌های نحویشان می‌تواند به ایجاد سیستم‌های کارتر در بازیابی اطلاعات منجر شود. به طور کلی می‌توان گفت وجود نداشتن واژگان مشابه میان دو متن دلیل بر نامربوط بودن آنها به یکدیگر نیست. حالت‌های متعددی ممکن است وجود داشته باشد که در آنها کاربر، نیاز اطلاعاتی خود را با لغات و واژگانی متفاوت از آنچه در مستندات سیستم وجود دارد بیان کند. و یا این که به علت کمبود اطلاعات قادر نباشد برخی کلید واژه‌های مناسب برای جستجوی خود را در اختیار سیستم قرار دهد. در تمامی این شرایط استفاده از روابط مفهومی میان مستندات و جستار کاربر، این امکان را به سیستم خواهد داد تا نتایج دقیق‌تر و کامل‌تری را به کاربران خود بازگرداند. هدف از این پایان نامه بررسی شیوه‌های مختلف به کارگیری دانش موجود در هستان‌شناسی‌ها در امر بازیابی اطلاعات و معرفی ایده‌های نوین در این رابطه می‌باشد. به طور خلاصه در این پایان نامه روش‌های جدیدی در زمینه کاربرد هستان‌شناسی در گسترش جستار، بازیابی معنایی مبتنی بر معیارهای شباهت مفهومی و بازیابی معنایی همراه با بازخورد ارتباطی خودکار مفهومی معرفی می‌شوند. دامنه مورد بررسی در این پایان نامه مستندات پزشکی می‌باشند. هستان‌شناسی به کار رفته در روش‌های معرفی شده MeSH یا Medical Subject Headings است که در برگیرنده در حدود ۲۴۰۰۰ مفهوم پرکاربرد در قلمرو پزشکی می‌باشد. همچنین برای مقایسه تاثیر به کار بردن هستان‌شناسی عمومی در مقابل هستان‌شناسی مختص دامنه در بخشی از تجربیات این پایان نامه از هستان‌شناسی Wordnet استفاده شده است. مجموعه آزمون‌های به کار رفته برای ارزیابی روش‌های معرفی شده نیز عبارتند از Medline و OHSUMED که از مجموعه آزمون‌های شناخته شده در ارزیابی سیستم‌های بازیابی اطلاعات می‌باشند.

کلمات کلیدی: بازخورد ارتباطی مفهومی، بازیابی اطلاعات مبتنی بر هستان‌شناسی، گسترش جستار، معیارهای شباهت معنایی.

۱- مقدمه.....	۸
۱-۱- تعاریف و مفاهیم اولیه.....	۸
۱-۱-۱- مقدمه‌ای بر وب‌معنایی و هستان‌شناسی.....	۸
۱-۱-۲- معرفی اجمالی ساختار وب‌معنایی.....	۸
۱-۱-۳- مقدمه‌ای بر بازیابی اطلاعات.....	۹
۱-۱-۴- مشکلات روش‌های سنتی بازیابی اطلاعات.....	۱۱
۲-۱- بازیابی معنایی اطلاعات.....	۱۲
۳-۱- معیارهای شباهت معنایی.....	۱۲
۴-۱- معیارهای شباهت معنایی در هستان‌شناسی MeSH.....	۱۸
۲- فعالیت‌های انجام شده در حیطه بازیابی معنایی اطلاعات.....	۲۲
۱-۲- گسترش مبتنی بر کلیدواژه جستار.....	۲۲
۲-۱-۱- استفاده از هستان‌شناسی‌های عمومی در گسترش جستار.....	۲۲
۲-۱-۲- استفاده از هستان‌شناسی‌های مختص دامنه در گسترش جستار.....	۲۷
۲-۲- استفاده از معیارهای شباهت معنایی در بازیابی اطلاعات.....	۳۷
۱-۲-۲- استفاده از معیارهای شباهت معنایی توسط هستان‌شناسی‌های عمومی.....	۳۸
۲-۲-۲- استفاده از معیارهای شباهت معنایی توسط هستان‌شناسی‌های مختص دامنه.....	۴۰
۳- معرفی تکنیک‌های جانبی به کار رفته.....	۴۵
۱-۳- تبدیل MeSH از قالب XML به پایگاه داده رابطه‌ای.....	۴۵
۲-۳- پردازش زبان طبیعی.....	۴۶
۳-۳- انطباق عبارات اسمی به مفاهیم MeSH.....	۴۷
۴-۳- معرفی مولفه بازیابی مبتنی بر کلیدواژه.....	۴۸
۴- معرفی روش‌های پیشنهادی.....	۵۱
۱-۴- روش‌های مبتنی بر استخراج خودکار مفاهیم.....	۵۱
۴-۱-۱- روش گسترش جستار.....	۵۱
۴-۱-۲- روش ترکیبی.....	۵۴
۴-۲- روش‌های مبتنی بر استخراج غیرخودکار مفاهیم.....	۵۹

- ۴-۲-۱- بازیابی معنایی اطلاعات با استفاده از مفاهیم استخراج شده توسط متخصص دامنه .. ۵۹
- ۴-۲-۲- روش بازخورد ارتباطی خودکار مفهومی ۶۱
- ۵- ارزیابی روش‌های پیشنهادی ۶۵
- ۵-۱- مجموعه آزمون Medline ۶۵
- ۵-۲- مجموعه آزمون OHSUMED ۶۵
- ۵-۳- نحوه ارزیابی روش‌های بازیابی اطلاعات معرفی شده ۶۶
- ۵-۴- معرفی معیارهای به کار رفته در ارزیابی روش‌ها ۶۸
- ۴-۱-۵- دقت ۶۸
- ۵-۲-۴- یادآوری ۶۸
- ۵-۳-۴- میانگین دقت کلی ۶۸
- ۵-۵- ارزیابی روش‌های پیشنهادی بر روی مجموعه آزمون Medline ۶۹
- ۵-۱-۵-۵- ارزیابی روش بازیابی اطلاعات مبتنی بر کلیدواژه ۶۹
- ۵-۲-۵-۵- ارزیابی روش بازیابی اطلاعات مبتنی بر مجموعه‌های مترادف WordNet ۶۹
- ۵-۳-۵-۵- ارزیابی روش گسترش جستار معرفی شده ۶۹
- ۵-۴-۵-۵- مقایسه روش‌های گسترش جستار و روش مبتنی بر کلیدواژه ۷۳
- ۵-۵-۵-۵- ارزیابی روش ترکیبی ۷۴
- ۵-۶-۵-۵- ارزیابی روش معنایی استفاده شده در روش ترکیبی ۷۴
- ۵-۷-۵-۵- مقایسه روش‌های مبتنی بر کلیدواژه، گسترش جستار و ترکیبی ۷۷
- ۵-۶-۵- ارزیابی روش‌های پیشنهادی بر روی مجموعه آزمون OHSUMED ۷۸
- ۵-۱-۶-۵- ارزیابی روش مبتنی بر کلیدواژه ۷۹
- ۵-۲-۶-۵- ارزیابی روش مبتنی بر گسترش جستار ۸۲
- ۵-۳-۶-۵- ارزیابی روش ترکیبی مبتنی بر استخراج خودکار مفاهیم ۸۵
- ۵-۴-۶-۵- ارزیابی روش ترکیبی مبتنی بر استخراج غیرخودکار مفاهیم ۸۵
- ۵-۶-۶-۵- ارزیابی روش بازخورد ارتباطی خودکار مفهومی ۸۸
- ۵-۷-۶-۵- ارزیابی روش بازخورد ارتباطی خودکار مفهومی برای حالت استخراج خودکار مفاهیم ۸۸
- ۵-۸-۶-۵- مقایسه روش‌های مبتنی بر کلیدواژه، ترکیبی و بازخورد ارتباطی خودکار مفهومی ... ۹۱

۹۳ ۷-۵- جمع بندی و نتیجه گیری

۱۰۰ مراجع

۱۰۵ واژه‌نامه

فهرست شکل ها

- شکل ۱-۱- معماری وب معنایی..... ۹
- شکل ۲-۱- دسته بندی روش های سنتی بازیابی اطلاعات و مدل های گسترش یافته آنها ۱۱
- شکل ۱-۲- جریان سیستمی ONTOSEARCH..... ۳۶
- شکل ۲-۲- ارزیابی روش های معرفی شده توسط MAO بر اساس دقت بر حسب یادآوری ۴۲
- شکل ۱-۵- نمودار دقت بر حسب یادآوری برای روش های گسترش جستار، WORDNET و مبتنی بر کلیدواژه ۷۳
- شکل ۲-۵- نمودار دقت بر حسب یادآوری برای روش های ترکیبی، گسترش جستار و مبتنی بر کلیدواژه ۷۷
- شکل ۳-۵- درصد بهبود یافت میانگین دقت کلی روش های آزمایش شده بر روی MEDLINE نسبت به روش مبتنی بر کلیدواژه..... ۷۸
- شکل ۴-۵- نمودار دقت بر حسب یادآوری برای روش های بازخورد ارتباطی خودکار مفهومی، ترکیبی و مبتنی بر کلیدواژه ۹۱
- شکل ۵-۵- درصد بهبود یافت میانگین دقت کلی روش های آزمایش شده بر روی OHSUMED نسبت به روش مبتنی بر کلیدواژه..... ۹۲
- شکل ۶-۵- نمودار دقت بر حسب یادآوری برای روش های آزمایش شده بر روی مجموعه یادگیری OHSUMED ۹۳

فهرست جدول ها

- جدول ۱-۱- معیارهای شباهت معنایی و خصوصیات آنها ۱۴
- جدول ۲-۱- امتیاز معیارهای شباهت معنایی در هستان‌شناسی MESH ۲۰
- جدول ۱-۲- ارزیابی روش‌های آزمایش شده توسط GONZALO ۲۴
- جدول ۲-۲- ارزیابی روش‌های آزمایش شده توسط GROOTJEN ۲۷
- جدول ۳-۲- ارزیابی روش‌های آزمایش شده توسط BODNER ۲۸
- جدول ۴-۲- ارزیابی روش‌های آزمایش شده توسط HERSH ۳۱
- جدول ۵-۲- ارزیابی روش‌های آزمایش شده توسط FU ۳۲
- جدول ۶-۲- ارزیابی روش‌های آزمایش شده با نام‌های مترادف ژنها توسط HERSH ۳۳
- جدول ۷-۲- ارزیابی روش‌های آزمایش شده توسط SRINIVASAN ۳۵
- جدول ۱-۵- قالب خروجی ابزار ارزیابی TREC-EVAL ۶۷
- جدول ۲-۵- ارزیابی روش مبتنی بر کلیدواژه بر روی MEDLINE ۷۰
- جدول ۳-۵- ارزیابی به کارگیری مجموعه‌های ترادف WORDNET بر روی MEDLINE ۷۱
- جدول ۴-۵- ارزیابی روش پیشنهادی گسترش جستار بر روی MEDLINE ۷۲
- جدول ۸-۲- مقایسه ONTOSEARCH با یک موتور جستجوی مبتنی بر کلیدواژه ۳۷
- جدول ۹-۲- نحوه تعیین ارتباط معنایی میان مفاهیم در روش RICHARDSON ۳۹
- جدول ۵-۵- ارزیابی روش ترکیبی بر روی MEDLINE ۷۵
- جدول ۶-۵- ارزیابی روش معنایی استفاده شده در روش ترکیبی بر روی MEDLINE ۷۶
- جدول ۸-۵- ارزیابی روش مبتنی بر کلیدواژه بر روی بخشی از OHSUMED ۸۰
- جدول ۹-۵- ارزیابی روش مبتنی بر کلیدواژه بر روی کل OHSUMED ۸۱
- جدول ۱۰-۵- ارزیابی روش گسترش جستار معرفی شده بر روی بخشی از OHSUMED ۸۳
- جدول ۱۱-۵- ارزیابی روش گسترش جستار معرفی شده بر روی کل OHSUMED ۸۴
- جدول ۱۲-۵- ارزیابی روش ترکیبی مبتنی بر استخراج خودکار مفاهیم بر روی OHSUMED ۸۶
- جدول ۱۳-۵- ارزیابی روش ترکیبی مبتنی بر استخراج غیر خودکار مفاهیم بر روی OHSUMED .. ۸۷
- جدول ۱۵-۵- ارزیابی روش بازخورد ارتباطی خودکار مفهومی بر روی OHSUMED ۸۹
- جدول ۱۶-۵- ارزیابی روش بازخورد ارتباطی خودکار مفهومی با مفاهیم استخراج شده خودکار بر روی OHSUMED ۹۰

جدول ۵-۱۷- مقایسه خصوصیات مجموعه آزمون MEDLINE و OHSUMED ۹۴

فصل اول

مقدمه و تعاریف

۱- مقدمه

در اثر افزایش حجم مستندات مربوط به زمینه‌های مختلف دانش بشری، نحوه ذخیره‌سازی و دسترسی به این مستندات تبدیل به مساله‌ای اساسی در دنیای امروز گشته‌است. چنانچه بتوان روش‌های نوینی ارائه نمود که نیاز روزافزون بشر به اطلاعات را با کیفیت بالاتری سیراب نمایند این روش‌ها از ارزش فوق العاده بالایی برخوردار خواهند بود. در این پایان نامه سعی شده است تا برای حوزه مستندات پزشکی، موتورهای جستجو با الگوریتم‌های معنایی نوینی معرفی شوند تا نیازهای اطلاعاتی کاربران را بهتر درک کرده و در نتیجه نتایج دقیق‌تری را برای جستجوهای کاربران به آنها بازگردانند.

۱-۱- تعاریف و مفاهیم اولیه

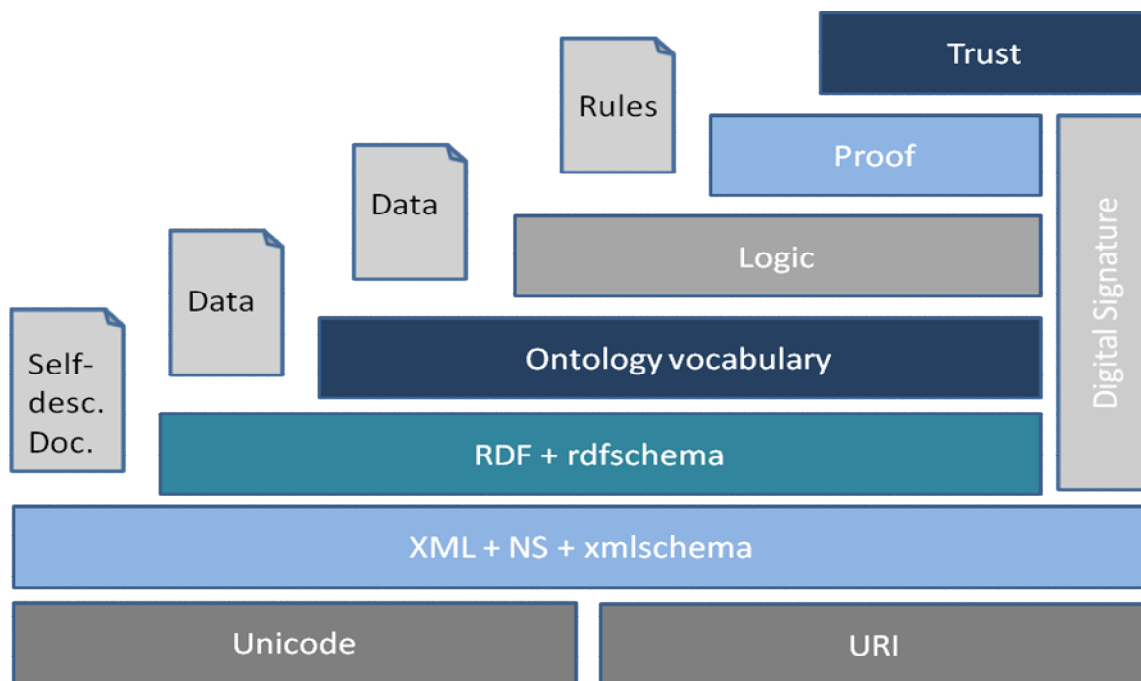
۱-۱-۱- مقدمه‌ای بر وب معنایی و هستان‌شناسی

از آنجا که بخش اعظم مستندات موجود در شبکه جهانی در حال حاضر در قالب HTML ذخیره شده اند این مستندات بیشتر مورد استفاده انسان قرار می گیرند تا عامل های نرم افزاری، به همین دلیل حیظه عملکرد و قدرت چنین عامل هایی در شبکه جهانی بسیار محدود و ناچیز است. با گسترش روز افزون محتوی شبکه جهانی و تعریف کاربرد های جدید برای آن به نظر می رسد که منابع و مستندات موجود در آن باید به گونه ای قابل فهم و تفسیر توسط ماشین ارائه شوند. علاوه بر این، روش‌های بازیابی اطلاعات در حال حاضر به علت گستردگی موضوعات موجود و حجم بالای مستندات، کارایی نسبتاً نامطلوبی دارند. وجود چنین مشکلاتی بود که موجب معرفی و ظهور وب معنایی در سال های پایانی هزاره دوم میلادی گردید. وب معنایی بر پایه هستان شناسی بنا شده است و کارایی آن تا حد زیادی وابسته به جامعیت، قدرت و ساختارمند بودن منابع هستان شناسی آن می باشد. از دیدگاه ساختاری وب معنایی را می توان به دو بخش موتور استنتاج و منابع هستان شناسی آن تقسیم کرد. آن چه در این جا بیشتر طرف توجه ما خواهد بود مسائل پیرامون هستان شناسی و موارد موجود در این راستاست.

۱-۱-۲- معرفی اجمالی ساختار وب معنایی

همان طور که در شکل ۱-۱ مشاهده می کنید [۴۲] وب معنایی بر پایه لایه های گوناگونی استوار است که یکی بر روی دیگری بنا شده است. هرچه از لایه های پایین این ساختار به سمت بالا حرکت کنیم حجم کار های انجام شده در آن زمینه کاهش یافته و مسائل حل نشده بیشتری باقی می ماند.

در سال‌های اخیر برای معرفی منابع موجود در وب استانداردهای فراوانی معرفی شده است. به عنوان نمونه می‌توان به استانداردهایی از قبیل (S) RDF، OIL، DAML و OWL اشاره کرد. استاندارد OWL که در حال حاضر از سوی W3C پیشنهاد می‌شود OWL است. OWL بر پایه (S) RDF بنا شده و از نوشتار XML سود می‌برد، ولی نسبت به (S) RDF امکانات بیشتری برای توصیف منابع دارد. OWL دارای سه نوع Lite، DL و Full می‌باشد. هر یک از این سه نوع دارای شرایط ویژه و دامنه استفاده مشخصی هستند. شایان ذکر است که هستان‌شناسی به کار رفته در این پژوهش به صورت استاندارد در قالب XML وجود دارد و چالش‌های مربوط به زبان‌های سطح بالاتر معرفی شده، مانند (S) RDF و یا OWL در مورد آن مطرح نمی‌شود، لیکن به دلیل معرفی حالت کلی وب‌معنایی در این بخش اشاره‌ای مختصر به ساختار کلی آن بی‌فایده نیست.



شکل ۱-۱- معماری وب‌معنایی

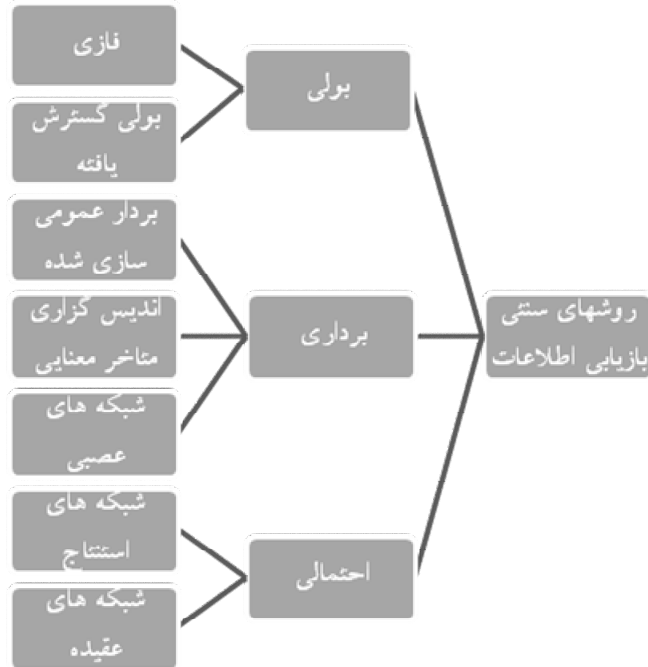
۱-۱-۳- مقدمه‌ای بر بازیابی اطلاعات

روش‌های سنتی بازیابی اطلاعات از طریق تطبیق کلیدواژه‌ها با یکدیگر مستندات را بازیابی می‌کنند. این نوع نگرش نسبت به بازیابی اطلاعات متضمن این مطلب است که معنای مستندات، توسط کلیدواژه‌های به کار رفته در آنها تعیین می‌شود. هرچند که چنین نگرشی به وضوح ساده سازی مساله

بازیابی اطلاعات است، کاربردهای بسیاری بر پایه همین دید در زمینه بازیابی اطلاعات شکل گرفته و طی سالیان متمادی مورد استفاده قرار گرفته‌اند. این ساده سازی در واقع به این علت رخ می‌دهد که جایگزینی یک مستند با تعدادی از کلیدواژه‌های به کار رفته در آن همواره باعث از دست دادن قسمتی از معنای مستند خواهد شد که توسط کلیدواژه‌های انتخاب شده قابل پوشش نمی‌باشد. در واقع به همین علت است که در اکثر موتورهای جستجوی سنتی همواره مقدار قابل توجهی از مستندات نامربوط در ازای جستار کاربر به او بازگردانده می‌شود و یا برخی از مستنداتی که به جستار کاربر مربوط هستند به طور کلی به عنوان مستندات مربوط، تشخیص داده نشده و به کاربر بازگردانده نمی‌شوند. در ادامه روش‌های اصلی بازیابی سنتی اطلاعات را مرور می‌کنیم.

سه روش اصلی بازیابی سنتی اطلاعات عبارتند از روش بولی، برداری و احتمالی. در روش بولی مستندات و جستارها به عنوان مجموعه‌هایی از کلیدواژه‌ها در نظر گرفته می‌شوند. بنابر این مدل بولی یک مدل از نوع تئوری مجموعه‌ها محسوب می‌شود. در روش برداری، مستندات و جستارها به عنوان بردارهایی در فضای t بعدی در نظر گرفته می‌شوند که بعدها این بردارها را کلیدواژه‌ها تعیین می‌نمایند. به همین دلیل روش برداری را از جمله روش‌های جبری در نظر می‌گیرند. در روش احتمالی نیز چارچوب مدل کردن مستندات و جستارها، یک چارچوب مبتنی بر احتمال است و به همین دلیل همان طور که از نام آن بر می‌آید یک روش احتمالی محسوب می‌شود.

با گذشت زمان مدل‌های جایگزینی برای هر یک از این سه روش سنتی معرفی شده است. بر طبق مدل‌های جایگزین مجموعه‌ای، مدل‌های فازی و گسترش یافته بولی معرفی شده‌اند. همچنین روش‌های برداری عمومی‌سازی شده، اندیس‌گذاری متاخر معنایی و مدل‌های شبکه عصبی بر اساس مدل‌های جایگزین جبری ارائه شده‌اند. و در پایان بر اساس مدل‌های جایگزین احتمالی، مدل‌های شبکه‌های استنتاج و شبکه‌های اعتقاد معرفی می‌شوند. در شکل ۱-۲ می‌توانید یک دسته‌بندی از این مدل‌های بازیابی اطلاعات را مشاهده نمایید [۴۳].



شکل ۱-۲- دسته بندی روش های سنتی بازیابی اطلاعات و مدل های گسترش یافته آنها

۱-۱-۴- مشکلات روش های سنتی بازیابی اطلاعات

از آنجا که در روش های سنتی بازیابی اطلاعات، مستندات و جستارها با کلیدواژه های موجود در آنها نمایش داده می شوند، حالت های متعددی به وجود می آید که جستار و مستند راجع به یک مفهوم خاص صحبت می کنند ولی به علت به کار بردن واژه های متفاوت برای اشاره به یک مفهوم، توسط موتورهای بازیابی سنتی اطلاعات، نامربوط تشخیص داده می شوند. همچنین موارد بسیاری وجود دارد که کاربر به علت عدم آگاهی از فرهنگ واژگان به کار رفته در یک حوزه، قادر نخواهد بود نیاز اطلاعاتی خود را به دقت بیان کند. برای حل چنین مشکلاتی بود که روش های بازیابی معنایی اطلاعات در مقابل روش های سنتی مطرح شدند.

جستجوی مفهومی به معنای جستجو بر اساس معنا به جای رشته ای از حروف، انگیزه حجم زیادی از تحقیقات و مطالعات در زمینه بازیابی اطلاعات بوده است. این دید نسبت به بازیابی اطلاعات در زمینه های شناخته شده ای مانند اندیس گذاری متاخر معنایی، مفهومی سازی زبانی و استفاده از فرهنگ های جامع و رده بندی ها دیده می شود.

۲-۱- بازیابی معنایی اطلاعات

بازیابی معنایی اطلاعات در برابر بازیابی سنتی اطلاعات مطرح می‌شود و هدف از آن، از میان برداشتن محدودیت‌ها و مشکلاتی است که در بازیابی سنتی اطلاعات وجود دارد. یکی از روش‌های به کار رفته در بازیابی اطلاعات استفاده از هستان‌شناسی برای بهبود دقت موتورهای جستجو است. استفاده از هستان‌شناسی در این راستا می‌تواند به دو صورت کلی انجام پذیرد. روش اول استفاده از هستان‌شناسی برای افزودن کلید واژه‌های مناسب به جستار اولیه کاربر است و روش دوم معرفی معیاری از شباهت معنایی میان مفاهیم برای استفاده در بازیابی اطلاعات در کنار شباهت‌های مبتنی بر کلیدواژه‌هاست.

۳-۱- معیارهای شباهت معنایی

برای تعیین شباهت‌های معنایی میان عبارات مختلف، روش‌های مختلفی وجود دارد که به چهار دسته زیر تقسیم می‌شوند:

- روش‌های شمارش یال: شباهت میان دو عبارت یا مفهوم را به عنوان تابعی از مسیر اتصال دهنده آن دو و موقعیت مفاهیم در یک رده بندی محاسبه می‌کند.
- روش‌های محتوای اطلاعات: تفاوت محتوای اطلاعات دو مفهوم یا عبارت را به عنوان تابعی از احتمال حضور آنها در دسته‌متن در نظر می‌گیرد. این احتمال حضور می‌تواند وابسته به دسته‌متن یا مستقل از آن باشد. به عنوان مثال چنان چه از نظر آماری تکرار یک عبارت در دسته‌متن مورد بررسی قرار بگیرد و احتمال بر پایه آن محاسبه شود، معیار حاصل وابسته به دسته‌متن خواهد بود و چنانچه این احتمال با توجه به موقعیت مفاهیم در یک سلسله مراتب که خود مستقل از دسته متن ایجاد شده است محاسبه گردد، معیار حاصل مستقل از دسته‌متن خواهد بود.
- روش‌های مبتنی بر ویژگی: شباهت دو عبارت را به عنوان تابعی از ویژگی‌های آنها (به عنوان مثال تعریف یا Glosses در Wordnet و یا scope note در MeSH) و یا مبتنی بر رابطه‌اشان با دیگر عبارات مشابه در رده‌بندی در نظر می‌گیرد. در این حالت تعداد بیشتر ویژگی‌های یکسان، شباهت معنایی را افزایش داده و برعکس تعداد ویژگی‌های غیریکسان، شباهت معنایی را کاهش می‌دهد.

- روش‌های ترکیبی: این روش‌ها ایده‌های موجود در روش‌های بالا را با یکدیگر ترکیب می‌کنند. به عنوان مثال شباهت عبارات را بر مبنای تطابق مترادف‌ها، عبارات همسایه و ویژگی‌های عبارات محاسبه می‌نمایند.

از منظری دیگر روش‌های تشخیص شباهت معنایی به دو دسته تقسیم می‌شوند:

- دسته اول روش‌هایی هستند که مفاهیم مورد نظر در یک هستان‌شناسی را با یکدیگر مقایسه می‌کنند.
- دسته دوم روش‌هایی هستند که عبارات مختلف از دو هستان‌شناسی متفاوت را با یکدیگر مقایسه می‌کنند.

یک نکته مهم در مورد معیارهای شباهت‌های معنایی مختلف این است که اکثر آن‌ها بالاترین شباهت را به مفاهیمی اختصاص می‌دهند که اولاً در هستان‌شناسی مربوطه، به یکدیگر نزدیک هستند و ثانياً نسبت به مفاهیمی با فاصله مشابه، در قسمت عمیق‌تری از سلسله مراتب واقع شده باشند (عبارات خاص‌تری باشند) [۳۸].

روش‌های مبتنی بر شمارش یال و محتوای اطلاعات از طریق استفاده از ساختارهای اطلاعاتی مانند موقعیت قرارگیری عبارات عمل می‌کنند و روش‌های مبتنی بر محتوای اطلاعات بیشتر برای عباراتی از یک هستان‌شناسی یکسان مناسب هستند. علت این مطلب آن است که ساختار و محتوای اطلاعات هستان‌شناسی‌های مختلف به طور مستقیم قابل مقایسه نیستند. از طرف دیگر روش‌هایی که از دو هستان‌شناسی برای معرفی معیار شباهت معنایی سود می‌برند معمولاً از روش‌های ترکیبی و یا مبتنی بر ویژگی استفاده می‌کنند.

شما می‌توانید در جدول ۱-۱ اطلاعاتی راجع به چند روش شناخته شده محاسبه شباهت معنایی را مشاهده کنید. در این جدول اطلاعاتی از قبیل تاثیر موقعیت مفاهیم در هستان‌شناسی، تاثیر ویژگی‌های یکسان بر افزایش یا کاهش شباهت، متقارن بودن معیار معرفی شده و نرمال بودن آن ذکر شده‌اند.

جدول ۱-۱- معیارهای شباهت معنایی و خصوصیات آنها

روش	نوع روش	افزایش با اشتراک	کاهش با تفاوت	تقارنی	نرمال شده در [۰۱]	موقعیت در سلسله مراتب
Rada [۱]	شمارش یال	آری	آری	آری	آری	خیر
Wu [۲]	شمارش یال	آری	آری	آری	آری	آری
Li [۳]	شمارش یال	آری	آری	آری	آری	آری
Leacock [۴]	شمارش یال	خیر	آری	آری	خیر	آری
Richardson [۵]	شمارش یال	آری	آری	آری	آری	آری
Resnik [۶]	محتوای اطلاعات	آری	خیر	آری	خیر	آری
Lin [۷]	محتوای اطلاعات	آری	آری	آری	آری	آری
Lord [۸]	محتوای اطلاعات	آری	خیر	آری	آری	آری
Jiang [۹]	محتوای اطلاعات	آری	آری	آری	خیر	آری
Tversky [۱۰]	ویژگی	آری	آری	خیر	آری	خیر
Rodriguez [۱۱]	ترکیبی	آری	آری	خیر	آری	خیر

در ادامه توضیحی مختصر راجع به هریک از این روش‌ها ارائه می‌شود.

- معیار Rada: یکی از قدیمی‌ترین معیارهای شباهت معنایی است که از روی فاصله بین دو مفهوم در یک سلسله مراتب میزان شباهت آنها به یکدیگر را محاسبه می‌کند.
- معیار Wu: Wu و Palmer معیار شباهت معنایی خود را برای دو مفهوم C_1 و C_2 به صورت زیر تعریف می‌کنند:

$$sim_{Wu\&Palmer}(C_1, C_2) = \frac{2 \times N_3}{N_1 + N_2 + 2 \times N_3} \quad (1-1)$$

که در آن N_1 تعداد گره‌های موجود در مسیر C_1 به خاص‌ترین مفهومی است که ابر مفهوم C_1 و C_2 است. N_2 نیز تعداد گره‌های موجود در مسیر C_2 به خاص‌ترین مفهومی است که ابر مفهوم C_1 و C_2 است و همچنین N_3 تعداد گره‌های موجود در مسیر خاص‌ترین ابر مفهوم C_1 و C_2 به ریشه هستان‌شناسی می‌باشد.

- معیار Li : این معیار از طریق رابطه‌ای که در آن فاصله دو مفهوم در هستان‌شناسی و عمق خاص‌ترین ابر مفهوم مشترکشان را در نظر می‌گیرد به صورت زیر محاسبه می‌شود:

$$sim_{Li}(c_1, c_2) = e^{-\alpha L} \times \frac{e^{\beta H} - e^{-\beta H}}{e^{\beta H} + e^{-\beta H}} \quad (2-1)$$

در این رابطه α و β پارامترهای غیر منفی هستند که نحوه تعامل طول مسیر میان دو مفهوم و عمق گره مشترک در برگیرنده آنها را مشخص می‌کند. L نشان‌دهنده طول مسیر میان دو مفهوم C_1 و C_2 است و H نیز عمق خاص‌ترین ابر مفهوم مشترک C_1 و C_2 را نشان می‌دهد.

- معیار Leacock: در این معیار که در رابطه (۳-۱) معرفی شده از سلسله مراتب IS-A و کوتاه‌ترین مسیر میان نام‌ها در این سلسله مراتب به عنوان معیاری برای محاسبه شباهت معنایی آنها استفاده می‌شود.

$$sim_{Leacock}(c_1, c_2) = -\log\left(\frac{N_p(c_1, c_2)}{2D}\right) \quad (3-1)$$

در این رابطه، $NP(c_1, c_2)$ نشان دهنده کوتاه‌ترین مسیر بر حسب گره بین مجموعه ترادف‌های C_1 و C_2 است. و D بیشینه عمق رده‌بندی را نشان می‌دهد.

- معیار Richardson: در این معیار شباهت معنایی به صورت شباهت کلاسی بیان می‌شود. و به طور کلی شباهت بین دو کلاس از روی محتوای اطلاعات اولین (خاص‌ترین) کلاسی که هر دو آنها را در بر می‌گیرد تخمین زده می‌شود. همچنین محتوای اطلاعات یک کلاس نیز از روی احتمال حضور آن در یک دسته متن بزرگ تخمین زده می‌شود. در رابطه (۴-۱) شما می‌توانید نحوه محاسبه معیار Richardson را مشاهده کنید:

$$sim_{Richardson}(c_1, c_2) = \max_{ci} [\log \frac{1}{P(c_i)}] \quad (4-1)$$

که در آن C_i مجموعه کلاس‌هایی است که هر دو کلاس c_1 و c_2 را در بر می‌گیرند و $P(c_i)$ کلاس احتمال کلاس c_i است و همچنین $\log 1/p(c_i)$ محتوای اطلاعات کلاس c_i را نشان می‌دهد.

- معیار Resnik: یک معیار شباهت معنایی برای سلسله مراتبی از نوع IS-A است که بر اساس محتوای اشتراکی اطلاعات معرفی می‌شود. در رابطه زیر شما می‌توانید نحوه محاسبه این معیار را مشاهده کنید:

$$sim_{Resnik}(c_1, c_2) = \max_{c \in S(c_1, c_2)} [-\log p(c)] \quad (5-1)$$

که در آن $S(c_1, c_2)$ مجموعه کلاس‌های در برگیرنده c_1 و c_2 هستند و $p(c)$ نیز بیانگر احتمال برخورد با نمونه‌ای از مفهوم c است.

- معیار Lin: در این معیار شباهت بین دو نمونه از مفاهیم بر حسب احتمال حضور آن مفاهیم در کل دسته متن و از روی رابطه زیر محاسبه می‌شود:

$$sim_{Lin}(x_1, x_2) = \frac{2 \times \log P(c_0)}{\log P(c_1) + \log P(c_2)} \quad (6-1)$$

در این رابطه $p(c_0)$ ، $p(c_1)$ و $p(c_2)$ از روی آمارهای مربوط به این مفاهیم در یک دسته‌متن برچسب‌گذاری شده بدست می‌آیند. به این صورت که تعداد دفعات تکرار این مفاهیم در دسته‌متن محاسبه شده و با توجه به کل مفاهیم موجود در متن نرمالایز می‌شوند. بدیهی است که در این رابطه c_1 و c_2 مفاهیم مورد نظر و c_0 خاص‌ترین مفهوم در برگیرنده c_1 و c_2 است.

- معیار Lord: این معیار از روی رابطه زیر محاسبه می‌شود:

$$sim_{Lord}(c_1, c_2) = -\ln Pms(c_1, c_2) \quad (7-1)$$