

بررسی روش های استخراج اطلاعات مبتنی بر یادگیری ماشین و مهندسی دانش

سمیه حیدری^۱، زهره بنائیان^۲، وحیده رشادت^{۳*}

۱- گروه کامپیوتر، دانشکده کامپیوتر، موسسه آموزش عالی پویش، قم، ایران

۲- گروه کامپیوتر، دانشکده کامپیوتر، موسسه آموزش عالی پویش، قم، ایران

۳- پژوهشکده فناوری اطلاعات، دانشگاه صنعتی مالک اشتر، تهران، ایران

چکیده

زبان طبیعی ابزار اصلی بیان مفاهیم برای انسان است و گفتار یا نوشتار را می توان رسانه انتقال آن دانست. در مقابل ماشین با این زبان نا آشنا است و نیازمند ترجمه آن به حوزه اطلاعات است. با توجه به اینکه در دهه های اخیر اطلاعات متنی در اینترنت رشد سریعی داشته و بخش قابل توجهی از این اطلاعات (اخبار آنلاین، مقالات علمی و کتب و...) به صورت غیر ساخت یافته و ناهمگن می باشد و اطلاعات غیر ساخت یافته قابل خواندن، سازماندهی و تحلیل توسط ماشین ها نیستند. برای اینکه بتوان از بین این حجم عظیم اطلاعات، انسان را در فهم و یافتن اطلاعات مورد نیاز یاری کرد باید بتوان متن غیر ساخت یافته را به اطلاعات ساخت یافته تبدیل کرد. در نتیجه وجود فناوری استخراج اطلاعات الزامی است. سیستم های استخراج اطلاعات با تبدیل اطلاعات به صورت ساخت یافته فهم آن را برای ماشین آسان و به انسان در درک بهتر این اطلاعات کمک می کنند. در این راستا استخراج اطلاعات دو روش کلی یادگیری ماشین و مهندسی دانش را معرفی می کند. یادگیری ماشین شامل روش های با ناظر که با مقادیر زیادی داده آموزشی کار می کنند، روش های بدون ناظر اغلب از خوشه بندی استفاده می کند و روش های نیمه ناظر با استفاده از هسته ها بوجود آمدند. روش استخراج آزاد اطلاعات نیز یکی دیگر از روش ها است که در مواردی که هدف کشف همه حقایق برجسته از متن بزرگ و متنوع است استفاده می شود. مهندسی دانش نیز شامل روش های استخراج الگو که از رویکرد تطبیق الگو برای شناسایی بخش لازم متن استفاده می کند و روش مبتنی بر قاب که در آن متون هم موضوع در یک خوشه قرار می گیرند تا برای هر کدام از آن ها قالبی با نقش های معنایی مشخص شود. هدف ما در این مقاله نقد و بررسی روش های استخراج اطلاعات (یادگیری ماشین و مهندسی دانش) می باشد.

کلمات کلیدی: استخراج اطلاعات، استخراج آزاد اطلاعات، یادگیری ماشین، مهندسی دانش

¹ Email: somayeheidary67@gmail.com

² Email: zbanayian@yahoo.com

* Email: reshadat@mut.ac.ir



۱. مقدمه

در دهه های اخیر اطلاعات متنی در اینترنت رشد سریعی داشته و بخش قابل توجهی از این اطلاعات (اخبار آنلاین، مقالات علمی و کتب و...) به صورت غیرساخت یافته و ناهمگن می باشد واطلاعات غیرساخت یافته قابل خواندن، سازماندهی و تحلیل توسط ماشین ها نیستند. برای اینکه بتوان از بین این حجم عظیم اطلاعات، انسان را در فهم و یافتن اطلاعات مورد نیاز یاری کرد باید بتوان متن غیرساخت یافته را به اطلاعات ساخت یافته تبدیل کرد. در نتیجه وجود فناوری استخراج اطلاعات^۱ الزامی است. سیستم های استخراج اطلاعات با تبدیل اطلاعات به صورت ساخت یافته فهم آن را برای ماشین آسان و به انسان در درک بهتر این اطلاعات کمک می کنند [۳، ۱].

برای مثال در جمله :

John is a graduate student at the University of Pennsylvania

سامانه استخراج اطلاعات John را بعنوان "person" و the University of Pennsylvania را بعنوان "University" برمی گرداند، student-of نیز رابطه بین John و University of Pennsylvania می باشد. student-of و Person می باشد. برچسب هایی است که از قبل توسط طراح سامانه تعیین شده اند.

با توجه به مثال بالا متوجه می شویم که دو وظیفه اصلی استخراج اطلاعات شامل استخراج موجودیت و استخراج رابطه می باشد [۲].

با استفاده از سامانه های استخراج اطلاعات می توان پایگاه دانشی ساخت یافته از متون ایجاد کرد. در این راستا استخراج اطلاعات با روش های مبتنی بر یادگیری ماشین (بناظر^۲، بدون ناظر^۳ و نیمه ناظر^۴) و مهندسی دانش تلاش در شناسایی حقایق دارند.

استخراج آزاد اطلاعات^۵ نیز روشی است که برای استخراج نمونه های رابطه در متون بزرگ مانند وب مورد استفاده می باشد و برخلاف روش های پیشین استخراج اطلاعات، استخراج همه روابط دلخواه از جملات موجود در متن را فراهم می کند [۱].

در دهه های اخیر شاهد تکثیر سریع اطلاعات متنی در منابع بی شمار روی اینترنت هستیم. بیشتر آن ها اطلاعات متنی غیر ساخت یافته اند و جستجو در آن ها مشکل است. این نیاز به رویکردهایی برای کشف دانش با ارزش از آن ها به صورت ساخت یافته را نشان می دهد که به ظهور فناوری استخراج اطلاعات منجر می شود. بنابراین استخراج اطلاعات، شناسایی مفاهیم از پیش تعریف شده و نادیده گرفتن اطلاعات بی ربط است [۴]. با تولید مجموعه عظیمی از اطلاعات متنی استخراج اطلاعات یک میدان مهم و در حال رشد است. مطالعات وسیعی در انواع تحقیقات مرتبط شامل پردازش زبان طبیعی و بازیابی اطلاعات انجام شده است [۵].

با توجه به اینکه بسیاری از اطلاعات از دست رفته در شکل متن باز روی صفحات وب در دسترس است پس برای استخراج رابطه روش های پردازش متن الزامی است [۶]. به همین منظور روش های کلی استخراج رابطه شامل مهندسی دانش و یادگیری ماشین معرفی می شود که روش های یادگیری ماشینی عبارتند از با ناظر، بدون ناظر و نیمه ناظر. روش های مهندسی دانش نیز شامل استخراج مبتنی بر قاب و استخراج مبتنی بر الگو می شود. استخراج اطلاعات آزاد نیز یکی دیگر از روش های استخراج رابطه است در مواردی که هدف کشف همه حقایق برجسته از متن بزرگ و متنوع است از این روش استفاده می شود.

- 1 information extraction
- 2 supervised
- 3 Un supervised
- 4 Semi supervised
- 5 Open information extraction

۲. استخراج اطلاعات

تشخیص و طبقه بندی روابط از پیش تعریف شده بین موجودیت های مشخص شده در متن [۴].

- رابطه بین یک فرد و یک سازمان

Steve Jobs works for Apple
Employee of (Steve Jobs ,Apple)

- رابطه بین یک فرد و محل

Mr. Smith gave a talk at the conference in New York
Located In (Smith ,New York)

- رابطه بین دو شرکت

Listed broadcaster TVN said its parent company, ITI Holdings, is considering various options for the potential sale. Subsidiary Of (TVN ,ITI Holding)

۲-۱. روش های یادگیری ماشین

یادگیری ماشینی^۱ یک بخش اساسی از هوش مصنوعی است که توانایی رفتار هوشمندانه بشر را به افزایش دانش و برای حل مشکلاتی که قبلا در اکثر سیستم های استخراج اطلاعات با آن مواجه می شد افزایش می دهد. بسیاری از سیستم های یادگیری ماشینی کنونی با توانایی شناختن خواص شناخته شده از داده مشخص شده اند و معمولا با رویکردهای یادگیری نظارتی هستند [۵].

۲-۱-۱. روش های باناظر

روش های باناظر با داده آموزشی کم کار می کنند و به دو دسته اصلی که در ادامه بیان می شود تقسیم می شود.

- دسته بندی براساس ویژگی^۲

این روش معمول استخراج رابطه است که مشکل دسته بندی را حل می کند. به طور خاص، هر جفت موجودیتی که در جمله اتفاق می افتد بعنوان نامزد مطرح می شود. هدف تخصیص برچسب کلاس به جفت موجودیت است که برچسب کلاس یک رابطه از پیش تعریف شده برای جفت موجودیت نامرتبط است. مهندسی ویژگی گام مهم در روش دسته بندی است.

- روش هسته^۳

مهم ترین کار در استخراج رابطه دسته بندی براساس هسته است. در یادگیری ماشینی، یک تابع هسته یا کرنل محصول داخلی نمونه های مشاهده شده را در بعضی زیر لایه های فضای برداری تعریف می کند. مزیت عمده استفاده از هسته این است که موارد مشاهده شده برای محاسبه شدن لازم نیست به صراحت به فضای برداری محصولات داخلی خود نگاشت شود [۷].

¹ Machine learning

² Feature-based classification

³ Kernel method

۲-۱-۲. روش های نیمه ناظر

روش های مبتنی بر هسته و مبتنی بر ویژگی از مقادیر زیاد داده آموزشی برای استخراج اطلاعات استفاده می کند. برای حل این مسئله روش یادگیری نیمه ناظر معرفی شد که با داده آموزشی کم کار می کند. در ادامه به بررسی روش های نیمه ناظر پرداخته خواهد شد.

• خودراه انداز

روش های مبتنی بر هسته و مبتنی بر ویژگی از مقادیر زیاد داده آموزشی برای استخراج اطلاعات استفاده می کند. برای حل این مسئله روش یادگیری نیمه ناظر معرفی شد که با داده آموزشی کم کار می کند. روش مهم در یادگیری نیمه ناظر روش خودراه انداز است که از مجموعه کوچک نمونه های رابطه شروع می شود و از الگوهای استخراج استفاده می کند. [۷].

سامانه اسنوبال توسط آگیتین^۱ و گراوانو^۲ [۸] برای استخراج اطلاعات در روش خودراه انداز عرضه شد. ایده این سامانه ساده است و با جفت موجودیت های مرتبط با رابطه هدف شروع می کند و در متن به جستجوی جفت موجودیت های مجاور هست، اگر جفت موجودیت ها بطور همزمان در متن رخ داده باشند، مفهوم همزمانی موجودیت ها احتمالاً به معنی الگویی برای رابطه هدف است. پس جفت موجودیت ها به نمونه رابطه اضافه می شوند و تا زمانی که شرایط دقیقی ایجاد شود پردازش ادامه دارد بطوریکه بیشتر الگوها و موجودیت ها به نتایج پردازش اضافه می شوند. یک گام مهم در شیوه خودراه انداز ارزیابی کیفیت الگوهای استخراج است در نتیجه فرآیند استخراج شامل الگوهای خراب نمی شود [۶].

• نظارت راه دور

با رشد وب اجتماعی بیشتر دانش انسان توسط کاربران زیادی در پایگاه های اطلاعات ذخیره می شود. نمونه کاملاً شناخته شده آن ویکی پدیا است. در این حالت ممکن است مجموعه بزرگ از موجودیت ها برای رابطه هدف باشند تا داده آموزشی تولید شود. روش نظارت راه دور ویژگی های استخراج شده از جملات متفاوت شامل هر جفت موجودیت را برای ایجاد بردار ویژگی غنی استفاده می کند [۱].

روش نظارت راه دور یا یادگیری خودنظارتی برای استخراج پایگاه های دانش بزرگ برای برچسب زدن خودکار موجودیت های در متن و استخراج ویژگی ها و آموزش دادن دسته بند بکار می رود. این روش فقط برای استخراج روابطی که از مرز جملات رد نشده اند و جملاتی که حاوی اشاره روشنی از فعل و فاعل رابطه است استفاده می شوند. استخراج نظارت راه دور بعنوان برچسب گذاری خودکار متن با خصوصیات و منابعی که منابع موجودیت ها از یک پایگاه دانش هستند استفاده می شوند. اگر دو موجودیت در یک رابطه باشند، هر جمله شامل این دو موجودیت ممکن است این رابطه را بیان کند. قبل از استفاده برچسب گذاری خودکار متن برای آموزش دسته بند، نمونه های حاوی واژگان مبهم کشف و دور انداخته می شوند. اولین رویکرد این است که اگر لغات مبهم هستند آن ها را از اشیا برای موجودیت هدف دور می اندازیم. در واقع اگر واژگان موضوع در سراسر کلاس مبهم باشد در نتیجه موضوع در کلاس خاص مبهم است.

۲-۱-۳. روش های بدون ناظر

این روش برای بهبود معایب روش های پیشین استخراج رابطه در متون بزرگ معرفی شد. در روش بدون ناظر هدف استخراج اطلاعات بدون نظارت، به طور خودکار القاء ساختار اطلاعات به عنوان انواع رابطه است، روش اصلی مورد

^۱Agichtein

^۲Gravano

استفاده در اینجا خوشه بندی می باشد [۱]. تا به اینجا از استخراج رابطه زمانی که نوع روابط استخراج شده شناخته شده هستند صحبت شد. مواردی هست که ما هیچ نوع رابطه خاصی در ذهن نداریم، اما می خواهیم انواع رابطه برجسته از یک مجموعه داده شده را کشف کنیم. شین یاما^۱ و همکاران در [۱۰] این مسئله را مطالعه کردند که به کشف رابطه محدود اشاره می کنند. آن ها ابتدا تعداد زیادی از مقالات خبری را از منابع مختلف در وب جمع آوری می کنند و سپس مقالات خبری در مورد یک رویداد یکسان را براساس شباهت لغوی خوشه بندی می کنند. در این روش آنها می توانند ویژگی های موجودیت را براساس وقوع متعدد آن در مقالات غنی کنند. سپس تجزیه نحوی را انجام می دهند و موجودیت های اسمی را از مقالات مختلف استخراج می کنند. در نهایت جفت موجودیت هایی که همزمان در همان مقاله رخ داده اند را براساس ویژگی هایشان خوشه بندی می کنند.

۲-۲. روش های مهندسی دانش

برای استخراج دانش از متن، نیاز داریم تا موجودیت های متن را استخراج نموده و نام های هم معنا که به یک موجودیت مربوط می شوند را باهم در یک گروه قرار دهیم. البته هر نام ممکن است در چند گروه قرار گیرد. به عنوان مثال نام «حسن روحانی» می تواند در گروه هایی متفاوت همراه بانام «رییس جمهور»، «سیاست مدار» و «استاد دانشگاه» قرار گیرد. این کار بسیار مشابه با مسئله «یادگیری هستان شناسی» خواهد بود. روش هایی وجود دارند که برای این کار از منابع دانش خارجی نظیر وردنت و یا ویکی پدیا استفاده می کنند؛ که به عنوان معروف ترین آن ها می توان از ویکی^۲ و یا گو^۳ نام برد. یاگو با استفاده از وردنت و همچنین دسته های موضوعی ویکی پدیا روشی را برای تولید پایگاه دانش به طور خودکار ارائه نموده است [۱۹].

۲-۲-۱. روش مبتنی بر قالب

منظور از قالب، نحوه بیان یک واقعه است، به همان شکلی که معمولاً بیان می شود. در بیان هر واقعه ای معمولاً تعدادی نقش معنایی در شکل های متنوع و البته محدود ظاهر می شوند. مثلاً در یک خبر مربوط به بمب گذاری، از نقش عامل بمب گذار، منطقه آسیب دیده و ... صحبت می شود. روشن است که هر قالب حجم زیادی از اطلاعات را در خود جای می دهد و پیشنهاد اقتباس و استفاده از آن به شکل بی ناظر کمی عجیب به نظر می رسد. استخراج اطلاعات به این نحو نیز تجربه شده است.

ویژگی های به دست آمده از متون بر اساس میزان باهم آیی آن ها خوشه بندی می شود تا به خوشه هایی که هر کدام در مورد موضوع مشخصی صحبت می کنند، برسیم. پس از این مرحله امیدواریم که مثلاً یک خوشه مربوط به اخبار بمب گذاری باشد و خوشه دیگر مربوط به اخبار آدم ربایی (در فضای متون خبری پلیسی صحبت می شود). متون هم موضوع در یک خوشه قرار می گیرند تا برای هر کدام از آن ها قالبی با نقش های معنایی مشخص، کشف شود. یادآوری می شود که فرایند کاملاً بی ناظر و مبتنی بر دانش است. توصیف از نقش معنایی به شکل آرگومان های ممکن آن ها در هر ویژگی استخراج می شود. مثلاً نقش «مکان آسیب دیده» می تواند «مفعول تخریب کردن» ظاهر شود. بعد از رسیدن به قالب های وقایع، استخراج اطلاعات آن ها بسیار آسان است. هر کدام از نقش های معنایی یکی از اطلاعات مورد نظر است که برای متون جدید به آسانی می توان آن ها را استخراج کرد [۱۷].

¹ Shinyama

² Wiki taxonomy

³yago



۲-۲-۲. روش مبتنی بر الگو

متن کاوی هم چنان یکی از موضوعات چالش برانگیز است که جهت استخراج دانش مفید از متن استفاده می شود تا به کاربران در یافتن الگوهای مفید و آنچه می خواهند، کمک می کند. از جمله مزایای استفاده از روش های مبتنی بر ترم عملکرد محاسباتی مناسب به علاوه نظریه های کامل جهت وزن دهی ترمها را می توان نام برد. اگرچه روش های مبتنی بر ترم از مشکلات معنایی^۱ و هم معنایی^۲ رنج می برند. منظور از معنایی این است که یک کلمه ممکن است چند مفهوم و معنا داشته باشد. و منظور از هم معنایی این است که چند کلمه ممکن است یک معنی و مفهوم را داشته باشند [۱۸].

در طول استخراج اطلاعات، این فرضیه وجود دارد که همواره رویکردهای مبتنی بر عبارت عملکرد بهتری از رویکردهای مبتنی بر ترم دارند. به عنوان مثال عبارت ممکن است معنای بیشتری را مانند اطلاعات با خود حمل کند. عبارات کمتر مبهم هستند و قابلیت متمایز کنندگی بیشتری نسبت به ترمهای فردی دارند. اما این فرضیه اکنون شانس چندانی ندارد به چند دلیل: (۱) عبارت ها خواص آماری نامرغوبی نسبت به ترمها دارند. (۲) تعدد وقوع در آنها بسیار پایین است. (۳) ممکن است تعداد زیادی عبارت تکراری و اضافه وجود داشته باشد.

پس الگوهای ترتیبی در حوزه متن کاوی به یک جایگزین امیدوارکننده برای «عبارت ها» تبدیل شد. زیرا الگوهای ترتیبی خواص آماری ترمها را به خوبی در نظر می گیرند. روش های تطبیق الگو، به طور وسیعی در قلمروی استخراج اطلاعات کاربرد دارند و به قلمروی یادگیری هستان شناس نیز به ارث رسیده اند در روش های مبتنی بر الگو، ورودی (معمولاً متن) به دنبال الگو یا کلمه کلیدی خاص که نشانگر رابطه مفهومی خاصی است جستجو می شود. این الگوها انواع مختلفی اعم از (نحوی یا معنایی، و عمومی یا خاص) دارند و برای استخراج عناصر مختلف هستان شناسی مثل روابط طبقه ای یا غیر طبقه ای و یا اصول بدیهی بکار می روند [۱۸].

۲-۳-۲. استخراج آزاد اطلاعات

در برخی موارد هدف کشف تمام حقایق مفید موجود در متن بزرگ و متنوع از جمله وب است که به استخراج اطلاعات آزاد اشاره دارد این روش اولین بار توسط بانکو^۴ [۱۱] معرفی شد. استخراج آزاد اطلاعات روشی است که برای کشف روابط متون بزرگ مانند وب استفاده می شود. در واقع در این روش به نوع رابطه خاص اشاره نمی شود و برخلاف روش های پیشین به مجموعه کوچک از روابط در متن محدود نمی شود و همه انواع وابستگی های دودویی موجود در متن را استخراج می کند و در این راستا از روش های بدون ناظر بهره می برد [۱].

چالش های سیستم های استخراج آزاد روابط شامل این است که این سیستمها نیز قادر به استخراج تمام روابط نیستند و از طرفی خروجی ناقص و نوفه دار^۵ دارند و نیز ممکن است استخراج اطلاعاتی را در بر نداشته باشند. از دیگر مشکلات این سیستمها این است که بدلیل ماهیت مقیاس پذیر بودن استخراج آزاد روابط، استفاده از ابزارهای عمیق پردازش زبان طبیعی نظیر تجزیه گر نحوی و معنایی که باعث بهبود قابل توجه نتایج و افزایش دقت می شود، ممکن نیست. از طرفی استفاده صرف از ابزارهای سطحی پردازش زبان طبیعی نظیر تجزیه گر سطحی، اجزای سخن و باعث کاهش چشم گیری در معیارهای کارایی استخراج گرها می شود.

¹ polysemy

² synonymy

³ ontology

⁴ Banko

⁵ noise



۲-۳-۱. سامانه های استخراج آزاد اطلاعات

در مسیر استخراج اطلاعات مجموعه ای از سامانه ها برای کمک به استخراج بخصوص در زمینه «استخراج اطلاعات آزاد» معرفی شدند که به شرح زیر است.

- سامانه تکست رانر^۱ [۱۰ و ۱۲]: اولین سامانه ای بود که با معرفی پارادایم «استخراج اطلاعات آزاد» عرضه شد. این سامانه با اعمال تعدادی قانون روی داده ها، برای خود تعدادی نمونه صحیح ایجاد کرده و سپس آن ها را یاد می گیرد. این روش را با نام روش خودناظر نامگذاری کرده اند. سپس از این ابزار برای استخراج رابطه از داده ها استفاده می شود. سامانه ای است که مجموعه بزرگ از سطرها را بدون نیاز به ورودی انسان استخراج می کند.
- سامانه ریورب^۲ [۱۴ و ۱۵]: این سامانه یکی دیگر از سامانه های استخراج آزاد اطلاعات است که فعل های موجود در متن را می یابد و سپس رابطه متناسب با هر فعل را استخراج می کند. تجزیه کننده نحوی را برای برچسب گذاری جملات استفاده می کند و محدودیت های واژگانی و نحوی را برای شناسایی واقعیات دودوئی بکار می برد.
- سامانه OLLIE^۳ [۶]: سامانه OLLIE برای استخراج آزاد اطلاعات ابتدا مجموعه سطرهایی از سامانه ریورب را با خود راه انداز مجموعه آموزشی بزرگ بکار می برد و قالب های الگوی باز را روی این مجموعه آموزشی یاد می دهد که این قالب های الگو در زمان استخراج بکار می روند. OLLIE بهترین شکل عبارت رابطه را بر مبنای قالب هایی روی عبارت رابطه ریورب تولید می کند.
- سامانه WOE [۶ و ۱۳]: این سامانه با استفاده از داده های ساخت یافته ای که در صفحات ویکی پدیا^۴ وجود دارد داده های مورد نیاز برای آموزش را ایجاد می کند. خودراه انداز مبتنی بر ویکی پدیا را استفاده می کند و دسترسی به عبارت رابطه ندارد. این سامانه محدودیت های معنایی - لغوی برای الگوها قرار نمی دهد و برای عبارت های رابطه که فعل میانجی شده دارد و شامل اسم نیست طراحی شده است.
- سامانه Kraken [۱۴]: سامانه استخراج آزاد اطلاعات Kraken برای گرفتن حقایق کامل از جملات عرضه شد و می تواند حقایق یک تایی، دو تایی تا چندتایی را استخراج کند.
- سامانه واندرلست^۵ [۱۴]: این سامانه با استفاده از گرامر سبک وابسته عمل می کند. مسیرهای وابسته را مطابق قواعد دستوری معتبر برای یافتن آرگومان های مرتبط با رابطه پیمایش می کند.
- سامانه اسنوبال^۶ [۸]: سامانه نیمه نظارتی برای استخراج اطلاعات است که با تعدادی داده آموزشی شروع به کار کرده و سعی می کند الگوهای مربوط به وقوع های متفاوت این نمونه ها را بیابد.
- سامانه Know-it-all [۱۶]: برخلاف اسنوبال که نیازی به داده ابتدایی برای شروع فرآیند استخراج اطلاعات ندارد. Know-it-all برای شروع کار خود نیاز به تعدادی الگو و شرح داده مورد نظر برای استخراج دارد. این الگوها وابسته به زبان و البته مستقل از رابطه هستند. سامانه با استفاده از الگوها و داده مورد نظر تعدادی عبارت تولید می کند و با استفاده از موتور جستجو صفحات وب مربوط به آن را بازیابی می کند و در نهایت اطلاعات از این صفحات بازیابی شده استخراج می گردد.

¹ Textrunner

² Reverb

³ Open language learning information extraction

⁴ Wikipedia

⁵ bootstrapping

⁶ Wanderlust

⁷ Snow ball

۳. مقایسه روش های استخراج اطلاعات

همانطور که در مطالب پیشین گفته شد روش های استخراج اطلاعات به دو دسته اصلی یادگیری ماشین و مهندسی دانش تقسیم می شود. رویکرد یادگیری ماشین در روش باناظر با داده آموزشی کم کار می کند و روش نیمه ناظر اغلب با داده آموزشی کم کار می کند و با نمونه رابطه های کوچک کار می کند و با استفاده از هسته ها بوجد آمدند و روش بدون ناظر از خوشه بندی استفاده می کند. روش مهندسی دانش نیز در دو روش کلی مبتنی بر قالب و الگو به استخراج اطلاعات می پردازد. در مبتنی بر الگو، کلمه کلیدی یا الگوی خاص از متن استخراج می شود و در مبتنی بر قالب، متون براساس باهم آیی خوشه بندی می شوند. روش استخراج آزاد اطلاعات نیز برای استخراج روابط از متون متنوع مانند صفحات وب مورد استفاده قرار می گیرد. در جدول زیر به تفکیک روش های استخراج اطلاعات مقایسه شده است و دسته بندی هر یک از این روش ها و سامانه های استخراج اطلاعات معرفی شده در این روش ها مشخص شده اند.

جدول ۱- مقایسه روش های استخراج اطلاعات

ویژگی های اصلی	نوع دسته بندی		روش استفاده شده	رویکردها
✓ استفاده از داده آموزشی زیاد	موجودیت متنی واژگانی متنی و نحوی	دسته بندی براساس ویژگی	باناظر	یادگیری ماشین
	هسته مبتنی بر ترتیب هسته مبتنی بر درخت هسته ترکیبی	دسته بندی براساس هسته		
✓ استفاده از نمونه رابطه و داده آموزشی کم	خودراه انداز نظارت راه دور		نیمه ناظر	مهندسی دانش
	کشف رابطه و القای الگو		بدون ناظر	
✓ خوشه بندی و استخراج قالب	خوشه بندی متون با موضوع یکسان و استخراج قالب مشخص		مبتنی بر قالب	مهندسی دانش
✓ الگوهای نحوی و معنایی	استخراج الگو یا کلمه کلیدی خاص از متن استخراج عناصر مختلف هستان شناسی		مبتنی بر الگو	

<ul style="list-style-type: none"> ✓ مقیاس پذیر ✓ بدون نیاز به تعریف روابط از پیش تعریف شده ✓ محدود به مجموعه کوچک از روابط در متن نیست 	<p>کشف همه روابط از متون بزرگ مانند وب</p>	<p>روش خودناظر</p>	<p>استخراج آزاد اطلاعات</p>
--	--	--------------------	-----------------------------

۴. نتیجه گیری

استخراج اطلاعات یک مسئله مهم در متن کاوی است و بطور گسترده در زمینه‌هایی مانند پردازش زبان طبیعی، بازیابی اطلاعات و وب‌کاوی مورد مطالعه قرار می‌گیرد. هدف از آن کشف اطلاعات ساخت یافته از متن نیمه ساخت یافته یا غیر ساخت یافته است. استخراج اطلاعات به دو دسته عمده تشخیص موجودیت اسمی و استخراج رابطه تقسیم می‌شود. به یافتن روابط معنایی بین موجودیت‌های متن استخراج رابطه گفته می‌شود [۱].

روش‌های یادگیری ماشینی همان‌طور که گفته شد بطور کلی به سه روش باناظر، بدون ناظر و نیمه ناظر تقسیم می‌شود. در روش باناظر یادگیری ماشینی همان‌طور که گفته شد بطور کلی به سه روش باناظر، بدون ناظر و نیمه ناظر تقسیم می‌شود. در روش باناظر یادگیری ماشینی، یک تابع هسته یا کرنل محصول داخلی نمونه‌های مشاهده شده را در بعضی زیر لایه‌های فضای برداری تعریف می‌کند [۱]. روش باناظر از مقادیر زیاد داده آموزشی برای استخراج اطلاعات استفاده می‌کند برای حل این مسئله روش یادگیری نیمه ناظر معرفی شد که با داده آموزشی کم کار می‌کند. روش نیمه ناظر به دو روش کلی خودراه‌انداز و نظارت راه دور تقسیم می‌شود. روش خودراه‌انداز از مجموعه کوچک از نمونه‌های رابطه شروع می‌کند و از الگوهای استخراج کمک می‌گیرد. در روش خودراه‌انداز فقط یک مجموعه کوچک از جفت موجودیت‌ها استفاده می‌شود. با رشد وب اجتماعی بیشتر دانش انسان توسط کاربران زیادی در پایگاه‌های اطلاعات ذخیره می‌شود، در این حالت ممکن است مجموعه بزرگ از موجودیت‌ها برای رابطه هدف باشند تا داده آموزشی تولید شود به همین دلیل روش نظارت راه دور معرفی شد. روش بدون ناظر نیز برای بهبود معایب روش‌های پیشین استخراج رابطه در متون بزرگ معرفی شد. و در نهایت استخراج آزاد اطلاعات که با هدف تسهیل کشف روابط از متون بزرگ و متنوع روی وب معرفی شد. در مواردی که هدف استخراج تمام حقایق مفید در متن است از این روش استفاده می‌شود [۴ و ۱].

روش‌های مهندسی دانش نیز به دو روش مبتنی بر قالب و مبتنی بر الگو کار می‌کنند. در مبتنی بر قالب ویژگی‌های بدست آمده از متن براساس باهم آیی آن‌ها خوشه بندی می‌شوند، در واقع متون با موضوع یکسان در یک خوشه قرار می‌گیرند و برای هر یک قالب مشخص استخراج می‌شود [۱۷]. در روش مبتنی بر الگو، متن بدنبال الگو یا کلمه خاص که نشان دهنده رابطه خاص است جستجو می‌شود. این الگوها برای استخراج عناصر مختلف هسته‌شناسی بکار می‌روند [۱۸].

مراجع

[1] C. C. Aggarwal and C. Zhai, Eds., Mining Text Data. Boston, MA: Springer US, 2012.

- [2] P. Z. Ives, "Acknowledgments," 2010.
- [3] M. Banko, "Open Information Extraction for the Web," 2009.
- [4] J. Piskorski and R. Yangarber, "Information Extraction: Past, Present and Future," 2013, pp. 23–49.
- [5] S. G. Small and L. Medsker, "Review of information extraction technologies and applications," *Neural Comput. Appl.*, vol. 25, no. 3–4, pp. 533–548, Sep. 2014.
- [6] Mausam, M. Schmitz, R. Bart, S. Soderland, and O. Etzioni, "emnlp12a.pdf," 2012.
- [7] C. C. Aggarwal and C. Zhai, Eds., *Mining Text Data*. Boston, MA: Springer US, 2012.
- [8] Eugene Agichtein and Luis Gravano. *Snowball: Extracting relations from large plain text collections*. In *Proceedings of the 5th ACM Conference on Digital Libraries*, pages 85–94, 2000.
- [9] B. B. Dalvi, W. W. Cohen, and J. Callan, "WebSets," in *Proceedings of the fifth ACM international conference on Web search and data mining - WSDM '12*, 2012, p. 243.
- [10] Yusuke Shinyama and Satoshi Sekine. *Preemptive information extraction using unrestricted relation discovery*. In *Proceedings of the Human Language Technology Conference of the North American Chapter of the Association for Computational Linguistics*, pages 304–311, 2006.
- [11] Michele Banko, Michael J. Cafarella, Stephen Soderland, Matthew Broadhead, and Oren Etzioni. *Open information extraction from the Web*. In *Proceedings of the 20th International Joint Conference on Artificial Intelligence*, pages 2670–2676, 2007.
- [12] M. banko, M. J. Cafarella, S. Soderland, M. Broadhead, and O. Etzioni, *Open Information Extraction for the Web*. University of Washington, 2009.
- [13] D. S. Weld and F. Wu, "Open Information Extraction using Wikipedia," *Proc. 48th Annu. Meet. Assoc. Comput. Linguist.*, no. July, pp. 118–127, 2010.
- [14] A. Akbik and A. Löser, "KrakeN: N-ary Facts in Open Information Extraction," *Proceedings of the Joint Workshop on Automatic Knowledge Base Construction and Web-scale Knowledge Extraction*. pp. 52–56, 2012.
- [15] O. Etzioni, M. Cafarella, D. Downey, S. Kok, A.-M. Popescu, T. Shaked, S. Soderland, D. S. Weld, and A. Yates, "Web-Scale Information Extraction in KnowItAll (Preliminary Results)," *WWW'04 Proc. 13th Int. Conf. World Wide Web*, pp. 100–110, 2004.
- [16] A. Fader, S. Soderland, and O. Etzioni, "Identifying relations for open information extraction," *Proc. Conf. ...*, pp. 1535–1545, 2011.
- [17] Chambers, N. and D. Jurafsky. *Template-based information extraction without the templates*. in *Proceedings of the 49th Annual Meeting of the Association for Computational Linguistics: Human Language Technologies-Volume 1*. 2011. Association for Computational Linguistics.
- [18] Zhong, N., Y. Li, and S.-T. Wu, *Effective pattern discovery for text mining*. *IEEE transactions on knowledge and data engineering*, 2012. 24(1): p. 30-44.
- [19] Hoffart, J., et al., *YAGO2: A spatially and temporally enhanced knowledge base from Wikipedia*. *Artificial Intelligence*, 2013. 194: p. 28-61