

## Organizational Data Security with Persian Content Document Management Model for News Media

Hamid Mirzaei Dehnavi<sup>1</sup> , Mohammad Ali Keramati<sup>2\*</sup> , Mohammad Ali Ameri Habibabadi<sup>3</sup>

<sup>1</sup> PhD student, Department of Information Technology Management, Central Tehran Branch, Islamic Azad University, Tehran, Iran.. Email: Msb66@chmail.ir)

<sup>2</sup> Professor, Department of Information Technology Management, Central Tehran Branch, Islamic Azad University, Tehran, Iran . Email: (\*Correspondence: Mohammadalikeramati@yahoo.com)

<sup>3</sup> Professor, Department of Information Technology Management, Central Tehran Branch, Islamic Azad University, Tehran, Iran) Email: M.ali.ameri.h@gmail.com

### ARTICLE INFO

#### Article history:

Article Type: Research paper

Received: 18 April 2025

Revised: 27 May 2025

Accepted: 08 June 2025

Available online: 28 June 2025

#### Keywords:

Data Security

Text Mining

Classification

Document Management

Persian Content

Information and Classification

### ABSTRACT

With the increasing growth of Persian news texts in the digital world and the Internet, an important and significant issue is the intelligent classification and management of news documents and our quick and cheap access to them. Given that there are various algorithms and methods for text mining and they are still being developed; but the main problem of these methods for organizations is to maintain the confidentiality of information, and also to maintain the principle of information confidentiality along with the use of modern science and technologies. Many existing programs in the field of data mining and text mining work on the Internet, which is practically a violation of the principle of information confidentiality for organizations. Therefore, our main question refers to how to ensure data security using the document management model, which is one of the main concerns of every organization. Existing programs in the field of data mining and text mining work on the Internet, which is practically useless for some organizations that have classified data; Fortunately, domestic companies have been working in the field of text mining for several years and have made significant progress in this field; but it is still not possible to use such software for organizations and military complexes. Therefore, it is necessary for organizations to implement the desired model natively in order to meet their needs in this field and have a suitable platform independent of the Internet network for analyzing information. News document management provides organizations with the opportunity to correctly classify received or existing news and documents, preserve, maintain and retrieve them. In this research, which is of an applied type and uses the experimental research method and text mining tools, the statistical population of the present research is Persian news for the first 6 months of 1402. The BBC news data is a collection of Persian news texts in a labeled form. BBC News includes 1700 news texts from this dataset, which only contains Persian news, Afghan news and global news. Also, this research reviews and describes the literature on document management, natural language processing and data analysis, and then introduces document management systems as a part of the knowledge management system that is responsible for acquiring, organizing and sharing knowledge in the organization. By reviewing, analyzing and processing in this research, we conclude that the Persian content document management model using the support vector machine model has a precision of 93.29, a precision of 93.32, a recall of 92.96 and an error of 6.71.

**Cite this article:** Mirzaei Dehnavi, Hamid<sup>o</sup>, Keramati, Mohammad Ali<sup>o</sup>, Ameri Habibabadi, Mohammad Ali<sup>o</sup> (2025). Organizational Data Security with Persian Content Document Management Model for News Media. Journal of Electronic and Cyber Defens. 2025; 13(2):89-101.

**DOR:** <https://dor.isc.ac/dor/20.1001.1.23224347.1404.13.2.8.5>

© Author(s) retain the copyright and full publishing rights

**Publisher:** Imam Hossein University.



## امنیت داده‌های سازمانی با مدل مدیریت مستندات محتوای فارسی رسانه‌های خبری

حمید میرزائی دهنوی<sup>۱</sup>، محمد علی کرامتی<sup>۲\*</sup>، محمدعلی عامری حبیب آبادی<sup>۳</sup>

<sup>۱</sup> دانشجوی دکتری، گروه مدیریت فناوری اطلاعات، واحد تهران مرکزی، دانشگاه آزاد اسلامی، تهران، ایران (Msb66@chmail.ir)  
<sup>۲</sup> استاد، گروه مدیریت فناوری اطلاعات، واحد تهران مرکزی، دانشگاه آزاد اسلامی، تهران، ایران (نویسنده مسئول: Mohammadalikeramati@yahoo.com)

<sup>۳</sup> استاد، گروه مدیریت فناوری اطلاعات، واحد تهران مرکزی، دانشگاه آزاد اسلامی، تهران، ایران (M.ali.ameri.h@gmail.com)

### مشخصات مقاله

#### تاریخچه مقاله:

نوع مقاله: علمی پژوهشی  
دریافت: ۱۴۰۴/۰۱/۱۷  
بازنگری: ۱۴۰۴/۰۲/۲۱  
پذیرش: ۱۴۰۴/۰۳/۱۸  
ارائه آنلاین: ۱۴۰۴/۰۴/۲۰

#### کلید واژه‌ها:

امنیت داده  
متن کلوی  
طبقه‌بندی  
مدیریت مستندات  
محتوای فارسی  
اطلاعات و دسته‌بندی

### چکیده (استایل عنوان چکیده)

با رشد روزافزون متون خبری فارسی در دنیای دیجیتال و اینترنت، موضوعی که مهم و حائز اهمیت است دست‌بندی و مدیریت هوشمند مستندات خبری و دسترسی سریع و ارزان مابین آن‌ها است. با توجه به اینکه الگوریتم‌ها و روش‌های مختلفی برای متن کلوی وجود دارد و همچنین نیز در حال توسعه هستند؛ اما مشکل اصلی این روش‌ها برای سازمان‌ها حفظ محرمانگی اطلاعات است. همچنین در کنار استفاده از علوم و فناوری‌های نوین اصل محرمانگی اطلاعات نیز حفظ شود. بسیاری از برنامه‌های موجود در موضوع داده کلوی و متن کلوی در بستر اینترنت کار می‌کنند که عملاً نقض اصل محرمانگی اطلاعات برای سازمان‌ها است. بنابراین سؤال اصلی ما اشاره به چگونگی تضمین امنیت داده با استفاده از مدل مدیریت مستندات که از دغدغه‌های اصلی هر سازمان به شمار می‌آید است. برنامه‌های موجود در موضوع داده کلوی و متن کلوی در بستر اینترنت کار می‌کنند که عملاً برای برخی از سازمان‌ها که داده‌های طبقه‌بندی شده دارند بدون استفاده از روش‌های مناسب چند سالی است که شرکت‌های داخلی شروع به فعالیت در زمینه متن کلوی نموده‌اند و پیشرفت‌های قابل توجهی هم در این زمینه داشته‌اند؛ اما کماکان نمی‌توان برای سازمان‌ها و مجموعه‌های نظمی از این گونه نرم‌افزارها استفاده نمود؛ لذا لازم است سازمان‌ها به منظور رفع نیاز خود در این حوزه، مدل موردنظر را به صورت بومی پیاده‌سازی کنند و یک بستر مناسب و مستقل از شبکه اینترنت برای تحلیل اطلاعات داشته باشند. مدیریت مستندات اخبار این امکان را برای سازمان‌ها فراهم می‌آورد که اخبار و اسناد دریافتی یا موجود را به رستی طبقه‌بندی نموده، آن‌ها را حفظ، نگهداری و بازیابی نماید. این پژوهش از نوع کاربردی و برای انجام آن از روش پژوهش آزمایشی و ابزار متن کلوی استفاده گردیده است. جمعاً آماری پژوهش حاضر اخبار فارسی عمده اول سال ۱۴۰۲ است. داده اخبار بی‌بی‌سی مجموعه‌ای از متون خبری فارسی به صورت برچسب خورده هستند. اخبار بی‌بی‌سی شامل ۱۷۰۰ عدد متن خبری از این مجموعه داده که فقط حاوی اخبار حوزه فارسی، اخبار حوزه افغانستان و سطح جهانی است. همچنین در این پژوهش به بررسی و تشریح ادبیات مدیریت مستندات، پردازش زبان طبیعی و تجزیه و تحلیل داده‌ها پرداخته شده است و سپس به مطالعه سیستم‌های مدیریت مستندات به عنوان بخشی از سیستم مدیریت دانش که وظیفه اکتساب سازمان‌دهی و اشتراک دانش در سازمان را بر عهده‌دارند معرفی می‌شود. با بررسی، تحلیل و پردازش در این پژوهش به این نتیجه می‌رسیم که مدل مدیریت مستندات محتوای فارسی با استفاده از مدل ماشین بردار پشتیبان دارای دقت ۹۳/۳۹٪، صحت ۹۳/۳۲٪، فراخوانی ۹۲/۹۶٪ و خطای ۶/۷۱٪ است.

**استناد:** میرزائی دهنوی، حمید<sup>۱</sup>، کرامتی، محمد علی<sup>۲</sup>، عامری حبیب آبادی، محمدعلی<sup>۳</sup>. امنیت داده‌های سازمانی با مدل مدیریت مستندات محتوای فارسی رسانه‌های خبری. پدافند الکترونیک و سایبری. (۱۴۰۴)؛ ۱۳ (۲): ۸۹-۱۰۱.

DOR: <https://dor.isc.ac/dor/20.1001.1.23224347.1404.13.2.8.5>

© نویسنده (گان) حق نشر و حقوق کامل انتشار را برای خود محفوظ می‌دارند.



ناشر: دانشگاه جامع امام حسین (ع).

OPEN ACCESS

## ۱- مقدمه

بدیهی است اتخاذ تصمیم در مقاطع زمانی مشخص و ارزیابی عملکرد پیشرفت برنامه‌ها، نیاز به اطلاعات دارد و تصمیمات باید مبتنی بر اطلاعات قابل اعتماد باشد. با در اختیار داشتن اطلاعات کامل، روشن، دقیق و به‌موقع، مدیران این امکان را خواهند داشت که تصمیماتی مطمئن و قابل اجرا اتخاذ نمایند. در غیر این صورت اطلاعات ناقص، مبهم و نارسا و توأم با حدس و گمان مسیر منطقی برنامه‌ها را تغییر داده و در نهایت سازمان را از اهداف اصلی خود دور می‌کند. با افزایش حجم اطلاعات الکترونیکی، نیاز به سیستم جامع برای مدیریت مستندات فزونی یافته است [۱]. شاید بسیاری از مدیران سازمان‌ها تصور کنند با صرف هزینه‌های بسیار و خرید بسته‌های نرم‌افزاری آرشو اسناد می‌توانند سیستم مدیریت مستندات هوشمند را در سازمان خود پیاده‌سازی و از آن بهره‌برداری نمایند. درحالی‌که لازمه پیاده‌سازی این سیستم در هر سازمان تنها خرید یک بسته نرم‌افزاری نبوده و مستلزم انجام اقدامات گوناگونی است [۲]. یکی از زیرمجموعه‌های اسناد متنی که هرروزه 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. به‌عبارت‌دیگر، پیش از انتشار اخبار در سایت‌های مختلف خبری نیاز است که دسته موضوعی هر متن خبری مشخص شود. با دسته‌بندی متون خبری که بر روی سایت‌های خبری منتشر می‌شوند، کاربران این سایت‌ها می‌توانند اخبار موردنظر خود را آسان‌تر و سریع‌تر پیدا کنند. از طرفی به دلیل حجم بالای اخباری که هر روزه منتشر می‌شوند، طبقه‌بندی اخبار به‌صورت دستی فرایندی 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 به‌صورت دستی بسیار زیاد است. این به این معنی است که هر سایت خبری تلاش می‌کند اخبار روزانه را سریع‌تر از دیگر سایت‌ها در صفحه خود بارگذاری کند. به همین دلیل است که طبقه‌بندی متون خبری به‌صورت 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.

سازمان‌ها داده‌های متنی بسیاری را رصد می‌کنند چگونه می‌توان این حجم از داده‌ها را مدیریت و امنیت آن را تضمین نمود؟ چگونه می‌توان اطلاعات مفیدی از این داده‌های متنی بی‌شمار به دست آورد؟ آیا این داده‌ها ارزشی دارند؟ الگوی مناسب

پیاده‌سازی مدیریت مستندات الکترونیکی در سازمان‌ها چیست؟ و آیا برای ایجاد آن نیاز به طی مراحل و فرایندهای خاصی است؟ این مراحل چیست و یک سیستم ایدئال طراحی شده چه ویژگی‌هایی را باید داشته باشد؟

متن کاوی فناوری مورد استفاده برای چنین مواردی است. با فناوری‌های متن کاوی و مدیریت مستندات می‌توان داده‌های متنی را بررسی و تحلیل کرده و از نتایج حاصل از این تحلیل، اطلاعات ارزشمندی کسب نمود [۳]. 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 [۴]. الگوریتم‌ها و روش‌های مختلفی برای متن کاوی وجود دارد و همچنان نیز درحال توسعه است اما به دلیل پیچیدگی‌ها موجود در زبان فارسی، سعی در بهینه نمودن روش متن کاوی در این حوزه را داشته‌ایم همچنین مشکل اصلی این روش‌ها برای سازمان‌ها 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 برخی از سازمان‌ها که داده‌های طبقه‌بندی شده دارند 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 سازمان‌ها و مجموعه‌های نظامی از این‌گونه نرم‌افزارها استفاده نمود. لذا لازم است سازمان‌های مذکور این موضوع را به‌صورت بومی پیاده‌سازی کنند و یک بستر مناسب و مستقل از شبکه اینترنت برای تحلیل اطلاعات داشته باشند؛ در این پژوهش به مدیریت مستندات محتوای فارسی رسانه‌های خبری پرداخته می‌شود که خروجی آن برای سازمان‌ها جمع‌آوری و دسته‌بندی متون مختلف از سایت‌های خبری و سپس دسته‌بندی و بازیابی (مدیریت مستندات) است. همچنین مدلی جامع و کارآمد جهت تأمین امنیت داده‌ها و پیاده‌سازی سیستم مدیریت مستندات آنلاین خبری در سازمان‌ها ارائه می‌گردد. ساختار مابقی مقاله شامل ادبیات موضوع، پیشینه تحقیق، ارائه روش پیشنهادی، ایجاد مدل طبقه‌بندی، ارزیابی و تحلیل داده‌ها و در نهایت نتیجه‌گیری است.

## ۲- ادبیات موضوع و پیشینه تحقیق

## ۲-۱- طبقه‌بندی مستندات فارسی

محوری‌ترین موضوع در حوزه متن کاوی، طبقه‌بندی مستندات است. اساس کار طبقه‌بندی مستندات بر پایه کلمات کلیدی و مهمی است که از مستندات استخراج می‌شود [۳]. برای دسته‌بندی متن از فناوری‌های استخراج اطلاعات<sup>۱</sup>، پردازش زبان

متن به اواخر دهه ۵۰ میلادی بازمی‌گردد [۸]. همچنین طبقه‌بندی متن از دهه ۶۰ میلادی آغاز شده و در دهه‌های ۸۰ و ۹۰ میلادی توجه محققان بیشتری را به خود جلب کرده است [۹].

### ۲-۳- الگوریتم بیز ساده

الگوریتم‌های بیزین ساده در چند دهه اخیر عملکرد بسیار خوبی داشته و به همین دلیل مورد توجه هر چه بیشتر محققان قرار گرفته‌اند [۱۰]. الگوریتم نایوبیز که مدلی احتمالی بر اساس نظریه بیز است، یکی از روش‌های ساده اما قدرتمند طبقه‌بندی متن محسوب می‌شود. الگوریتم طبقه‌بندی نایوبیز یکی از ساده‌ترین و پرکاربردترین الگوریتم‌ها در میان طبقه‌بندی‌ها است. الگوریتم نایوبیز گاوسی<sup>۵</sup>، برنولی<sup>۶</sup> و چندجمله‌ای<sup>۷</sup> از زیرمجموعه‌های الگوریتم نایوبیز هستند. الگوریتم نایوبیز گاوسی برای داده‌های گسسته استفاده می‌شود و دو الگوریتم نایوبیز برنولی و چندجمله‌ای برای داده‌های پیوسته استفاده می‌گردد. این الگوریتم از یک مدل احتمالی بر پایه قضیه بیز استفاده می‌کند و با مستقل در نظر گرفتن ویژگی‌ها، احتمال وقوع داده جدید را در هر یک از دسته‌ها پیش‌بینی می‌کند. در نهایت سند به دسته‌ای تعلق می‌گیرد که بیشترین احتمال را دارد. اگرچه فرض مستقل در نظر گرفتن ویژگی‌ها در واقعیت درست نیست، اما الگوریتم نایوبیز با در نظر گرفتن این فرض نتایج خیلی خوبی ارائه داده است [۱۰].

این روش متداول‌ترین روش در دسته‌بندی متن است. در این روش متن به صورت مجموعه‌ای از کلمات مستقل از یکدیگر و مستقل از محل قرار گرفتن در متن در نظر گرفته می‌شود. تعریف تابع احتمال هر متن از حاصل ضرب احتمال کلمات آن و احتمال رخداد متنی با آن طول به دست می‌آید. عمل یادگیری در این بخش تخمین پارامترهایی برای تولید مدلی است که فقط از متون بر چسب‌دار استفاده می‌کند [۱۱].

### ۲-۴- الگوریتم ماشین بردار پشتیبان<sup>۸</sup>

ماشین بردار پشتیبان یک روش یادگیری ماشین است که اولین بار در سال ۱۹۹۵ توسط واپنیک<sup>۹</sup> برای حل مسئله تشخیص الگوی دو کلاسه مورد استفاده قرار گرفت. این روش در یک فضای برداری تعریف می‌شود و هدف مسئله یافتن یک جداکننده است که به بهترین شکل بتواند داده‌های دو کلاس داده‌شده را از هم جدا کند. مجموعه داده‌های آموزش داده‌شده به سیستم در صورتی جدایی‌پذیر به صورت بهینه خواهند بود که بدون خطا تعیین کلاس شده و فاصله بین نزدیک‌ترین نمونه آموزشی به ابر صفحه حداکثر باشد. وقتی که داده‌ها جدایی‌پذیر خطی نباشند،

طبیعی<sup>۱</sup> و یادگیری ماشین<sup>۲</sup> به طور وسیع استفاده می‌شود. به طور کلی هدف یک دسته‌بند متون، دسته‌بندی اسناد در قالب تعداد معینی از دسته‌های از پیش تعیین شده است. هر سند می‌تواند در یک، چند و یا هیچ دسته‌ای قرار بگیرد. در مورد هر سند به این سؤال پاسخ داده خواهد شد که این سند در کدام یک از دسته‌ها قرار می‌گیرد. این موضوع می‌تواند در قالب یک یادگیری خودکار قرار گیرد تا با استفاده از آن بتوان هر سند را به طور خودکار به دسته‌ای نسبت داد [۳]. بعد از انتخاب مجموعه داده و پاک‌سازی متون به کمک روش نرمال شده فرکانس کلمه - معکوس فرکانس سند به ویژگی‌ها وزن داده می‌شود و در طی دو مرحله ویژگی‌ها با استفاده از روش فرکانس سند و مربع‌کاری<sup>۳</sup> انتخاب می‌شوند [۴].

دسته‌بندی متون کاربردهای فراوانی می‌تواند داشته باشد که از جمله آن می‌توان به دسته‌بندی گفتاری که ترکیبی از دسته‌بندی متون و تشخیص گفتار است، دسته‌بندی متون چندرسانه‌ای از طریق عنوان‌های متنی، تشخیص نویسنده برای متون نامشخص با مورد بحث، تشخیص زبان برای متونی که زبان آن‌ها نامشخص است، تشخیص خودکار جنس متن، بایگانی اسناد، دسته‌بندی صفحات وب، یادگیری خودکار علایق مطالعاتی و پژوهشی کاربران، فیلتر کردن خودکار پست الکترونیکی بر اساس محتوا و... اشاره نمود [۲].

### ۲-۲- پردازش زبان طبیعی<sup>۴</sup>

پردازش زبان طبیعی در دهه ۵۰ میلادی به عنوان فصل مشترک علوم هوش مصنوعی و زبان‌شناسی به وجود آمد و در دهه ۸۰ دچار تحولات و تغییرات اساسی شد [۵]. پردازش متن تا دهه ۸۰ میلادی با استفاده از روش‌های قانون محور انجام می‌گرفت. اگرچه روش‌های قانون محور ویژگی‌های مثبتی از قبیل قابل درک بودن برای انسان و قابلیت ارتقا کیفیت در طول زمان را دارا هستند، مشکلات و محدودیت‌هایی نیز دارند. لازمه نوشتن قوانین در روش‌های قانون محور داشتن دانش عمیقی از آن حوزه است و از طرفی این روش‌ها بسیار زمان‌بر بوده و از نظر محاسباتی پیچیده هستند [۶]. این مشکلات و محدودیت‌ها باعث شد در دهه ۹۰ میلادی روش‌های یادگیری ماشین یا به عبارتی روش‌های آماری محور در پردازش متون طبیعی محبوب‌تر از روش قانون محور ظاهر شوند. در واقع در دهه ۹۰ میلادی بود که پژوهشگران به روش‌های یادگیری ماشین توجه بسیاری نشان دادند که این به دلیل توسعه سخت‌افزاری و نرم‌افزاری علم رایانه اتفاق افتاد [۷]. پژوهش‌ها نشان می‌دهند که طبقه‌بندی متن از نخستین تحقیقات انجام گرفته در حوزه متن‌کاوی است و نقش مهمی در آن ایفا می‌کند. به گفته فرانسویس تاریخچه پردازش

<sup>5</sup> Gaussian Naïve Bayes

<sup>6</sup> Bernoulli Naïve Bayes

<sup>7</sup> Multinomial Naïve Bayes

<sup>8</sup> Support Vector Machine

<sup>9</sup> Vapnik

<sup>1</sup> Natural Language Processing (NLP)

<sup>2</sup> Machine Learning

<sup>3</sup> Chi Square

<sup>4</sup> Natural Language Processing (NLP)

کرده اند. برای استخراج کلمات کلیدی از ساختار درهم‌ریزی Trie و دادگان مورد استفاده، دادگان همشهری بوده است. مقایسه دسته‌بندی متون فارسی با استفاده از الگوریتم‌های KNN و FKNN و انتخاب ویژگی‌ها بر اساس بهره اطلاعات و فرکانس سند [۱۶]. در این مقاله دو الگوریتم KNN و FKNN را 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.

در پژوهشی دیگر با عنوان بهبود عملکرد طبقه‌بندی متون فارسی با استفاده از تجزیه و تحلیل مؤلفه‌های اصلی با کمک معیار میانگین یادآوری و دقت [۱۷]. 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.

همچنین در پژوهش کاوش متون فارسی در وب با استفاده از تحلیل معنایی و روش طبقه‌بندی چندلایه‌ای اطلاعات [۱۸]. 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.

ارائه روش دسته‌بندی متون با تکنیک کاهش ویژگی فیلتری و یادگیری ماشین [۱۹]. 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.

دسته‌بندی متون لاتین با روش‌های انتخاب ویژگی فیلتری، پوششی و الگوریتم‌های بیز ساده [۲۰]. 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.

بررسی تطبیقی پرکاربردترین الگوریتم‌های طبقه‌بندی در پیش‌بینی و اعتبارسنجی مشتری [۲۱]. این مقاله به بررسی و مقایسه تطبیقی عملکرد آن‌ها در اعتبارسنجی و شناسایی مشتریان خوش حساب و بدحساب با استفاده از یک مجموعه داده واقعی آلمانی پرداخته است. نتایج این بررسی عملکرد بهتر درخت تصمیم را گزارش کرده است.

در کارهای مرتبط انجام شده 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 است که تحقیقات برای متون انگلیسی توسعه داده شده و به کار گرفته می‌شود؛ اما از آنجایی که زبان فارسی پیچیدگی‌های ساختاری نسبت به سایر زبان‌ها دارد و همچنین تحقیقات کمتری در این زمینه انجام گرفته است بر آن شدیم که ابتدا روش

ورودی‌های اولیه به صورت غیرخطی به فضای با ابعاد بالاتر نگاشت می‌شوند. به این ترتیب در فضای جدید داده‌ها با استفاده از تابع هسته می‌توانند با ابر صفحه مناسب، بدون افزایش پیچیدگی محاسبات و به صورت خطی از هم جدا شوند. در حقیقت تابع هسته از شباهت بین داده‌ها در فضای اولیه برای یافتن شباهت بین بردارها در فضای با ابعاد بالاتر استفاده می‌کند. تابع هسته می‌تواند تابع چند جمله‌ای، تابع تانژانت هیدرولیک و یا توابع مناسب دیگر انتخاب شود [۱۲]. این الگوریتم در برابر نویز داده‌ها مقاوم است، می‌تواند با تعداد زیادی ویژگی کار کند و در طبقه‌بندی متون عملکرد خوبی دارد [۱۳].

## ۵-۲- الگوریتم نزدیک‌ترین همسایه

الگوریتم طبقه‌بند نزدیک‌ترین همسایه، برخلاف الگوریتم طبقه‌بند نایبیز، فرایند طبقه‌بندی را تا حضور داده آزمون به تعویق می‌اندازد. به همین دلیل است که این الگوریتم در گروه یادگیرنده‌های تنبل طبقه‌بندی می‌شود [۱۴]. در الگوریتم طبقه‌بند نزدیک‌ترین همسایه، میزان شباهت داده آزمون با داده‌های آموزش با استفاده از یک معیار شباهت اندازه‌گیری می‌شود. سپس از کلاس همسایه‌های آن داده برای تشخیص کلاس داده مورد نظر استفاده می‌شود. برای اندازه‌گیری میزان شباهت بین دو سند (متن) می‌توان از معیارهای مختلفی از جمله فاصله اقلیدسی، ضریب شباهت جا کارد و شباهت کسینوسی استفاده کرد. فاصله اقلیدسی روش مناسبی برای اندازه‌گیری میزان شباهت بین دو متن است. پس از اندازه‌گیری میزان شباهت داده آزمون با هر یک از داده‌های آموزش و مقایسه این مقادیر، تعداد K تا از نزدیک‌ترین همسایه‌های (مشابه‌ترین اسناد به) سند آزمون انتخاب می‌شوند و در نهایت سند آزمون به دسته‌ای تعلق می‌گیرد که بیشترین همسایه‌هایش به آن دسته تعلق دارند؛ بنابراین می‌توان گفت که روش نزدیک‌ترین همسایه یک روش مبتنی بر مورد است چراکه از همه داده‌های آموزشی استفاده می‌کند. مقدار K می‌تواند هر عدد صحیح باشد. اگر مقدار K برای ۱ در نظر گرفته شود، دسته شبیه‌ترین همسایه به داده مورد نظر نسبت داده می‌شود. به دلیل وجود داده‌های پرت مقدار ۱ مقدار مناسبی برای K نیست و بهتر است از روش تجربی برای تخمین بهترین مقدار K بهره گرفته شود. در پژوهش حاضر از روش تجربی ۱ به منظور تعیین مقدار بهینه K استفاده شده است. در این روش با تکرار الگوریتم با مقادیر مختلف K، مقداری که باعث کمترین نرخ خطا می‌شود به دست آمده و بهترین مقدار K مشخص می‌گردد.

## ۶-۲- پیشینه پژوهش

در پژوهش استخراج کلمات کلیدی جهت طبقه‌بندی متون فارسی [۱۵]. از ترکیب روش‌های Wordnet و Porter استفاده

است که شرکت‌های داخلی شروع به فعالیت در زمینه متن‌کاوی نموده‌اند و پیشرفت‌های قابل توجهی هم در این زمینه داشته‌اند؛ اما کماکان نمی‌توان برای سازمان‌ها و مجموعه‌های نظامی از این گونه نرم‌افزارها استفاده نمود؛ لذا لازم است سازمان‌های مذکور این موضوع را به صورت بومی پیاده‌سازی کنند و یک بستر مناسب و مستقل از شبکه اینترنت برای تحلیل اطلاعات داشته باشند؛ لذا بر آن شدیم تا در این پژوهش روشی برای دسته‌بندی هوشمند متون ارائه نمائیم که بتواند از سایت‌های خبری، اخبار مختلف را دسته‌بندی نماید. امروزه سامانه‌های مدیریت متون فارسی که بهینه باشد بسیار احساس می‌شود. مدیریت متون، یکی از زیرمجموعه‌های متن‌کاوی است. متن‌کاوی 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 دودسته برای طبقه‌بندی اخبار را تشکیل می‌دهند. طبقه‌بندی متون قدمت بسیار زیادی در حوزه متن‌کاوی دارد این موضوع از سال ۱۹۶۰ میلادی مورد توجه محققین بوده است ولی بارونق کامپیوتر و نرم‌افزارها، مانند دیگر موضوعات هوش مصنوعی و داده‌کاوی در دهه ۹۰ میلادی توجه به این موضوع رشد چشمگیری پیدا کرد و مورد توجه قرار گرفت. از آنجایی که در انجام پژوهش‌های مرتبط با هوش مصنوعی و زبان‌شناسی رایانشی، از علوم زبان‌شناسی و علوم رایانه استفاده می‌شود، ماهیت بین‌رشته‌ای علم زبان‌شناسی رایانشی در تعیین روش و ابزار انجام پژوهش حاضر نقش به‌سزایی دارد. در این بخش در پی آشنایی با گام‌های پیشنهادشده برای انجام پژوهش هستیم. برای انجام تحقیق باید گام‌های ذیل مشخص شود که در شکل ۱ سیر مراحل مشخص و گویا است.

مراحل روش پیشنهادی به شرح ذیل است.

- الف- جمع‌آوری متون خبری فارسی (اخبار حوزه ایران، اخبار حوزه افغانستان و اخبار سطح جهانی)
  - ب- پیش‌پردازش و آماده‌سازی اخبار
  - ج- انتخاب ویژگی
  - د- مقایسه الگوریتم‌های یادگیری ماشین برای دسته‌بندی (استفاده از مدل اتوماتیک<sup>۱</sup>)
  - ه- ایجاد مدل موردنظر برای دسته‌بندی<sup>۲</sup>
  - و- ارزیابی مدل ارائه شده
- به طور خلاصه اجزای اساسی سیستم پیشنهادی در شکل ۲ قابل مشاهده است.

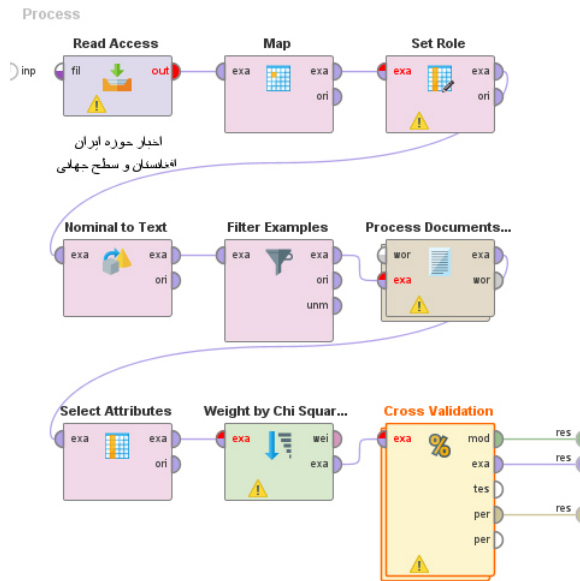
دسته‌بندی متون فارسی را بهبود ببخشیم و در نهایت با بومی‌سازی روش پیشنهادی و همچنین عدم وابستگی به پراکسی‌ها و شبکه اینترنت امنیت را برای سازمان‌ها تا حدودی در این زمینه تأمین نماییم. در تحقیقات اشاره شده در پیشینه تحقیق تلاش برای بهینه نمودن پارامترهای روش دسته‌بندی، بالابردن سرعت اجرای دسته‌بندی و همچنین کاهش درصد خطا است که در این پژوهش با استفاده از روش پیشنهادی موفق شدیم پارامترهای مذکور را تعمیم و درصد خطا را نیز کاهش بدهیم.

### ۳- روش پیشنهادی

امروزه استفاده از اینترنت و فناوری اطلاعات، در واقع در کنار مزایایی که دارد، خطر برملا شدن اسرار خصوصی را به همراه دارد. اتصال به اینترنت در سیستم‌ها و شبکه‌های کامپیوتری موجب شده است تا متخصصان امنیت، روزه‌روز شاهد تهدیدات امنیتی جدیدی باشند. مکانیسم‌های رمزنگاری و احراز هویت در سیستم‌های مدیریت پایگاه داده به معنای در امان بودن داده‌ها از آسیب‌پذیری‌های احتمالی نبوده و این مکانیزم‌ها توانایی مقابله با انواع حملات را ندارند. باید سعی شود که امنیت داده‌ها و حفظ حریم خصوصی در هنگام داده‌کاوی، کشف روابط بین داده‌ها و ذخیره در پایگاه داده‌ها حفظ شود. هدف از داده‌کاوی عمومی‌سازی اطلاعات است، نه اینکه اطلاعات شخصی را در اختیار عموم قرار دهیم. لازم به ذکر است که استفاده از داده‌کاوی جهت ایجاد مدل و استخراج الگو، امری ضروری است. حفظ حریم خصوصی در داده‌کاوی یک موضوع مهم در زمینه داده‌کاوی و امنیت پایگاه داده است. صاحبان داده‌ها به علت ترس از افشای اطلاعات شخصی و محرمانه خود توسط دیگران، چندان تمایلی جهت انجام داده‌کاوی روی داده‌های خود نداشته؛ ولی این مطلب را نیز می‌دانند که بدون انجام داده‌کاوی به نتایج و دانش مفید از داده‌های یکدیگر دسترسی پیدا نمی‌کنند. داده‌ها مادامی که پردازش نشوند و اطلاعات مفید از آن‌ها جهت تصمیم‌گیری‌های راهبردی، مباحث رقابتی و رسیدن به سود بیشتر استخراج نگردد، برای مدیران و برنامه‌ریزان چندان کارساز نیست. داده‌کاوی روی داده‌های مشترک باید طوری انجام شود که حریم شخصی صاحبان داده‌ها حفظ گردد؛ بنابراین مشکل اصلی بحث حفظ حریم شخصی در استخراج دانش، داده‌کاوی نیست، بلکه نحوه انجام داده‌کاوی است. الگوریتم‌ها و روش‌های مختلفی برای متن‌کاوی وجود دارد و همچنان نیز در حال توسعه است؛ اما مشکل اصلی این روش‌ها برای سازمان‌ها 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 داده‌کاوی و متن‌کاوی در بستر اینترنت کار می‌کنند که عملاً برای برخی از سازمان‌ها که داده‌های طبقه‌بندی شده دارند بدون استفاده است؛ خوشبختانه چند سالی

<sup>1</sup> Auto Model

<sup>2</sup> Classification



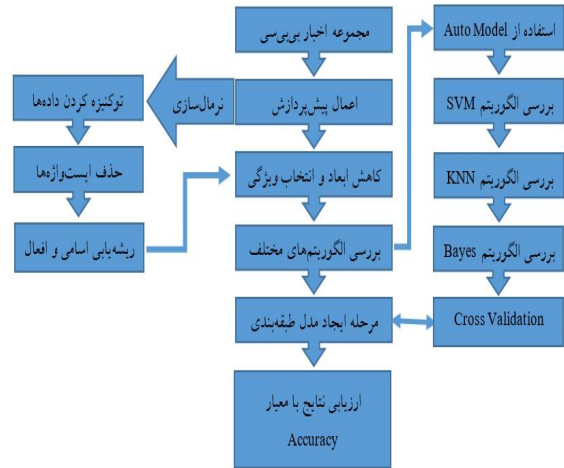
شکل (۳): نمایی کلی از فرآیند پیشنهادی در نرم افزار رپیدماینر

### ۳-۲- جمع آوری اخبار

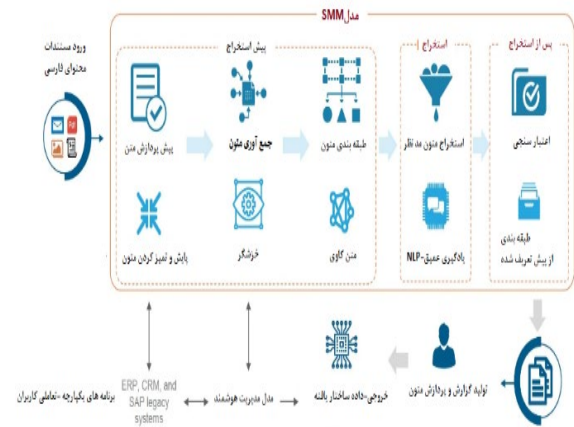
پایگاه خبری موردنظر، آرشیو اخبار بی بی سی است که این مجموعه داده به وسیله یک خزنده<sup>۱</sup> که به زبان پایتون نوشته شده است جمع آوری شده است جامعه آماری پژوهش حاضر اخبار فارسی ۶ماهه اول سال ۱۴۰۲ است. داده اخبار بی بی سی مجموعه ای از متون خبری فارسی به صورت برچسب خورده است. اخبار بی بی سی شامل ۱۷۰۰ عدد متن خبری از این مجموعه داده که فقط حاوی اخبار حوزه فارسی، اخبار حوزه افغانستان و سطح جهانی است. این خزنده با یک حلقه while مدام در حال جمع آوری اخبار بوده و اخبار مربوطه را در محل از قبل تعیین شده ذخیره می نماید. اخبار جمع آوری 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 را ارائه می دهد. در شکل ۵ صفحه مربوط به داشبورد مدیریتی از نحوه جمع آوری اخبار را از بستر اینترنت و فضای مجازی ارائه می دهد.

### ۳-۳- آماده سازی و پیش پردازش

در فاز آماده سازی متون، فرمت فایل از جیسون<sup>۲</sup> به اکسس<sup>۳</sup> تبدیل می شود، سپس متن سند که شامل کاراکترهای پشت سرهم است به نمایی که برای الگوریتم های یادگیری و طبقه بندی مناسب باشد تبدیل می شود. این کار توسط پلاگین های<sup>۴</sup> رپید ماینر<sup>۵</sup> و همچنین قطعه کدی که به زبان



شکل (۱): روش پیشنهادی مدل طبقه بندی



شکل (۲): مراحل مختلف ایجاد مدل طبقه بندی متون

### ۳-۱- تجزیه و تحلیل داده ها

جامعه آماری پژوهش حاضر اخبار فارسی ۶ماهه اول سال ۱۴۰۲ است. داده اخبار بی بی سی مجموعه ای از متون خبری فارسی به صورت برچسب خورده است. اخبار بی بی سی شامل ۱۷۰۰ عدد متن خبری از این مجموعه داده که فقط حاوی اخبار حوزه فارسی، اخبار حوزه افغانستان و سطح جهانی است. همچنین برای افزایش دقت، فراخوانی و همچنین کارایی کلی الگوریتم های طبقه بندی از روش های متعدد پیش پردازش از قبیل حذف کلمات اضافه و اصلاح فاصله به نیم فاصله در ساخت مجموعه آموزش استفاده کرده و از سه روش طبقه بندی ماشین بردار پشتیبان، بیز ساده و کی - نزدیک ترین همسایه برای طبقه بندی اسناد بهره برده ایم. برای انجام مراحل پیش پردازش، وزن دهی، انتخاب ویژگی و طبقه بندی از نرم افزار رپید ماینر استفاده کرده ایم.

ابتدا مجموعه داده اکسس خوانده می شود. سپس برای ستون های آن نقش تعریف می شود. محتوای خوانده شده به صورت متن تبدیل می شود تا قابل نوک نیزه شدن توسط رپید ماینر شود. در ادامه مراحل انجام کار شرح داده می شود.

<sup>1</sup> Crawler

<sup>2</sup> Json

<sup>3</sup> Access(\*.acddb)

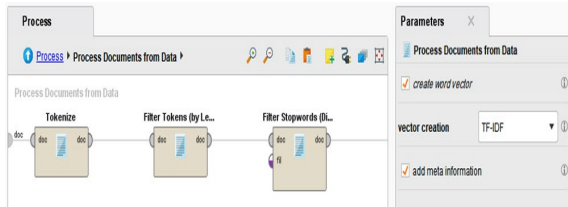
<sup>4</sup> Plugin

<sup>5</sup> Rapid Miner

قابل مشاهده است از الگوریتم پردازش متن معکوس فراوانی متن استفاده شده است.

در جدول ۱ فهرستی از فراوانی هر لغت در دسته مربوطه پس از عملیات توکنیزه کردن نمایش داده شده است.

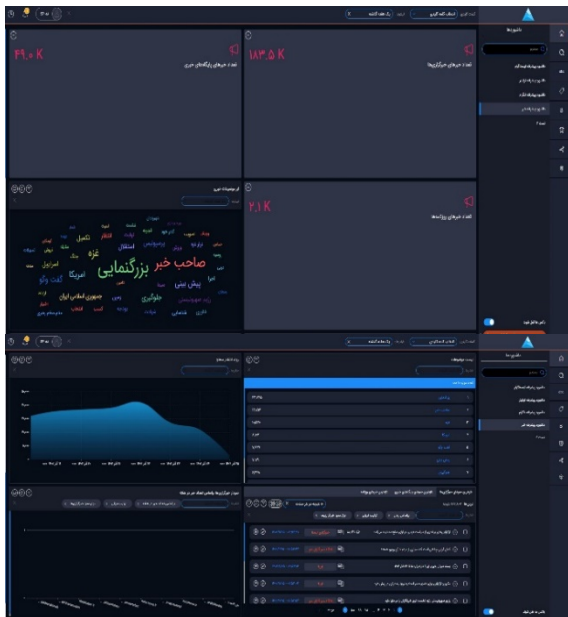
در شکل ۷ نمایی از داشبورد مدیریتی فراوانی کلمه‌ها را مشاهده می‌کنید.



شکل (۶): کاربرد الگوریتم پردازش متن TF-IDF

جدول (۱): نمونه‌ای از فراوانی هر کلمه در سه دسته خبری

کلمه توکن شده	اخبار سطح جهانی	اخبار حوزه ایران	اخبار حوزه افغانستان
تراکم	۳۷	۳۱	۷
تورویستی	۷	۱۷	۱۱
تورویست	۹	۹	۱۳
تُرک	۳۲	۱۴	۱۹



شکل (۷): نمایی از داشبورد مدیریتی فراوانی کلمه

### ۳-۵- کاهش ابعاد

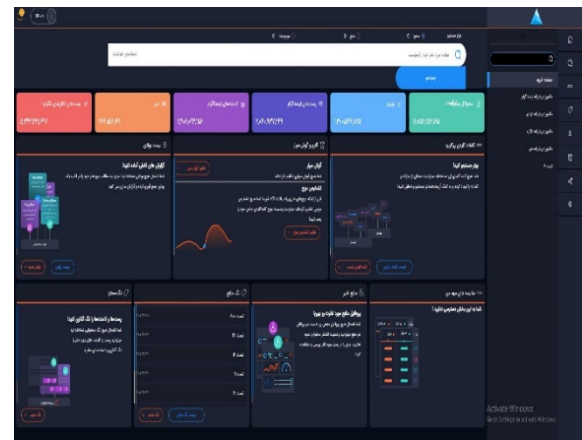
به علت سنگینی پردازش‌های موجود بر روی متون؛ و حذف حشویات و مواردی که باعث عدم خوانایی برنامه و کاهش دقت کار می‌شود، از روش کاهش ابعاد استفاده می‌شود. برای این منظور از ابزار فیلتر توکن استفاده می‌شود. این فیلتر یک ورودی متن می‌گیرد و خروجی آن هم فیلتر شده همان متن است. برای

پایتون نوشته شده است انجام می‌گیرد. در این کد از کتابخانه هم استفاده شده است. این فرایند شامل موارد ذیل است:

- حذف تگ‌های اچ‌تی‌ام‌ال<sup>۱</sup> و یا ایکسم<sup>۲</sup>
- کدگذاری متون به یونی‌کد<sup>۳</sup>
- به دست آوردن ریشه کلمات و حذف پیشوندها و پسوندها

category	body
1	یک نظرسنجی تازه جهانی نشان دادند صدها میلیون نفر از بسیاری جوانان دربار تغییر نگرش احساس می‌کنند. تقریباً ۶۰ درصد جوانان که در این بررسی مورد پرسش قرار گرفته گفتند که خیلی ناخوشایند احساس نگرانی می‌کنند.
2	امین شاد حدادی ۱۰ هزار پانجمین در افراسیاب بل مرزی میان ترکیه و آمریکا به شدت گرفتن بحران انسانی در این منطقه منجر شده است. این بل شهر بل درین نگران را به شهر آرتینا در ترکیه متصل می‌کند و اردوگاه موقت که در کنار آن شکل گرفته، در روزهای اخیر به شدت گسترش یافته است.
3	انحرفان ما در بهبهان انظار از سیاست‌های ایران در شهری کرمانجک تاوه ترین انحراف ما در ایران در شهری کرمانجک در استان خرمساران نمود پیدا کرده. در بهبهان اتفاقا یکی از مبداهای مرکزی بهبهان در انحراف‌های همین ایام سال گذشته هر کتابی منتخب بود.
4	انتخاب حوزه ایران
5	انتخاب حوزه ایران
6	انتخاب حوزه ایران
7	انتخاب حوزه ایران
8	انتخاب حوزه ایران
9	انتخاب حوزه افغانستان
10	انتخاب حوزه ایران
11	انتخاب حوزه ایران
12	انتخاب حوزه افغانستان
13	انتخاب حوزه ایران

شکل (۴): نمایی از مجموعه داده اخبار بی‌بی‌سی



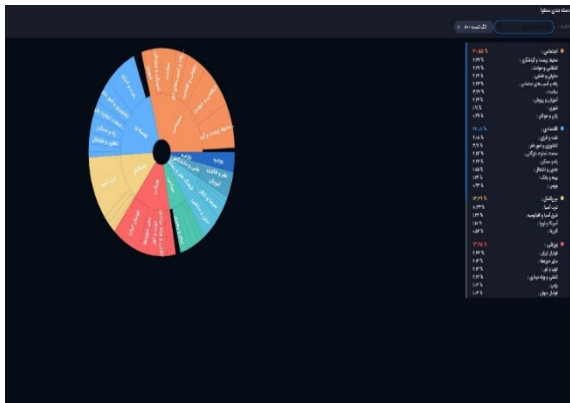
شکل (۵): نمایی از داشبورد مدیریتی جمع‌آوری خبر

### ۳-۴- نشان‌دار کردن اطلاعات (توکنیزه)

به منظور نوک نیزه کردن اطلاعات و تولید بردارهای قابل پردازش؛ از روش وزن دهی کلمات معکوس فراوانی متن استفاده شده است. این کار به کمک ابزار پردازش متن ریپید ماینر به نام توکنایز انجام می‌شود. تصویر زیر خروجی موردنظر را نمایش می‌دهد. دقت کنید که تمامی کنترل‌های موجود در تصویر در کامپوننت پردازش سند روی داده هستند. همان‌طور که در تصویر

<sup>1</sup> html  
<sup>2</sup> xml  
<sup>3</sup> UTF-8

ریشه کلمات قابل استفاده نیستند [۱]. اغلب روش های استفاده شده در این زمینه، روش های مبتنی بر یادگیری هستند. از مهم ترین روش های استفاده شده در طبقه بندی متون که در این پژوهش نیز استفاده شده است می توان به بیزین ساده، ماشین بردار پشتیبان و کی نزدیک ترین همسایه اشاره کرد. شکل زیر داشبورد مدیریتی مربوط به طبقه بندی ها را نمایش می دهد.



شکل (۸): نمایی از داشبورد مدیریتی طبقه بندی خبر

#### ۵- ارزیابی روش پیشنهادی و تحلیل نتایج

معمولاً برای ارزیابی صحت دسته بندی اسناد از معیار صحت<sup>۲</sup>، معیار فراخوانی<sup>۳</sup>، معیار دقت<sup>۴</sup> و معیار ارزیابی<sup>۵</sup> F استفاده می شود. دقت مربوط به درستی نسبت دادن داده های آزمون به کلاس مربوطه است و معیار ارزیابی<sup>۵</sup> F ترکیب هماهنگ از دقت و فراخوانی است. معادلات مربوطه در جدول زیر که موسوم به ماتریس درهم ریختگی<sup>۶</sup> آست آمده است. برای درک دو معیار ارزیابی دقت و فراخوانی و همچنین معیار صحت، ابتدا باید با عبارات مثبت واقعی<sup>۷</sup>، مثبت کاذب<sup>۸</sup>، منفی واقعی<sup>۹</sup> و منفی کاذب<sup>۱۰</sup> آشنا بود. این چهار عبارت برای مقایسه برچسب کلاس ها در ماتریس درهم ریختگی استفاده می شوند. جدول ۳ این عناوین را نشان می دهد.

صحت: نسبت مقدار موارد صحیح طبقه بندی 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 طبقه بندی کرده است که به صورت زیر محاسبه می شود:

$$Precision = \frac{TP}{TP+FP} \quad (1)$$

کاهش ابعاد کار؛ بهامهرهای کمتر از ۴ و بیشتر از ۲۵ را حذف کردیم. این کار به افزایش صحت کار کمک بیشتری می کند.

همان طور که در جدول ۲ مشاهده می کنید تعداد بار تکرار هر یک از کلمات پردازش شده توسط الگوریتم مربوطه در هر سه دسته ایران، افغانستان و سطح جهانی قابل مشاهده است.

همچنین برای کاهش اندازه ویژگی ها و بهینه تر شدن انتخاب ویژگی ها از عملگر قدرتمند وزن دهی با مربع کاری<sup>۱</sup> استفاده شد. خروجی این عملگر یک مجموعه از اوزان کلمات است که با توجه به فناوری های مربع کاری تولید و مورد استفاده قرار می گیرد.

جدول (۲): نمونه ای از کلمات پس از وزن دهی با روش مربع کاری

دسته	آزاد	آذربایجان	آخرین	آثار	آنتن	آبان
اخبار سطح جهانی	.	.	.	.	۰.۰۰۷	.
اخبار سطح جهانی	.	.	.	.	.	.
اخبار سطح جهانی	.	.	.	.	.	.
اخبار حوزه ایران	.	.	.	.	۰.۰۳۹	۰.۰۳۹
اخبار حوزه ایران	.	.	.	.	.	.
اخبار سطح جهانی	.	.	۰.۰۰۸	.	.	.

#### ۴- ایجاد مدل طبقه بندی به کمک الگوریتم های

##### یادگیری ماشین

بیشتر سامانه های طبقه بندی خودکار متون، برای متون زبان انگلیسی طراحی شده اند و معمولاً برای متون فارسی به دلیل ماهیت زبان فارسی و در دسترس نبودن مجموعه

<sup>۱</sup> Weight By Chi Squared

<sup>۲</sup> Precision

<sup>۳</sup> Recall

<sup>۴</sup> Accuracy

<sup>۵</sup> F-measure (F-score)

<sup>۶</sup> Confusion Matrix

<sup>۷</sup> True Positive

<sup>۸</sup> False Positive

<sup>۹</sup> True Negative

<sup>۱۰</sup> False Negative

## ۷- استفاده از ابزار مدل اتوماتیک

روش استفاده از مدل اتوماتیک به این صورت است که ابتدا داده موردنظر 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.

### ۷-۱- بررسی معیار دقت الگوریتم بیز ساده

در این بخش تأثیر روش مربع‌کاری و نرمال‌سازی کلمات بر روی عملکرد این الگوریتم را به صورت جداگانه بررسی می‌کنیم. فرایند انتخاب ویژگی تأثیر زیادی در بهبود عملکرد الگوریتم طبقه‌بند به خصوص در الگوریتم طبقه‌بند بیز ساده که به فرایند انتخاب ویژگی بسیار وابسته است، دارد بنابراین پیش‌بینی می‌شود که عملکرد الگوریتم طبقه‌بند بیز ساده با اعمال روش مربع‌کاری بهبود یابد. جدول ۴ میانگین معیار دقت را برای الگوریتم طبقه‌بند بیز ساده در سه حالت بدون اعمال نرمال‌سازی کلمات، اعمال نرمال‌سازی کلمات و با اعمال فرایند انتخاب ویژگی مربع‌کاری نشان می‌دهد.

در شکل ۹ نیز سایر معیارهای ارزیابی برای ترکیب الگوریتم بیز ساده و مربع‌کاری آمده است.

### جدول (۴): بررسی نرمال‌سازی کلمات و روش مربع‌کاری بر الگوریتم

بیز ساده

معیار ارزیابی	بیز ساده و عدم نرمال‌سازی	بیز ساده و اعمال نرمال‌سازی کلمات	بیز ساده و مربع‌کاری
میانگین دقت	۵۲٫۴٪	۵۳٫۷٪	۵۶٫۲٪

### Performances

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

### Confusion Matrix

	خبر حوزہ افغانستان true	خبر حوزہ ایران true	خبر حوزہ چغالی true	class precision
خبر حوزہ افغانستان pred.	12	14	68	72.34%
خبر حوزہ ایران pred.	44	203	143	52.05%
خبر حوزہ چغالی pred.	2	0	0	100.00%
class recall	3.45%	93.55%	32.23%	

شکل (۹): سایر معیارها برای الگوریتم بیز ساده

### جدول (۳): ماتریس درهم‌ریختگی

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

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:

$$Recall = \frac{TP}{TP+FN} \quad (2)$$

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.

$$Accuracy = \frac{TP+TN}{TP+TN+FP+FN} \quad (3)$$

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

$$F - Measure = 2 \times \frac{precision \times recall}{precision + recall} \quad (4)$$

## ۶- ابزار مورد استفاده

در این تحقیق از مجموعه داده اخبار فارسی به عنوان مجموعه آموزش و مجموعه آزمون استفاده شده است. این پیکره شامل ۱۷۰۰ خبر بین سال‌های ۱۴۰۰ تا ۱۴۰۲ است. 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، وزن دهی، انتخاب ویژگی و طبقه‌بندی از نرم‌افزار رپید مایند استفاده کرده‌ایم.

## ۷-۲- بررسی معیار دقت الگوریتم کی-نزدیکترین

### همسایه

در این بخش تأثیر نرمال سازی کلمات و روش مربع کاری بر روی عملکرد الگوریتم کی - نزدیک ترین همسایه را به صورت جداگانه بررسی می کنیم. جدول ۵ میانگین معیار دقت را برای الگوریتم کی - نزدیک ترین همسایه در سه حالت بدون اعمال نرمال سازی کلمات، اعمال نرمال سازی کلمات و با اعمال فرایند انتخاب ویژگی مربع کاری نشان می دهد.

در شکل ۱۰ نیز سایر معیارهای ارزیابی برای ترکیب الگوریتم کی-نزدیکترین همسایه و مربع کاری آمده است.

### جدول (۵): بررسی نرمال سازی کلمات و روش مربع کاری بر الگوریتم

کی-نزدیکترین همسایه

معیار ارزیابی	کی - نزدیک ترین همسایه و عدم نرمال سازی کلمات	کی - نزدیک ترین همسایه و اعمال نرمال سازی کلمات	کی - نزدیک ترین همسایه و مربع کاری
میانگین دقت	٪۸۶،۱	٪۸۷،۳	٪۸۹،۱

### Performances

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

### Confusion Matrix

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

شکل (۱۰): سایر معیارها برای الگوریتم کی-نزدیکترین همسایه

## ۷-۳- بررسی معیار دقت الگوریتم ماشین بردار

### پشتیبان

در این بخش تأثیر نرمال سازی کلمات و روش مربع کاری بر روی عملکرد الگوریتم ماشین بردار پشتیبان را به صورت جداگانه بررسی می کنیم. جدول ۶ میانگین معیار دقت را برای الگوریتم ماشین بردار پشتیبان در سه حالت بدون اعمال نرمال سازی کلمات، اعمال نرمال سازی کلمات و با اعمال فرایند انتخاب ویژگی مربع کاری نشان می دهد.

در شکل ۱۱ نیز سایر معیارهای ارزیابی برای ترکیب الگوریتم ماشین بردار پشتیبان و مربع کاری آمده است.

### جدول (۶): بررسی تأثیر نرمال سازی کلمات و روش مربع کاری بر

الگوریتم ماشین بردار پشتیبان

معیار ارزیابی	ماشین بردار پشتیبان و نرمال سازی کلمات	ماشین بردار پشتیبان و نرمال سازی کلمات	ماشین بردار پشتیبان و عدم نرمال سازی کلمات
میانگین دقت	٪۸۸،۸	٪۸۸،۸	٪۸۷،۵

### Performances

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

### Confusion Matrix

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

شکل (۱۱). سایر معیارها برای الگوریتم ماشین بردار پشتیبان

## ۸- 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 باتوجه به تحریم روی سایت های فارسی، بدون فیلتر اعمال نمی شوند و استفاده از پراکسی باعث کند شدن روند فعالیت خزش گر می شود. همچنین یک مجموعه واحد از ریشه لغات زبان فارسی وجود ندارد به همین دلیل عملیات ریشه یابی تأثیر زیادی در نتیجه عملکرد ما ندارد.

## ۹- نتیجه گیری

متن کاوی به عنوان یک زمینه در حال رشد و پرکاربرد، به دنبال کشف دانش از متون غیر ساخت یافته است. به دلیل مشکلات ساختاری زبان فارسی، تحقیقات کمتری در این زمینه صورت گرفته است. باتوجه 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 طبقه بندی متون انجام گرفته است بر روی متون زبان انگلیسی و چینی بوده است. در این پژوهش روشی برای طبقه بندی متون فارسی ارائه شده است. روش کلی بر مبنای یادگیری ماشینی استوار است که در دوفاز آموزش و آزمون ارائه شده است. باتوجه به تجزیه و تحلیل پیچیده متون فارسی از اکستنشن های<sup>۱</sup> متن

<sup>1</sup> Extension

- [7] S.Eyheramendy, A.GenKin, W.Ju, H.Lewis, D.Madigan, "Sparse bayesian classifiers for text categorization", *Journal of Intelligence Community Research and Development*, 13, 2003.
- [8] L.A.Francis, "Taming Text: An Introduction to Text Mining", Paper presented at the Casualty Actuarial Society Forum, 2006.
- [9] I.A.BolshaKov, A.GelbuKh, "Computational linguistics: models, resources, applications: Instituto Politecnico Nacional", 2004.
- [10] A.McCallum, K.Nigam, "A comparison of event models for naive bayes text classification", Paper presented at the AAAI-98 workshop on learning for text categorization, 1998.
- [11] Kardan.A, M.Keyhani Nejad, "Presenting a Model for Extracting Information from Text Documents Based on Text Mining in the Field of E-Learning", *Iranian Information and Communication Technology Scientific and Research Quarterly*, Year 4, Issues 11 and 12, pp. 47-54, 2012. (In Persian)
- [12] M.Sheikhi, S.Akbarpour, A.Farzan, "Text Mining of Persian Texts for Its Classification", *Fourth Iranian Conference on Electrical and Electronic Engineering*, 2012. (In Persian)
- [13] M.Bazghandi, Q.Tabrizi, J.Vafaei, "First International Conference on Natural Language and Handwriting Processing", Semnan University, 2012. (In Persian)
- [14] P. MulaK, N. Talhar, "Analysis of distance measures using K-nearest neighbor algorithm on KDD dataset", *International Journal of Science and Research*, 4(7), 2101-2104, 2015.
- [15] S.Arabi-Narei, M.Vahidi-Asl, B.Minaei-Bidgoli, "Extracting Keywords for Persian Text Classification", *First Iranian Data Mining Conference*, Tehran, <https://civilica.com/doc/33094>, 2007. (In Persian)
- [16] M.E.Basiri, S.Nemati, N.Ghasem Aghaei, "Comparison of Persian Text Classification Using K-Nearest Neighbor and fkNN Algorithms and Selecting Features Based on Information Gain and Document Frequency", *Thirteenth Annual Conference of the Iranian Computer Association*, Kish Island, <https://civilica.com/doc/41786>, 2007. (In Persian)
- [17] A.Ghanbari-Sorkhi, F.Ebrahimi, "Improving the performance of Persian text classification using principal component analysis with the help of average recall and precision criteria", *14th National Electrical Engineering Student nConference*, Kermanshah, <https://civilica.com/doc/121558>, 2011. (In Persian)
- [18] M.Abedininia, L.Allahdadi, F.Sheikhi, "Exploring Persian texts on the web using semantic analysis and multi-layer information classification method", *First National Conference on New Approaches in Computer Engineering and Information Retrieval*, Rasht, <https://civilica.com/doc/225377>, 2013. (In Persian)
- [19] P.Jalali, S.M.Hashemi, A.Siyahi, A.Heydari, "Presentation of a text classification method using filter feature reduction and machine learning techniques", *International Conference on New Research Findings in Electrical Engineering and*

موجود در ریپید ماینر استفاده شده که نتایج خوبی نیز در برداشته است. پس از ساخته شدن بردار ویژگی توسط ابزار قدرتمند پردازش متن از داده‌ها<sup>۱</sup> به کمک ابزار مدل اتوماتیک الگوریتم-های مختلف بررسی شد و مشخص گردید بهترین نتیجه مربوط به الگوریتم ماشین بردار پشتیبان است سپس با استفاده از ابزار اعتبارسنجی متقابل طبقه‌بندی اسناد صورت گرفت و نتایج این ابزار نیز نشان می‌دهد مدل ماشین بردار پشتیبان بهره‌وری<sup>۲</sup> بهتری داشته است و استفاده از روش‌های بی‌زین ساده و نزدیک‌ترین همسایه نتوانست نتیجه بهتری به دست آورد با اینکه ما در مدل پیشنهادی موضوع زمان اجرا<sup>۳</sup> را نیز به‌عنوان معیار مهم در اولویت قرار ندادیم. همچنین به‌عنوان کارهای آتی می‌توان به استفاده از روش‌های متفاوت وزن دهی مثل کلمه به بردار<sup>۴</sup>، استفاده از روش‌های جدید انتخاب ویژگی مثل الگوریتم کرم شب‌تاب<sup>۵</sup>، بهبود روش پردازش زبان طبیعی، بهبود و تقویت الگوریتم‌های ریشه‌یابی و نوک نیزه کردن و بررسی روش‌های تقویت بردار ویژگی مناسب برای زبان فارسی، مورد بررسی قرارداد.

## ۱۰- مراجع

- [1] B.Bina, M.Rahgozar, A.DehMobad, "Automatic Classification of Persian Texts", *13th Annual Conference of the Iranian Computer Association*, Kish Island, Computer Association, Sharif University of Technology, 2007. (In Persian)
- [2] I.Jamali, S.J.Mir-Abedini, A.Haroonabadi, "Presenting a model for Persian text classification using a combination of classification methods", *Journal of Telecommunications Engineering*, Volume 7, Issue 23, 2017. (In Persian)
- [3] M.Karimianesh, H.Shirazi, "Comparison of feature weighting methods in the document classification process", *First National Conference on New Approaches in Computer Engineering and Information Retrieval in Iran*, 2013. (In Persian)
- [4] A.Siyahi, S.M.Hashemi, S.Mazraa, "Improving a Feature Selection Method for Text Classification with Text Mining Algorithms", *Second National Congress of Iranian New Technologies for Sustainable Development*, Tehran, Center for Strategies for Sustainable Development, Mehr Arvand Institute of Higher Education, 2015. (In Persian)
- [5] P.NadKarni, M.Ohno-Machado, L.Chapman, "Natural language processing: an introduction", *Journal of the American Medical Informatics Association*, 18(5), 544-551, 2011.
- [6] I.Moulinier, J.G.Ganascia, "Applying an existing machine learning algorithm to text categorization", Paper presented at the *International Joint Conference on Artificial Intelligence*, 1995.

<sup>1</sup> Process Document From Data

<sup>2</sup> Performance

<sup>3</sup> Run Time

<sup>4</sup> Word2vec

<sup>5</sup> Firefly Algorithm

- [21] M.Golparvar, S.J.Mousavi-Rad, "A Comparative Study of the Most Used Classification Algorithms in Customer Prediction and Validation", Fifth National Conference on Applied Research in Management and Accounting, Tehran, Iranian Management Association., (2018). (In Persian)
- Computer Science, Tehran, Nikan Institute of Higher Educationhttps, 2015.(In Persian)
- [20] H.PourAbbasi, M.Abbasi-Dezfuli, "Latin Text Classification with Filtering Covering and Naive Bayes Algorithms", First National Conference on Computer and Information Technology, Sepidan, Islamic Azad University, Sepidan Branch, 2017.(In Persian)