

## ارائه روش ترکیبی تحلیل پایگاه داده های خبری با استفاده از RapidMine

### مطالعه موردی: متون خبری فارسی

سعید احمدیان<sup>۱\*</sup>، کریم صابری<sup>۲</sup>

۱- کارشناسی ارشد دانشگاه امام حسین (ع)، ۲- کارشناسی ارشد دانشگاه خواجه نصیرالدین طوسی

(دریافت: ۱۳۹۸/۱۱/۰۲، پذیرش: ۱۳۹۸/۱۲/۱۳)

#### چکیده

یکی از مهم ترین و پرکاربردترین رویکردها در مدیریت سامانه های تحت وب، استفاده از خبرگزاری های آنلاین و پایگاه داده های خبری هست. خبرگزاری های آنلاین و تحت وب به مانند قلب تپنده یک جامعه، جزء حیاتی ترین منابع اطلاعاتی یک جامعه به حساب می آیند. اهمیت این موضوع در مورد اخبار سازمان های نظامی دوچندان خواهد بود؛ بنابراین، داشتن روشی منحصر به فرد، صحیح و مفهومی جهت تحلیل و دسته بندی این منابع، می تواند مزایای متعددی فراهم نماید و به تصمیم گیران سازمان ها به ویژه سازمان های نظامی در انجام تصمیمات کمک کند. در این پژوهش، در مرحله اول به بررسی و شناخت مفاهیم انواع روش های داده کاوی می پردازیم. در مرحله دوم انواع روش های پردازش متن، از منظرهای مختلف باید صورت گیرد تا بتوان راه کارهای مثبت و نقاط قوت انواع روش ها را شناسایی کرده و برای تحلیل پایگاه داده های خبری بهترین روش را شناسایی کرد. در انتها پایگاه داده های خبری را مورد تحلیل قرار داده تا بتوان اطلاعات پنهان موجود در این داده ها را کشف و مورد استفاده قرار داد.

#### کلیدواژه ها:

داده کاوی، متن کاوی، پایگاه داده های خبر

#### ۱- مقدمه

اخبار موجود در دنیای امروز غالباً به عنوان یک پایگاه داده ثبت اخبار و به صورت آرشیوی انجام می شود در حالی که در میان آن ها دانش ها و روابط پنهان بسیار مفیدی قرار دارند. کشف این دانش ها می تواند برای تصمیم گیران و تحلیل گران در حوزه های مختلف و به ویژه در مسائل مرتبط با پدافند سایبری و تحلیل خبرها به صورت آئی بسیار کاربردی و جالب توجه باشد. با توجه به حجم عظیم این داده ها بایستی با استفاده از روش ها و فن های علمی، این داده ها تجزیه و تحلیل شوند و روابط و دانش های پنهان میان آن ها استخراج و در دسترس مدیران و تصمیم گیران قرار گیرند. با این کار بازدهی فرایند تحلیل و شناسایی خبرها به صورت هدفمند افزایش می یابد، به این علت که تحلیل خبرها از حالت تحلیل فرد کاربر خارج شده و توسط روش های داده کاوی با بازدهی بسیار بالاتری نسبت به تحلیل دستی کاربر، به انجام می رسد. استفاده از این روش ها می تواند به کلیه سازمان های داخل کشور اعم از نظامی و غیر نظامی این توانایی را بدهد که اخبار رسیده به سازمان متناسب با الگوی تعیین شده همان

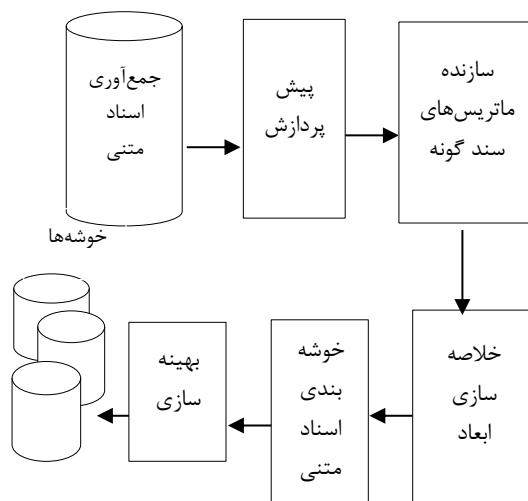
سازمان، تحلیل شود و در اختیار تصمیم گیران آن سازمان قرار گیرد. تصمیم گیران بدین وسیله توانایی تصمیماتی صحیح، مناسب با موقعیت موجود، سریع و با دانش بیشتری پیدا خواهند نمود.

لذا در این تحقیق با تشریح دقیق کاربرد روش های داده کاوی و متن کاوی در تحلیل خبرهای منتشر شده و کشف روابط پنهان میان آن ها، چگونگی استخراج روابط و نتایج کارامل در سامانه پایگاه خبری ارائه خواهد شد. در نهایت با دانش های پنهان در پایگاه داده در دسترس قرار گرفته که تحلیل و تفسیر خواهند شد، یک چارچوب برای تحلیل خبرها و چگونگی روابط پنهان میان خبرها بر اساس روش های پیش بینانه داده کاوی ارائه خواهد شد. به طور کلی اهدافی که از ارائه این روش دنبال می شود را می توان بدین گونه ذکر کرد:

۱. تعیین کلمات کلیدی و ماهیت خبرها بر اساس الگوهای متن کاوی و تهیه گزارش ها تحلیلی از آن ها.
۲. کشف دانش های پنهان از خبرهای تولید شده بر اساس روند سری های زمانی موجود در بین آن ها.
۳. کشف و تحلیل روابط پنهان میان کلمات کلیدی و ماهیت خبرها در گذشته.
۴. کشف کلمات مهم در خبرهای منتشر شده که تأثیر گذاری بسیاری در ذهنیت مردم و خوانندگان خبرها دارند و تعیین جهت گیری های ماهیت خبرها بر اساس روش های داده کاوی.

<sup>۱</sup> رایانامه نویسنده پاسخگو: s.ahmadian@ihu.ac.ir

### ۳- تعاریف متن کاوی



شکل (۲): مراحل متن کاوی.

می توان گفت که متن کاوی از تکنیک های بازیابی اطلاعات، استخراج اطلاعات همچنین پردازش کردن زبان طبیعی<sup>۱</sup> استفاده می کند و آن ها را به الگوریتم ها و متدهای KDD، داده کاوی، یادگیری ماشین و آماری مرتبط می کنند. با توجه به ناحیه های تحقیق گوناگون، بر هر یک از آن ها می توان تعاریف مختلفی از متن کاوی در نظر گرفت در ادامه برخی از این تعاریف بیان می شوند:

- متن کاوی = استخراج اطلاعات: در این تعریف متن کاوی متناظر با استخراج اطلاعات در نظر گرفته می شود (استخراج واقعیت ها<sup>۲</sup> از متن).
- متن کاوی = کشف داده متنی: متن کاوی را می توان به عنوان متدها و الگوریتم هایی از فیلهای یادگیری ماشین و آماری برای متن ها باهدف پیدا کردن الگوهای مفید در نظر گرفت. برای این هدف پیش پردازش کردن متون ضروری است. در بسیاری از روش ها، متدهای استخراج اطلاعات، پردازش کردن زبان طبیعی یا برخی پیش پردازش های ساده برای استخراج داده از متون استفاده می شود. سپس می توان الگوریتم های داده کاوی را بر روی داده های استخراج شده اعمال کرد.
- متن کاوی = فرایند KDD: که در بخش بعد به طور کامل توضیح داده شده است.

۵. دسته بندی و خوشه بندی صحیح اطلاعات و داده ها که موجب جلوگیری از نشر اخبار محرمانه نظامی و غیرنظامی خواهد شد.

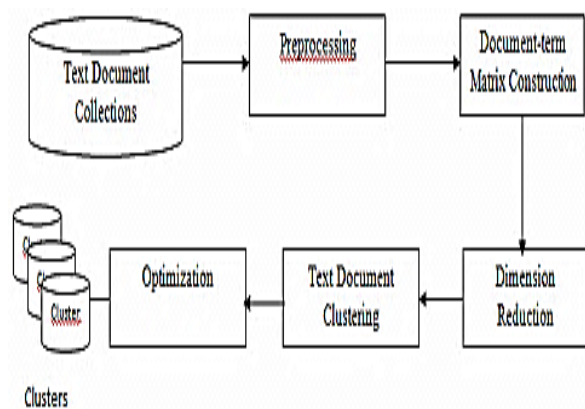
### ۲- تعریف مفاهیم

متن کاوی: استخراج اطلاعات از متن شامل ارائه یک قالب ساختارمند مانند یک پایگاه داده از اطلاعات دلخواه موجود در متن هست [۱۴]. متن کاوی به مفهوم جستجوی الگوها در متون بدون ساختار است. متن کاوی ترکیبی از فن های کاوش داده و استفاده آن ها با کمک لغت شناسی موضوعی هست [۱۵].

دو معیار اصلی برای ارزیابی کارایی سامانه های اطلاعات وجود دارد. اول این که، چه درصدی از اطلاعات استخراجی صحیح هستند و دوم این که، چه درصدی از اطلاعات موجود در متن استخراج یافته اند [۱۷].

یکی از مراحل متداول در تحلیل اسناد متنی خوشه بندی این اسناد هست. با خوشه بندی، اسناد مشابه را می توان با سرعت کاوش کرد، به آسانی می توان موضوعات و زیر موضوعات آن ها را درک کرد و به طور مؤثری بین آن ها برای بسیاری از کاوش ها پرس و جو انجام داد. خوشه بندی اسناد متنی، نقش مهمی را در سازمان دهی مؤثر سند، خلاصه سازی، استخراج موضوع و بازیابی اطلاعات ایفا می کند و یک ابزار کارا برای مدیریت سربار اطلاعات است [۱۶].

می توان متن کاوی را گونه ای اختصاصی از خوشه بندی دانست با این تفاوت که به دلیل داشتن ابعاد بالاتر کمی دشوارتر هست. قدم های اساسی برای متن کاوی در شکل (۱) به نمایش درآمده است.



شکل (۱): مراحل متن کاوی .

<sup>۱</sup> Natural language processing (NLP)

<sup>۲</sup> Facts

#### ۴- کشف دانش و ارتباط آن با متن کاوی

کشف دانش در پایگاه داده‌ها (KDD)، این لغت به همه شیوه‌هایی اشاره دارد که هدف آن‌ها پی‌بردن به ارتباط و نظم بین اطلاعات قابل مشاهده است. لغت KDD برای توصیف همه مراحل استخراج اطلاعات از پایگاه داده و نیز بیان اهداف کارهای اولیه کاربرد قوانین تصمیم‌گیری است. در کل KDD فرآیند یافتن اطلاعات و الگوهای مفید از داده را گویند و داده‌کاوی بهره‌گیری از الگوریتم‌هایی برای یافتن اطلاعات مفید در فرآیند KDD است.

از جمله خصوصیتی که برای اندازه‌گیری کیفیت الگوهای پیداشده در داده می‌توان استفاده کرد عبارت‌اند از: قابلیت فهم انسان، اعتبارسنجی با معیارهای آماری و تازگی و مفید بودن. کشف دانش در پایگاه داده را می‌توان به‌عنوان یک فرایند که به‌وسیله چندین گام پردازش کردن<sup>۱</sup> تعریف می‌شود، در نظر گرفت. این گام‌ها به‌منظور استخراج الگوهای مفید باید بر روی مجموعه داده‌ای اعمال شوند. این گام‌ها به‌صورت تکراری اجرا می‌شوند و برخی گام‌ها نیاز به بازخورد از کاربر دارند. یک کاربر سامانه KDD به‌منظور انتخاب زیرمجموعه صحیحی از داده‌ها باید درک بالایی از قلمرو داده‌ها، رده مناسبی از الگوها و معیار خوبی برای الگوهای جالب داشته باشد؛ بنابراین سامانه KDD باید ابزارهایی با اثر تعاملی داشته باشد نه سامانه‌های تجزیه‌وتحلیل خودکار. طبق [۱۹] گام‌ها را می‌توان به‌صورت زیر بیان کرد: (۱) درک کردن کسب‌وکار (۲) درک کردن داده (۳) آماده‌سازی داده (۴) مدل کردن (۵) ارزیابی (۶) deployment. مرحله پیش‌پردازش غالباً یکی از مراحل زمان‌بر و درعین حال بسیار مهم در کسب نتیجه مطلوب است. مخصوصاً در متن‌کاوی که نیاز به متدهای پیش‌پردازش کردن خاصی برای تبدیل داده متنی به فرمتی که برای الگوریتم‌های داده‌کاوی مناسب است، داریم.

#### ۵- ناحیه‌های جستجوی مرتبط

سه روش اساسی در مواجهه با این حجم وسیع از اطلاعات غیر ساخت‌یافته وجود دارد که عبارت‌اند از: بازیابی اطلاعات<sup>۲</sup>، استخراج اطلاعات<sup>۳</sup> و پردازش زبان طبیعی. بازیابی اطلاعات: اصولاً مرتبط است با بازیابی مستندات و مدارک. کار معمول در بازیابی اطلاعات این است که با توجه به

نیاز مطرح‌شده از سوی کاربر، مرتبط‌ترین متون و مستندات و یا درواقع کلمه را از میان دیگر مستندات یک مجموعه بیرون بکشد. این عمل یافتن دانش نیست بلکه تنها آن مجموعه‌ای از کلمات را که به نظرش مرتبط‌تر به نیاز اطلاعاتی جستجوگر است را به او تحویل می‌دهد. این روش به‌واقع هیچ‌دانشی و حتی هیچ اطلاعاتی را برایمان به ارمغان نمی‌آورد.

پردازش زبان طبیعی: هدف کلی NLP<sup>۴</sup> رسیدن به یک درک بهتر از زبان طبیعی توسط کامپیوترهاست. تکنیک‌های مستحکم و ساده‌ای برای پردازش کردن سریع متن به کار می‌روند. همچنین از تکنیک‌های تحلیل زبان‌شناسی نیز برای پردازش کردن متن استفاده می‌شود.

استخراج اطلاعات: هدف روش‌های استخراج اطلاعات، استخراج اطلاعات خاص از سندهای متنی است. استخراج اطلاعات می‌تواند به‌عنوان یک‌فاز پیش‌پردازش در متن‌کاوی بکار برود. استخراج اطلاعات عبارت‌اند از نگاشت کردن متن‌های زبان طبیعی (مثلاً گزارش‌ها، مقالات journal، روزنامه‌ها، ایمیل‌ها، صفحات وب، هر پایگاه داده متنی و ...) به یک نمایش ساخت‌یافته و از پیش تعریف‌شده یا قالب‌هایی که وقتی پر می‌شوند، منتخبی از اطلاعات کلیدی از متن اصلی را نشان می‌دهند. یک‌بار اطلاعات استخراج‌شده و سپس اطلاعات می‌توانند در پایگاه داده برای استفاده‌های آینده، ذخیره شوند.

#### ۶- روش‌های پیش‌پردازش کردن متون

برای کاوش کردن<sup>۵</sup> مجموعه بزرگی از اسناد ضروری است که اسناد پیش‌پردازش شوند و اطلاعات در یک ساختار داده‌ای مناسب برای پردازش‌های بعدی ذخیره شوند. روش‌های اصلی پیش‌پردازش متون عبارت‌اند از: مدل فضای بردار<sup>۶</sup> [۲۲]، مدل احتمالی [۲۳] و مدل منطقی [۲۴].

#### ۷- روش‌های متن‌کاوی

دلیل اصلی به کار بردن روش‌های داده‌کاوی برای اسناد متنی، ساختاربندی کردن آن‌هاست. ساختارهای معروف عبارت‌اند از: کاتالوگ‌های کتابخانه یا نمایه‌های کتاب. مشکل نمایه‌های طراحی‌شده به‌صورت دستی، زمان موردنیاز برای نگهداری آن‌ها است؛ بنابراین برای منابع اطلاعاتی که خیلی تغییر می‌کنند مثل

<sup>۴</sup> Natural language processing

<sup>۵</sup> Mining

<sup>۶</sup> Vector space

<sup>۱</sup> Processing

<sup>۲</sup> Information Retrieval

<sup>۳</sup> Information Extraction

تقسیم می‌شوند؛ اما هر نوع شباهتی را نمی‌توان به صورت عددی تخمین زد. برای مثال فاصله بین مفاهیم UMIST و OLSE. حتی اگر این فواصل بتوانند به اعداد انتقال داده شوند، ممکنه انتقال مشکل باشد؛ بنابراین، تجزیه و تحلیل خوشه مبتنی بر فاصله قدیمی در زمان پردازش کردن مفاهیم ممکنه نتایج بامعنی تولید نکند.

در خوشه‌بندی مفهومی، خوشه‌ها تنها مجموعه‌ای از اشیاء با شباهت عددی نیستند. خوشه‌ها به‌عنوان گروهی از اشیاء که با یکدیگر یک مفهوم را نشان می‌دهند، هستند. برای خوشه‌بندی مفهومی ما به مجموعه‌ای از صفات برخی اشیاء نیاز داریم (یک‌زبان توصیف برای مشخص کردن خوشه‌های چنین اشیایی) و یک معیار کیفیت خوشه‌بندی است. هدف، تقسیم‌بندی کردن اشیاء در خوشه‌ها به‌گونه‌ای است که معیار کیفیت بیشینه شود و درعین حال تعیین کردن توصیفات عمومی از این خوشه‌هاست. توجه شود که در خوشه‌های مفهومی خصوصیات خوشه با بررسی و دقت به فرآیند مشخص کردن خوشه‌ها به وجود می‌آید. این یک تفاوت اصلی بین خوشه‌بندی مفهومی و قدیمی است. در روش خوشه‌بندی قدیمی خوشه‌ها مطابق با یک معیار شباهت مشخص می‌شوند. این معیار شباهت یک تابع است که فقط خصوصیات اشیاء مقایسه می‌شوند. در مقابل در خوشه‌بندی مفهومی به شرح یا توصیف هم توجه می‌شود.

### ۳-۷- روش‌های ترکیبی

تعداد زیادی روش در فاز استخراج دانش وجود دارد. درعین حال تمام این روش‌ها را شاید بتوان به دودسته اصلی تقسیم کرد. این دودسته اصلی، روش‌های مبتنی بر کارایی و روش‌های مبتنی بر دانش هستند. در روش اول، طراحان نگران کارایی سامانه هستند و به‌گونه‌ای سامانه را طراحی می‌کنند که بهترین کارایی و سرعت را داشته باشد. روش‌های رایج‌تر در این نوع نگرش، روش‌های آماری و نیز شبکه‌های عصبی هستند. روش‌های آماری بر پایه هر نوع اطلاعات آماری است که از متون قابل استخراج است. مواردی چون تکرار لغات به‌تنهایی، تکرار لغات باهم و چیزهایی شبیه آن.

در سوی دیگر روش‌های مبتنی بر دانش قرار دارند که از زاویه دید دیگری به این مسئله نگاه می‌کنند. آن‌ها سعی می‌کنند اولاً تا حد ممکن مفاهیم موجود را از داخل مجموعه متون استخراج کنند و ثانیاً بین این مفاهیم روابطی برقرار کنند. استفاده از این روش بسیار وابسته به NLP است. در حقیقت این هدفی است که NLP نیز آن را دنبال می‌کند و آن درک متن است. دستگاه‌هایی

وب مناسب نیستند. متدهای موجود برای ساختاربندی کردن مجموعه‌ها عبارت‌انداز: روش‌های رده‌بندی و روش‌های خوشه‌بندی.

### ۷-۱- رده‌بندی

هدف از رده‌بندی متون نسبت دادن کلاس‌های از پیش تعریف‌شده به اسناد متنی است؛ مثلاً یک خبر جدید که وارد می‌شود بگوییم متعلق به کلاس ورزشی یا سیاسی یا هنری. برای رده‌بندی اسناد روش‌های گوناگونی به کار می‌روند. در رده‌بندی یک مجموعه آموزش از اسناد وجود دارد که برای این مجموعه کلاس‌ها مشخص است. با استفاده از این مجموعه مدل رده‌بندی مشخص می‌شود، سپس با استفاده از آن کلاس سند جدید که وارد می‌شود، مشخص می‌گردد. برای اندازه‌گیری کارایی مدل رده‌بندی، یک مجموعه آزمودن که مستقل از مجموعه آموزش است در نظر گرفته می‌شود؛ و برچسب‌هایی که برای این اسناد توسط مدل تخمین زده می‌شود با برچسب واقعی اسناد مقایسه می‌شود. روش‌های رده‌بندی:

- انتخاب‌ترم ایندکس<sup>۱</sup>
- رده‌بندی کننده Naive Bayes
- رده‌بندی نزدیک‌ترین همسایه
- درخت تصمیم‌گیری

### ۷-۲- الگوریتم خوشه‌بندی

انتخاب الگوریتم بیشتر وابسته به مجموعه داده و وظیفه‌ای است که باید انجام شود. خوشه‌بندی سلسله‌مراتبی و خوشه‌بندی رابطه‌ای باینری از معروف‌ترین الگوریتم‌های خوشه‌بندی هستند. در روش‌های قبلی، خوشه‌ها در درخت‌های خوشه‌ای (سلسله‌مراتبی) قرار می‌گرفتند. به طوری که خوشه‌های مرتبط در یک شاخه از درخت بودند.

این تحلیل به‌طور خودکار می‌تواند برای وظایف زیادی مفید واقع شود. همچنین می‌تواند یک دیدی از محتویات یک مجموعه بزرگ اسناد فراهم کند. وظیفه دیگر، شناسایی ساختارهای پنهان در گروه‌های اشیاء است. فرآیند پیدا کردن اطلاعات مربوط می‌تواند آسان و دقیق شود. به‌علاوه، تمایلات جدید (سیاست‌های جدید) که بیان‌نشده‌اند در اسناد دیگر، می‌تواند در اسناد کشف گردد و سرانجام کپی اسناد در یک مجموعه می‌تواند حذف شود.

تحلیل خوشه‌بندی، یک فن آماری قدیمی است. اشیاء بر اساس معیار شباهت و فاصله عددی بین اشیاء به گروه‌هایی

<sup>1</sup> Index term selection

اطلاعات استخراج می‌شوند و سپس اطلاعات می‌توانند در پایگاه داده ذخیره شوند و برای پرس‌وجو کاوش گردند و در زبان طبیعی خلاصه شوند. اولین گام ضروری پیش‌پردازش کردن زبانی (زبان‌شناختی) است. این گام شامل تعدادی فن‌های زبان‌شناختی مثلاً tokenization، بخشی از برچسب‌گذاری گفتار و... است. اهداف کلی این مقاله را می‌توان به دودسته تقسیم نمود: (۱) مدیریت کردن اطلاعات ذخیره‌شده در پایگاه داده متنی (مجموعه‌ای اسناد) (۲) استخراج کردن دانش مفید.

این روش از دو کامپوننت تحلیل متن<sup>۱</sup> و داده‌کاوی تشکیل شده است. کامپوننت اول داده‌های نیمه ساخت‌یافته را به داده‌های ساخت‌یافته‌تر ذخیره‌شده در پایگاه داده تبدیل می‌کند؛ و کامپوننت دوم تکنیک‌های داده‌کاوی را بر روی خروجی کامپوننت اول اعمال می‌کند. هدف این روش مدیریت کردن اطلاعات (طبقه‌بندی کردن سندها در دسته‌های مناسب) و کاوش داده برای کشف کردن دانش مفید است؛ بنابراین در این روش ابتدا ترم‌ها و رویدادها استخراج شده و در پایگاه داده ذخیره می‌شوند. سپس الگوریتم خوشه‌بندی مناسبی (استفاده از الگوریتم Rock و مفهوم لنیک) بر روی پایگاه داده حاصل شده اعمال شده و اسناد گروه‌بندی می‌شوند به طوری که اسناد مشابه در یک گروه قرار گیرند. سپس یک الگوریتم رده‌بندی مناسب (درخت تصمیم‌گیری) برای معتبر سازی بیشتر نتایج حاصل از خوشه‌بندی و بهره‌برداری بهتر از دانش کشف‌شده اعمال می‌شود. برای جزئیات بیشتر این روش می‌توان به برای جزئیات بیشتر این روش می‌توان به مراجعه کرد.

## ۸- روش تحقیق و چارچوب پیشنهادی

در تحقیقاتی که امروزه صورت می‌گیرد از انواع روش‌ها استفاده می‌شود. این روش‌ها یا به صورت کمی و با استفاده از آمار و ارقام به نتیجه موردنظر می‌رسد، یا با استفاده از روش‌های کیفی. روش‌های تحقیق را در صورتی که برداری دوطرفه در نظر بگیریم، می‌توان گفت که یک‌طرف از بردار را روش‌های کمی و طرف دیگر را روش‌های کیفی تشکیل می‌دهند. روش‌های کمی با تکیه بر آمار و ارقام انجام می‌شود و روش‌های کیفی به کمک انواع فن‌های مشاهده و مصاحبه، به جمع‌آوری اطلاعات می‌پردازد. ویژگی مشترک این دو روش، مجهز نمودن محققین در تمام مراحل تحقیق به انواع فن‌های جمع‌آوری اطلاعات و

که از این روش‌ها استفاده می‌کنند در حال حاضر زیاد نیستند ولی DR-LINK [۲۵] از دانشگاه Syracuse یکی از آنهاست.

### ۷-۳-۱- روش Discotex<sup>۱</sup>

این روش توسط Kanya در سال ۲۰۰۷ ارائه شده است [۲۰]. این روش یک چارچوب جدید برای متن‌کاوی بر اساس یکپارچه کردن سامانه استخراج اطلاعات (IE) و ماژول استنتاج کردن قوانین استاندارد (KDD) ارائه می‌کند. IE، اسناد متنی را به داده ساخت‌یافته‌تر تبدیل می‌کند. در واقع تکه‌های خاصی از داده را در سندها به زبان طبیعی جستجو می‌کند و مجموعه‌ای از اسناد متنی نیمه ساخت‌یافته را به پایگاه داده ساخت‌یافته‌تری تبدیل می‌کند. در این روش RAPIER و BWI<sup>۲</sup> برای ساختن IE استفاده شده است. سپس پایگاه داده ساخته شده به وسیله ماژول IE توسط ماژول KDD برای کاوش بیشتر دانش استفاده می‌شود. در یک نسخه بهبودیافته از این روش از قوانین به دست آمده از ماژول KDD برای پیش‌بینی کردن اطلاعات از دست‌رفته و بهبود دقت ماژول IE استفاده می‌شود. برای ساختن ماژول KDD از RIPPER و APRIORI استفاده شده است.

### ۷-۳-۲- روش Textminer

در این روش ابتدا داده نیمه‌ساخت‌یافته تغییر می‌یابد مثلاً سندها به داده ساخت‌یافته ذخیره‌شده در یک پایگاه داده تبدیل می‌شوند. کامپوننت دوم تکنیک‌های داده‌کاوی را روی خروجی کامپوننت اول اعمال می‌کند بیشتر روش‌ها برای متن‌کاوی، الگوریتم‌های کاوش را روی برچسب‌های نسبت داده‌شده به هر سند اعمال می‌کنند. این برچسب‌ها ممکن است کلمات کلیدی استخراج شده از سند یا فقط فهرستی از کلمات در سند موردنظر باشند. در روش textminer الگوریتم‌های کاوش بر روی ترم‌ها (دنباله معنی‌دار از کلمات مثلاً (department of computation) ترکیب شده با رویدادهای (مجموعه معنی‌دار از ترم‌ها، مثلاً در یک دامنه مالی، خرید بین شرکت A و B) استخراج شده از سندها اعمال می‌شود. نویسندگان این مقاله معتقد هستند که مهم‌ترین فاکتورهای مشخصه که یک سند را توصیف می‌کنند ترم‌ها و رویدادهای بیان شده در سند هستند. این اطلاعات در یک جدول به نام EvantType نگه‌داری می‌شوند. استخراج اطلاعات فناوری مهمی است که یک گام پیش‌پردازش دارد. در این روش یک‌بار

<sup>۱</sup> Discovery from text extraction

<sup>۲</sup> Boosted Wrapper induction

<sup>۳</sup> Text Analysis Component

موجود را به خوبی بشناسیم. پس از جمع‌آوری داده لازم است تا داده‌ها پاک‌سازی شوند. یکی از مشکلات شایع داده‌ها پایین بودن کیفیت آن‌ها است، عملیاتی که مشکل کیفیت را برطرف می‌کند، پاک‌سازی داده نامیده می‌شود [۲۵]. حذف رکوردهای پرت یکی از اقدام‌ها برای پاک‌سازی داده است. داده‌های تکراری یا دو نسخه‌ای باید شناسایی و حذف شوند. همچنین داده‌های مغشوش را نیز می‌بایست شناسایی و ترمیم گردند.

مرحله بعد مدل‌سازی است. در این مرحله دو فاز عملیاتی صورت گرفته است. فاز اول مربوط به استفاده از الگوریتم k-means جهت خوشه‌بندی انواع متن‌های خبری است بدین منظور که مشخص شود هر متن خبر مربوط به کدام حوزه خبری هست. فاز دوم تحلیل دسته‌بندی اخبار موجود به کمک روش‌هایی مانند درخت تصمیم است که ارائه‌دهنده الگوریتمی جهت مشخص کردن نوع اخبار و دسته‌بندی آن‌هاست. فاز اول یعنی استفاده از الگوریتم خوشه‌بندی k-means به کمک نرم‌افزار Rapid Miner و فاز دوم نیز به کمک نرم‌همین نرم‌افزار صورت گرفته است. در این مرحله به دنبال توسعه الگوریتم تحلیلی و بهبود نتایج حاصله هستیم.

## ۹- پیاده‌سازی طرح و نتیجه‌گیری

داده‌های موجود جهت تحلیل و پیاده‌سازی در پروژه مربوط به خبرهای مختلف از خبرگزاری‌های اینترنتی هست. پایگاه داده‌ها شامل اطلاعات جمع‌آوری شده از حدود ۱۰۰ پایگاه خبری اینترنتی با سه زبان مختلف فارسی، انگلیسی و عربی بود. در مجموع دو جدول مختلف در این پایگاه داده وجود دارد.

یک جدول مربوط به اطلاعات خبر منتشره که حاوی اطلاعاتی مانند شناسه گروه، شناسه منبع، شناسه گزارش‌گیر، شناسه نوع، شناسه خبر، آدرس خبر، تاریخ انتشار، تاریخ ایجاد، شماره فایله، مسیر فایله، آدرس فایله و وضعیت خبر هست. جدول دوم با نام جدول متن نیز شامل اطلاعاتی مانند شناسه متن، شناسه خبر، شناسه فایله، شناسه زبان، متن و نوع متن می‌باشد. تعداد کل داده‌ها در جدول خبر شامل ۱۶۳۴۸۶ سطر یکتا و در جدول متن ۴۴۷۴۰۷ سطر یکتا هست. دلیل تمایز بین تعداد موجود در دو جدول نوع داده‌های<sup>۱</sup> مختلف در جدول متن هست، به این صورت که هر شناسه خبر یکتا در جدول خبر، انواع مختلفی متن ذخیره شده در جدول متن را داراست، مانند عنوان خبر، متن خبر و خلاصه خبر مربوط به یک خبر خاص.

تجزیه و تحلیل‌های متنوع است. در واقع روش‌های کمی در مجموع با شمارش و اندازه‌گیری جنبه‌هایی از زندگی اجتماعی سروکار دارند، در حالی که روش‌های کیفی بیشتر با تولید توصیف‌های استدلالی و کشف معناها و تغییرهای کنشگران اجتماعی سروکار دارند.

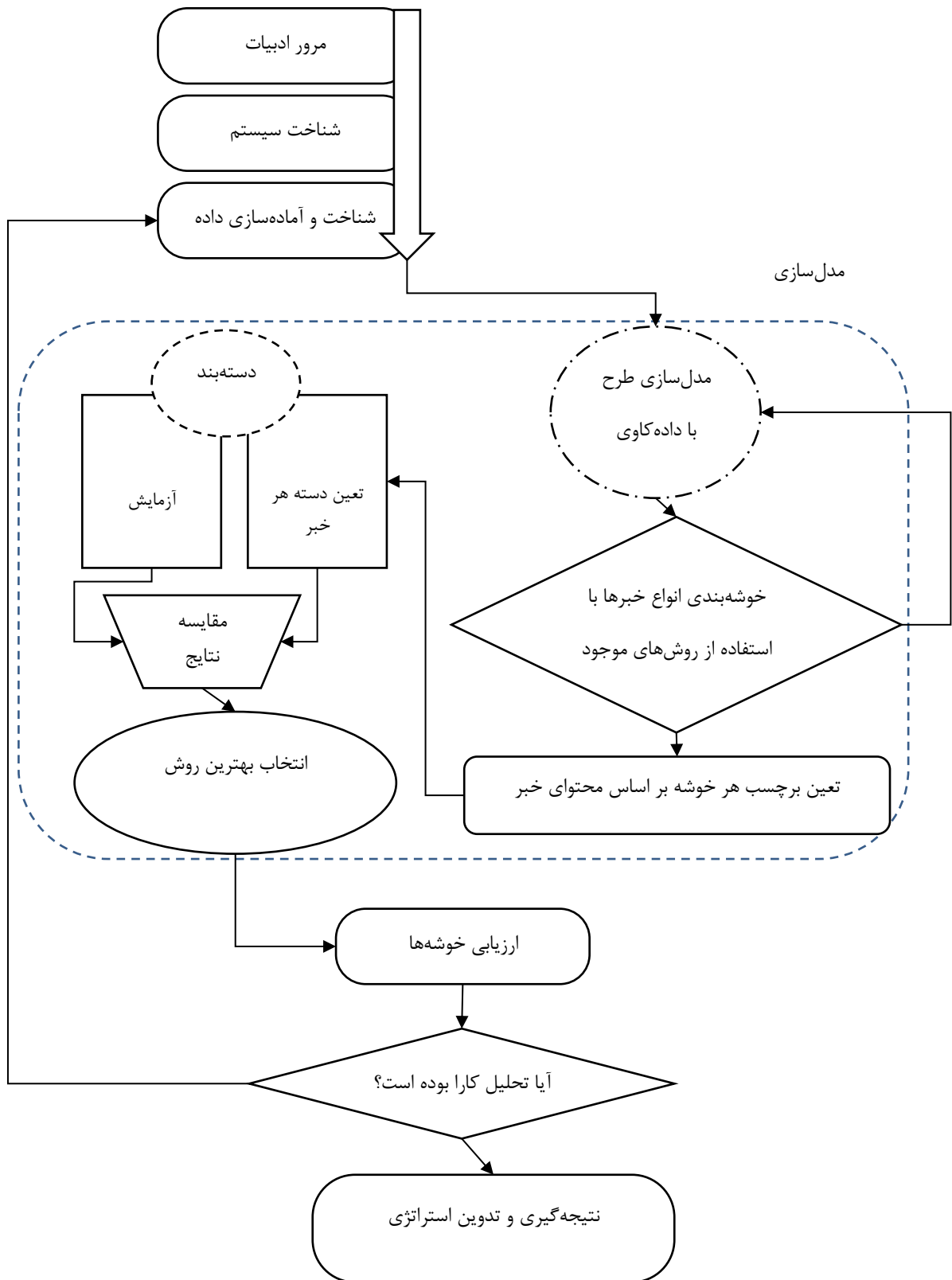
باین حال تمام روش‌های تحقیق، به استفاده از کمیت یا کیفیت خلاصه نمی‌شوند، بلکه می‌توان رویکردی را اتخاذ نمود که بسته به نوع تحقیق، از هر دو نوع روش استفاده کند. به این نوع از روش‌ها، ترکیبی می‌گویند [۲۲-۲۰]. در واقع ظهور و استفاده از روش‌های ترکیبی به منظر قوت دادن به تحقیقات صورت می‌گیرد [۱۸].

در تحقیق حاضر نیز از روش ترکیبی استفاده شده است. دلیل این امر استفاده از پایگاه داده‌های متنی متشکل از خبر است. این پایگاه داده‌ها به صورت کمی و کیفی، از انواع خبرها گردآوری شده است. همچنین به دلیل استفاده از روش‌های داده‌کاوی از استاندارد و متدولوژی CRISP برای نشان دادن چارچوب تحقیق استفاده می‌کنیم. شکل زیر چارچوب پیشنهادی تحقیق را نشان می‌دهد. به منظور تحلیل موردنظر در چارچوب ذکر شده، مرور ادبیات مباحثی مانند داده‌کاوی، تحلیل خبر، الگوریتم k-means و متن‌کاوی و انواع روش‌های آن به‌طور مفصل مورد بررسی قرار گرفته است. این بررسی به منظور احاطه کامل تر بر مباحث و به دست آوردن ایده صورت گرفته است. فاز بعدی از تحقیق شامل بررسی سامانه موجود برای پیاده‌سازی تحلیل است. برای این منظور تحلیل‌های مختلف متن و خبر به‌عنوان مورد مطالعاتی مدنظر قرار گرفت. دلیل این انتخاب وجود داده‌های با حجم بالا در این حوزه و ارتباط تنگاتنگ با مباحث تحلیلی است.

ما در طول این تحقیق در راستای شناخت فضای موردبحث از انواع صفحات وب در زمینه اخبار و مرور ادبیات این حوزه استفاده کردیم. مرحله بعد مربوط به فاز شناخت داده و آماده‌سازی آن جهت استفاده در مرحله پیاده‌سازی است. در این مرحله بررسی می‌شود که چه داده‌هایی موردنیاز است و این داده‌های موردنیاز در چه قالبی، مناسب فاز پیاده‌سازی است. در این تحقیق از دو پایگاه داده مختلف استفاده شده است. اولین آن‌ها پایگاه داده اطلاعات نوع خبر است.

دومین پایگاه داده مربوط به متن خبرها هست که در نهایت بتوان تحلیلی دقیق‌تری به کمک الگوریتم‌های داده‌کاوی برای این بخش ارائه داد. برای اینکه بدانیم چه داده‌هایی موردنیاز است، ابتدا باید مرور ادبیات را به‌طور کامل انجام دهیم و سامانه

<sup>۱</sup> Type ID



شکل (۳): روش انجام تحقیق

انگلیسی و همچنین ۲۹۴۸ رکورد و متن عربی شناسایی و حذف گردید.

```
DELETE
FROM DBO.News_Text
WHERE LanguageId = '2'
```

الگوریتم (۳): واکشی خبر بر اساس نوع زبان

```
DELETE
FROM News_Text
WHERE LanguageId = '3'
```

الگوریتم (۴): واکشی خبر بر اساس نوع زبان

ستون «شناسه زبان» که نشان‌دهنده نوع زبان خبر بود، حذف گردید. این ستون شامل مقادیر متمایزکننده برای زبان‌های فارسی، عربی و انگلیسی بود. حال که تنها یک مقدار برای کلیه متون فارسی وجود دارد حذف این ستون اختلالی در پایگاه داده‌ای نمی‌کند.

ستون‌های «آدرس خبر<sup>۲</sup>» و «شناسه منبع» با یکدیگر ارتباط نظیر به نظیر دارند؛ بنابراین ستون «آدرس خبر» حذف گردید.

ستون شناسه نوع<sup>۳</sup> در جدول خبر دارای شش مقدار متمایز بوده و تنها درباره دو مقدار آن (۰ و ۱) اطلاع داریم. با توجه به کم بودن تعداد رکوردهای باقیمانده، دو راه کار داریم: الف) حذف رکوردهای باقیمانده ب) در نظر نگرفتن این ستون در تحلیل داده‌ها.

همچنین در جدول متن تمام رکوردهایی که مقدار ستون «شناسه فایل<sup>۴</sup>» در آن‌ها مخالف صفر است حذف شدند؛ زیرا این رکوردها صرفاً برچسب تصاویر خبرها هستند و حاوی محتوای معناداری برای تحلیل نخواهند بود.

```
DELETE
FROM News_Text
WHERE FileId != '0'
```

الگوریتم (۵): حذف خبری که شناسه فیلد آن صفر است

برای برقراری ارتباط بین جداول خبر و متن، از فیلد مشترک آن‌ها (شناسه خبر) در جدول پیوند استفاده شد. در جدول پیوند تمام رکوردهایی ظاهر می‌شوند که شناسه خبر مشترک بین جداول خبر و متن را داشته باشند؛ بنابراین «تعداد رکوردهای پیوند دو جدول» با «تعداد رکوردهای پیوند دو جدول بعد از حذف» با یکدیگر برابر هستند.

برای شناسایی خبرهای معتبر بین دو جدول از عملیات پیوند<sup>۱</sup> استفاده کردیم که تعداد سطرهای حاصله از این عملیات ۴۴۷۳۹۲ است. به این معنی که سطرهای مختلف بین دو جدول شناسایی شده و تعداد کل به دست می‌آید. با انجام این عمل مشخص می‌شود که تعدادی داده اضافی در هر دو جدول وجود دارد. تعداد شناسه خبرهایی که در جدول خبر هستند ولی در جدول متن پیدا نشده‌اند برابر ۳۳۶۷ بود، به این معنی که شناسه خبر موجود ولی متنی برای خبر وجود ندارد. همچنین تعداد خبرهایی که در جدول خبر هستند ولی در جدول (۱) متن به آن‌ها ارجاعی داده نشده ۱۵ عدد است.

همان‌گونه که در بالا بیان شد در جدول متن ۱۵ رکورد وجود داشتند که شناسه خبر آن‌ها در جدول خبر یافت نشد. این ۱۵ رکورد در واقع ۶ شناسه خبر متمایز بودند؛ بنابراین، بایستی این شناسه خبرها در جدول متن و در نهایت جدول پیوند حذف گردند.

```
delete
from dbo.News_Text
where NewsId in (select t2.NewsId
                 from dbo.News_Text t2
                 where t2.NewsId not in (select t1.NewsId
                                         from dbo.News_News t1))
```

الگوریتم (۱): حذف خبری که شناسه خبر آن‌ها در جدول متن یافت نشد

در جدول خبر ۳۳۶۷ رکورد وجود داشتند که شناسه خبر آن‌ها در جدول متن یافت نشد؛ بنابراین بایستی این شناسه خبرها در جدول خبر و نهایتاً پیوند حذف گردند.

```
delete
from dbo.News_News
where NewsId in (select t1.NewsId
                 from dbo.News_News t1
                 where t1.NewsId not in (select t2.NewsId
                                         from dbo.News_Text t2))
```

الگوریتم (۲): حذف خبری که شناسه خبر آن‌ها در جدول متن یافت نشد

با توجه به این که هدف اصلی پروژه متن‌کاوی اطلاعات موجود از داده‌های خبری فارسی حاصل از خبرگزاری‌های جمهوری اسلامی ایران است، باید کلیه رکوردهای غیرفارسی از جدول متن پاک گردد. در این راستا ۷۴۶۲ رکورد و متن

<sup>۲</sup> URL

<sup>۳</sup> Type ID

<sup>۴</sup> File ID

<sup>۱</sup> Join



## ۱۰- نتیجه‌گیری

در این مقاله، روش ترکیبی‌ای جهت تحلیل داده‌های خبری خبرگزاری‌های فارسی ارائه گردیده است. روش پیشنهادی از استاندارد CRISP تبعیت کرده و در طی انجام مراحل آن از الگوریتم‌های شرح داده‌شده کیفی و کمی داده‌کاوی و الگوریتم Kmeans استفاده می‌کند.

در طی انجام مراحل این تحقیق ابتدا مفاهیم موردنیاز که در انجام پروژه کاربرد دارند به‌صورت کامل شرح داده‌شده‌اند و در مرحله بعد به بررسی سامانه موجود، جهت پیاده‌سازی تحلیل پرداخته می‌شود. فاز بعد به شناخت داده و آماده‌سازی در جهت استفاده در مرحله پیاده‌سازی پرداخته می‌شود. درنهایت در مرحله پیاده‌سازی سعی گردیده است که با کدهای دستوری ساده و قابل‌فهم بر روی پایگاه داده‌ها و روش پیشنهادی پیاده‌سازی این تحقیق که با استفاده از نرم‌افزار RapidMine انجام گردیده به‌صورت کامل شرح داده شود.

استفاده از این روش ترکیبی تحلیل، علاوه بر فراهم کردن تحلیلی مناسب بر روی خبرها و پیدا کردن روابط پنهان بین اخبار می‌تواند مزایای فرعی متعددی همچون پایه‌گذاری روشی بومی جهت تحلیل اخبار، روشی سازگار با زبان فارسی، پایه‌گذاری تحلیل‌های معنایی بیشتر بر روی واژگان رایج در اخبار فارسی و درنهایت کمک به تصمیم‌گیران سازمان‌های نظامی و غیرنظامی جهت انجام تصمیم‌گیری کامل و صحیح را داشته باشد.

## ۱۱- مراجع

- [1] U. Fayyad, G. Piatetsky-Shapiro, and P. Smyth, "From data mining to knowledge discovery in databases," AI magazine, vol. 17, p. 37, 1996.
- [2] J. Han and M. Kamber, "Data Mining," Southeast Asia Edition: Concepts and Techniques: Morgan kaufmann, 2006.
- [3] M. J. Berry and G. S. Linoff, "Data mining," techniques: for marketing, sales, and customer relationship management: John Wiley & Sons, 2004.
- [4] H. C. Koh and C. K. Low, "Going concern prediction using data mining techniques," Managerial Auditing Journal, vol. 19, pp. 462-476, 2004.
- [5] N. Aggarwal, A. Kumar, H. Khatteer, and V. Aggarwal, "Analysis the effect of data mining techniques on database," Advances in Engineering Software, vol. 47, pp. 164-169, 2012.
- [6] D. T. Larose, "k- Nearest Neighbor Algorithm," Discovering Knowledge in Data: An Introduction to Data Mining, pp. 90-106, 2005.

```
SELECT *
FROM dbo.News_Text t1
INNER JOIN dbo.News_News t2
ON t1.NewsId = t2.NewsId
```

الگوریتم (۶): واکشی متن خبر با شناسه مشترک

تمامی تغییرات صورت گرفته جهت شناسایی و پیش‌پردازش داده‌های موجود در جدول زیر به نمایش درآمده است.

جدول (۱): نتایج حاصله از پردازش صورت گرفته اخبار

تعداد	توضیح
۱۶۳۴۸۶	تعداد رکوردهای جدول خبر (سطر یکتا): شناسه خبر)
۴۴۷۴۰۷	تعداد رکوردهای جدول متن (سطر یکتا): شناسه متن)
۴۴۷۳۹۲	تعداد رکوردهای پیوند دو جدول
۳۳۶۷	تعداد شناسه خبرهایی که در جدول خبر هستند ولی در جدول متن پیدا نشدند
۱۵	تعداد شناسه خبرهایی که در جدول متن هستند ولی در جدول خبر به آن‌ها ارجاعی داده نشده است.
۱۶۰۱۲۵	تعداد رکوردهای جدول متن که شناسه خبر یکتا دارند.
۱۶۰۱۱۹	تعداد رکوردهای جدول خبر بعد از حذف شناسه خبر ناموجود در جدول متن
۴۴۷۳۹۲	تعداد رکوردهای جدول متن بعد از حذف شناسه خبر ناموجود در جدول خبر
۴۴۷۳۹۲	تعداد رکوردهای جدول متن که حاوی متن انگلیسی می‌باشند.
۱۸۱۰۹	تعداد رکوردهای جدول متن که حاوی برچسب تصاویر بودند.
۴۱۸۸۷۳	تعداد رکوردهای نهایی تا پایان مرحله ۹

جدول حاصل از پیوند دو جدول خبر و متن شامل اطلاعات، شناسه خبر، شناسه متن، نوع متن، متن، شناسه گروه، شناسه منبع، وضعیت و شناسه نوع هست. بدین ترتیب تعداد ستون‌های داده از ۲۰ ستون به ۸ ستون کاهش یافت. دلیل عمده این کاهش بعد نداشتن اطلاعات کافی از ستون موردنظر و یا افزونگی داده‌ای بود.

- [16] S. E. Robertson, "The probability ranking principle," *Journal of Documentation*, vol. 33, pp. 294–304, 1977.
- [17] C. J. Van Rijsbergen, "A non-classical logic for information retrieval," *The Computer Journal*, vol. 29(6), pp. 481–485, 1986.
- [18] H. Zhuge et al., "An Automatic Semantic Relationships Discovery Approach," *The 13th International World Wide Web Conference (WWW2004)*, New York, USA, May 2004.
- [19] T. Joachims, "Text categorization with support vector machines: Learning with many relevant features," In C. Nedellec and C. Rouveirol, editors, *European Conf. on Machine Learning (ECML)*, 1998.
- [20] S. Dumais, J. Platt, D. Heckerman, and M. Sahami, "Inductive learning algorithms and representations for text categorization," In *7th Int. Conf. on Information and Knowledge I*, 1998.
- [21] K. Nigam, A. McCallum, S. Thrun, and T. Mitchell, "Text classification from labeled and unlabeled documents using em," *Machine Learning*, vol. 39, pp. 103–134, 2000.
- [22] F. Sebastiani, "Machine learning in automated text categorization," *ACM Computing Surveys*, vol. 34, pp. 1–47, 2002.
- [23] N. Kanya and S. Geetha, "Information Extraction," *A Text Mining Proach*, produced IEEE, 2007.
- [24] H. Karanikas, C. Tjortjis, and B. theodoulidis, "An approach to text mining information extraction," 2001.
- [25] M. Rajman, "Text Mining, Knowledge extraction from unstructured textual data," *Proc. of Eurostat Conference*, Francfort (Deutschland), May 1997.
- [7] M. S. Deshpande and D. V. Thakare, "Data mining system and applications: A review," *International Journal of Distributed and Parallel systems (IJDPDS)*, vol. 1, pp. 32–44, 2010.
- [8] Y. M. Chae, S. H. Ho, K. W. Cho, D. H. Lee, and S. H. Ji, "Data mining approach to policy analysis in a health insurance domain," *International journal of medical informatics*, vol. 62, pp. 103–111, 2001.
- [9] R. Alguliev and R. Aliguliyev, "Experimental investigating the F-Measure as similarity measure for automatic text summarization," *Applied and Computational Mathematics*, vol. 6, no. 2, pp. 278–287, 2007.
- [10] M. E. Califf and R. J. Mooney, "Bottom-up relational learning of pattern matching rules for information extraction," *Journal of Machine Learning Research*, vol. 4, pp. 177–210, 2003.
- [11] M. A. Hearst, "Untangling text data mining," In *Proceedings of the 37th annual meeting of the Association for Computational Linguistics on Computational Linguistics*, Association for Computational Linguistics, pp. 3–10, June 1999.
- [12] U. M. Fayyad, G. Piatetsky-Shapiro, and P. Smyth, "Knowledge discovery and data mining: Towards a unifying framework," In *Knowledge Discovery and Data Mining*, pp. 82–88, 1996.
- [13] V. Kumar and M. Joshi, *What is datamining?*, 2003. <http://wwwusers.cs.umn.edu/~mjoshi/hpdmntut/sld004.htm>
- [14] R. Feldman and I. Dagan, "Kdt – knowledge discovery in texts," In *Proc. Of the First Int. Conf. on Knowledge Discovery (KDD)*, pp. 112–117, 1995.
- [15] G. Salton, A. Wong, and C. S. Yang, "A vector space model for automatic indexing," *Communications of the ACM*, vol. 18(11), pp. 613–620, 1975. (see also TR74-218, Cornell University, NY, USA).

---

## **Provides a Hybrid Method of Analyzing News Databases Using RapidMine Case Study: Persian News Texts**

**S. Ahmadian \*, K. Saberi**

Imam Hossein Comprehensive University

### **Abstract**

One of the most widely used approaches in the management of Web-based systems, use the online agency and Data related to them . Online web agency, like the heart of a community and is one of the most critical information resources in a society. The importance of this topic in the news military organizations will be doubled. So having a right way and a unique conceptual approach to analyze and categorize the resources, Can provide numerous benefits and assistance to decision-makers within the organization, especially military organizations in carrying out decisions. In this study, firstly we contribute to the study and understanding of the concepts of data mining methods. Secondly, analysis of different methods of text processing must be done from different perspectives So as to identify the positive strategies and strengths of a variety of methods and The Finally, the analysis of news database, so as to uncover information hidden in the data and how to use it. Keywords: Data mining, text mining, News Database, Mixed-Method.

**Keywords:** Data Mining, Text Mining, News Databases

---

\* Corresponding author E-mail: s.ahmadian@ihu.ac.ir