

مدیریت مستندات محتوای فارسی رسانه‌های آنلاین خبری در جامعه اطلاعاتی

● حمید میرزائی دهنوی ●

گروه مدیریت فناوری اطلاعات، واحد تهران مرکزی، دانشگاه آزاد اسلامی، تهران، ایران.

● محمد علی کرامتی ●

گروه مدیریت صنعتی، واحد تهران مرکزی، دانشگاه آزاد اسلامی، تهران، ایران. (نویسنده مسئول)

● محمد علی عامری حبیب آبادی ●

گروه مدیریت، پژوهشگاه علوم انتظامی و مطالعات اجتماعی، تهران، ایران.

تاریخ پذیرش: ۱۴۰۱/۰۲/۲۱

تاریخ دریافت: ۱۴۰۱/۱۱/۱۵

چکیده

عصری که ما در آن زندگی می‌کنیم، عصر اطلاعات است و برای سازمان‌ها مهم‌ترین مسئله، اشراف بر همین اطلاعات است. با رشد روزافزون اخبار در دنیای دیجیتال و اینترنت، موضوعی که مهم و حائز اهمیت می‌شود، دسته‌بندی این اطلاعات و دسترسی سریع و ارزان ما به آن‌ها می‌باشد. این مهم به دست نمی‌آید مگر با انجام روش‌هایی که از آن به‌عنوان دسته‌بندی متون یاد شده است. هدف این پژوهش دسته‌بندی متون خبری در دسته‌های از قبل مشخص شده می‌باشد که با استفاده از ابزار مدل اتوماتیک که یکی از زیرمجموعه‌های متن‌کاوی محسوب می‌شود صورت می‌پذیرد. با توجه به اهمیت موضوع و کاری که در این زمینه برای زبان‌های دیگر دنیا انجام گرفته است، نیاز به طبقه‌بندی متون فارسی به خوبی احساس می‌شود. این نکته قابل توجه می‌باشد که تحقیقات برای متون انگلیسی توسعه داده شده و به کارگرفته می‌شود اما از آنجایی که زبان فارسی پیچیدگی‌های ساختاری نسبت به سایر زبان‌ها دارد و همچنین تحقیقات کمتری در این زمینه انجام گرفته است، این پژوهش از نوع کاربردی، توسعه‌ای می‌باشد که برای انجام آن می‌توان به روش پژوهش آزمایشی و استفاده از ابزار متن‌کاوی اشاره کرد، به این صورت که در یک محیط کاملاً تحت کنترل و با توانایی ثابت نگه داشتن سایر متغیرها انجام می‌گردد. در جامعه اطلاعاتی دسته‌بندی متون به‌وسیله افراد نخبه و به‌صورت دستی انجام می‌گیرد. دسته‌بندی متون آن‌هم با این حجم و به‌صورت دستی غیرممکن به نظر می‌رسد، بنابراین ما ناگزیریم که به دنبال روش‌هایی برای دسته‌بندی خودکار متون باشیم. از سوی دیگر ذخیره‌سازی، پردازش و تحلیل این حجم از اطلاعات تبدیل به چالشی جدی شده است. با توجه به حجم بالای اخبار، داده‌ها، اطلاعات، اسناد و پیچیدگی حفظ و نگهداری آنها، لازم است سیستمی جهت مدیریت دریافت، حفظ و نگهداری اخبار موجود، بکار گرفته شود. پیچیدگی سازمان‌ها نیاز به متمرکز بودن اخبار، اسناد، طبقه‌بندی درست، گردش صحیح اخبار و سهولت در دسترسی به آنها را ایجاد می‌نماید. مدیریت مستندات این امکان را برای سازمان‌های اطلاعاتی فراهم می‌آورد که اخبار و اسناد دریافتی با موجود را به درستی طبقه‌بندی نموده، آنها را حفظ، نگهداری و بازیابی نماید. با بررسی، تحلیل و پردازش در این پژوهش به این نتیجه می‌رسیم که دقت و نتایج حاصل روش پیشنهادی روی متون خبری برخط نشان می‌دهد: مدل ماشین بردار پشتیبان دارای دقت ۹۳/۲۹، صحت ۹۳/۳۲، فراخوانی ۹۲/۹۶ و خطای ۶/۷۱ است.

واژگان کلیدی: طبقه‌بندی، مدیریت مستندات، محتوای فارسی، متن‌کاوی، اطلاعات و دسته‌بندی

یکی از زیرمجموعه‌های اسناد متنی که هر روزه با آن روبه‌رو هستیم، اخبار است. برخلاف چند دهه گذشته امروزه دیگر حجم غالب اخبار در روزنامه‌ها و مجلات چاپ نمی‌شود، بلکه بیشتر حجم اخبار منتشرشده در سایت‌های خبری و شبکه‌های اجتماعی دیده می‌شود. وجود حجم بسیار زیاد اخبار در این پایگاه‌های خبری، نیاز به طبقه‌بندی موضوعی اخبار را ایجاد می‌کند. به عبارت دیگر، پیش از انتشار اخبار در سایت‌های مختلف خبری نیاز است که دسته‌ی موضوعی هر متن خبری مشخص شود. با دسته‌بندی متون خبری که بر روی سایت‌های خبری منتشر می‌شوند، کاربران این سایت‌ها می‌توانند اخبار مورد نظر خود را آسان‌تر و سریع‌تر پیدا کنند. از طرفی به دلیل حجم بالای اخباری که هر روزه منتشر می‌شوند، طبقه‌بندی اخبار به صورت دستی فرآیندی پرهزینه و نیازمند تعداد زیادی نیروی انسانی متخصص است و علاوه بر این، معیار زمان که در انتشار اخبار، عامل بسیار مهمی است، در طبقه‌بندی اخبار به صورت دستی بسیار زیاد است. این به آن معنی است که هر سایت خبری تلاش می‌کند اخبار روزانه را سریع‌تر از دیگر سایت‌ها در صفحه‌ی خود بارگذاری کند. به همین دلیل است که طبقه‌بندی متون خبری به صورت خودکار امری بسیار کاربردی و لازم است. با طبقه‌بندی خودکار متون خبری، نویسندگان سایت‌ها و شبکه‌های خبری می‌توانند صدها متن خبری را سریع‌تر، آسان‌تر و با صرف هزینه و نیروی متخصص کمتری منتشر نمایند.

جامعه اطلاعاتی داده‌های متنی بسیاری را رصد می‌کنند، چگونه می‌توان این حجم از داده‌ها را مدیریت کرد؟ چگونه می‌توان اطلاعات مفیدی از این داده‌های متنی بی‌شمار بدست آورد؟ آیا این داده‌ها ارزشی دارند؟ متن کاوی فناوری مورد استفاده برای چنین مواردی است. با فناوری‌های متن کاوی می‌توان داده‌های متنی را بررسی و تحلیل کرده و از نتایج حاصل از این تحلیل، اطلاعات ارزشمندی کسب نمود. داده‌های متنی هیچ‌گونه ارزشی ندارند مگر اینکه متن کاوی شوند. متن کاوی به شناسایی الگوها، کلمات کلیدی، موضوعات و دیگر ویژگی‌های موجود در متن می‌پردازد. الگوریتم‌ها و روش‌های مختلفی برای متن کاوی وجود دارد و همچنان نیز در حال توسعه می‌باشد، اما مشکل اصلی این روش‌ها برای دستگاه‌های اطلاعاتی حفظ محرمانگی اطلاعات است. باید در کنار استفاده از علوم و فناوری‌های نوین اصل محرمانگی اطلاعات نیز حفظ شود. بسیاری از برنامه‌های موجود در موضوع داده کاوی و متن کاوی در بستر اینترنت کار می‌کنند که عملاً برای جامعه اطلاعاتی بدون استفاده است. خوشبختانه چند سالی است که شرکت‌های داخلی شروع به فعالیت در

زمینه متن کاوی نموده‌اند و پیشرفت‌های قابل توجهی هم در این زمینه داشته‌اند، اما کماکان نمی‌توان برای مجموعه‌های امنیتی و نظامی از این گونه نرم‌افزارها استفاده نمود. لذا لازم است مجموعه‌های اطلاعاتی امنیتی این موضوع را به صورت بومی پیاده‌سازی کنند و یک بستر مناسب برای تحلیل اطلاعات داشته باشند. لذا بر آن شدیم تا در این مقاله به مدیریت مستندات محتوای فارسی رسانه های آنلاین خبری پردازیم که خروجی آن برای سازمان های اطلاعاتی این است که بتوانند از سایت‌های خبری، متون مختلف را جمع آوری و سپس دسته‌بندی و بازیابی نمایند.

روش شناسی پژوهش

امروزه سامانه‌های طبقه‌بند متون گسترده‌تری پیدا کرده‌اند و وجود یک سیستم طبقه‌بند متون فارسی که بهینه باشد بسیار احساس می‌شود. طبقه‌بندی متون، یکی از زیرمجموعه‌های متن کاوی است. متن کاوی بر روی پردازش متون تمرکز دارد. در پردازش متون سعی می‌شود دانشی از متون خام استخراج گردد. در طبقه‌بندی متون هدف مشخص کردن دسته متن است. به طور مثال: خبری که در یک پایگاه خبری درج شده است یک خبر سیاسی یا ورزشی است. در این مثال سیاسی و ورزشی بودن خبر دو دسته برای طبقه‌بندی اخبار را تشکیل می‌دهند. طبقه‌بندی متون قدمت بسیار زیادی در حوزه متن کاوی دارد. این موضوع از سال ۱۹۶۰ میلادی مورد توجه محققین بوده است ولی با رونق کامپیوتر و نرم‌افزارها، مانند دیگر موضوعات هوش مصنوعی و داده کاوی در دهه ۹۰ میلادی توجه به این موضوع رشد چشمگیری پیدا کرد و مورد توجه قرار گرفت. از آن جایی که در انجام پژوهش‌های مرتبط با هوش مصنوعی و زبان شناسی رایانشی، از علوم زبان شناسی و علوم رایانه استفاده می‌شود، ماهیت بین رشته‌ای علم زبان شناسی رایانشی در تعیین روش و ابزار انجام پژوهش حاضر نقش به سزایی دارد. در این بخش در پی آشنایی با گام‌های پیشنهاد شده برای انجام پژوهش می‌باشیم. برای انجام تحقیق باید گام‌های ذیل مشخص شود:

الف) جمع آوری متون خبری فارسی (اخبار حوزه ایران، اخبار حوزه افغانستان و اخبار سطح جهانی)

ب) پیش پردازش و آماده سازی اخبار

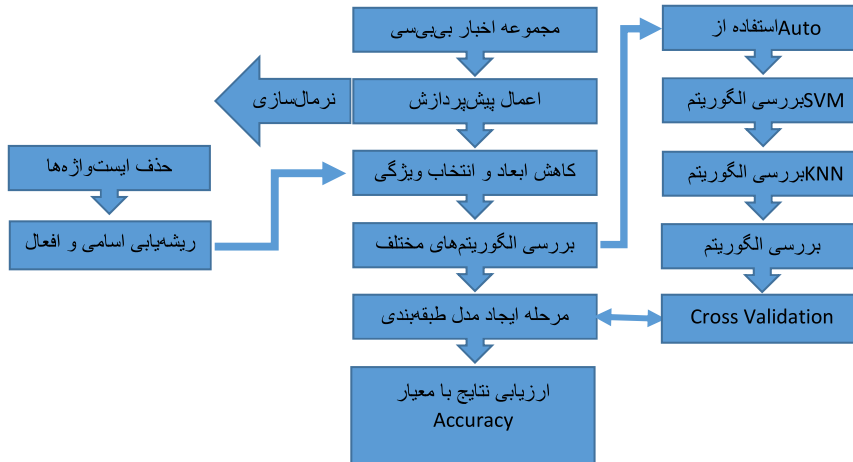
ج) انتخاب ویژگی

د) مقایسه الگوریتم‌های یادگیری ماشین برای دسته‌بندی (استفاده از مدل اتوماتیک^۱)

ه) ایجاد مدل موردنظر برای دسته‌بندی^۱

و) ارزیابی مدل ارائه شده

به طور خلاصه اجزای اساسی سیستم پیشنهادی در تصویر زیر قابل مشاهده می‌باشد.



شکل (۱): مراحل مختلف ایجاد مدل طبقه‌بندی متون خبری

با توجه به اینکه این پژوهش از داده‌های جمع‌آوری شده پایگاه خبری فارسی بی بی سی استفاده می‌کند، در گروه پژوهش‌های موردی قرار می‌گیرد. همچنین به دلیل اینکه به بررسی و طبقه‌بندی متون خبری فارسی می‌پردازد، یک تحقیق کاربردی به حساب می‌آید. به علاوه چون از داده‌های بی‌بی‌سی آرشیو شده استفاده می‌کند، مقطعی است. طبقه‌بندی اسناد در این پژوهش، مبتنی بر الگوریتم‌های یادگیری ماشین است. الگوریتم‌های مورداستفاده در این پژوهش عبارت است از: ماشین بردار پشتیبان، بیزین ساده و کی-نزدیک‌ترین همسایه. شکل زیر نمایی کلی از فرآیند پیشنهادی پژوهش را ارائه می‌دهد.

مبانی نظری

متن کاوی

متن کاوی شاخه‌ای از داده کاوی یا همان کشف دانش است. از نظر فیاد و همکارانش، کشف دانش، فرآیند غیربديهی تشخیص الگوهای معتبر، نو، مفید و درنهایت قابل درک در داده‌هاست. متن کاوی از

1 Classification.

فناوری‌های بازیابی اطلاعات و استخراج اطلاعات استفاده می‌کند (Usama, et.al ۱۹۹۶).

آماده‌سازی متون سیستم طبقه‌بندی اخبار فارسی

داده‌های ورودی به برنامه طبقه‌بند اخبار، متونی هستند که از سایت بی‌بی‌سی جمع‌آوری شده‌اند، با توجه به اینکه این متون در ویرایشگرهای فارسی متفاوتی تایپ شده است و نیز هنگام بارگذاری روی اینترنت ممکن است دچار تغییراتی شده باشد لازم است در ابتدا بازبینی کلی روی متون انجام شود. قسمتی از بازبینی به شیوه دستی انجام می‌شود، ولی در حجم تعداد بالا ویرایش دستی ممکن نیست و باید از روش‌های مرتبط استفاده نمود.

پیش‌پردازش متن

پس از آماده‌سازی اولیه متون، پیش‌پردازش انجام می‌شود. در واقع پیش‌پردازش، اولین گام برای تطابق مستندات متنی با نمایش آن‌ها در یک قالب مناسب می‌باشد. ثابت شده است که تنها ۳۳ درصد کلمات در متن مفید هستند و می‌توان از آن‌ها برای استخراج اطلاعات استفاده نمود (هاشمی، ۱۳۹۴).

با اینکه مجموعه گسترده‌ای از روش‌ها در پردازش زبان طبیعی به کار می‌روند، فناوری‌های به کار رفته را می‌توان به سه دسته کلی تقسیم نمود. روش‌های آماری، روش‌های ساختاری مبتنی بر الگو و روش‌های مبتنی بر استنتاج. باید توجه داشت که این راهکارها لزوماً از هم جدا نیستند. در واقع، جامع‌ترین مدل‌ها از ترکیب هر سه این روش‌ها استفاده می‌کنند. تفاوت این راهکارها در نوع عملیات پردازشی است که قادر به انجام آن هستند و میزان قواعدی که در مقابل آموزش یادگیری خودکار از روی داده‌های زبانی نیاز دارند (بهرام پور و همکاران، ۱۳۹۴).

برای کاوش کردن حجم قابل توجهی از اسناد، ضروری است که اسناد پیش‌پردازش شوند و اطلاعات در یک ساختار داده‌ای مناسب برای پردازش‌های بعدی ذخیره گردد. زمانی که داده‌های ورودی در دسترس است باید آن‌ها را برای ورود به الگوریتم‌های یادگیری ماشین آماده کنیم. این مرحله برای داده‌های متنی به این معناست که آن‌ها را از حالت غیر ساخت یافته به فرمت ساختاریافته و قابل تشخیص برای ماشین تبدیل کنیم (پرئی و همکاران، ۱۳۹۵).

پیش‌پردازش

شامل پیش‌پردازش متون و پیش‌پردازش ادبی است. پیش‌پردازش متون، بر روی مجموعه اسناد

جمع آوری شده، نشانه گذاری^۱، حذف کلمه توقف یافته^۲ و ریشه یابی^۳ انجام می شود. پیش پردازش ادبی برای افزایش اطلاعات عبارات به کار می رود. بدین منظور دیدگاه ها، برچسب گذاری ادات سخن، تکه کردن متن، ابهام زدایی حس کلام^۴ و تجزیه کردن^۵ را مکرر اعمال می کند (آقا کاردان و همکاران، ۱۳۹۱).

طبقه بندی متن

محوری ترین موضوع در حوزه متن کاوی، طبقه بندی متون است. اساس کار طبقه بندی متون بر پایه کلمات کلیدی و مهمی است که از مستندات استخراج می شود (کریمی منش و همکاران، ۱۳۹۲). برای دسته بندی متن از فناوری های استخراج اطلاعات^۶، پردازش زبان طبیعی^۷ و یادگیری ماشین^۸ به طور وسیع استفاده می شود. به طور کلی هدف یک دسته بند متون، دسته بندی اسناد در قالب تعداد معینی از دسته های از پیش تعیین شده می باشد. هر سند می تواند در یک، چند و یا هیچ دسته ای قرار بگیرد. در مورد هر سند به این سؤال پاسخ داده خواهد شد که این سند در کدام یک از دسته ها قرار می گیرد. این موضوع می تواند در قالب یک یادگیری خودکار قرار گیرد تا با استفاده از آن بتوان هر سند را به طور خودکار به دسته ای نسبت داد. بعد از انتخاب مجموعه داده و پاک سازی متون به کمک روش نرمال شده فرکانس کلمه-معکوس فرکانس سند به ویژگی ها وزن داده می شود و در طی دو مرحله ویژگی ها با استفاده از روش فرکانس سند و مربع کای^۹ انتخاب می شوند (سیاحی و همکاران، ۱۳۹۴).

هدف از طبقه بندی متون نسبت دادن کلاس های از پیش تعریف شده به اسناد متنی است. در طبقه بندی، یک مجموعه آموزشی از اسناد با کلاس های معین وجود دارد. با استفاده از این مجموعه، مدل طبقه بندی معین شده و کلاس سند جدید مشخص می گردد. برای کارایی مدل طبقه بندی، یک مجموعه آزمودن، مستقل از مجموعه آموزشی در نظر گرفته می شود. برچسب های تخمین زده شده

1. tokenization
2. Stop-word removal
3. stemming
4. Word Sense Disambiguation(WSD)
5. Text Chunking
6. Information Extraction
7. Natural Language Processing(NLP)
8. Machine Learning
9. Chi Square

با برچسب واقعی اسناد مقایسه می شود. نسبت اسنادی که به درستی طبقه بندی شده اند به تعداد کل اسناد، دقت^۱ نامیده می شود.

مراحل اصلی طبقه بندی متن

مراحل اصلی طبقه بندی متن در دو فاز که فاز اول شامل دو مرحله انتخاب ویژگی و وزن دهی و فاز دوم شامل تخمین عملکرد و انتخاب یک الگوریتم طبقه بندی می باشد دسته بندی می شود. (کریمی منش و همکاران، ۱۳۹۲).

کاربردهای طبقه بندی متن

دسته بندی متون کاربردهای فراوانی می تواند داشته باشد که از جمله آن می توان به دسته بندی گفتاری که ترکیبی از دسته بندی متون و تشخیص گفتار است، دسته بندی متون چند رسانه ای از طریق عنوان های متنی، تشخیص نویسنده برای متون نامشخص با مورد بحث، تشخیص زبان برای متونی که زبان آن ها نامشخص است، تشخیص خودکار جنس متن، بایگانی اسناد، دسته بندی صفحات وب، یادگیری خودکار علایق مطالعاتی و پژوهشی کاربران، فیلتر کردن خودکار پست الکترونیکی بر اساس محتوا و... اشاره نمود (جمالی و همکاران، ۱۳۹۶).

پردازش زبان طبیعی^۲

پردازش زبان طبیعی در دهه ۵۰ میلادی به عنوان فصل مشترک علوم هوش مصنوعی و زبان شناسی به وجود آمد و در دهه ۸۰ دچار تحولات و تغییرات اساسی شد (Nadkarni, et.al, ۲۰۱۱). پردازش متن تا دهه ۸۰ میلادی با استفاده از روش های قانون محور انجام می گرفت. اگرچه روش های قانون محور ویژگی های مثبتی از قبیل قابل درک بودن برای انسان و قابلیت ارتقا کیفیت در طول زمان را دارا هستند، مشکلات و محدودیت هایی نیز دارند. لازمی نوشتن قوانین در روش های قانون محور داشتن دانش عمیقی از آن حوزه است و از طرفی این روش ها بسیار زمان بر بوده و از نظر محاسباتی پیچیده هستند (Moulinier, et.al, ۱۹۹۵).

این مشکلات و محدودیت ها باعث شد در دهه ۹۰ میلادی روش های یادگیری ماشین یا به عبارتی روش های آماری محور در پردازش متون طبیعی محبوب تر از روش قانون محور ظاهر شوند. در واقع در دهه ۹۰ میلادی بود که پژوهشگران به روش های یادگیری ماشین توجه بسیاری نشان دادند که این

1.precision
2.Natural language processing(NLP)

به دلیل توسعه‌ی سخت‌افزاری و نرم‌افزاری علم رایانه اتفاق افتاد (Eyheramendy, et.al, ۲۰۰۳). پژوهش‌ها نشان می‌دهند که طبقه‌بندی متن از نخستین تحقیقات انجام گرفته در حوزه‌ی متن‌کاوی است و نقش مهمی در آن ایفا می‌کند. به گفته‌ی فرانسیس (Francis, ۲۰۰۶). همچنین به گفته‌ی گلباخ و بلشاکو طبقه‌بندی متن از دهه ۶۰ میلادی آغاز شده و در دهه‌های ۸۰ و ۹۰ میلادی توجه محققان بیشتری را به خود جلب کرده است (BolshaKov, et.al, ۲۰۰۴).

هدف کلی پردازش زبان طبیعی رسیدن به یک درک بهتر از زبان طبیعی توسط کامپیوترهاست. کامپیوترها فناوری‌های مستحکم و ساده‌ای برای پردازش سریع متن به کار می‌برند. هم‌چنین از فناوری‌های آنالیز زبان‌شناسی نیز برای پردازش متن استفاده می‌کنند. نقش پردازش زبان طبیعی در متن‌کاوی فراهم کردن یک سیستم در مرحله استخراج اطلاعات با داده‌های زبانی است.

پیشینه پژوهش

پیشینه پژوهش در قالب جدول شماره ۱، شامل نام پژوهشگر، عنوان پژوهش، یافته‌های پژوهشی و سال پژوهشی قابل دسترسی می‌باشد.

جدول شماره (۱) پیشینه پژوهش

پژوهشگر	سال	عنوان پژوهش	یافته‌های پژوهش
سیده عربی زری و همکاران	۱۳۸۱	استخراج کلمات کلیدی جهت طبقه‌بندی متون فارسی	از ترکیب روش‌های Wordnet و Porter استفاده کرده‌اند. برای استخراج کلمات کلیدی از ساختار درهم‌ریزی Trie و دادگان مورد استفاده، دادگان همشهری بوده است.
محمد احسان بصیری	۱۳۸۴	مقایسه دسته‌بندی متون فارسی با استفاده از الگوریتم‌های KNN و FKNN و انتخاب ویژگی‌ها بر اساس بهره اطلاعات و فرکانس سند	در این مقاله دو الگوریتم KNN و FKNN با دو روش استخراج ویژگی IG و DF آزموده می‌شود که بهترین ترکیب استفاده از روش FKNN و IG است. با این ترکیب به میانگین دقت ۸۰٪ می‌رسیم.
مسلم محمدی و همکاران	۱۳۸۷	استفاده از شبکه‌های عصبی CC4 برای رده‌بندی اسناد فارسی	در مرحله پیش‌پردازش کلمات عمومی حذف و سایر کلمات ریشه‌یاب می‌شوند. سپس با استفاده از روش مبتنی بر فراوانی کلمه، اسناد با اندازه‌های مختلف به یک فضای K بعدی با اندازه‌ی ثابت نگاشت می‌شوند. دقت روی دادگان ایسنا حدود ۹۰٪ گزارش شده است.
ایوب باقری و همکاران	۱۳۸۱	دسته‌بندی متون خبری فارسی با استفاده از الگوریتم بیز ساده	در این مقاله با استفاده از بیز ساده و کاهش فضای حالت با روش واریانس فرکانس مدت دسته‌بندی قابل قبولی ارائه شد.
ایمان بیبا و همکاران	۱۳۸۱	طبقه‌بندی خودکار متون فارسی	از روش‌های tri-gram و quad-gram با معیارهای اندازه‌گیری فاصله منتهن، اندازه‌ی دایس و ضرب نقطه‌ای، به همراه روش یادگیری KNN مورد بررسی قرار گرفته و بهترین نتیجه از ترکیب quad-gram و ضرب نقطه‌ای به دست آمد.
پوشین تقی‌زاده و همکاران	۱۳۸۱	ارائه روشی جدید در طبقه‌بندی متون فارسی با استفاده از دانش معنایی	با استفاده از معیارهای ICF و Uni، ویژگی‌ها را انتخاب کرده و با استفاده از گنج‌وازه و وزن‌دهی TFIDF و SVM دسته‌بندی متون پرداخته شده و به نتایج پایداری رسیده است.

پژوهشگر	سال	عنوان پژوهش	یافته‌های پژوهش
علی قنبری سرنی	۱۳۹۰	بهبود عملکرد طبقه‌بندی متون فارسی با استفاده از تجزیه و تحلیل مؤلفه‌های اصلی با کمک معیار میانگین یادآوری و دقت	نتایج به دست آمده نشان داده است که با در نظر گرفتن روش‌های دسته‌بندی KNN و Bayesian در روش پیشنهادی بهبود قابل توجهی در طبقه‌بندی متون فارسی و کاهش مدت زمان آزمون با ویژگی‌های استخراج شده به دست خواهد آمد.
حمید حسن پور	۱۳۹۱	استخراج بهترین ویژگی از متون فارسی با استفاده از تجزیه و تحلیل مؤلفه‌های اصلی با کمک میانگین یادآوری و الگوریتم ژنتیک	روش وزن‌دهی TFCRF را با استفاده از معیارهای ارتباط مثبت و منفی، به کار برده‌اند. آن‌ها با دو دسته‌بند KNN و بیزین به دسته‌بندی متون پرداختند و نسبت به کارهای قبلی نتایج بهتری گرفتند
محسن طاهری نیا	۱۳۹۱	دسته‌بندی متون فارسی با استفاده از یادگیری نیمه نظارت شده	این پژوهش از ترکیبی از مثال‌های آموزشی برجسب دار و بدون برجسب برای یادگیری استفاده می‌کند و از این فن برای دسته‌بندی متون فارسی استفاده کرده و نتایج خوبی بدست آمد.
محسن زلفی و همکاران	۱۳۹۲	دسته‌بندی متون فارسی با استفاده از روش آنالیز معنایی پنهان احتمالاتی	از روش فاصله‌یابی اقلیدسی در فضای ماتریس‌های کاهش بعد استفاده کرده است. نتایج حاصل شده نشان می‌دهد که این روش در بهبود عملکرد سیستم نقش مؤثری دارد.
مهتابی برافانی و همکار	۱۳۹۲	استفاده از ترکیب شبکه‌های عصبی جهت دسته‌بندی متون فارسی مبتنی بر الگوریتم‌های GA، کی-نزدیکترین همسایه، PCA جهت کاهش ویژگی	از ترکیب دو شبکه عصبی پرسپترون MLP در دسته‌بندی مستندات XML بر روی پایگاه داده روزنامه همشهری استفاده کرده‌اند. از روش TFCRF برای وزن‌دهی و PCA برای ارزش‌دهی ویژگی‌ها و از الگوریتم KNN برای دسته‌بندی استفاده شده است.
نیلا عبدالله عابدینی نیا و همکاران	۱۳۹۲	کاوش متون فارسی در وب با استفاده از تحلیل معنایی و روش طبقه‌بندی چندلایه‌ای اطلاعات	از مدل فضای برداری برای دسته‌بندی متون فارسی در وب استفاده شده است. روش پیشنهادی این امکان را می‌دهد که با دقت و سرعت بیشتری متن خود را در وب بیابند.
حمید حسن پور و همکار	۱۳۹۳	بهبود دقت سیستم دسته‌بندی خودکار اسناد فارسی به کمک هستان‌شناسی فارسی نت	معیار χ^2 برای انتخاب ویژگی و روش وزن‌دهی TFIDF در وزن‌دهی به کار گرفته شده است. نتایج نشان داده است که این روش یک گام مؤثر در بهبود کارایی دسته‌بند است.
پیمان جلالی و همکاران	۱۳۹۴	ارائه‌ی روش دسته‌بندی متون با تکنیک کاهش ویژگی فیلتری و یادگیری ماشین	یک روش جدید انتخاب ویژگی بر اساس الگوریتم ژنتیک ارائه می‌شود. اثربخشی روش پیشنهادی با استفاده از الگوریتم بیزین ساده و طبقه‌بندی انجمنی در سه مجموعه مختلف از داده‌های متن عربی ارزیابی می‌شود.

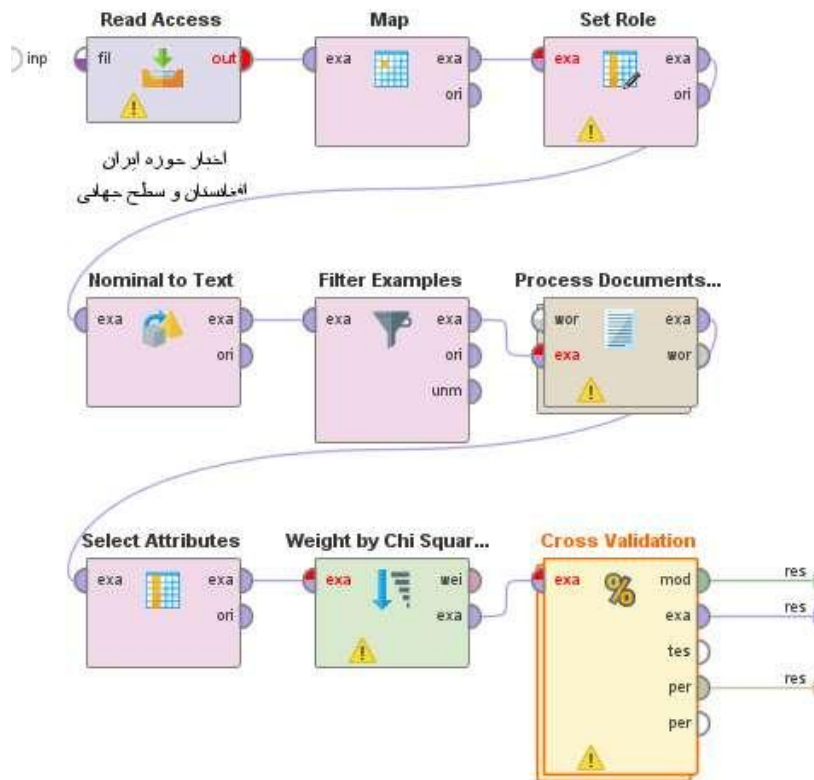
پژوهشگر	سال	عنوان پژوهش	یافته‌های پژوهش
عارف، ساسی و همکاران	۱۳۹۴	بهبود یک روش مبتنی بر انتخاب ویژگی به منظور دسته‌بندی متون با الگوریتم‌های متن کاوی	از الگوریتم انتخاب ویژگی فیلتری، برای کاهش پیچیدگی دسته‌بندی و از الگوریتم‌های یادگیری بیز ساده، درخت تصمیم و ماشین بردار پشتیبان برای ارزیابی و بهبود کارایی استفاده شده است.
ایمان ابراهیمی و همکاران	۱۳۹۴	رده‌بندی متون فارسی با استفاده از ماشین بردار پشتیبان مبتنی بر روش‌های انتخاب ویژگی PCA و الگوریتم ژنتیک	ابتدا با استفاده از روش TF-CRF به کلمات وزن‌دهی انجام می‌شود سپس از روش‌های PCA و الگوریتم ژنتیک استفاده شده است. نتایج تجربی نشان می‌دهد که این روش می‌تواند با دقت ۷۸٪ عمل رده‌بندی را انجام دهد.
محبوبه صبایی	۱۳۹۴	الگوریتمی جدید برای طبقه‌بندی مستندات، مبتنی بر وزن‌دهی به ویژگی‌ها و فایل‌ها	با استفاده از روش PCA ابعاد ویژگی‌ها کاهش داده می‌شود سپس با استفاده از SVM به پیاده‌سازی مدل پیشنهادی می‌پردازیم و در نهایت صحت مدل با روش اعتبار سنجی ۱۰ مرحله‌ای به صحت ۹۱٫۸۶٪ می‌رسد که نسبت به کارهای پیشین انجام گرفته صحت بالاتری دارد.
عبدالصمد کهرآزمی	۱۳۹۷	بهینه‌سازی نتایج موتورهای جستجو برای مفاهیم تخصصی حوزه پزشکی با به کارگیری روش‌های پردازش متن و داده کاوی	ا هزار مورد استفاده تانگرا و متلب بوده و دقت روش پیشنهادی نسبت به الگوریتم KNN ۰۰۰۵۱٪، نسبت به شبکه عصبی ۰۰۱۳۱٪ و نسبت به الگوریتم نایویز ۰۰۳۱۱٪ بهبود داشته است.
محمد بهروزیان، زیاد و همکاران	۱۳۹۳	استفاده از تکنیک‌های داده کاوی در دسته‌بندی خودکار اسناد متنی	روش پیشنهادی از روش فیلتری بهره می‌برد و با استفاده از دسته‌بندی بیز ساده و درخت تصمیم پیاده‌سازی شده است. نتایج نشان‌دهنده برتری روش ترکیبی نسبت به دسته‌بندی منفرد می‌باشد.
سیده زهرا سجادی و همکاران	۱۳۹۷	دسته‌بندی روش‌های توصیف متن در عقیده کاوی از شبکه‌ی اجتماعی توییتر	برای استفاده از روش‌های یادگیری ماشین نیازمند استخراج ویژگی از متون می‌باشیم که در این مقاله به دسته‌بندی انواع روش‌های مختلف برای توصیف متن در عقیده کاوی از شبکه‌ی اجتماعی توییتر پرداخته شده است.
سعید جمالی، اسکویی و همکار	۱۳۹۷	نظر کاو، کاوش نظرات و تحلیل احساسات در زبان فارسی	با استفاده از برجسب گذاری POS برای جملات زبان فارسی و سپس دادن امتیاز به هر ویژگی و تعیین مثبت، منفی و یا خنثی بودن ویژگی و در نهایت تعیین جهت‌گیری معنایی کل نظر با استفاده از امتیاز تمام ویژگی‌ها ارائه شده است.
سید محمد جواد مقدم و همکار	۱۳۹۷	وب کاوی و متن کاوی با رویکرد ماشین یادگیری	در ابتدا با مفاهیم اساسی و کاربردی وب کاوی و متن کاوی و آنالیز احساس پرداخته شده و در ادامه این مفاهیم را با رویکرد و فناوری‌های یادگیری ماشین بررسی می‌نماید.

پژوهشگر	سال	عنوان پژوهش	یافته‌های پژوهش
ملیکا تجاری و همکار	۱۳۹۷	ارائه یک روش جدید مبتنی بر ماشین بردار پشتیبان در جهت طبقه‌بندی تراکنش‌های موبایل	در این مقاله تمرکز بر روی استفاده از SVM در طبقه‌بندی تراکنش‌های موبایل است. برای این منظور عمل انتخاب ویژگی با استفاده از الگوریتم PCA انجام می‌شود. برای مقایسه روش پیشنهادی، از شبکه عصبی پرسپترون چندلایه استفاده شده است. نتایج نشان می‌دهد ماشین بردار پشتیبان با میانگین مربع خطای ۰.۲۳۶ و دقت ۹۴.۱ درصد به ازای کل داده‌ها طبقه‌بندی تراکنش‌های موبایل را انجام می‌دهد.
شفاق سوادگر طاقتی و همکار	۱۳۹۷	تشخیص ویژگی‌های مؤثر برای شناسایی بیماری قلبی با استفاده از فیلترهای انتخاب ویژگی و طبقه‌بندی آن‌ها به کمک تکنیک‌های داده‌کاوی	در این مقاله از دسته‌بند Naive Bayes استفاده شده است. هم‌چنین از دو فیلتر انتخاب ویژگی Relief و Wrapper استفاده شده است.
اله‌ام صدری و همکار	۱۳۹۷	ارائه الگوریتمی برای بازشناسی ارقام دست‌نویس فارسی با استفاده از ترکیب طبقه‌بندها	نتایج آزمایش‌ها بر روی پایگاه داده‌ای شامل ۶۰۰۰۰ تصویر اعداد نشان می‌دهد که روش همجوشی بر اساس این الگو بهتر از سایر طرح‌ها می‌باشد.
حدیث پورعاسی و همکار	۱۳۹۶	دسته‌بندی متون لاتین با روش‌های انتخاب ویژگی فیلتری، پوششی و الگوریتم‌های بیز ساده	در مرحله یادگیری از الگوریتم‌های خانواده بیز ساده استفاده شده است. روش پیشنهادی در مقایسه با کارهای گذشته در این زمینه بهبود بسیار خوبی داشته است.

تجزیه و تحلیل داده‌ها

به منظور فهم بیشتر روش پیشنهادی در ابتدا نمای کلی از فرآیند پیشنهادی در نرم افزار ریپد ماینر را مشاهده می‌نمایید که ابتدا مجموعه داده اکسس خوانده، سپس برای ستون‌های آن نقش تعریف می‌گردد و محتوای خوانده شده به صورت متن تبدیل می‌شود تا قابل توکنیزه شدن توسط ریپد ماینر باشد. در ادامه مراحل انجام کار شرح داده می‌شود.

۱۷۸



شکل (۲): نمایی کلی از فرآیند پیشنهادی در نرم‌افزار ریپد ماینر

جمع‌آوری اخبار

پایگاه خبری مورد نظر، آرشیو اخبار بی‌بی‌سی می‌باشد که این مجموعه داده به وسیله یک خزنده که به زبان پایتون نوشته شده است جمع‌آوری شده است. این خزنده با یک حلقه `while` مدام در حال جمع‌آوری اخبار از سایت BBC بوده و اخبار مربوطه را در محل از قبل تعیین شده ذخیره می‌نماید. اخبار جمع‌آوری شده در سه دسته خبری حوزه ایران، حوزه افغانستان و سطح جهانی توسط برچسب مشخص شده است و آموزش و آزمایش روی آن صورت خواهد گرفت. شکل زیر نمایی از مجموعه داده اخبار بی‌بی‌سی را ارائه می‌دهد.

B body	A category	№
یک نظرسنجی تازه جهانی نشان دهنده عمق اضطرابی است که بسیاری جوانان درباره تغییر اقلیم احساس می کنند. تقریباً ۶۰ درصد جوانانی که در این بررسی مورد پرسش قرار گرفتند گفتند که خیلی یا شدیداً احساس نگرانی می کنند.	اخبار سطح جهانی	1
جمع شدن حدود ۱۰ هزار پناهجو در اطراف پل مرزی میان مکزیک و آمریکا به شدت گرفتن بحران انسانی در این منطقه منجر شده است. این پل شهر دل ریو، تکراس را به شهر آکوتیا در مکزیک متصل می کند و اردوگاه موقت که در کنار آن شکل گرفته، در روزهای اخیر به شدت گسترش یافته است.	اخبار سطح جهانی	3
اعتراض‌ها در بهبهان، انتقاد از سیاست منطقه‌ای ایران در شهری کوچک تازه ترین اعتراض‌ها در ایران در شهری کوچک در استان خوزستان نمود پیدا کرده، در بهبهان اتفاقا یکی از میدانهای مرکزی بهبهان در اعتراض‌های خونین آبان سال گذشته هم کانونی ملتهب بود.	اخبار حوزه ایران	4
آیت الله علی خامنه ای با صدور حکمی سرتیب خلیان حمید واحدی را به عنوان فرمانده جدید نیروی هوایی ارتش منصوب کرده است. او در این حکم خواستار "ارتقای توان و آمادگی‌های رزمی نیروی هوایی شده است."	اخبار حوزه ایران	5
بیم قرن پس از آن که آمریکا به عنوان نخستین کشور جهان، پرچم خود را روی کره ماه نصب کرد، پرچم چین به عنوان دومین کشور جهان روی این تنها قمر کره زمین به اهتزاز در آمد. تصاویری که اداره ملی فضانوردی چین منتشر کرده پرچم سرخ و ۵ ستاره این کشور را نشان می دهد.	اخبار سطح جهانی	6
دادگامی در لندن می گوید شیخ محمد آل مکتوم، حاکم دبی، در جریان دعوای حقوقی با هیا بنت حسین، همسر سابق خود، برای گرفتن حق حضانت فرزندانشان دستور حک کردن تلفن او را داده است. به گفته دادگاه برای این کار از نرم افزار اسرائیلی "پکاسوس" استفاده شده است.	اخبار سطح جهانی	7
در چارچوب طرح کوواکس سازمان بهداشت جهانی، شب گذشته یک پرواز ویژه "سامان ایر" اولین محموله واکسن های "استرازنکا" را از هند وارد دوشنبه، پایتخت تاجیکستان، کرده به گفته مسئولان وزارت بهداشت و تأمین اجتماعی تاجیکستان، در مرحله نخست ۱۹۲ هزار دوز از این واکسن وارد شده است.	اخبار سطح جهانی	8
صهران خان، نخست وزیر پاکستان، در یک مصاحبه اختصاصی با بی بی سی گفته است که پاکستان در تعامل با همسایگان افغانستان و بطور "جمعی" در مورد به رسمیت شناختن یا نشناختن دولت طالبان تصمیم خواهد گرفت.	اخبار حوزه افغانستان	9
سید خطیب زاده سخنگوی وزارت خارجه ایران در نشست هفتگی با رسانه‌ها گفته است که مذاکرات ایران و هریستان "به مسیر جدی تری وارد شده" و "بقشه راه و مسیری در حال گفت و گو است." به گفته وی "توافقات حیات اهرامی ایران با هریستان این بوده که در فضای عمومی این مذاکرات منتشر نشود."	اخبار حوزه ایران	10
وزارت خزانهداری آمریکا اعلام کرد که سه نفر به نام‌های دانیل بهزاد فردوس، ماتوئل مهرزاد فردوس و محمد رضا دوفیلیان را از لیست افراد تحریم شده توسط این کشور خارج شده اند، و در عین حال تاکید کرده است که تغییر در لیست تحریم‌ها ارتباطی با مذاکرات وین ندارد.	اخبار حوزه ایران	11
کمیسیون دسترسی به اطلاعات افغانستان از ریاست جمهوری این کشور خواسته است که اسدالله خالد، سرپرست وزارت دفاع و احمد جاوید رسولی، رئیس عمومی اداره ملی احصائیه و معلومات (مرکز آمار) این کشور را به دلیل عدم همکاری با رسانه‌ها برای یک ماه کسر حقوق کند.	اخبار حوزه افغانستان	12
ده ماه بعد از اعدام خضر قوبیل در زندان ارومیه، دیوان عالی کشور حکم اعدام او را لغو کرده است. یکی از اعضای خانواده قوبیل به بی بی سی فارسی گفت از شنیدن خبر لغو حکم او چند ماه پس از اعدامش جا خورده اند و غمشان بسیار افزون تر شده و با شکایت از قاضی شیخولو اعتراض خود را شروع کرده اند.	اخبار حوزه ایران	13

شکل (۳) نمایی از مجموعه داده اخبار بی بی سی

آماده سازی و پیش پردازش

در فاز آماده سازی متون، فرمت فایل از جیسون^۱ به اکسس^۲ تبدیل می شود، سپس متن سند که شامل کاراکترهای پشت سر هم است به نمایشی که برای الگوریتم های یادگیری و طبقه بندی مناسب باشد تبدیل می شود. این کار توسط پلاگین های^۳ رپید ماینر^۴ و همچنین قطعه کدی که به زبان پایتون نوشته شده است انجام می گیرد. در این کد از کتابخانه هضم استفاده شده است. این فرآیند شامل موارد ذیل است:

حذف تگ های اچ تی ام ال^۵ و یا ایکس ام ال^۶

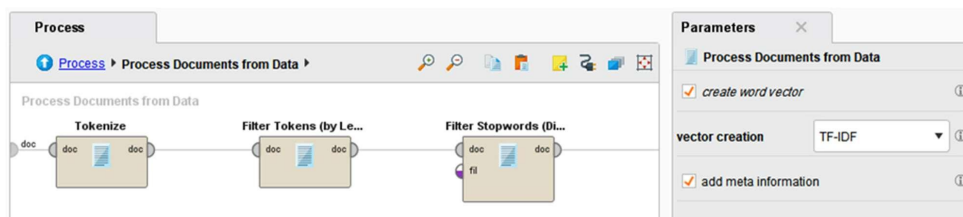
کد گذاری متون به یونی کد^۷

به دست آوردن ریشه کلمات و حذف پیشوندها و پسوندها

- 1.json
- 2.Access(*.accdb)
- 3.Plugin
- 4.Rapid Miner
- 5.html
- 6.xml
- 7 UTF-8.

توکنیزه کردن اطلاعات

به منظور توکنیزه کردن اطلاعات و تولید بردارهای قابل پردازش، از روش وزن‌دهی کلمات معکوس فراوانی متن استفاده شده است. این کار به کمک ابزار پردازش متن ریپدمااینر به نام توکنایز انجام می‌شود. تصویر زیر خروجی مورد نظر را نمایش می‌دهد. دقت کنید که تمامی کنترل‌های موجود در تصویر در کامپوننت پردازش سند روی داده می‌باشند. همان‌طور که در تصویر قابل مشاهده می‌باشد از الگوریتم پردازش متن معکوس فراوانی متن استفاده شده است.



شکل (۴): کاربرد الگوریتم پردازش متن TF-IDF

در جدول زیر فهرستی از فراوانی هر لغت در دسته مربوطه پس از عملیات توکنیزه کردن نمایش داده شده است.

جدول (۲) نمونه‌ای از فراوانی هر کلمه در سه دسته خبری

کلمه توکن شده	اخبار حوزه افغانستان	اخبار حوزه ایران	اخبار سطح جهانی
ترامپ	۷	۳۱	۳۷
تروریستی	۱۱	۱۷	۷
تروریست	۱۳	۹	۹
ترک	۱۹	۱۴	۳۲

کاهش ابعاد

به علت سنگینی پردازش‌های موجود بر روی متون و حذف حشویات و مواردی که باعث عدم خوانایی برنامه و کاهش دقت کار می‌شود، از روش کاهش ابعاد استفاده می‌شود. برای این منظور از ابزار فیلتر توکنز استفاده می‌شود. این فیلتر یک ورودی متن می‌گیرد و خروجی آن هم فیلتر شده همان متن می‌باشد. برای کاهش ابعاد کار، توکن‌های کمتر از ۴ و بیشتر از ۲۵ را حذف کردیم. این کار به افزایش صحت کار کمک بیشتری می‌کند.

همان‌طور که در جدول ۳ مشاهده می‌کنید تعداد بار تکرار هر یک از کلمات پردازش شده توسط

الگوریتم مربوطه در هر سه دسته ایران، افغانستان و سطح جهانی قابل مشاهده می‌باشد. همچنین برای کاهش اندازه ویژگی‌ها و بهینه‌تر شدن انتخاب ویژگی‌ها از عملگر قدرتمند وزن‌دهی با مربع کای استفاده شد. خروجی این عملگر یک مجموعه از اوزان کلمات می‌باشد که با توجه به فناوری‌های مربع کای^۱ تولید و مورد استفاده قرار می‌گیرد.

جدول (۳): نمونه‌ای از کلمات پس از وزن‌دهی با روش مربع کای

دسته	آبان	آتنن	آثار	آخرین	آذربایجان	آزاد
اخبار سطح جهانی	۰	۰,۰۰۷	۰	۰	۰	۰
اخبار سطح جهانی	۰	۰	۰	۰	۰	۰
اخبار سطح جهانی	۰	۰	۰	۰	۰	۰
اخبار حوزه ایران	۰,۰۳۹	۰,۰۳۹	۰	۰	۰	۰
اخبار حوزه ایران	۰	۰	۰	۰	۰	۰
اخبار سطح جهانی	۰	۰	۰	۰,۰۰۸	۰	۰

ایجاد مدل طبقه‌بندی به کمک الگوریتم‌های یادگیری ماشین

بیشتر سامانه‌های طبقه‌بندی خودکار متون، برای متون زبان انگلیسی طراحی شده‌اند و معمولاً برای متون فارسی به دلیل ماهیت زبان فارسی و در دسترس نبودن مجموعه ریشه کلمات قابل استفاده نیستند (بینا و همکاران، ۱۳۸۶). اغلب روش‌های استفاده شده در این زمینه، روش‌های مبتنی بر یادگیری هستند. از مهم‌ترین روش‌های استفاده شده در طبقه‌بندی متون که در این پژوهش نیز استفاده شده است می‌توان به بیزین ساده، ماشین بردار پشتیبان و کی نزدیک‌ترین همسایه اشاره کرد.

ارزیابی

معمولاً برای ارزیابی صحت دسته‌بندی کننده اسناد از معیار صحت^۲، معیار فراخوانی^۳، معیار دقت^۴ و معیار ارزیابی^۵ F استفاده می‌شود. دقت مربوط به درستی نسبت دادن داده‌های آزمون به کلاس

1 Weight By Chi Squared.

2 Precision.

3 Recall.

4. Accuracy

5 F-measure (F-score)

مربوطه است و معیار ارزیابی F ترکیب هماهنگ از دقت و فراخوانی است. معادلات مربوطه در جدول زیر که موسوم به ماتریس درهم‌ریختگی^۱ است آمده است. برای درک دو معیار ارزیابی دقت و فراخوانی و همچنین معیار صحت، ابتدا باید با عبارات مثبت واقعی^۲، مثبت کاذب^۳، منفی واقعی^۴ و منفی کاذب^۵ آشنا بود. این چهار عبارت برای مقایسه‌ی برچسب کلاس‌ها در ماتریس درهم‌ریختگی استفاده می‌شوند. جدول ۴ این عناوین را نشان می‌دهد.

جدول (۴) ماتریس درهم‌ریختگی

برچسب واقعی		برچسب	
منفی	مثبت	مثبت	برچسب پیش‌بینی شده
مثبت کاذب (FP)	مثبت واقعی (TP)	مثبت	
منفی واقعی (TN)	منفی کاذب (FN)	منفی	

صحت: نسبت مقدار موارد صحیح طبقه‌بندی شده توسط الگوریتم از یک کلاس مشخص، به کل تعداد مواردی که الگوریتم چه به صورت صحیح و چه به صورت غلط، در آن کلاس طبقه‌بندی کرده است که به صورت زیر محاسبه می‌شود:

$$\text{Precision} = \frac{TP}{TP + FP} \quad \text{فرمول ۱:}$$

فراخوانی: نسبت مقداری موارد صحیح طبقه‌بندی شده توسط الگوریتم از یک کلاس به تعداد موارد حاضر در کلاس مذکور که به صورت زیر محاسبه می‌شود:

$$\text{Recall} = \frac{TP}{TP + FN} \quad \text{فرمول ۲:}$$

دقت: در واقع این معیار عمومی‌ترین معیار محاسبه کارایی الگوریتم‌های طبقه‌بندی است که نشان می‌دهد، طبقه‌بند طراحی شده چند درصد از کل مجموعه اسناد را به درستی دسته‌بندی کرده است.

- 1 Confusion Matrix
2. True Positive
3. False Positive
4. True Negative
5. False Negative

فرمول ۳:

$$Accuracy = \frac{TP + TN}{TP + TN + FP + FN}$$

معیار ارزیابی F: با توجه به محاسبات انجام گرفته برای معیارهای دقت و فراخوانی، در این مرحله می‌توان مقدار کمیت وزن‌دار معیار ارزیابی F را محاسبه نمود. معیار ارزیابی F، پارامتر مناسبی برای ارزیابی کیفیت کلاس‌بندی می‌باشد و همچنین توصیف‌کننده میانگین وزن‌دار مابین دو کمیت دقت و فراخوانی می‌باشد. برای یک الگوریتم کلاس‌بندی کننده در شرایط ایده آل، مقدار این کمیت برابر با ۱ می‌باشد و در بدترین وضعیت برابر با صفر می‌باشد. این پارامتر با توجه به رابطه زیر محاسبه می‌شود:

$$F - Measure = 2 \times \frac{precision \times recall}{precision + recall}$$

فرمول ۴:

ابزار مورد استفاده

در این تحقیق از مجموعه داده اخبار بی‌بی‌سی فارسی به‌عنوان مجموعه آموزش و مجموعه آزمون استفاده شده است. این پیکره شامل ۱۷۰۰ خبر بین سال‌های ۱۴۰۰ تا ۱۴۰۱ می‌باشد. برای افزایش دقت، فراخوانی و همچنین کارایی کلی الگوریتم‌های طبقه‌بندی از روش‌های متعدد پیش‌پردازش از قبیل حذف کلمات اضافه و اصلاح فاصله به نیم‌فاصله در ساخت مجموعه آموزش استفاده کرده و از سه روش طبقه‌بندی ماشین بردار پشتیبان، بیز ساده و کی-نزدیکترین همسایه برای طبقه‌بندی اسناد بهره برده‌ایم. برای انجام مراحل پیش‌پردازش، وزن‌دهی، انتخاب ویژگی و طبقه‌بندی از نرم‌افزار رپیدماینر نسخه ۹،۶ استفاده کرده‌ایم.

استفاده از ابزار مدل اتوماتیک

روش استفاده از مدل اتوماتیک به این صورت است که ابتدا داده مورد نظر را به عنوان ورودی معرفی می‌کنیم، سپس از پلاگین پیش‌پردازش برای تمیز کردن داده‌ها استفاده می‌شود و معیارهای ارزیابی هر یک از الگوریتم‌های یادگیری ماشین را به تفکیک ارائه کرده و بهترین الگوریتم را به لحاظ دقت، صحت، زمان اجرا و ... پیشنهاد می‌دهد. در ادامه نتایج هر یک از الگوریتم‌ها آورده شده است.

بررسی معیار دقت الگوریتم بیز ساده

در این بخش تاثیر روش مربع کای و نرمال سازی کلمات بر روی عملکرد این الگوریتم را به صورت جداگانه بررسی می کنیم. فرآیند انتخاب ویژگی در بهبود عملکرد الگوریتم طبقه بند به خصوص در الگوریتم طبقه بند بیز ساده که به فرآیند انتخاب ویژگی بسیار وابسته است، تاثیر زیادی دارد بنابراین پیش بینی می شود که عملکرد الگوریتم طبقه بند بیز ساده با اعمال روش مربع کای بهبود یابد. جدول ۵ میانگین معیار دقت را برای الگوریتم طبقه بند بیز ساده در سه حالت بدون اعمال نرمال سازی کلمات، اعمال نرمال سازی کلمات و با اعمال فرآیند انتخاب ویژگی مربع کای نشان می دهد.

جدول (۵): بررسی تاثیر نرمال سازی کلمات و روش مربع کای بر الگوریتم بیز ساده

معیار ارزیابی	بیز ساده و عدم نرمال سازی	بیز ساده و اعمال نرمال سازی کلمات	بیز ساده و مربع کای
میانگین دقت	۵۲,۴٪	۵۳,۷٪	۰/۵۶,۲

در شکل زیر نیز سایر معیارهای ارزیابی برای ترکیب الگوریتم بیز ساده و مربع کای آمده است.

Performances

Criterion	Value	Standard Deviation
Accuracy	56.2%	± 3.0%
Classification Error	43.8%	± 3.0%

Confusion Matrix

	true class precision	خيار حوزه افغانستان true	خيار حوزه ايران true	خيار سطح جهانی true
pred. class precision	72.34%	12	14	68
pred. class recall	52.05%	44	203	143
pred. class precision	100.00%	2	0	0
class recall		3.45%	93.55%	32.23%

شکل (۵) سایر معیارها برای الگوریتم بیز ساده

بررسی معیار دقت الگوریتم کی-نزدیکترین همسایه

در این بخش تاثیر نرمال سازی کلمات و روش مربع کای بر روی عملکرد الگوریتم کی-نزدیکترین همسایه را به صورت جداگانه بررسی می کنیم. جدول ۶ میانگین معیار دقت را برای الگوریتم کی-نزدیکترین همسایه در سه حالت بدون اعمال نرمال سازی کلمات، اعمال نرمال سازی کلمات و با اعمال فرآیند انتخاب ویژگی مربع کای نشان می دهد.

جدول (۶): بررسی تاثیر نرمال سازی کلمات و روش مربع کای بر الگوریتم کی-نزدیکترین همسایه

معیار ارزیابی	کی-نزدیکترین همسایه و عدم نرمال سازی کلمات	کی-نزدیکترین همسایه و اعمال نرمال سازی کلمات	کی-نزدیکترین همسایه و مربع کای
میانگین دقت	۸۶,۱٪	۸۷,۳٪	۸۹,۱٪

در شکل زیر نیز سایر معیارهای ارزیابی برای ترکیب الگوریتم کی-نزدیکترین همسایه و مربع کای آمده است.

Criterion	Value	Standard Deviation
Accuracy	89.1%	± 2.7%
Classification Error	10.9%	± 2.7%

Confusion Matrix

	true	اختیار حوزه ایران true	اختیار حوزه افغانستان true	class precision
اختیار سطح جهانی pred.	182	24	9	84.65%
اختیار حوزه ایران pred.	15	188	1	92.16%
اختیار حوزه افغانستان pred.	2	2	64	94.12%
class recall	91.46%	87.85%	86.49%	

بررسی معیار دقت الگوریتم ماشین بردار پشتیبان

در این بخش تاثیر نرمال سازی کلمات و روش مربع کای بر روی عملکرد الگوریتم ماشین بردار پشتیبان را به صورت جداگانه بررسی می کنیم. جدول ۷ میانگین معیار دقت را برای الگوریتم ماشین بردار پشتیبان در سه حالت بدون اعمال نرمال سازی کلمات، اعمال نرمال سازی کلمات و با اعمال فرآیند انتخاب ویژگی مربع کای نشان می دهد.

جدول (۷): بررسی تاثیر نرمال سازی کلمات و روش مربع کای بر الگوریتم ماشین بردار پشتیبان

معیار ارزیابی	ماشین بردار پشتیبان و عدم نرمال سازی کلمات	ماشین بردار پشتیبان و اعمال نرمال سازی کلمات	ماشین بردار پشتیبان و اعمال نرمال سازی کلمات و مربع کای
میانگین دقت	۸۷,۵٪	۸۸,۸٪	۹۰,۸٪

در شکل زیر نیز سایر معیارهای ارزیابی برای ترکیب الگوریتم ماشین بردار پشتیبان و مربع کای آمده است.

Performances

Criterion	Value	Standard Deviation
Accuracy	90.8%	± 2.6%
Classification Error	9.2%	± 2.6%

Confusion Matrix

	اخبار سطح جهانی true	اخبار حوزه ایران true	اخبار حوزه افغانستان true	class precision
اخبار سطح جهانی pred.	183	18	8	87.56%
اخبار حوزه ایران pred.	14	195	2	92.42%
اخبار حوزه افغانستان pred.	2	1	64	95.52%
class recall	91.96%	91.12%	86.49%	

شکل (۷): سایر معیارها برای الگوریتم ماشین بردار پشتیبان

محدودیت های تحقیق

جمع آوری اخبار از سایت های مختلف یا سایت واحد نیاز به یک خزنده با قابلیت جمع آوری متون فارسی دارد که خزنده های آنلاین نیز با توجه به تحریم روی سایت های فارسی، بدون فیلتر اعمال نمی شوند و استفاده از پراکسی باعث کند شدن روند فعالیت خزشگر می شود. همچنین یک مجموعه واحد از ریشه لغات زبان فارسی وجود ندارد به همین دلیل عملیات ریشه یابی تاثیر زیادی در نتیجه عملکرد ما ندارد.

نتیجه گیری و پیشنهادات

متن کاوی به عنوان یک زمینه در حال رشد و پر کاربرد، به دنبال کشف دانش از متون غیر ساخت یافته است. به دلیل مشکلات ساختاری زبان فارسی، تحقیقات کمتری در این زمینه صورت گرفته است. با توجه به اهمیت پردازش داده های متنی فارسی در کشور و در مجموعه های اطلاعاتی، در این پژوهش به مسئله طبقه بندی اخبار که قابلیت بسط به حوزه های مختلف متن کاوی را دارد پرداخته شده است. در

گذشته بیشتر کاری که برای طبقه‌بندی متون انجام گرفته است بر روی متون زبان انگلیسی و چینی بوده است. در این پژوهش روشی برای طبقه‌بندی متون فارسی ارائه شده است. روش کلی بر مبنای یادگیری ماشینی استوار است که در دو فاز آموزش و آزمون ارائه شده است. با توجه به تجزیه و تحلیل پیچیده متون فارسی از اکستنشن‌های^۱ متن موجود در ریدمایندر استفاده شده که نتایج خوبی نیز در برداشته است. پس از ساخته شدن بردار ویژگی توسط ابزار قدرتمند پردازش متن از داده‌ها^۲ به کمک ابزار مدل اتوماتیک الگوریتم‌های مختلف بررسی شد و مشخص گردید بهترین نتیجه مربوط به الگوریتم ماشین بردار پشتیبان می‌باشد سپس با استفاده از ابزار اعتبارسنجی متقابل طبقه‌بندی اسناد صورت گرفت و نتایج این ابزار نیز نشان می‌دهد مدل ماشین بردار پشتیبان بهره‌وری^۳ بهتری داشته است و استفاده از روش‌های بیزین ساده و نزدیک‌ترین همسایه نتوانسته نتیجه بهتری به دست آورد با اینکه ما در مدل پیشنهادی موضوع زمان اجرا^۴ را نیز به عنوان معیار مهم در اولویت قرار ندادیم. همچنین به عنوان کارهای آتی می‌توان به استفاده از روش‌های متفاوت وزن‌دهی مثل کلمه به بردار^۵، استفاده از روش‌های جدید انتخاب ویژگی مثل الگوریتم کرم شب‌تاب^۶، بهبود روش پردازش زبان طبیعی، بهبود و تقویت الگوریتم‌های ریشه‌یابی و توکنیزه کردن و بررسی روش‌های تقویت بردار ویژگی مناسب برای زبان فارسی، مورد بررسی قرار داد.

-
- 1.Extension
 - 2.process document from data
 - 3.performance
 - 4.Run time
 - 5.Word2vec
 - 6.Firefly algorithm

۱. باقری، ایوب؛ فرزانه فر، حامد؛ سرایی، محمدحسین و احمدزاده، محمدرضا (۱۳۸۷). دسته‌بندی متون خبری فارسی با استفاده از الگوریتم بیز ساده، دومین کنفرانس داده‌کاوی ایران، تهران، <https://civilica.com/doc/70524>
۲. برفامی، مهدی و فاطری، سهیل (۱۳۹۲). استفاده از ترکیب شبکه‌های عصبی جهت دسته‌بندی متون فارسی مبتنی بر الگوریتم‌های GA، کی-نزدیکترین همسایه، PCA جهت کاهش ویژگی، اولین همایش ملی رویکردهای نوین در مهندسی کامپیوتر و بازیابی اطلاعات، رشت، [225887/https://civilica.com/doc](https://civilica.com/doc/225887)
۳. بصیری، محمد احسان؛ نعمتی، شهلا و قاسم آقایی، ناصر (۱۳۸۶). مقایسه دسته‌بندی متون فارسی با استفاده از الگوریتم‌های کی-نزدیکترین همسایه و fkNN و انتخاب ویژگی‌ها بر اساس بهره اطلاعات و فرکانس سند، سیزدهمین کنفرانس سالانه انجمن کامپیوتر ایران، جزیره کیش، [41786/https://civilica.com/doc](https://civilica.com/doc/41786)
۴. بینا، بهاره؛ رهگذر، مسعود و ده موبد، آذین (۱۳۸۶). طبقه‌بندی خودکار متون فارسی سیزدهمین کنفرانس سالانه انجمن کامپیوتر ایران، جزیره کیش، انجمن کامپیوتر، دانشگاه صنعتی شریف.
۵. حسن‌پور، حمید؛ قنبری سرخی، علی و پارسی، اشکان (۱۳۹۱). استخراج بهترین ویژگی از متون فارسی با استفاده از تجزیه و تحلیل مؤلفه‌های اصلی با کمک میانگین یادآوری و الگوریتم ژنتیک، نخستین کنفرانس بین‌المللی پردازش خط و زبان فارسی.
۶. حسن‌پور، حمید و مدنی، صبا سادات (۱۳۹۳). بهبود دقت سیستم دسته‌بندی خودکار اسناد فارسی به کمک هستان‌شناسی فارسی، مجله علمی پژوهشی، رایانش نرم و فناوری اطلاعات، جلد ۳، شماره ۱
۷. زمانی، محسن؛ دیانت، روح‌الله و صادق زاده، مهدی (۱۳۹۲). دسته‌بندی متون فارسی با استفاده از روش آنالیز معنایی پنهان احتمالاتی، همایش ملی کاربرد سیستم‌های هوشمند (محاسبات نرم) در علوم و صنایع، قوچان، [206251/https://civilica.com/doc](https://civilica.com/doc/206251)
۸. طاهری نیا، محسن (۱۳۹۱). دسته‌بندی متون فارسی با استفاده از یادگیری نیمه نظارت‌شده، چهارمین کنفرانس مهندسی برق و الکترونیک ایران، گناباد، [164226/https://civilica.com/doc](https://civilica.com/doc/164226)
۹. عابدینی نیا، مائده؛ الله دادی، لاله و شیخی، فاطمه (۱۳۹۲). کاوش متون فارسی در وب با استفاده از تحلیل معنایی و روش طبقه‌بندی چندلایه‌ای اطلاعات، اولین همایش ملی رویکردهای نوین در مهندسی کامپیوتر و بازیابی اطلاعات، رشت، [225377/https://civilica.com/doc](https://civilica.com/doc/225377)
۱۰. عربی نرئی، سمیه؛ وحیدی اصل، مجتبی و مینایی بیدگلی، بهروز (۱۳۸۶). استخراج کلمات کلیدی جهت

- طبقه-بندی متون فارسی، اولین کنفرانس داده‌کاوی ایران، تهران، <https://civilica.com/doc/33094>
۱۱. قنبری سرخی، علی و ابراهیمی، فاطمه (۱۳۹۰). بهبود عملکرد طبقه‌بندی متون فارسی با استفاده از تجزیه و تحلیل مؤلفه‌های اصلی با کمک معیار میانگین یادآوری و دقت، چهاردهمین کنفرانس دانشجویی مهندسی برق کشور، کرمانشاه، <https://civilica.com/doc/121558>
۱۲. مقصودی، نوشین و همایون‌پور، محمدمهدی (۱۳۸۸). ارائه روشی جدید در طبقه‌بندی متون فارسی با استفاده از دانش معنایی "، پانزدهمین کنفرانس بین‌المللی سالانه انجمن کامپیوتر ایران.
۱۳. آقا کاردان، احمد و کیهانی نژاد، مینا (۱۳۹۱). ارائه مدلی برای استخراج اطلاعات از مستندات متنی، مبتنی بر متن‌کاوی در حوزه یادگیری الکترونیکی. فصلنامه علمی-پژوهشی فناوری اطلاعات و ارتباطات ایران، سال چهارم، شماره‌های ۱۱ و ۱۲، ص ۴۷ تا ۵۴
۱۴. بازقندی، مهدی؛ تدین تبریزی، قمرناز و وفایی جهان، مجید (۱۳۹۱). نخستین کنفرانس بین‌المللی پردازش خط و زبان طبیعی، دانشگاه سمنان.
۱۵. بهرام‌پور، اکبر؛ بهشتی، همایون و لاکتراشی، طیبه (۱۳۹۴). بررسی روش‌ها و الگوریتم‌های دسته‌بندی اخبار با استفاده از پردازش زبان طبیعی (NLP)، دومین کنفرانس ملی توسعه علوم مهندسی، تنکابن، موسسه آموزش عالی آیندگان.
۱۶. پرئی، اعظم‌السادات و حمیدی، حجت اله (۱۳۹۵). ارائه رویکردی برای مدیریت و سازماندهی اسناد متنی با استفاده از تجزیه تحلیل هوشمند متن. فصل‌نامه علمی- پژوهشی پژوهشگاه علوم و فناوری اطلاعات ایران، دوره ۳۲، شماره ۴، ص ۱۱۷۱ تا ۱۲۰۲
۱۷. جمالی، ایمان؛ میرعابدینی، سید جواد و هارون‌آبادی، علی (۱۳۹۶). ارائه‌ی یک مدل جهت دسته‌بندی متون فارسی با استفاده از ترکیب روش‌های دسته‌بندی، مجله مهندسی مخابرات، سال هفتم، شماره ۲۳
۱۸. سیاحی، عارف؛ هاشمی، سید محسن و مزرعه، سعید (۱۳۹۴). بهبود یک روش مبتنی بر انتخاب ویژگی به منظور دسته‌بندی متون با الگوریتم‌های متن‌کاوی، دومین کنگره سراسری فناوری‌های نوین ایران با هدف دستیابی به توسعه پایدار، تهران، مرکز راهکارهای دستیابی به توسعه پایدار، موسسه آموزش عالی مهر اروند.
۱۹. شیخی، مریم؛ اکبرپور، شاهین و فرزانه، علی (۱۳۹۱). متن‌کاوی متون فارسی در راستای طبقه‌بندی آن. چهارمین کنفرانس مهندسی برق و الکترونیکی ایران.
۲۰. کریمی منش، مصطفی و شیرازی، حسین (۱۳۹۲). مقایسه روش‌های وزن‌دهی ویژگی در فرایند

طبقه‌بندی مستندات. اولین کنفرانس ملی رویکردهای نوین در مهندسی کامپیوتر و بازیابی اطلاعات ایران ۲۱. هاشمی، سید محسن (۱۳۹۴). بهبود دسته‌بندی متون فارسی با ترکی روش دو مرحله‌ای انتخاب ویژگی و الگوریتم‌های یادگیری ماشین، کنفرانس بین‌المللی یافته‌های نوین پژوهشی در مهندسی برق و علوم کامپیوتر، تهران، موسسه آموزش عالی نیکان.

22. BolshaKov, I. A. & GelbuKh, A. (2004). Computational linguistics: models, resources, applications: Instituto Politecnico Nacional.
23. Eyheramendy, S. GenKin, A. Ju, W.H. Lewis, D. D. & Madigan, D. (2003). Sparse bayesian classifiers for text categorization. Journal of Intelligence Community Research and Development, 13.
24. Francis, L. A. (2006). Taming Text: An Introduction to Text Mining. Paper presented at the Casualty Actuarial Society Forum
25. McCallum, A. & Nigam, K. (1998). A comparison of event models for naive bayes text classification. Paper presented at the AAAI-98 workshop on learning for text categorization.
26. Moulinier, I. & Ganascia, J. G. (1995). Applying an existing machine learning algorithm to text categorization. Paper presented at the International Joint Conference on Artificial Intelligence.
27. P. MulaK, & N. Talhar (2015). Analysis of distance measures using K-nearest neighbor algorithm on KDD dataset. International Journal of Science and Research, 4(7), 2101-2104.
28. NadKarni, P. M. Ohno-Machado, L. & Chapman, W. W. (2011). Natural language processing: an introduction. Journal of the American Medical Informatics Association, 18(5), 544-551.