

تاریخ دریافت: ۹۷/۰۶/۲۱

تاریخ پذیرش: ۹۷/۱۲/۰۱

مروری بر روش‌ها و پژوهش‌های دسته‌بندی خودکار متون فارسی

هادی ویسی*

استادیار دانشکده علوم و فنون نوین، دانشگاه تهران

پست الکترونیکی: h.veisi@ut.ac.ir

پویان پارسافرد

کارشناسی ارشد پردیس بین‌المللی، دانشگاه تهران

پست الکترونیکی: pooyanparsafard@ut.ac.ir

چکیده

می‌شود، بیانگر انتخاب برچسب یا موضوع برای اسناد متنی است. دسته‌بندی اسناد امروزه جزو موضوع‌های مهم و حائز اهمیت در پردازش زبان طبیعی^۲ محسوب می‌شود، از آن جهت که نقش ویژه‌ای را در کنترل و اداره حجم روز افزون داده و اطلاعات متنی ایفا می‌کند [۱]. اسناد فاقد این نوع دسته‌بندی به نوعی دچار ابهام مفهوم اولیه در حجم بالای متون خواهند شد و سیستم را وادار به بررسی و خواندن تک‌تک این اسناد فاقد برچسب، توسط نیروی انسانی خواهد کرد. پتانسیل مصرف شده توسط نیروی انسانی به خدمت گرفته شده را می‌توان با الگوریتم‌ها و روش‌های دسته‌بندی اسناد جایگزین کرد تا بتوان از نیروی انسانی به شکل کارا تر و بهینه‌تر در دیگر عرصه‌های پردازش داده استفاده کرد. در راستای این عمل، زمان مفید سیستم و بازدهی آن نیز به طور چشم‌گیری افزایش خواهد یافت. دسته‌بندی اسناد همچنان که حجم اطلاعات متنی را کاهش می‌دهد، مدیریت و جهت‌دهی اسناد در راستای پردازش آن‌ها را آسان‌تر و بهینه‌تر می‌سازد تا در ادامه، بازیابی اسناد نیز با سهولت و سرعت بیشتری انجام بگیرد.

در عصر فناوری ارتباطات که بخش بزرگی از آن را

دسته‌بندی اسناد متنی یا تشخیص عنوان به فرآیند شناسایی خودکار موضوع یک سند متنی (مانند هنری، ورزشی، سیاسی، علمی و ...) گفته می‌شود که در کاربردهای مختلف پردازش زبان طبیعی مانند بازیابی اطلاعات و تحلیل متون مورد استفاده است. یک سامانه دسته‌بندی‌کننده خودکار متون، مشابه اغلب سامانه‌های بازنمایی الگو، از دو گام مهم استخراج ویژگی و دسته‌بندی تشکیل شده است. در این مقاله، مروری بر روش‌های رایج برای استخراج ویژگی و دسته‌بندی در این سامانه‌ها صورت گرفته و پژوهش‌هایی که در این حوزه برای زبان فارسی انجام شده است، مرور شده‌اند. همچنین، تحلیلی از نقاط قوت و ضعف روش‌های موجود و مقایسه کارهای صورت گرفته با همدیگر ارائه شده است. واژه‌های کلیدی: دسته‌بندی متون فارسی؛ پردازش زبان طبیعی؛ مرور روش‌ها؛ استخراج ویژگی و دسته‌بندی

مقدمه

دسته‌بندی اسناد^۱ که به آن تشخیص عنوان^۲ نیز گفته

* نویسنده مسئول

1- Document Classification

2- Topic Identification

3- Natural Language Processing (NLP)

فضای داده‌های متنی و نوشتاری تشکیل می‌دهد، فرآیند تولید اسناد با سهولت و سرعت بالا انجام می‌گیرد و همین امر باعث ایجاد روزافزون نمونه‌های فاقد برچسب از این نوع داده‌ها شده است. این نوع داده‌ها که تحت عنوان متون دسته‌بندی نشده^۴ از آن‌ها یاد می‌شود، عمل بررسی و تحلیل محتوای سند را دشوارتر می‌کند. از این رو با توجه به رشد سریع این نوع داده، دسته‌بندی اسناد روز به روز از اهمیت بالاتری برخوردار خواهد شد. توسعه و گسترش اینترنت که در پی آن داده‌های دیجیتال از جمله متن استخراج می‌شوند بستر پردازش داده را برای فناوری‌های جدیدتر فراهم می‌کند. در مورد مسئله کلان داده^۵، که همواره موضوعی بحرانی و جدی در زمینه داده است، یکی از راهکارهای ابتدایی متداول که در مقدمه هر نوع پردازش داده خام صورت می‌گیرد، دسته‌بندی اسناد است.

کاربردهای بسیار متنوعی برای دسته‌بندی اسناد وجود دارد [۱، ۲]. برای مثال می‌توان شناسایی هرزنامه الکترونیکی، جویشرها، تحلیل محتوای دیجیتال، تعیین اعتبار، مشخص نمودن گروه‌هایی از مشتری‌ها که خصوصیات و علایق مشترکی دارند، تشخیص میزان تاثیر داروها و موثر بودن را نام برد.

مراحل و بخش‌های (اصولی) انجام یک فرآیند دسته‌بندی معمولاً بدین ترتیب است [۱]:

- پیش‌پردازش: پیش‌پردازش اولین مرحله فرآیند دسته‌بندی به حساب می‌آید. عملیاتی همچون حذف ایست‌واژه‌ها^۶، ریشه‌یابی^۷ کلمات، و رفتارهایی از این قبیل در این مرحله صورت می‌گیرد. از آنجایی که اعمال پیش‌پردازش در اکثر مواقع نتایج دسته‌بندی را بهبود می‌بخشد، این مرحله از اهمیت بالایی برخوردار است به طوری که تحقیقات زیادی در زمینه تاثیرگذاری عملیات پیش‌پردازش بر کیفیت دسته‌بندی، انجام گرفته است.

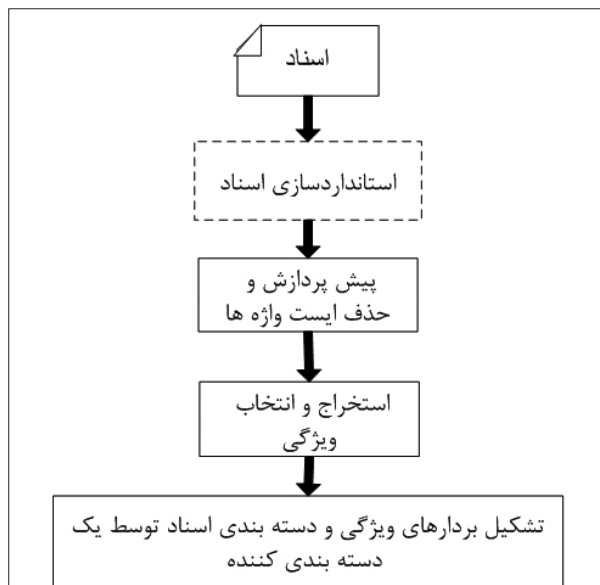
- استخراج ویژگی: بعد از پیش‌پردازش، برای اجرای فعالیت‌های آتی نیازمند استخراج عبارت‌های پردازش هستیم. در این مقاله، منظور از عبارت، همان عبارت پردازش^۸ است. هر سند حاوی متونی می‌باشد که می‌توان آن را در عبارت‌های بزرگ‌تر به صورت پاراگراف یا جمله دید، اما نهایتاً هر جمله خود متشکل از عبارت‌های کوچک‌تری چون کلمه می‌باشد، پس اگر کلمه را کوچک‌ترین عبارت پردازش در نظر بگیریم، عبارت یا عبارت‌های مورد پردازش همان کلمه یا کلمات استخراجی هستند. هر عبارت استخراجی یک ویژگی محسوب می‌شود، استخراج این ویژگی‌ها ما را با مجموعه‌ای از آن‌ها مواجه می‌کند، به طوری که اگر بخواهیم بردار را نماینده سند در نظر بگیریم هر سند با برداری از مجموعه ویژگی‌های آن سند برابری می‌کند. بنابراین اسناد را می‌توان مجموعه‌ای متشکل از بردارهای ویژگی فرض کرد.

- انتخاب ویژگی^۹: این مرحله که متناسب با نیاز روش دسته‌بندی، به صورت اختیاری انجام می‌گیرد، عمل انتخاب مجموعه ویژگی‌های مفید از مجموعه ویژگی‌های اولیه را انجام می‌دهد که طی این روند ویژگی‌هایی را که همبستگی کمتری با بقیه ویژگی‌ها دارند حذف می‌نماید و با انتخاب زیرمجموعه‌های کاراتر از ویژگی، سرعت دسته‌بندی افزایش پیدا می‌کند. این عمل کاهش بعد یا ابعاد، که در راستای حل مسئله افزایش ابعاد^{۱۰} (چندبعدی) انجام می‌گیرد منجر به کاهش فضای ویژگی و همچنین فضای مورد نیاز برای پردازش داده‌ها می‌گردد.

- دسته‌بندی: در پایان، فرآیند دسته‌بندی با انتخاب یک دسته‌بند مناسب تحت یادگیری انجام می‌گیرد. یادگیری در دسته‌بندی اسناد متنی، در اکثر حالات به این دو صورت انجام می‌گیرد: یادگیری با نظارت^{۱۱} و یادگیری بدون نظارت^{۱۲}، که در دسته‌بندی با نظارت سیستم یا

8- Term
9- Feature Selection
10- High Dimensionality
11- Supervised
12- Unsupervised

4- Uncategorized
5- Big Data
6- Stopwords
7- Stemming



شکل ۱: مراحل کلی فرآیند دسته‌بندی اسناد

مراحل سعی شده است که مهم‌ترین و برجسته‌ترین روش‌های انتخاب و ویژگی، انواع روش‌های نمایش اسناد و دسته‌بندی‌کننده‌های مورد استفاده شده در دسته‌بندی اسناد فارسی معرفی شود و توصیفاتی در راستای چگونگی روند هر کدام ارائه شود. همچنین، معیارهای ارزیابی متداول دسته‌بندی به وضوح تشریح خواهند شد.

۲-۱ معماری کلی روش‌های دسته‌بندی اسناد

در فرآیند دسته‌بندی اسناد، در اولین گام اسناد مورد پردازش تحت یک سری عملیات استانداردسازی یک‌دست می‌شوند. در گام بعدی پس از اعمال مرحله پیش‌پردازش و حذف ایست‌واژه‌ها، ویژگی‌های مورد نظر انتخاب و استخراج می‌شوند تا در نهایت با به نمایش درآوردن اسناد به شکل بردارهایی از ویژگی، دسته‌بندی اسناد توسط یک دسته‌بندی‌کننده مطلوب صورت گیرد. روند و مراحل کلی فرآیند دسته‌بندی اسناد در شکل ۱ قابل مشاهده می‌باشد. اسناد دریافتی بعد از نرمال کردن و استاندارد شدن، پیش‌پردازش شده و سپس ویژگی‌های مناسب از سند استخراج شده و پس از آن کار دسته‌بندی انجام می‌شود. هرکدام از این مراحل در ادامه به صورت دقیق‌تری تشریح شده است.

دسته‌بندی‌کننده، تحت داده‌هایی آموزش داده می‌شود که دارای برچسب‌های از قبل تعیین شده باشند. به بیانی دیگر اسناد به همراه برچسب یا دسته تعلق گرفته به آن، به منظور آموزش^{۱۳} سیستم مورد استفاده قرار می‌گیرند. اما در یادگیری بدون نظارت، معمولاً الگوریتم‌های خوشه‌بندی^{۱۴} هستند که نقش مهمی را ایفا می‌کنند. دسته دیگری از روش‌ها نیز وجود دارند که به آن‌ها نیمه‌نظارتی^{۱۵} گفته می‌شود و معمولاً استفاده همزمان از روش‌های بانظارت و بدون نظارت است.

در این مقاله، روش‌های رایج برای استخراج ویژگی و دسته‌بندی در سامانه‌های دسته‌بندی متون مرور شده است و پژوهش‌هایی که در این حوزه برای زبان فارسی انجام شده بررسی شده‌اند. برای این منظور، ساختار ارائه مطالب در این مقاله بدین شکل می‌باشد که در بخش دوم، به معرفی مراحل و گام‌های یک فرآیند دسته‌بندی پرداخته که در راستای آن انواع دسته‌بندی‌کننده‌ها و روش‌های دسته‌بندی عنوان شده و سپس، روش‌های گوناگون انتخاب ویژگی شرح و توضیح داده خواهد شد. همچنین در ادامه همین بخش، انواع معیارهای ارزیابی و روش‌های مقایسه عملکرد فرآیندهای دسته‌بندی همراه با معادله و نحوه محاسبه آن‌ها بیان خواهد شد. در بخش سوم، مروری جامع و موشکافانه به روی کارهای انجام شده و پژوهش‌های گزارش شده در زمینه دسته‌بندی اسناد فارسی صورت گرفته است که با مطالعه نتایج گردآوری شده حاصل از بررسی‌های انجام گرفته و همچنین یک دید کلی و مشاهده نتایج در یک نگاه بتوان در بخش چهارم، به نتیجه‌گیری و یک جمع‌بندی مفید، که هدف و ایده‌نهایی نوشتن این مقاله است، رسید.

۲- دسته‌بندی اسناد متنی

معماری و مراحل انجام یک فرآیند دسته‌بندی اسناد، عناوین اصلی این بخش را تشکیل می‌دهند. در کنار توضیح

13- Train
14- Clustering
15- Semi-Supervised

۲-۲ استانداردسازی اسناد و پیش پردازش

متون پیش از ورود به اولین مرحله پردازش (یعنی پیش پردازش) نیازمند یک سری اصلاحات اولیه می باشند [۱]. به بیان دیگر، اسناد قبل از این که وارد هرگونه پردازش و محاسبات اصلی شوند، بررسی می شوند که آیا استانداردها و معیارهای ورود به پردازش متن رعایت شده است یا خیر؟ بدین مفهوم که هر سند ممکن است بسته به منبع^{۱۶} آن، ویژگی های نوشتاری یا متنی خاصی را دارا باشد، مانند یک سند خبری از یک خبرگزاری که بسته به ضوابط و شیوه نگارش خبرگزاری، از یک سری نویسه های ویژه و شکل مانند، در متن استفاده شده است. استانداردسازی متن، کمک بزرگی را به پردازش های متنی آتی می کند از آن جهت که استخراج اطلاعات ارزشمند از داده خام، هزینه و بار محاسباتی را به نحو مطلوبی کاهش داده و در نتیجه سرعت کار بر روی داده ها نیز افزایش می یابد. منظور از استانداردسازی، شامل موارد مختلفی است از جمله:

- یکدست کردن نحوه نوشتن کلماتی که دارای صورت نوشتاری مختلف دارد، مانند «مسئول و مسوول»، «میشود، می شود و می شود» یا «کتابها و کتاب ها و کتاب ها»
 - یکسان کردن کدهای نویسه ها (کاراکترها) مانند «ک فارسی و عربی» و «ی فارسی و عربی»
 - حذف نویسه های اضافی مانند بولت
 - اصلاح فاصله گذاری علائم سجاوندی
- ایستواژه ها به دسته ای از کلمات اطلاق می گردد که قرارگیری آنها در فرآیند دسته بندی اسناد، بی اهمیت و بی ارزش می باشد، مثل حروف اضافه. حذف این عناصر در کنار جایگزینی کلمات با کلمه ریشه معادل آنها (ریشه یابی)، فرآیند پیش پردازش را تشکیل می دهد.

۲-۳ استخراج و انتخاب ویژگی از اسناد متنی

در پردازش زبان طبیعی به صورت عمومی و در دسته بندی اسناد به صورت مشخص، متن اسناد برای

پردازش باید به مدل فضای برداری (VSM^{۱۷}) برده شود و یا اصطلاحاً از آن ویژگی استخراج شود. ویژگی های استخراج شده از متن را می توان شامل دو دسته کلی سبدها کلمات^{۱۸} و جاسازی کلمات^{۱۹} دانست. در نمایش سبدها کلمات، خود کلمات و یا نمایشی عددی از آنها (مثلاً به صورت دودویی و یا بر اساس شمارش تعداد آنها) به عنوان بردار ویژگی متن در نظر گرفته می شود. این رویکرد به دلیل عدم در نظر گرفتن وابستگی بین کلمات و همچنین دارا بودن ابعاد بالای بردارهای ویژگی دارای ضعف هستند. در رویکرد جاسازی کلمات برای استخراج ویژگی، از یادگیری و برداری کردن هر کلمه یا هر سند استفاده می شود. برخی از روش های این دو گروه از استخراج ویژگی به صورت زیر هستند:

- روش های مبتنی بر سبدها کلمات [۳]
 - ♦ بردارهای دودویی
 - ♦ فراوانی عبارت^{۲۰} (TF)
 - ♦ فراوانی عبارت - معکوس فراوانی سند^{۲۱} (TF-IDF)
 - روش های مبتنی بر جاسازی کلمات
 - ♦ تحلیل معنایی پنهان^{۲۲} (LSA) [۴]
 - ♦ تحلیل آماری معنایی پنهان^{۲۳} (PLSA) [۵]
 - ♦ بردار کلمات مبتنی بر شبکه عصبی^{۲۴} [۶، ۷]
- در روش های مبتنی بر سبدها کلمات، برای به نمایش در آوردن هر سند، تعداد n کلمه (معمولاً بر اساس فراوانی تکرار در یک پیکره مانند داده های آموزش) را انتخاب کرده و به ازای هر سند یک وزن به آن کلمه داده می شود تا این بردار n بعدی نمایشی از آن سند باشد. این موضوع در رابطه (۱) برای سند d_k نشان داده شده است.

$$d_k = (w_{1,k}, w_{2,k}, \dots, w_{i,k}, \dots, w_{n,k}) \quad (1)$$

در این رابطه $w_{i,k}$ بیانگر کلمه i ام در سند d_k است و یا

17- Vector Space Model
18- Bag of Words
19- Word Embedding
20- Term Frequency
21- Term Frequency-Inverse Document Frequency
22- Latent Semantic Analysis
23- Probabilistic Latent Semantic Analysis
24- Neural Network based Word Vectors

16- Source

در صورت استفاده از N تایی^{۲۶}ها، بیانگر N تایی نام است. در نمایش برداری دودویی، مقادیر $w_{i,k}$ یک (در صورت حضور کلمه k در سند d_k) و یا صفر (در صورت عدم حضور کلمه k در سند d_k) است. در فراوانی عبارت، مقدار $w_{i,k}$ برابر با تعداد تکرار کلمه k در سند d_k است که معمولاً از مقدار نرمال شده آن (روی جمع تکرار کلمات و یا روی مقدار بیشینه تکرارها مانند رابطه (۲) استفاده می‌شود.

(۲)

$$tf_{i,k} = \frac{freq_{i,k}}{\max_l freq_{l,k}} \rightarrow tf_{i,k} = freq(w_{i,k})$$
 که در این معادله، $tf_{i,k}$ نمایانگر فراوانی عبارت k در سند i و $\max_l freq_{l,k}$ فراوانی پرتکرارترین عبارت k در سند k می‌باشد.

در روش فراوانی عبارت - معکوس فراوانی سند (TF-IDF) از ترکیب دو مقدار فراوانی عبارت (معادله ۲) و معکوس فراوانی سند (معادله ۳) به مقدار وزن حاصل شده در رابطه ۴ می‌رسیم [۳]. معادله ۳ بیان می‌کند idf_i معکوس (نرمال شده با لگاریتم) فراوانی تعداد اسنادی است که عبارت k در آن رخ داده است، N تعداد همه اسنادها n_i و تعداد اسنادهایی که شامل امین عبارت هستند. ایده IDF این هست که اگر عبارتی (کلمه‌ای) در همه اسنادها تکرار شده باشد، آن عبارت (کلمه) کمکی به ما در تمایز بین دسته‌های مختلف نمی‌کند و عبارت‌هایی (کلماتی) مفیدتر هستند که فقط در تعداد محدودی سند تکرار شده باشند. به عنوان مثال کلمه «تورینگ» در متن‌های علمی و فناوری و کلمه «رونالدو» در متن‌های ورزشی بیشتر از متن‌های دیگر رخ می‌دهند و از این رو، این دو کلمه به ترتیب برای تمایز دسته‌های علمی و ورزشی مناسب‌تر هستند.

$$idf_i = \log \frac{N}{n_i} \quad (۳)$$

$$TF-IDF_{i,k} = tf_{i,k} \times idf_i = \frac{freq_{i,k}}{\max_l freq_{l,k}} \times \log \frac{N}{n_i} \quad (۴)$$
 روش‌های دیگری مانند آنتروپی نرمال شده برای

استخراج ویژگی از متن وجود دارند که دارای منطقی مشابه TF-IDF هستند.

در دسته دوم روش‌ها، سعی می‌شود هر کلمه یا هر سند به نحوی به یک بردار عددی تبدیل شود که هم تعداد ویژگی‌ها کمتر باشند و هم بردارها حاوی اطلاعاتی دیگر (معنایی، نحوی و ...) در مورد آن کلمه یا سند باشند به نحوی که کلمات مشابه (مثلاً از نظر معنایی) دارای بردارهای نزدیک به هم باشند. به عنوان مثال بردار کلمه «خانه» با بردار کلمه «منزل» نزدیک به همدیگر باشند.

در میان روش‌های جاسازی کلمات برای استخراج ویژگی، روش تحلیل معنایی پنهان (LSA) [۴] که آن را نمایه‌سازی معنایی پنهان (LSI^{۲۶}) می‌نامند به عنوان یکی از روش‌های جبری-آماری برای تبدیل کلمات و اسناد به بردار ویژگی مورد استفاده قرار می‌گیرد. این روش سعی در یافتن معنی پنهان کلمات (عنوان یا موضوع سند) در اسنادها دارد. در این روش، فرض بر این است که کلماتی که به طور همزمان در یک سند با موضوع مشخص رخ می‌دهند، از نظر معنایی به هم مرتبط هستند و سندهایی که دارای موضوع مشابهی هستند، حاوی کلمات مشابهی هستند. در این روش، در یک سند، کلمات قابل مشاهده هستند ولی عنوان/موضوع آن پنهان (Latent) است. این روش با بهره‌گیری از تجزیه مقادیر تکین (SVD^{۲۷})، ماتریس عبارت-سند ساخته شده از فراوانی تکرار عبارت‌ها در اسناد را تجزیه می‌کند و با نگاشت بردار فراوانی سند یا عبارت به زیرفضاهای کوچک‌تر، بردار سند یا عبارت را پیدا می‌کند. اگر A ماتریس عبارت-سند دادگان موجود باشد، تجزیه SVD آن به صورت رابطه (۵) خواهد بود که در آن T ماتریس بردارهای کلمات، D ماتریس بردارهای مستندات و S یک ماتریس قطری حاوی مقادیر ویژه است که جهت کاهش بعد استفاده می‌شود.

$$A = TSD^T \quad (۵)$$

با انتخاب k مقدار از بیشترین مقادیر ویژه (از ماتریس)

26- Latent Semantic Indexing (LSI)
27- Singular Value Decomposition

25- N-gram

قطری این ماتریس استخراج شده و ابعاد ماتریس A با توجه به رابطه ۶ کاهش می‌یابد.

$$A \approx A_k = T_k S_k D_k^T \quad (6)$$

در واقع، روش LSA متون و کلمات را به یک فضای معنایی مشترک تصویر می‌کند و آن را می‌توان یک روش کاهش ابعاد بردارهای ویژگی دانست که از ایده‌ای مشابه به تحلیل اجزای اصلی^{۲۸} [۸] بهره می‌برد. با این روش، بردارهای حاصل (در حدود صد ویژگی) از بردارهای سبب کلمات (در حدود چند هزار ویژگی) حاصل می‌شود که بردارهای عبارات حاوی اطلاعات معنایی پنهان آن عبارات هستند.

نسخه احتمالی این روش PLSA [۵] نام دارد که متغیرهای پنهان (عنوان/موضوع) را با متغیرهای قابل مشاهده (سندها و کلمات) مرتبط می‌کند. مشکل اصلی استفاده از روش‌هایی همچون LSA این است که این روش‌ها بدون نظارت هستند یعنی نسبت به توزیع دسته‌ها کور هستند و آن را در نظر نمی‌گیرند. بنابراین ویژگی‌ای که با روش LSA انتخاب می‌شود، لزوماً در راستای توزیع دسته اسنادی که بهترین جداکننده هستند، قرار نمی‌گیرد [۱]. استفاده از روش‌های پیش‌برنده^{۲۹} در کنار ویژگی‌های مفهومی تهیه شده از روش PLSA بدون نظارت باعث بهبود دقت این نوع دسته‌بندی شده‌اند [۹].

تحلیل معنایی پنهان احتمالاتی تجزیه را توسط یک مدل آماری به نام مدل وجه^{۳۰} آغاز می‌کند. مشاهدات به صورت (w, d) از کلمات و مستندات موجود هستند. احتمال رخداد کلمه در سند با استفاده از رابطه ۷ محاسبه می‌گردد که در آن C موضوعات موجود بوده و احتمالات شرطی توسط الگوریتم درست‌نمایی بیشینه^{۳۱} محاسبه می‌شوند [۵].

(۷)

$$P(w, d) = \sum_c P(c) P(d|c) P(w|c) = P(d) \sum_c P(c|d) P(w|c)$$

روش بردار کلمات، یکی از روش‌های نوین بازنمایی

متون برای تبدیل کلمات به بردار ویژگی با استفاده از توان یادگیری شبکه‌های عصبی مصنوعی^{۳۲} است که موفقیت آن در مدل‌سازی مختلف پردازش زبان طبیعی، آن‌ها را به روشی اصلی تبدیل کرده است. در این رویکرد، با استفاده از پیکره‌های متنی^{۳۳} برای هر کلمه یک بردار (در حدود یکصد تا سیصد ویژگی) به صورت بی‌نظارت یاد گرفته می‌شود که به دلیل استفاده از اطلاعات بافت^{۳۴}، بردارهای حاصل دارای اطلاعات معنایی و نحوی هستند [۱۰]. منظور از اطلاعات بافت برای یک کلمه، اطلاعات کناری آن کلمه در متن مانند کلمات اطراف آن است. هرچند روش‌های مختلفی برای ایجاد بردار کلمات وجود دارد اما از میان آن‌ها سبب کلمه پیوسته (CBOW)^{۳۵} و پرس-چندتایی^{۳۶} [۱۱] رایج‌تر بوده و در ابزارهای متن‌بازی مانند Word2Vec^{۳۷} در دسترس هستند.

مدل سبب پیوسته کلمات از بافت متنی $[w(t-2), w(t-1), w(t+1), w(t+2), \dots]$ ، برای پیش‌بینی کلمه $w(t)$ استفاده می‌شود. این مدل از مدل زبانی شبکه عصبی پیش‌خور گرفته شده است اما در آن لایه غیر خطی پنهان حذف شده است. ورودی شبکه، بردارهای کدگذاری شده کلمات موجود در بافت با کدگذاری یکی فعال^{۳۸} هستند. وزن بین لایه ورودی و لایه خروجی به صورت ماتریس W با ابعاد $N \times V$ است که در آن V تعداد واژگان و N ابعاد (تعداد ویژگی) هر کلمه است. با استفاده از این ماتریس وزن‌دار می‌توان برای هر کلمه موجود در واژگان یک وزن محاسبه کرد و هر سطر ماتریس W بازنمایی N بعدی از بردارهای کلماتی است که در ورودی دریافت شده‌اند. سپس در لایه پنهان میانگین بردارهای کلمات محاسبه می‌شود. از لایه پنهان به لایه خروجی نیز ماتریس وزن‌دار W' با ابعاد $N \times V$ وجود دارد. بردار حاصل از خروجی این ماتریس وزن‌دار به عنوان توزیع احتمالاتی کلمه هدف

32- Artificial Neural Network

33- Text Corpus

34- Context Information

35- Continuous Bag of Words

36- Skip-gram

37- <https://code.google.com/p/word2vec/>

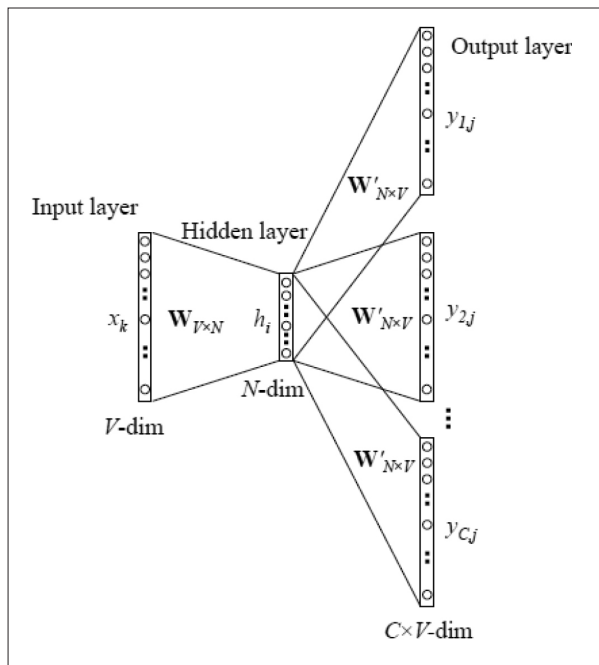
38- One-Hot Encoding

28- Principal Component Analysis

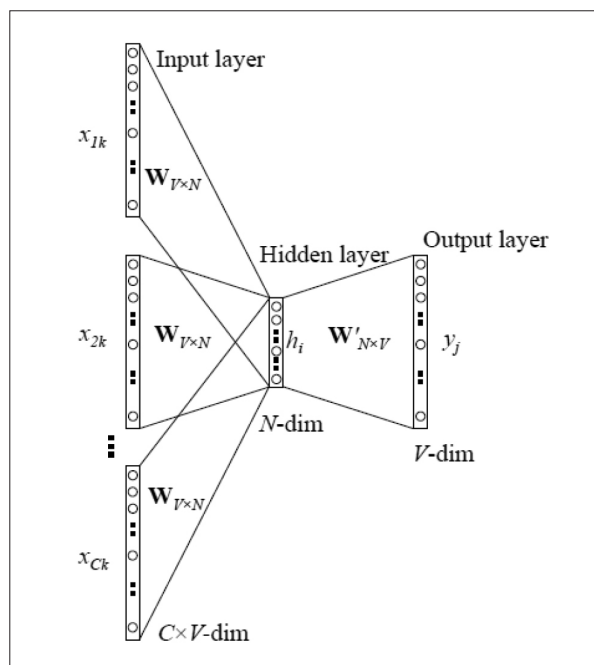
29- Boosting

30- Aspect Model

31- Maximum Likelihood



شکل ۳: معماری شبکه عصبی پرش چندتایی [۱۲]



شکل ۲: معماری شبکه عصبی سبب پیوسته کلمات [۱۲]

دیریکله (LDA)^{۴۱} [۱۵]، روش مدل‌سازی عنوان تنک (STC)^{۴۲} [۱۶] و روش بردار جهانی^{۴۳} [۱۷] اشاره کرد.

۲-۳-۱ انتخاب ویژگی

در حالتی که از رویکرد سبب کلمات به‌عنوان ویژگی استفاده می‌کنیم، با ابعاد بالای بردارهای ویژگی روبه‌رو هستیم که با افزایش بردارهای ویژگی مشکل افزایش ابعاد (چندبعدی) بردارهای ویژگی پدیدار می‌گردد که حل این مسئله خود به بهبود کیفیت دسته‌بندی اسناد نیز کمک شایانی خواهد نمود. راه حل مورد استفاده برای این مشکل استفاده از روش‌های انتخاب ویژگی است. در کل روش‌های انتخاب ویژگی به دو دسته تقسیم‌بندی می‌شوند، روش‌های پوشه‌ای^{۴۴} و روش‌های تصفیه‌ای (فیلتری)^{۴۵}، که روش‌های پوشه‌ای از الگوریتم‌های فراگیری اکتشافی^{۴۶} استفاده می‌کنند و با تخمین کارایی سیستم نتیجه‌گیری می‌کنند که آیا مجموعه ویژگی‌های انتخاب شده سودمند هستند یا خیر. روش‌های پوشه‌ای جواب‌های بهتری را برای مسائل مقیاس کوچک ارائه می‌دهند اگر چه

41- Latent Dirichlet Allocation (LDA)

42- Sparse Topical Model (STC)

43- Global Vector For Word Representaion

44- Wrapper Methods

45- Filter Methods

46- Heuristic Learning Algorithms

برگردانده می‌شود. شکل ۲ معماری مدل سبب پیوسته کلمات را نشان می‌دهد.

روش پرش چندتایی مشابه با مدل سبب پیوسته کلمات است اما هدف آن تعیین پارامترهایی است که احتمال کلمه جاری را در بافت متن به حداکثر برساند. معماری این مدل در شکل ۳ آورده شده است که در آن با استفاده از کلمه ورودی $w(t)$ ، بافت متن $[w(t-2), w(t-1), w(t+1), w(t+2), \dots]$ را پیش‌بینی می‌شود. به طور دقیقتر، از کلمه جاری به‌عنوان ورودی یک دسته‌بند خطی-لگاریتمی با لایه بازتاب پیوسته برای پیش‌بینی کلمات بافت متن در یک اندازه پنجره قبل و بعد از کلمه جاری بهره‌برده شود. ورودی شبکه، بردار کد گذاری شده کلمه ورودی است و از لایه ورودی به لایه پنهان ماتریس وزن‌دار W مرتبط با بردار کلمه ورودی وجود دارد.

روش‌های مرتبط دیگری جهت ایجاد بازنمایی برداری کلمات با رویکرد جاسازی کلمات وجود دارد که از آن‌ها می‌توان به روش‌های مبتنی بر خوشه‌بندی مانند خوشه‌بندی براون^{۳۹} [۱۳]، روش‌های مبتنی بر شبکه عصبی برای ایجاد مدل زبانی^{۴۰} [۱۴]، روش تخصیص پنهان

39- Brown Clustering

40- Neural Network Language Modeling

زمان محاسباتی بالایی را می‌طلبند تا آن‌جا که در مسائلی با تعداد ویژگی‌های بالا تقریباً ناکارآمد به نظر می‌آیند. در مقابل روش‌های تصفیه‌ای زمان کمتری را بابت مسائل با مقیاس‌های بزرگ‌تر مصرف می‌کنند پس می‌توانند نقش مفیدتری را نسبت به روش‌های پوشه‌ای، در دسته‌بندی اسناد بازی کنند. روش‌های تصفیه‌ای به جای به کارگیری الگوریتم‌های فراگیری اکتشافی، از قدرت تمایز ویژگی‌ها و دیگر اندازه‌گیری‌های آماری استفاده می‌کنند.

همان‌گونه که اشاره شد، روش‌های متنوعی برای انتخاب ویژگی وجود دارند که هر کدام بسته به نوع روابط و ضوابط مدون شده، اقدام به انتخاب ویژگی می‌نمایند. در زیر به چند نمونه از بارزترین آن‌ها اشاره شده است:

- خی‌دو آماری^{۴۷} (χ^2)
- بهره اطلاعات^{۴۸} (IG)
- اطلاعات متقابل^{۴۹} (MI)
- فراوانی سند^{۵۰} (DF)
- اهمیت عبارت^{۵۱} (TS)

در فرایند انتخاب ویژگی به روش خی‌دو آماری (χ^2)، یک ویژگی بر اساس میزان همبستگی‌اش با یک دسته یا طبقه انتخاب می‌شود [۱۸]. نحوه محاسبه این معیار به صورت رابطه ۸ تعریف می‌شود [۱۹]:

(۸)

$$\chi^2(t_i, c_j) = \frac{N \times (a_{ij}d_{ij} - b_{ij}c_{ij})^2}{(a_{ij} + b_{ij}) \times (a_{ij} + c_{ij}) \times (b_{ij} + d_{ij}) \times (c_{ij} + d_{ij})}$$

که در معادله ۸، N معرف تعداد اسناد است و سایر متغیرها به صورت زیر هستند:

a_{ij} تعداد سندهایی است که حاوی ویژگی t_i و متعلق به دسته c_j می‌باشند،

b_{ij} تعداد اسنادی است که حاوی ویژگی t_i نیست اما متعلق به دسته c_j می‌باشند،

c_{ij} تعداد سندهایی است که حاوی ویژگی t_i هست اما متعلق به دسته c_j نیست،

d_{ij} تعداد اسنادی است که نه حاوی ویژگی t_i است، و نه متعلق به دسته c_j .

انتخاب ویژگی به روش بهره اطلاعات (IG)، تعداد بیت‌های اطلاعاتی را می‌شمرد که از آن‌ها برای پیش‌بینی دسته در زمان حضور یا عدم حضور یک ویژگی در سند، استفاده می‌کند [۲۰]. رابطه آن به صورت زیر است:

$$IG(t_i) = - \sum_{j=1}^c P(c_j) \times \log P(c_j) + P(t_i) \times \sum_{j=1}^c P(c_j|t_i) \times \log P(c_j|t_i) + P(\bar{t}_i) \times \sum_{j=1}^c P(c_j|\bar{t}_i) \times \log P(c_j|\bar{t}_i) \quad (9)$$

با توجه به رابطه ۹، C نمایانگر تعداد دسته‌های سند و $P(c_j)$ احتمال رخداد دسته c_j ، $P(t_i)$ احتمال رخداد ویژگی t_i ، $P(\bar{t}_i)$ احتمال رخ ندادن ویژگی شرطی رخداد دسته c_j وقتی که ویژگی t_i هم رخ دهد، $P(c_j|\bar{t}_i)$ احتمال شرطی رخداد دسته c_j وقتی که ویژگی t_i هم رخ ندهد، است.

روش اطلاعات متقابل (رابطه ۱۰)، یک استاندارد همبستگی بین ویژگی و دسته است و می‌تواند برای تشخیص ارتباط ویژگی‌ها با یک دسته خاص به کار رود. اندازه اطلاعات متقابل با استفاده از نظریه اطلاعات حاصل شده که آن‌ها را بین دسته‌ها و ویژگی‌ها مدل‌سازی می‌کند و باید متذکر شد که این روش، ضعیفی را همراه دارد که آن، تمایل به انتخاب ویژگی‌های تنک^{۵۲} است که $P(t_i)$ پایینی دارند [۱۹].

$$MI(t_i) = \sum_{j=1}^c P(c_j) \times \log \frac{P(t_i|c_j)}{P(t_i)} \quad (10)$$

در رابطه ۱۰، $P(t_i|c_j)$ احتمال رخداد ویژگی t_i در احتمال شرطی رخداد ویژگی t_i وقتی که دسته c_j هم رخ دهد و $P(c_j)$ احتمال رخداد دسته c_j است.

در انتخاب ویژگی به روش اهمیت عبارت (TS)، اهمیت یک عبارت t با یک دسته c [۲۱] این‌گونه تعریف می‌شود:

(۱۱)

$$TS(t, c) = \begin{cases} \frac{\log(\max\{P(t), P(c)\})}{1 - \log(\min\{P(t), P(c)\})}, & \text{if } P(t, c) = 0 \\ \frac{\log(\max\{P(t), P(c)\}) - \log P(t, c)}{1 - \log(\min\{P(t), P(c)\})}, & \text{Otherwise} \end{cases}$$

47- Chi Square Statistic
48- Information Gain
49- Mutual Information
50- Document Frequency
51- Term Significance

که در اینجا، $P(t)$ بیانگر احتمال این که یک سند حاوی عبارت t باشد، $P(t, c)$ احتمال این که یک سند متعلق به دسته c باشد و حاوی عبارت t باشد، است.

برای اندازه‌گیری ارتباط یک عبارت در فضای سراسری ویژگی، همه امتیازات معین دسته در نظر گرفته می‌شوند، با فرض این که تعداد دسته‌ها باشد، وزن یک عبارت در همه دسته‌ها این‌گونه محاسبه می‌شود (رابطه ۱۲).

$$TS_{max}(t) = \max_{i=1, \dots, m} \{TS(t, c_i)\} \quad (12)$$

همه ویژگی‌ها، با توجه به وزن اهمیت عبارت $TS_{max}(t)$ به ترتیب کاهشی رتبه‌بندی خواهند شد و سپس به تعداد از پیش تعیین شده، عبارت‌هایی که بالاترین وزن‌ها را دارند برای فرآیند دسته‌بندی انتخاب می‌شوند. الگوریتم زیر فرآیند انتخاب ویژگی اهمیت عبارت را برای دسته‌بندی بیان می‌کند و همچنین این الگوریتم، دارای ملاحظات است که در [۲۱]، بیان شده‌اند.

۲-۴ روش‌های دسته‌بندی برای استفاده در دسته‌بندی اسناد متنی

به‌صورت کلی همه دسته‌بندی‌کننده‌های موجود را می‌توان برای کاربرد دسته‌بندی متن هم به کار برد که هر کدام به تنهایی پارامترهای قابل تنظیم مختص به خود را دارند که تغییر در هر پارامتر الگوریتم، ممکن است نتایج دسته‌بندی را کاملاً متفاوت سازد. گاهی اوقات چندین روش دسته‌بندی (دسته‌بندی‌کننده‌ها) به روی یک مجموعه داده مشترک نتایج خوبی را ارائه می‌دهند، در این موقعیت زمان مورد نیاز برای آموزش سیستم است که تعیین کننده عملکرد بهتر است. تا به حال تلاش‌های بسیاری در زمینه دسته‌بندی اسناد انجام گرفته است که دسته‌بندی‌کننده‌های متفاوت با درصدهای موفقیت متفاوت گزارش شده است، از جمله برجسته‌ترین این دسته‌بندی‌کننده‌ها می‌توان به موارد زیر اشاره کرد:

- بیز ساده^{۵۴}

- K نزدیک‌ترین همسایه (KNN)^{۵۵}
- ماشین بردار پشتیبان (SVM)^{۵۶}
- شبکه عصبی مصنوعی (ANN)
- مدل مخفی مارکوف^{۵۷} (HMM)
- درخت تصمیم^{۵۸}
- N تایی آماری^{۵۹}
- قوانین فازی^{۶۰}

مبنای دسته‌بندی به روش بیز ساده، احتمال‌هاست و احتمال حضور کلمات در یک سند، تعیین کننده دسته آن سند هستند. توجه دسته‌بندی به روش K نزدیک‌ترین همسایه به اسناد آموزش است و روش آن بدین ترتیب است که ابتدا شباهت اسناد آموزش با سند ورودی محاسبه می‌شود، سپس به تعداد K تا سند از اسناد آموزش که به سند ورودی شبیه‌تر هستند، همسایگان سند مورد نظر محسوب می‌شوند. روش ماشین بردار پشتیبان یکی از موفق‌ترین روش‌های دسته‌بندی است که تا به حال برای دسته‌بندی متون گزارش شده است و با پیدا کردن بهترین مرز تصمیم‌گیری بین دسته‌های مختلف (با شروع از دسته‌بندی دو دسته‌ای) از روی مجموعه داده آموزش، یک ابرصفحه^{۶۱} بین این هر دو دسته به دست می‌آورد. در این روش، سعی می‌شود فاصله بین صفحات مرزی دو دسته و ابرصفحه دسته‌بندی بیشینه^{۶۲} باشد. شبکه عصبی مصنوعی یک مدل ریاضی است که ساختار و پردازش‌های مغز انسان را شبیه‌سازی می‌کند و به‌صورت جلورو^{۶۳} یا بازگشتی^{۶۴} استفاده می‌شود. یک شبکه عصبی معمولاً شامل لایه ورودی (که ویژگی‌های سند یا کلمات را دریافت می‌کند)، لایه خروجی (که بیانگر دسته‌ها یا موضوعات مختلف است) و تعدادی لایه‌های مخفی است [۲۲]. از شبکه‌های موفق مورد استفاده در

54- K-Nearest Neighbors
55- Support Vector Machines
56- Hidden Markov Model
57- Decision Tree
58- Statistic N-Gram
59- Fuzzy Rules
60- Hyperplane
61- Maximum
62- Forward
63- Recurrent

53- Naive Bayes

- دقت دسته که همان $Precision(c_j)$ است.
- معیار دسته که همان $F - Measure(c_j)$ است.
- درستی دسته که همان $Accuracy(c_j)$ است.

$$Recall(c_j) = \frac{A_j}{A_j + B_j} \quad (۱۳)$$

در این معادله، معیار فراخوانی دسته از درصد سندهای درست دسته‌بندی شده، در میان همه سندهایی است که باید به دسته هدف^{۶۸} تخصیص می‌یافتند.

$$Precision(c_j) = \frac{A_j}{A_j + C_j} \quad (۱۴)$$

با توجه به معادله ۱۴، معیار دقت دسته از درصد سندهای درست دسته‌بندی شده، در میان همه سندهایی است که به دسته هدف تخصیص یافته‌اند.

$$F - Measure(c_j) = \frac{2P(c_j) \times R(c_j)}{P(c_j) + R(c_j)} \quad (۱۵)$$

و در معادله ۱۵، معیار اندازه F دسته از یک میانگین هارمونیک^{۶۹} از $Precision(c_j)$ و $Recall(c_j)$ می‌باشد.

$$Accuracy(c_j) = \frac{A_j + D_j}{A_j + B_j + C_j + D_j} \quad (۱۶)$$

در معادله ۱۶، درستی، بیانگر درصد تعداد سندهای درست دسته‌بندی شده است.

برای تجمیع نتایج روی همه نمونه‌ها، از میانگین‌گیری استفاده می‌شود که این تجمیع می‌تواند میانگین‌گیری ریز^{۷۰} یا میانگین‌گیری درشت^{۷۱} باشد. معیارهای ارزیابی ریز که از معادلات ۱۳، ۱۴ و ۱۵ مشتق می‌شوند، این‌گونه تعریف می‌شوند:

$$R_{micro} = \frac{\sum_{j=1}^c A_j}{\sum_{j=1}^c (A_j + B_j)} \quad (۱۷)$$

$$P_{micro} = \frac{\sum_{j=1}^c A_j}{\sum_{j=1}^c (A_j + C_j)} \quad (۱۸)$$

$$F_{micro} = \frac{2P_{micro} \times R_{micro}}{P_{micro} + R_{micro}} \quad (۱۹)$$

68- Target Category
69- Harmonic Mean
70- Micro Averaging
71- Macro Averaging

دسته‌بندی متون می‌تواند به شبکه جلورو پرسپترون چند لایه (MLP)^{۶۴} و شبکه بازگشتی حافظه کوتاه مدت ماندگار (LSTM)^{۶۵} اشاره کرد. روش مدل مخفی مارکوف به‌عنوان روشی موفق در مدل‌سازی داده‌های متوالی مانند گفتار، می‌تواند با مدل‌سازی احتمالاتی دسته‌ها نسبت به دسته‌بندی متون از روی مشاهدات (کلمات) اقدام کند [۲۳]. تصمیم و غلبه^{۶۶}، روند کاری دسته‌بندی به روش درخت تصمیم را تشکیل می‌دهد به طوری که داده‌های آموزش از قبل دسته‌بندی شده توسط عبارت تقسیم‌بندی می‌شوند. انتخاب عبارت را می‌توان بر عهده معیار انتخاب ویژگی بهره اطلاعات، یا معیارهای گوناگون دیگر گذاشت. پس در ادامه مجموعه آموزش به دو زیرمجموعه تقسیم می‌شود، یکی زیرمجموعه‌ای که شامل عبارت می‌شود و دیگری زیرمجموعه‌ای که عبارت را شامل نمی‌شود. فرآیند مذکور را دوباره برای آن دو زیر مجموعه تکرار می‌کنیم. به همین شکل، ادامه اسناد را تقسیم‌بندی می‌نماییم و این روند تا زمانی تکرار می‌شود و ادامه پیدا می‌کند که همه اسناد موجود در یک زیرمجموعه متعلق به یک دسته باشند. گره^{۶۷} آخر تبدیل به برگ می‌شود به عبارتی دیگر فرآیند تقسیم و غلبه را تا زمانی ادامه می‌دهیم که دیگر عبارتی برای تقسیم‌بندی زیرمجموعه‌ها موجود نباشد [۱].

۲-۵ معیارهای ارزیابی سامانه‌های دسته‌بندی اسناد متنی

برای معرفی معیارهای ارزیابی دسته‌بندی اسناد، ابتدا باید پارامترهای مورد استفاده در این معیارها را بشناسیم. این معیارها شامل برچسب‌هایی است که نشان می‌دهد آیا سند مورد آزمون در دسته مورد نظر برچسب‌گذاری شده است یا خیر، و آیا سند مورد آزمون در دسته مورد نظر دسته‌بندی شده است یا خیر. با این توصیفات، می‌توان چهار معیار ارزیابی دسته را این‌گونه تعریف کرد:

- فراخوانی دسته که همان $Recall(c_j)$ است.

64- Multi-Layer Perceptron
65- Long Short-Term Memory
66- Decision and Conquer
67- Node

در میانگین‌گیری درشت وزن همه دسته‌ها برابر در نظر گرفته می‌شود و به صورت رابطه‌های زیر محاسبه می‌شوند.

$$R_{macro} = \frac{1}{C} \times \sum_{j=1}^c R(c_j) \quad (20)$$

$$P_{macro} = \frac{1}{C} \times \sum_{j=1}^c P(c_j) \quad (21)$$

$$F_{macro} = \frac{2P_{macro} \times R_{macro}}{P_{macro} + R_{macro}} \quad (22)$$

$$Accuracy_{macro} = \frac{1}{C} \times \sum_{j=1}^c Accuracy(c_j) \quad (23)$$

ذکر نتایج آورده‌ایم. آزمایش‌های دسته‌بندی به روی زبان فارسی با سیستم‌های دسته‌بندی متفاوت انجام گرفته است، که با توجه به تنوع حضور دسته‌بندی‌کننده‌های مختلف و همچنین استفاده از انواع گوناگون روش‌های انتخاب ویژگی و ترکیب و تلفیق الگوریتم‌های کاربردی دیگر، نتایج گسترده‌ای به دست آمده است. مشابه کار جاری، برخی کارهای مروری مختصر به زبان فارسی وجود دارد [۲۴] که فقط به مرور روش‌های یادگیری بسنده کرده‌اند و کارهای فارسی را بررسی نکرده‌اند.

مجموعه آزمایش‌های دسته‌بندی مرتبط با هر دسته‌بندی‌کننده در بخش مربوطه بررسی شده‌اند. سپس در انتهای فرآیند بررسی همه این آزمایش‌ها به روی اسناد فارسی، نگاهی کلی به تاثیرات روش‌های مختلف مورد استفاده و نحوه استفاده از آن‌ها بر کیفیت دسته‌بندی اسناد فارسی خواهیم داشت.

۳-۱ k نزدیک‌ترین همسایه (KNN)

روش KNN که یکی از دسته‌بندی‌کننده‌های معمول است، درصد بزرگی از آزمون‌های انجام شده در زمینه دسته‌بندی را به خود اختصاص داده است. این روش دسته‌بندی که به تنهایی قادر به ارائه یک دسته‌بندی خوب و مطلوب نیست، معمولاً با روش‌ها و الگوریتم‌های دیگر به کار می‌رود. آزمون KNN به روش‌های مختلف بر روی اسناد فارسی صورت گرفته است. در ادامه به ذکر شرایط، روش‌ها (روش انتخاب ویژگی و ...) و پارامترهای به کار گرفته در هر آزمون می‌پردازیم.

آزمایش انجام گرفته بر روی مجموعه داده فارسی همشهری، دسته‌بندی به روش k نزدیک‌ترین همسایه و انتخاب ویژگی TF-IDF می‌باشد. طی بررسی [۲۵] این روش دسته‌بندی نیز به مانند روش ماشین بردار پشتیبان، که هر دو بر روی یک مجموعه داده مشترک آزمایش شده‌اند، نتایج خوبی را در بر داشته است و بیانگر تاثیر بسزای تعداد ویژگی در روند نتایج است. روش دیگر استفاده از مقدار kهای متفاوت است.

۳- مروری بر پژوهش‌های انجام شده در دسته‌بندی اسناد فارسی

توضیحات مراحل و همچنین روند کارها و پژوهش‌های برجسته انجام شده در زمینه دسته‌بندی اسناد فارسی، همراه با گزارش بهترین نتایج به دست آمده توسط آن‌ها، در این بخش گردآوری شده‌اند تا با هدف تحلیل و بررسی روش‌ها و الگوریتم‌های استفاده شده، بتوان در انتهای این مقاله یک نتیجه‌گیری مناسب از کارهای انجام شده در زمینه زبان فارسی داشته باشیم.

هر زبان طبیعی، پیچیدگی‌ها و قواعد خاص خود را دارد و زبان فارسی نیز از این قاعده مستثنی نیست. با توجه به ساختار زبان فارسی، ویژگی‌هایی در این زبان وجود دارد که کار دسته‌بندی را (نسبت به زبان انگلیسی) کمی دشوارتر می‌سازند. به عنوان مثال، برای اشاره به چند نمونه از این دشواری‌ها می‌توان بحث تفکیک کردن عبارت‌ها (عبارت‌های پردازش) و یا تشخیص فضاهای خالی (واقعی) را مطرح کرد. در زمینه دسته‌بندی اسناد تا به امروز، مطالعات مختلفی به روی اسناد فارسی صورت گرفته است. اما تا آنجایی که سعی شده بیشتر تحقیقات و بررسی‌های برجسته‌ای (در محدوده دسته‌بندی اسناد فارسی) که تا الان صورت گرفته‌اند، در این پژوهش با

بررسی نتایج این روش نشان می‌دهد که هر چقدر مقدار k به سمت اعداد پایین‌تری میل کند نتایج بهتری را در پیش خواهیم داشت. به عبارتی دیگر، اجرای الگوریتم به منظور پیش‌بینی دسته نمونه جدید، هر چقدر با تعداد همسایگان کمتری صورت بگیرد نتایج بهتری را در پی خواهد داشت. فاکتور سوم مورد بررسی معیار فاصله می‌باشد، این معیار با استفاده از روابط شباهت، میزان شباهت یا عدم شباهت نمونه‌ها را نسبت به یکدیگر می‌سنجد. بررسی‌ها نشان می‌دهد که معیار فاصله کازین، بهترین عملکرد را در محاسبه میزان شباهت بین نمونه‌های آموزش و آزمون ارائه می‌دهد و بعد از آن، معیار فاصله همبستگی نیز معیار مفیدی برای اندازه‌گیری شباهت مورد نظر است.

روش دسته‌بندی k نزدیک‌ترین همسایه و روش ماشین بردار پشتیبان، با توجه به این که به روی یک مجموعه داده مشترک فارسی، تحت یک تحقیق مشترک [۲۵] مورد آزمون و بررسی قرار گرفتند، هر دو نتایج خوبی را گزارش دادند با این تفاوت که، روش k نزدیک‌ترین همسایه عملکرد بهتری را نسبت به روش ماشین بردار از خود نشان داد، آن چه که به صورت مشترک بین این دو الگوریتم تاثیر صعودی از خود به جای می‌گذارد، تعداد ویژگی‌ها است که با افزایش آن‌ها روند صعودی دقت دسته‌بندی را مشاهده خواهیم کرد. انتخاب یک معیار شباهت درست و درخور در روش k نزدیک‌ترین همسایه، عملکرد دسته‌بندی را بهبود خواهد داد.

یک بررسی دیگر دسته‌بندی اسناد فارسی، بر روی مجموعه داده‌ای با نام همشهری انجام گرفته است [۲۶]. در این مطالعه کارایی دسته‌بندی متاثر از استفاده N تایی حروف برای تبدیل اسناد به بردارهای عددی مورد بررسی قرار گرفته است. همچنین اثر دسته‌بندی‌کننده بهبود یافته k نزدیک‌ترین همسایه نیز مورد مطالعه قرار گرفته است. در دسته‌بندی به روش k نزدیک‌ترین همسایه، توزیع نامتقارن اسناد آموزش در دسته‌های مختلف به روی دسته‌بندی‌کننده تاثیر می‌گذارد. از آنجایی که در

الگوریتم k نزدیک‌ترین همسایه، دسته برنده دسته‌ای است که بیشترین همسایه را داراست، یک دسته با اسناد آموزش بیشتر نیز شانس بیشتری برای برنده شدن خواهد داشت.

روش N تایی حروف، یک روش متداول برای نمایش بردارها به شکل عددی است. در این روش، هر سند متنی به تکه‌هایی به طول N از حروف همجوار تقسیم می‌شود. بردار مرتبط با هر سند شامل یک لیست از N تایی‌های غیر تکراری، حاوی عدد تکرار هر N تایی است. در این بررسی اثر طول N تایی‌های متفاوت بین ۲ و ۱۰ بر کارایی دسته‌بندی‌کننده تحلیل می‌شود.

در این مطالعه آزمایش‌های متعددی به منظور ارزیابی عملکرد روش جدید انجام گرفته، دقت دسته‌بندی‌کننده با رویه اعتبارسنجی ۵ تایی اندازه‌گیری شده، به این معنا که متن به پنج بخش منحصر تقسیم می‌شود و الگوریتم برای هر بخش یک بار اجرا می‌شود که هر بار یک بخش به عنوان مجموعه آزمون، و مابقی ۴ بخش دیگر به عنوان مجموعه آموزش مورد پردازش قرار می‌گیرند. نهایتاً از نتایج ۵ بار اجرا میانگین گرفته می‌شود.

در آزمایش اول، دسته‌بندی‌کننده KNN بهبود یافته بر روی مجموعه داده اعمال شده است، که عملکرد دسته‌بندی‌کننده k نزدیک‌ترین همسایه بهبود یافته در مقایسه با حالت بدون تغییر الگوریتم بهتر بوده است. در آزمایش دوم، از طول N تایی‌های متفاوت در عملیات پیش پردازش استفاده شده و سپس دسته‌بندی‌کننده نوین و بهینه شده اعمال شده است. در این حالت بهترین عملکرد با طول N تایی ۸ بوده است. در آزمایش سوم، حذف ایست‌واژه‌ها و کلمات با فراوانی سند کم در دست اقدام قرار گرفته که در این آزمایش عدد مورد نظر برای پارامتر فراوانی سند کم برابر با ۱ است، نتیجه آزمایش نشان می‌دهد که حذف کلمات با فراوانی سند ۱، موجب تاثیر منفی در این روش می‌شود. در آزمایش چهارم، یک مجموعه داده جدید حاوی ۱۰ دسته از موضوعات مختلف جهت مقایسه دو روش

دسته‌بندی اسناد فارسی، یعنی دسته‌بندی با استفاده از الگوریتم k نزدیک‌ترین همسایه بهبود یافته (روش مورد بحث) و دسته‌بندی با استفاده از الگوریتم ماشین بردار پشتیبان با در نظر گرفتن گنجینه واژگان، روش مطرح شده در [۲۷]، به کار رفته است. نتیجه این آزمایش بیانگر نتایج بهتر در مورد روش اول است.

در یک مطالعه دیگر [۲۸]، مجموعه داده همشهری مورد بررسی قرار گرفته که اسناد مورد بررسی شامل ۱۶۰۰ سند آموزش و ۴۰۰ سند آزمون می‌باشد. از آن جایی که روش K نزدیک‌ترین همسایه به تنهایی نتایج مطلوبی را ارائه نمی‌دهد، اعمال روش‌ها و الگوریتم‌های دیگر به KNN به منظور بالا بردن دقت و کارایی الگوریتم در دستور اجرا قرار می‌گیرند. در این تحلیل، از شبکه واژگان استفاده شده است تا سندهای قرار گرفته در یک دسته شباهت بیشتری را به هم داشته باشند. با استفاده از این شبکه (واژگان)، فراوانی کلماتی که با هم در ارتباطند از دقت بیشتری برخوردار می‌شود.

در شبکه واژگان، نحوه نمایش بردارهای اسناد به صورت مجموعه‌ای از کلمات است. فرآیند آغازین پس از استخراج ویژگی، ریشه‌یابی کلمات است. بعضی از کلمات به اشکال دستوری متفاوتی استفاده می‌شوند اما کماکان یک معنی را دارند. عمل ریشه‌یابی می‌تواند به حل مشکل چندشکلی بودن کمک کند. همچنین کلمات پر تکراری مانند "در" و "بر" مناسب برای تشخیص دسته اسناد نیستند، بنابراین باید حذف شوند. اما بعضی از این کلمات پیشوند کلمات دیگر هستند و می‌توانند در معنای کلمه تاثیرگذار باشند و حتی باعث تغییر در دسته اسناد شوند. بنابراین به جای قرار دادن آن‌ها در لیست ایست‌واژه‌ها و حذف مستقیم آن‌ها، از روش بهره اطلاعات استفاده می‌شود و سپس کلمات با فراوانی بالا و بهره اطلاعات پایین حذف می‌شوند. این عملیات این اطمینان را می‌دهد که کلمات با اثرگذاری کم در دسته‌بندی حذف خواهند شد. شبکه واژگان [۲۹، ۳۰]، یک پایگاه داده لغوی از هر زبان است.

هر کلمه موجود در سند به شبکه واژگان داده می‌شود و فراوانی آن محاسبه می‌شود.

۳-۲ N تایی آماری

دسته‌بندی به روش N تایی، بر روی مجموعه داده همشهری صورت گرفته است [۳۱]. سودمندی استفاده از این روش این است که در مقایسه با روش‌های یادگیری دسته‌بندی سنتی در حین انجام دسته‌بندی، ویژگی‌های با فراوانی کم را نادیده نمی‌گیرد. روش مورد نظر یک نوع مدل‌سازی آماری زبان (به صورت متنی و نوشتاری) محسوب می‌شود. مدل‌سازی آماری زبان، احتمالات رخ داده‌های دنباله کلمات را در مجموعه آموزش محاسبه می‌کند، سپس با استخراج دنباله کلمات از مجموعه آزمون بیشترین احتمال را برای دسته مربوط به سند مورد پرسش پیش‌بینی می‌کند [۳۱]. در واقع هدف مدل‌سازی، قرار دادن دنباله‌های کلماتی که به واقع رخ داده‌اند در احتمالات بالا و همچنین قرار دادن دنباله‌هایی از کلمات که هرگز رخ نداده‌اند در احتمالات پایین است.

در واقع، اصل استفاده از مدل زبانی N تایی به عنوان دسته‌بندی‌کننده متن بیان این مطلب است که دسته سند مورد پرسش باید به گونه‌ای تعریف شود که سند پرسش بیشترین شباهت را با مدل دسته تعیین شده داشته باشد. بنابراین یک مدل زبانی مجزا برای هر دسته، آموزش داده می‌شود و دسته‌بندی یک سند جدید با ارزیابی میزان درست‌نمایش تحت هر دسته، و سپس انتخاب دسته صورت می‌گیرد. پارامتر n یک فاکتور کلیدی در مدل‌سازی زبانی N تایی به حساب می‌آید. n کوچک اطلاعات کافی را برای یک مدل وابستگی‌های کلمات با دقت بالا ارائه نمی‌دهد و همچنین n خیلی بزرگ نیز مشکلات پراکندگی داده را در آموزش سیستم به وجود می‌آورد. یک چند جمله‌ای خالص بیز ساده، نام یک حالت خاص از دسته‌بندی‌کننده (N تایی) است که در آن $n=1$ است و در این حالت از هموارکننده افزودن یک^{۷۲} استفاده می‌شود.

منظور از هموارکننده، تکنیک‌هایی هستند که برای تنظیم بیشینه درست‌نمایی تخمین زده شده از میان احتمالات به منظور تولید احتمالات دقیق‌تر، به کار می‌روند. تکنیک‌های هموارسازی نه تنها احتمالات صفر را از بین می‌برند بلکه در تلاشند تا دقت مدل را نیز بهبود ببخشند [۳۲]. چند نمونه از مهم‌ترین روش‌های هموارسازی از این قرار است [۳۳]:

[۳۴]: هموارکننده افزودن یک برای جلوگیری از احتمالات صفر وانمود می‌کند که رخ دادن هر N تایی ناچیزتر از آن چیزی است که در واقعیت هست، این روش به هر شمارش، 1 را اضافه می‌کند. ایده هموار کننده تخفیف قطعی^{۷۳} این است که، احتمالات کلمات دیده شده را با تفریق کردن یک مقدار ثابت از تعدادشان می‌کاهد. رفتار هموارکننده عقب‌گرد^{۷۴} بدین گونه است که به دنبال کوچک‌تر شدن N تایی‌هاست، به این معنا که در مدل‌سازی زبانی سه‌تایی، تنها از احتمالات سه‌تایی استفاده نمی‌شود بلکه احتمالات دو تایی و یک تایی نیز در نظر گرفته می‌شود، اما در هموارکننده مذکور اگر سه‌تایی پیدا نشد به سراغ دو تایی می‌رود و اگر دو تایی هم پیدا نشد به سراغ یک تایی می‌رود. طبق نتایج حاصل شده، هموارکننده عقب‌گرد بهترین دقت را در مجموعه داده مورد بررسی و در مدل با $n=3$ داشته است. با توجه به نتایج بررسی شده و در نظر گرفتن n متعادل برای مقابله با مشکلات پراکندگی داده، مدل‌سازی زبانی سه‌تایی بهترین دقت را در دسته‌بندی اسناد فارسی ارائه می‌دهد و حتی این نتایج با اعمال روش هموارسازی عقب‌گرد بهتر نیز خواهند شد.

۳-۳ شبکه عصبی مصنوعی

فرآیند آموزش دسته‌بندی به این روش، بر روی ۱۰۵۰ سند خبری از مجموعه اسناد فارسی همشهری ۲ انجام گرفته است که روش دسته‌بندی یادگیری چندی‌سازی برداری، یک روش دسته‌بندی الگویی است که بر اساس شبکه‌های نگاشت خودسازمان‌ده می‌باشد [۳۵، ۳۶]. در این روش فضای ورودی به چند ناحیه منحصر به فرد

73- Absolute Discounting
74- Back-Off

تقسیم می‌شود سپس برای هر ناحیه یک بردار ساخته می‌شود که نمایانگر یک دسته خاص تعریف شده است. بردار وزنی برای یک واحد خروجی، معمولاً با یک بردار کتاب کد منتسب می‌شود که نمایانگر دسته‌ای است که واحد نمایش می‌دهد. طی فرآیند آموزش، محل قرارگیری واحدهای خروجی با تنظیم کردن وزن‌هایشان از میان یک آموزش با نظارت با مبنای دسته‌بندی‌کننده بیز، صورت می‌گیرد. روش دسته‌بندی جلو رو به روی مجموعه داده فارسی تحت بررسی [۳۷] اعمال شده که در پی آن نتایج خوبی گزارش شده است.

۳-۴ بیز ساده

بیز ساده یکی از روش‌های آماری برای دسته‌بندی به شمار می‌آید، بدین گونه که دسته‌های مختلف، هر کدام به عنوان یک فرضیه دارای احتمال در نظر گرفته می‌شوند. استفاده از تکنیک‌های کمکی و روش‌های انتخاب ویژگی مختلف در کنار بیز ساده، قالب آزمون‌های صورت گرفته برای این روش از دسته‌بندی را (در این پژوهش) تشکیل می‌دهند. یکی از این تکنیک‌ها احتساب بردار نماینده است. در این روش از مجموعه داده همشهری شامل ۱۶۰۰۰۰ خبر فارسی بین سال‌های ۱۹۹۷ تا ۲۰۰۲ میلادی استفاده شده است. تاثیر بردار نماینده، بردار متشکل از کلمات مرتبط با هم همراه با درجه ارتباطشان، بر بهبود کیفیت دسته‌بندی به روی اسناد فارسی در این بررسی [۳۸] مورد مطالعه قرار گرفته است. در این فرآیند، بردار نماینده حکم یک پایه و اساس برای ارزیابی بهتر روش بیز ساده را دارد.

تکنیک دیگر استفاده از روش‌های انتخاب ویژگی فراوانی سند (DF)، واریانس فراوانی عبارت (TFV)، اطلاعات متقابل (MI)، اطلاعات متقابل تغییر یافته (MMI) است. این بررسی [۳۹] که به روی مجموعه داده‌ای متشکل از ۸۲۹ نظردهی فارسی صورت گرفته، تنها مطالعه‌ای است که با این حجم کم داده (فقط حاوی دو دسته، نظرات مثبت و نظرات منفی) صورت گرفته است. این روش‌ها عملکردهای بهتری نسبت

به روش MI سنتی دارند. دلیل عملکرد ضعیف MI این است که این روش فقط از اطلاعات بین ویژگی و دسته متناظر با هم (جاری) استفاده می‌کند و با بقیه ویژگی‌ها و دسته‌های دیگر کاری ندارد. وقتی سه روش TFV، DF و MMI را با هم مقایسه می‌کنیم، مشاهده می‌کنیم که MMI در رقابت بر سر دسته منفی، با بهبودهای ۱/۲ و ۱/۷۶ درصدی بر آن دو روش دیگر برتر است اما، وقتی که به سراغ نتایج F دسته مثبت می‌رویم دو روش DF و TFV برتری‌های ۰/۰۴ و ۰/۳ درصدی را نسبت به عملکرد روش MMI دارا هستند. از طرفی مزیت روش MI این است که از تمام اطلاعات وابسته به یک ویژگی استفاده می‌کند و همچنین تمام فاکتورهای بین ویژگی‌ها و دسته‌ها را در نظر می‌گیرد نه بخشی از آن‌ها را. روش MMI علاوه بر دسته‌بندی نظرات فارسی (که حاوی دو دسته بودند) می‌تواند به دیگر حوزه‌های دسته‌بندی اسناد (با دسته‌های بیشتر) نیز کمک کند.

در [۲۴] با محوریت دسته‌بندی متون کوتاه فارسی (متون توثیت‌ها)، از ویژگی TF-IDF و چهار روش بیز ساده، نزدیک‌ترین همسایه، درخت تصمیم و SVM استفاده شده است و نتیجه‌گیری شده است که روش‌های بیز ساده و نزدیک‌ترین همسایه به ترتیب بهترین و ضعیف‌ترین روش برای دسته‌بندی ۶ دسته از متون بوده است.

۳-۵ منطق فازی

دسته‌بندی اسناد فارسی به روش منطق فازی در [۴۰] بررسی شده است که به روی ۹۰۰ سند استخراج شده از وبگاه‌های همشهری، ویکی‌پدیا، دیگر وبگاه‌های فارسی، انجام گرفته است. دو سوم این مجموعه به‌عنوان مجموعه داده آموزش و مابقی تحت عنوان مجموعه داده آزمون مورد استفاده قرار گرفته است. در این روش دسته‌بندی روابط فازی که مبنای دسته‌بندی هستند از نظریه فازی مشتق می‌شوند.

با هدف دسته‌بندی اسناد، با توجه به ویژگی‌های خاص زبان فارسی می‌توان از مدل فازی استفاده کرد.

ایده اصلی این روش، از پیوندهای بین دسته‌های از پیش تعریف شده‌شان مشتق می‌شود، که این پیوندها از تابع عضویت فازی به‌دست می‌آیند. مقدار تابع مذکور بین صفر و یک است که مقدار صفر مربوط به حالتی است که عنصر انتخاب شده عضو هیچ دسته‌ای نیست و مقدار یک بیانگر این حالت است که عنصر انتخاب شده کاملاً عضو یک دسته است، همچنین مقادیر بین صفر و یک نیز معرف عضویت و درجه شباهت عنصر با دسته مربوطه است [۴۰].

یکی از روش‌های تعیین کننده مقادیر تابع عضویت استفاده از عوامل هوشمند، بر اساس قوانین فازی (If ... Then) است. برای تعریف این عوامل استاندارد خاصی وجود ندارد بلکه عموماً اندیشه‌هایی که به ارتباطات مهم می‌پردازند در پیش‌زمینه این عوامل قرار می‌گیرند. عوامل هوشمند اعمال شده در این بررسی، به منظور تعیین میزان درجه عضویت عبارتهای اسناد آموزش به دسته‌ها استفاده می‌شوند. در پی آن، برای دریافت پیوندهای فازی از میان کلمات، عوامل هوشمند سازوکارهای شناسایی دسته اسناد مورد پرسش را فراهم می‌کنند [۴۰].

۳-۶ ماشین بردار پشتیبان

روش ماشین بردار پشتیبان، یکی از بهترین و محبوب‌ترین روش‌های دسته‌بندی متون به حساب می‌آید. استفاده از این روش در مسائل دسته‌بندی روند جدیدی است که در سالیان اخیر مورد توجه بسیاری قرار گرفته است. همچنین توابع (هسته) گوناگونی که می‌توانند در دل محاسبات این روش قرار بگیرند، نتایج جالبی را در شرایط متفاوت از خود برجای می‌گذارند. یکی از دلایل محبوبیت این روش، انعطاف‌پذیری بالای آن در شرایط متفاوتی از همکاری با دیگر الگوریتم‌هاست. روش انتخاب ویژگی، استفاده از توابع هسته گوناگون، استفاده از گنجینه واژگان و در نظر گرفتن روابط معنایی و ... همگی شرایطی هستند که در آزمون‌های دسته‌بندی اسناد فارسی به روش ماشین بردار پشتیبان پیش رو مورد بررسی قرار گرفته‌اند.

روش ماشین بردار پشتیبان، به روی مجموعه داده فارسی همشهری که جزو مجموعه داده‌های معتبر اسناد فارسی محسوب می‌شود، نتایج خوبی را در بر داشته است [۲۵]. از روش انتخاب ویژگی TF-IDF برای انتخاب برترین ویژگی‌های هر دسته استفاده شده است و همچنین از ۷۰ درصد مجموعه داده همشهری به عنوان مجموعه آموزش و مابقی که ۳۰ درصد می‌باشد به عنوان مجموعه داده آزمون استفاده گردیده است.

عامل دیگر تاثیرگذار بر کیفیت عملکرد ماشین بردار پشتیبان تابع هسته می‌باشد [۲۵]. انتخاب یک تابع هسته مناسب از اهمیت بالایی برخوردار است، زیرا تابع هسته تعیین کننده فضای ویژگی است که نمونه‌های مجموعه داده آموزش دسته‌بندی می‌شوند.

دسته‌بندی اسناد فارسی با استفاده از گنجینه واژگان و روش SVM، به هدف بررسی تاثیر آن در نتایج دسته‌بندی صورت گرفته است [۲۷]. به کارگیری و عدم به کارگیری گنجینه واژگان، موردهای مورد آزمایش هستند که در پی آن گزارش‌هایی حاکی از بهبود نتایج در شرایط حضور گنجینه واژگان ارائه شده‌اند. لازم به ذکر است مجموعه داده مورد استفاده متشکل از اسناد موجود در وبگاه ویکی‌پدیا فارسی، آرشیو روزنامه همشهری، و نشریات سروش و رشد [۴۱] می‌باشند.

گسترش بردار ویژگی با اضافه کردن کلمات استخراجی از گنجینه واژگان، عملی است که هدفش کمک به دسته‌بندی‌کننده در زمان‌هایی است که مجموعه داده آموزش کامل و مناسب همه دسته‌ها نیست. فرهنگ طایفی [۴۲] نام گنجینه واژگانی است که در این بررسی استفاده می‌شود. با توجه به ساختار گنجینه واژگان، گنجینه برای هر دسته به هدف یافتن کلمات مرتبطی که نمایندگان بهتری برای آن دسته هستند جستجو می‌شود. انتخاب یک کلمه از گنجینه مبتنی بر همان ویژگی‌هایی است که از فاز انتخاب ویژگی اولیه حاصل شده‌اند به این ترتیب که، هر کلمه انتخاب شده (از فاز انتخاب ویژگی) در گنجینه

جستجو می‌شود، همه کلمات مرتبط (در گنجینه) با کلمه انتخاب شده، نامزد اضافه شدن به بردار ویژگی می‌شوند. سرانجام اگر فراوانی کلمه نامزد شده در اسناد متعلق به دسته مورد نظر در مجموعه داده آموزش، از یک آستانه از قبل تعیین شده بزرگ‌تر باشند و کلمه قبلاً در بردار ویژگی موجود نباشد، به بردار ویژگی اضافه خواهد شد. نکته قابل ذکر این است که اگر آستانه مقداری کوچک در نظر گرفته شود کلمات غیر مرتبط در بردار ویژگی به تناسب زیاد خواهند شد. همچنین انتخاب یک مقدار آستانه بزرگ، کلمات مرتبط بیشتری را به بردار ویژگی اضافه خواهد کرد. مقدار آستانه نقش مهمی را در اضافه نشدن کلمات غیر مرتبط به بردار ویژگی ایفا می‌کند.

در بررسی [۴۳]، اسناد متشکل از ۳ مجموعه داده همشهری نسخه ۲ [۴۴]، پرسیکا [۴۵] و تابناک با برچسب دسته‌های مشترک به عنوان مجموعه داده آموزش و سه مجموعه داده متفاوت D1، D2، و D3 با ۶ دسته مشترک که به ترتیب حاوی ۵۸۵، ۹۷۶ و ۱۲۲۵ سند به تصادف انتخاب شده می‌باشد به عنوان مجموعه داده‌های آزمون استفاده شده است. در دسته‌بندی اسناد فارسی به روش ارتباطات معنایی بین کلمات، آن چه که باعث تغییرات مهم و کارایی بهتر نسبت به روش‌های دیگر آماری می‌شود، نحوه وزندهی عبارتهاست. روش وزندهی عبارتها بر اساس روابط معنایی بین کلمات برتری است که روش‌های دیگر وزندهی از به کارگیری آن معذورند و به جای آن از اطلاعات آماری استفاده می‌کنند. این روش ارتباط بین عبارتها را اندازه‌ای از وابستگی در نظر می‌گیرد، همچنین برای تعیین معناهای (ارتباطات معنایی) دسته‌ها بر اساس قرارگیری عبارتها در (برچسب‌های) دسته‌ها، بردار ویژگی هر دسته را با استفاده از گنجینه واژگان گسترش می‌دهد.

سیستم دسته‌بندی اسناد فارسی مذکور (مورد بررسی) با سیستم‌های طراحی شده دیگر برای اسناد فارسی مورد مقایسه قرار گرفته‌اند. برای این منظور بهترین نتایج سه سیستم که شبیه سیستم مورد بررسی بوده‌اند، مورد

ساده و SVM)، به دسته‌بندی اسناد پرداخته است. این روش با به‌کارگیری مجموعه کلمات اسمی به‌عنوان، کلمات کلیدی و استفاده از دسته‌بندی‌کننده SVM بهترین مقدار F را برای نه دسته، ۸۸ درصد گزارش داده است.

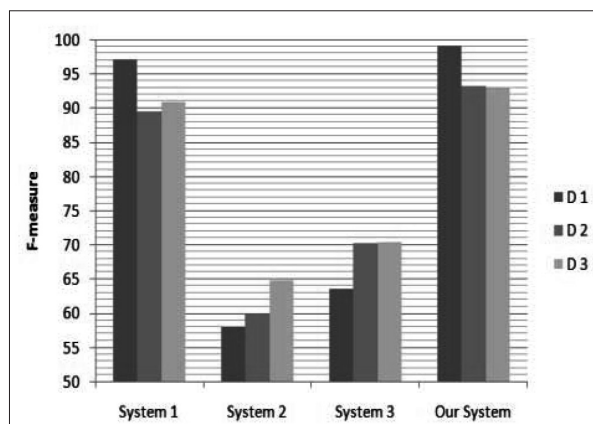
در یکی از کارهای اخیر در [۴۸] از روش‌های مدل‌سازی عنوان برای استخراج ویژگی استفاده شده و روش SVM برای دسته‌بندی متون فارسی روی پیکره بیجن‌خان به کار گرفته شده است که در آن روش‌های مختلف مدل‌سازی عنوان شامل LDA، PLSA و STC با روش TF-IDF روی تعداد مختلف دسته (عنوان) بررسی شده و نشان داده شده است که روش‌های مدل‌سازی عنوان به طور چشمگیری کارایی بالاتری را نتیجه می‌دهند.

۳-۷ مدل مخفی مارکوف (HMM)

این روش مبتنی بر مدل مخفی مارکوف که جزو روش‌های آماری است، به مدل‌سازی توالی دنباله کلمات یک سند می‌پردازد. در بررسی [۲۳]، که به روی مجموعه داده فارسی بی‌جن‌خان انجام شده است، اسناد مورد دسته‌بندی شامل هشت دسته ادبی، مذهبی، اقتصادی، هنری، پزشکی، تاریخی، سیاسی و ورزشی می‌باشند. در این روش، هر سند به صورت رشته‌ای از کلمات متوالی در نظر گرفته می‌شود، همچنین از قابلیت این روش در مدل‌سازی داده‌های ترتیبی، بهره گرفته می‌شود. مدل مخفی مارکوف در میان سایر دسته‌بندی‌کننده‌ها (شبکه عصبی، بیز ساده، شبکه بیز، ماشین بردار پشتیبان، درخت تصمیم و k نزدیک‌ترین همسایه)، بیشترین مقدار میانگین دقت (۸۳ درصد)، فراخوانی (۷۶ درصد) و F (۷۹ درصد) را به روی هشت دسته، گزارش داده است. پس از مدل مخفی مارکوف، شبکه عصبی است که برترین نتایج را ارائه داده است.

خلاصه و جمع‌بندی

در این بخش سعی شده که نتایج همه آزمایش‌های دسته‌بندی اسناد فارسی که در بخش قبل بیان شدند در جدول ۱ گردآوری شود. در هر سطر از جدول به ترتیب



شکل ۴: مقایسه سیستم‌های دسته‌بندی اسناد فارسی [۴۳]

مطالعه قرار گرفتند. نیاز به توضیح است که سه سیستم با شرایط یکسان و با آزمایش به روی هر سه مجموعه داده D1، D2، D3 پیاده‌سازی شده‌اند (شکل ۴). این سیستم‌ها عبارتند از:

سیستم ۱: دسته‌بندی اسناد فارسی با روش انتخاب ویژگی DF و استفاده از SVM و گنجینه واژگان توسط [۲۷].
سیستم ۲: دسته‌بندی اسناد فارسی با روش انتخاب ویژگی IG و استفاده از SVM (به جای KNN) و گنجینه واژگان توسط [۲۸].

سیستم ۳: دسته‌بندی اسناد فارسی با نمایش عددی بردارها به روش N تایی و احتساب ضریب دسته با استفاده از KNN (معیار شباهت دایس) توسط [۲۶].

با این توضیح که سیستم ۱، فقط از گنجینه واژگان به منظور گسترش بردارهای ویژگی و TF-IDF به هدف وزن‌دهی استفاده کرده است و اطلاعاتی مانند ارتباطات معنایی را در نظر نگرفته است، همچنین سیستم‌های ۲ و ۳ از اطلاعات آماری برای وزن‌دهی استفاده کرده‌اند.

در بررسی [۴۶] که به روی مجموعه داده فارسی بی‌جن‌خان [۴۷] انجام شده است، به منظور نمایش برداری کلمات کلیدی اسناد، از روش TF-IDF استفاده شده است. در گام بعدی، به وسیله این بردارها و با استفاده از الگوریتم‌های تحلیل آماری معنایی پنهان، به بردارهای ویژگی نمایش‌دهنده اسناد رسیده است. در ادامه با دسته‌بندی‌کننده‌های گوناگون (k نزدیک‌ترین همسایه، بیز

نوع دسته‌بندی‌کننده، جزئیات دسته‌بندی (روش انتخاب ویژگی، شرایط، پارامترهای الگوریتم، الگوریتم‌های مکمل و کمکی) و بهترین نتیجه حاصل از ارزیابی دسته‌بندی آورده شده است. لازم به ذکر است که هر خط افقی، یک جدا کننده بین دو مجموعه داده متفاوت می‌باشد به این معنا که نتایج قرار گرفته بین دو خط افقی حاصل از اعمال آزمایش‌های دسته‌بندی به روی یک مجموعه داده مشترک است. بنابراین مقایسه نتایج قرار گرفته بین خطوط افقی از اعتبار بیشتری (نسبت به مقایسه به روی کل نتایج جدول) برخوردار است و دلیل بدیهی این امر نیز، یکسان بودن شرایط استفاده از مجموعه داده در آزمایش‌های دسته‌بندی (قرار گرفته بین خطوط افقی) است.

زمانی که طبیعت یک مجموعه داده با فرضیات و تدابیر مورد نظر یک دسته‌بندی تطبیق داده شود، دقت آن دسته‌بندی افزایش پیدا می‌کند. اما عدم تطبیق و ضعف تناسب یک مدل نیز می‌تواند ناشی از انتخاب نوع ویژگی‌های استفاده شده در مدل، انتخاب روش‌های امتیازدهی، انتخاب توابع شباهت و خیلی از عوامل درگیر دیگر باشد. حال با یک دید کلی می‌توان به بحث در مورد دسته‌بندی‌کننده‌ها پرداخت.

دسته‌بندی‌کننده بیز ساده، یک دسته‌بندی‌کننده خیلی ساده است که خیلی خوب روی داده‌های متنی و عددی عمل می‌کند، ساده پیاده‌سازی می‌شود و از نظر بار محاسباتی بسیار ارزان‌تر از هر الگوریتم دسته‌بندی دیگر است. یکی از محدودیت‌های این دسته‌بندی‌کننده این است که زمانی که ویژگی‌ها با هم ارتباط بسیار قوی دارند، ضعیف عمل می‌کند. همچنین با توجه به قواعد دسته‌بندی اسناد، این نوع دسته‌بندی‌کننده ناتوان از شمارش تعداد رخداد کلمات در یک بردار ویژگی است و اما در مقابل، دسته‌بندی‌کننده k نزدیک‌ترین همسایه، مشخصات محلی یک سند را حفظ می‌کند ولی زمان زیادی را بابت دسته‌بندی می‌طلبد و همچنین پیدا کردن مقدار بهینه k همواره یک مشکل جدی است چنان‌که اگر نقاط، توزیع یکنواختی را نداشته باشند این عمل دشوارتر نیز خواهد شد. همان‌گونه که انتخاب

یک مقدار k مناسب در کیفیت دسته‌بندی تاثیرگذار است، عملکرد این نوع دسته‌بندی‌کننده به معیار فاصله اعمال شده نیز وابسته است. عموماً انتخاب یک k کوچک، در جایی که اسناد اشتراکات (روی هم‌افتادگی) زیادی را با هم دارند و به دسته‌های مختلف تعلق دارند بهتر نتیجه می‌دهد چرا که، دسته‌های فشرده‌تری را حاصل می‌نماید. در حالی که مقدار k های بزرگ‌تر در مقابل نوفه‌ها ایمن‌تر و مقاوم‌تر است و مرزهای متعادل‌تری را بین دسته‌ها قائل می‌گردد. اگرچه دسته‌بندی‌کننده مرکز ثقل ساده و خطی است، اما اغلب مدل آن از عدم تطبیق داده‌ها با فرضیاتش رنج می‌برد. این دسته‌بندی‌کننده همیشه با این فرض اقدام به دسته‌بندی می‌کند که اگر شباهت یک سند با مرکز ثقل یک دسته نسبت به دسته‌های دیگر بیشترین مقدار را داشته باشد، آن سند متعلق به همان دسته مورد نظر (با بیشینه شباهت) است. قطعاً این فرضیه در موارد خاص و استثناً همیشه خوب جواب نمی‌دهد. یک دسته‌بندی‌کننده مبتنی بر قانون مانند درخت تصمیم، به راحتی قابل فهم است و پیچیدگی مسئله را می‌کاهد اما زمان زیادی را بابت آموزش اسناد مصرف می‌کند. همچنین با توجه به ساختار درختی شکل آن، در سطوح بالاتر یک زیردرخت اشتباهاتی رخ می‌دهد که گاهی قابل چشم‌پوشی نمی‌باشد. در اداره متغیرهای پیوسته ناتوان است و از مشکلات عدم تناسب نیز رنج می‌برد. دسته‌بندی‌کننده‌های بر پایه و اساس شبکه عصبی، نتایج بسیار خوبی را در دامنه‌های پیچیده ارائه می‌دهند و همچنین قادر به اداره داده‌های پیوسته و گسسته می‌باشد. از معایب آن می‌توان زمان زیاد آموزش و دشوار بودن تفسیر نتایج یادگیری توسط کاربران را نام برد. در سال‌های اخیر استفاده همه‌گیر از روش‌های یادگیری عمیق ۷۵ و کارایی بالاتر آن‌ها لزوم توجه و به کارگیری این روش‌ها را برای دسته‌بندی متون فارسی پررنگ‌تر می‌کند، به ویژه پیش‌بینی می‌شود استفاده از شبکه‌های بازگشتی مانند LSTM در این موضوع منجر به بهبود کارایی این سیستم‌ها شود.

جدول ۱: خلاصه‌ای از پژوهش‌های صورت گرفته در دسته‌بندی اسناد فارسی

| مرجع | دسته‌بندی‌کننده | جزئیات (روش انتخاب ویژگی، شرایط، پارامترهای الگوریتم، الگوریتم‌های مکمل و کمکی) | | | | نتیجه ارزیابی |
|------|-----------------|---|----------------------|------------|---------------------------|----------------|
| [۲۵] | KNN | انتخاب ویژگی = TFIDF | تعداد ویژگی = ۴۰۰۰ | K=3 | معیار فاصله = Cosine | ریز P=97.3 |
| [۲۵] | KNN | انتخاب ویژگی = TFIDF | تعداد ویژگی = ۱۰۰۰ | K=1 | معیار فاصله = Cosine | ریز P=98.8 |
| [۲۵] | KNN | انتخاب ویژگی = TFIDF | تعداد ویژگی = ۱۰۰۰ | K=3 | معیار فاصله = Euclidean | ریز P=67 |
| [۲۵] | KNN | انتخاب ویژگی = TFIDF | تعداد ویژگی = ۱۰۰۰ | K=3 | معیار فاصله = Cityblock | ریز P=45.2 |
| [۲۵] | KNN | انتخاب ویژگی = TFIDF | تعداد ویژگی = ۱۰۰۰ | K=3 | معیار فاصله = Correlation | ریز P=94.3 |
| [۲۶] | KNN | N تایی نمایش عددی بردار = | احتساب ضریب دسته | K=3 | Dice = معیار فاصله | ریز F=90.3 |
| [۲۶] | KNN | N تایی نمایش عددی بردار = | احتساب ضریب دسته | K=8 | Dice = معیار فاصله | ریز F=93.9 |
| [۲۸] | KNN | انتخاب ویژگی = IG | احتساب گنجینه واژگان | K=10 | - | میانگین F=91.2 |
| [۲۶] | KNN | N تایی نمایش عددی بردار = | احتساب ضریب دسته | K=8 | Dice = معیار فاصله | ریز F=91 |
| [۲۷] | SVM | انتخاب ویژگی = DF | احتساب گنجینه واژگان | - | - | ریز F=89 |
| [۳۱] | N تایی آماری | تعداد نمونه در train=9000 | هموار کننده Back-Off | N=3 | - | Etp&Ppx=98 |
| [۳۷] | LVQ 1 | - | - | $\alpha=1$ | - | درشت F=80.3 |
| [۳۷] | OLVQ 1 | - | - | $\alpha=1$ | - | درشت F=84.4 |
| [۳۷] | LVQ 2.1 | - | - | $\alpha=1$ | - | درشت F=88.4 |
| [۳۷] | LVQ 3 | - | - | $\alpha=1$ | - | درشت F=88.8 |
| [۳۷] | OLVQ 3 | - | - | $\alpha=1$ | - | درشت F=89.9 |
| [۳۷] | SVM | - | - | - | - | درشت F=79.8 |
| [۴۹] | KNN | - | - | - | - | درشت F=68.6 |
| [۳۹] | بیز ساده | انتخاب ویژگی = MI | - | - | - | درشت F=57.5 |
| [۳۹] | بیز ساده | انتخاب ویژگی = DF | - | - | - | درشت F=83 |
| [۳۹] | بیز ساده | انتخاب ویژگی = TFV | - | - | - | درشت F=83.5 |
| [۳۹] | بیز ساده | انتخاب ویژگی = MMI | - | - | - | درشت F=84 |
| [۳۸] | بیز ساده | انتخاب ویژگی = MI | احتساب بردار نماینده | - | - | درشت F=78.4 |
| [۳۸] | بیز ساده | انتخاب ویژگی = MI | - | - | - | درشت F=75.3 |
| [۴۰] | قوانین فازی | - | - | - | - | درشت F=91 |
| [۲۵] | SVM | انتخاب ویژگی = TFIDF | تعداد ویژگی = ۱۰۰۰ | - | تابع هسته = Polynomial | ریز F=93.8 |
| [۲۵] | SVM | انتخاب ویژگی = TFIDF | تعداد ویژگی = ۱۰۰۰ | - | تابع هسته = Linear | ریز F=82.2 |
| [۲۵] | SVM | انتخاب ویژگی = TFIDF | تعداد ویژگی = ۱۰۰۰ | - | تابع هسته = RBF | ریز F=93.1 |
| [۲۵] | SVM | انتخاب ویژگی = TFIDF | تعداد ویژگی = ۱۰۰۰ | - | تابع هسته = MLP | ریز F=82.8 |
| [۲۵] | SVM | انتخاب ویژگی = TFIDF | تعداد ویژگی = ۱۰۰۰ | - | تابع هسته = Quadtrac | ریز F=83.8 |

| | | | | | | |
|------|------------|-------------------------|----------------------|------|-------------------------|----------------|
| [۲۷] | SVM | انتخاب ویژگی=DF | - | - | - | ریز F=90 |
| [۲۷] | SVM | انتخاب ویژگی=DF | احتساب گنجینه واژگان | - | - | ریز F=94 |
| [۲۷] | SVM | انتخاب ویژگی=DF | احتساب گنجینه واژگان | - | احتساب وزن معنایی | میانگین F=99 |
| [۲۷] | SVM | انتخاب ویژگی=DF | احتساب گنجینه واژگان | - | - | میانگین F=97.1 |
| [۲۷] | SVM | انتخاب ویژگی=DF | - | - | - | میانگین F=71.3 |
| [۲۸] | SVM | انتخاب ویژگی=IG | احتساب گنجینه واژگان | - | - | میانگین F=58 |
| [۲۶] | KNN | نمایش عددی بردار N=تایی | احتساب ضریب دسته | K=8 | معیار فاصله = Dice | میانگین F=64 |
| [۲۷] | SVM | انتخاب ویژگی=DF | احتساب گنجینه واژگان | - | احتساب وزن معنایی | میانگین F=93.1 |
| [۲۷] | SVM | انتخاب ویژگی=DF | احتساب گنجینه واژگان | - | - | میانگین F=89.5 |
| [۲۷] | SVM | انتخاب ویژگی=DF | - | - | - | میانگین F=78 |
| [۲۸] | SVM | انتخاب ویژگی=IG | احتساب گنجینه واژگان | - | - | میانگین F=60 |
| [۲۶] | KNN | نمایش عددی بردار N=تایی | احتساب ضریب دسته | K=8 | معیار فاصله = Dice | میانگین F=70 |
| [۲۷] | SVM | انتخاب ویژگی=DF | احتساب گنجینه واژگان | - | احتساب وزن معنایی | میانگین F=92.8 |
| [۲۷] | SVM | انتخاب ویژگی=DF | احتساب گنجینه واژگان | - | - | میانگین F=90.8 |
| [۲۷] | SVM | انتخاب ویژگی=DF | - | - | - | میانگین F=78.6 |
| [۲۸] | SVM | انتخاب ویژگی=IG | احتساب گنجینه واژگان | - | - | میانگین F=65 |
| [۲۶] | KNN | نمایش عددی بردار N=تایی | احتساب ضریب دسته | K=8 | معیار فاصله = Dice | میانگین F=70 |
| [۴۶] | SVM | PLSA & TFIDF | - | - | - | میانگین F=88 |
| [۲۳] | HMM | هر عنوان ۵ سند | States ۹ | - | - | میانگین F=79 |
| [۲۳] | NN | MLP | 2900 in, 46 hi, 8 ou | - | HID-OUT=sigmoid | میانگین F=73 |
| [۲۳] | SVM (SMO1) | - | - | - | تابع هسته=Polynomial | میانگین F=71 |
| [۲۳] | KNN | - | - | K=50 | معیار فاصله = Euclidean | میانگین F=67 |
| [۵۰] | بیز ساده | TFIDF | - | - | - | درستی=74 |
| [۵۰] | SVM | TFIDF | - | - | - | درستی=67 |
| [۵۰] | درخت تصمیم | TFIDF | - | C4.5 | - | درستی=66 |
| [۵۰] | KNN | TFIDF | - | K=3 | - | درستی=46 |
| [۴۸] | SVM | PLSA | - | - | - | درستی=67 |
| [۴۸] | SVM | LDA | - | - | - | درستی=68 |
| [۴۸] | SVM | MedSTC ² | - | - | - | درستی=87 |

- formation Retrieval, 50-57.
6. Mikolov, T., Yih, W. T., & Zweig, G., 2013. "Linguistic regularities in continuous space word representations". In Proceedings of the 2013 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies (pp. 746-751).
 7. Pennington, J., Socher, R. and Manning, C., 2014. "Glove: Global vectors for word representation". In Proceedings of the 2014 conference on empirical methods in natural language processing (EMNLP) (pp. 1532-1543).
 8. Jolliffe, I.T., 2002. Principal Component Analysis. Springer.
 9. Cai, L., Hofmann, T., 2003. Text categorization by Boosting Automatically Extracted Concepts. Proceedings of the 26th Annual International ACM SIGIR Conference on Research and Development in Informaion Retrieval, 182-189.
 10. Bengio Y., Ducharme R., Vincent P., and Janvin C., 2003, "A neural probabilistic language model," The Journal of Machine Learning Research, vol. 3, pp. 1137-1155.
 11. Mikolov T., Chen K., Corrado G., and Dean J., 2013, "Efficient estimation of word representations in vector space," In ICLR.
 12. Rong X., 2014, "word2vec Parameter Learning Explained".
 13. Brown P. F., Desouza P. V., Mercer R. L., Pietra V. J. D., and Lai J. C., 1992, "Class-based n-gram models of natural language," Computational linguistics, vol. 18, pp. 467-479.
 14. Bengio Y., Schwenk H., Senécal J.-S., Morin F., and Gauvain J.-L., 2006, "Neural probabilistic language models," in Innovations in Machine Learning, pp. 137-186.
 15. Blei D. M., Ng A. Y., and Jordan M. I., 2003, "Latent dirichlet allocation," the Journal of machine Learning research, vol. 3, pp. 993-1022.
 16. Zhu, J. and Xing, E.P., 2012. Sparse topical coding. arXiv preprint arXiv:1202.3778.
 17. Pennington J., Socher R., and Manning C. D., 2014, "Glove: Global vectors for word representation," Proceedings of the Empirical Methods in Natural Language Processing (EMNLP 2014), vol. 12, pp. 1532-1543.
 18. Yang, Y., Pedersen, J.O., 1997. A Comparative Study on Feature Selection in Text Categorization. Proceedings of the 14th International Conference on Machine Learning 97, 412-420.
 19. Zong, W., Wu, F., Chu, L., Sculli, D., 2015. A Discriminative and Semantic Feature Selection Method for Text Categorization. International Journal of Production Economics, 215-222.
 20. Aghdam, M.H., Ghasem-Aghaee, N., Basiri, M.E., 2009. Text Feature Selection Using Ant Colony Optimization. Expert Syst. Appl. 36 (3), 6843-6853.
 21. Basu, T., Murthy, C.A., 2012. Effective Text Classification by a Supervised Feature Selection Approach. Data Mining Workshops (ICDMW), IEEE 12th International Conference, 918-925.

در قیاس با بقیه دسته‌بندی‌کننده‌ها ماشین بردار پشتیبان، عمل تشخیص مشخصات ذاتی داده را بهتر و دقیق‌تر انجام می‌دهد. این روش دسته‌بندی قابلیت یادگیری مستقل در فضای ویژگی چندبعدی را داراست. پیچیدگی دسته‌بندی‌کننده در تنظیم پارامترها و تعیین هسته محاسبات است. نهایتاً آن چه را که به طور کلی می‌توان نتیجه گرفت این است که دسته‌بندی‌کننده‌های مختلف با توجه به خصوصیات داده به روی مجموعه داده‌های متفاوت، عملکردهای متفاوتی را از خود نشان می‌دهند و هیچ‌گاه نمی‌توان به طور قطع، این ادعا را داشت که یکی از آن‌ها در همه موارد برترین دسته‌بندی‌کننده است، هر چند که ماشین بردار پشتیبان با مدل نمایشی VSM، در اکثر مسائل دسته‌بندی اسناد نتایج خوبی را گزارش می‌دهد.

در مورد روش‌های استخراج ویژگی، می‌توان گفت که در میان روش‌های کلاسیک، TF-IDF به صورت کلی کارایی قابل قبولی داشته است و روش‌های مبتنی بر جاسازی کلمات و مدل‌سازی عنوان به دلیل بهره‌گیری از اطلاعات بیشتر نقش موثری در بهبود کارایی دارند. در میان کارهای مرور شده، تحقیقی که از روش‌های جاسازی کلمات مبتنی بر شبکه عصبی که در سال‌های اخیر بسیار رایج شده، استفاده کند، یافت نشد که برآورد می‌شود این نمایش از متون، بتواند منجر به بهبود کارایی سیستم‌های دسته‌بندی متون شود.

مراجع

1. Aggarwal, C.C., Zhai, C.X., 2012. A Survey of Text Classification Algorithms. Mining Text Data Book, 163-222.
2. Sebastiani, F., 2002. "Machine learning in automated text categorization". ACM computing surveys (CSUR), 34(1), pp.1-47.
3. Salton, G., Yang, C.S., 1999. On the Specification of Term Variants in Automatic Indexing. Journal of Documentation, 29(2), 351-372.
4. Deerwester, S., Dumais, S., Harshman, H., 1990. Indexing by Latent Semantic Analysis. Journal of the American Society for Information Science, 391.
5. Hofmann, T., 1999. Probabilistic Latent Semantic Indexing. Proceedings of the 22nd Annual International ACM SIGIR Conference on Research and Development in In-

- of Persian Textual Documents Using Learning Vector Quantization. In Proceedings of the International Conference on Natural Language Processing and Knowledge Engineering, IEEE, 1-6.
38. Jafari, A., Hosseinejad, M., Amiri, A., 2011. Improvement in Automatic Classification of Persian Documents by Means of Naïve Bayes and Representative Vector. 1st International eConference on Computer and Knowledge Engineering (ICCKE), IEEE, 226-229.
 39. Sarace, M., Bagheri, A., 2013. Feature Selection Methods in Persian Sentiment Analysis. Natural Language Processing and Information Systems, Springer-Verlag Berlin Heidelberg, 7934, 303-308.
 40. Yari, A., Abbasi, A., MomenBellah, S., 2010. Presenting a fuzzy relation to classify the Persian Web documents. Intelligent Computing and Intelligent Systems (ICIS), IEEE International Conference, 2, 220-223.
 41. Bijankhan, M., 2008. 100 Millions Word Farsi Corpus. Technical Report, Research Center for Intelligent Signal Processing.
 42. Fararoy, J., 2008. Farhang-e Teyfi. Tehran: Hermes Press.
 43. Parseh, S., Baraani, A., 2014. Improving Persian Document Classification Using Semantic Relations between Words. International Journal of Advanced Studies in Computers, Science and Engineering, ARXIV, 16-22.
 44. AleAhmad, A., Amiri, H., Darrudi, E., Rahgozar, M., Oroumchian, F., 2009. Hamshahri: A Standard Persian Text Collection. Knowledge-Based Systems, 22, 382-387.
 45. Eghbalzadeh, H., Hosseini, B., Khadivi, S., Khodabakhsh, A., 2012. Persica: A Persian Corpus for Multi-purpose Text Mining and Natural Language Processing. Sixth International Symposium on Telecommunications (IST), IEEE, 1207-1214.
 ۴۶. خاکی اردکانی بصیرا، ۱۳۹۲. به کارگیری مدل‌های موضوع پنهان جهت دسته بندی مستندات فارسی و ارائه راهکارهای مناسب جهت بهبود آن. پایان نامه، دانشگاه صنعتی شریف.
 47. Bijankhan, M., Sheykhzadegan, J., Bahrani, M., Ghayoomi, M., 2011. Lessons from Building a Persian Written Corpus: Peykare. Language Resources and Evaluation, 45, 143-164.
 48. Ahmadi, P., Tabandeh, M. and Gholampour, I., 2016, May. Persian text classification based on topic models. In Electrical Engineering (ICEE), 2016 24th Iranian Conference on (pp. 86-91). IEEE.
 49. Bina, B., Ahmadi, M.H., Rahgozar, M., 2008. Farsi Text Classification Using N-gram and Knn Algorithm A Comparative Study. DMIN, 385-390.
 ۵۰. انامرادنژاد عیسی، حبیبی جعفر. ۱۳۹۶. بررسی کارایی الگوریتم‌های کلاس بندی متن بر روی متون کوتاه فارسی. کنفرانس بین المللی فناوری اطلاعات، مهندسی کامپیوتر و مخابرات.
 22. Chen, J., Pan, H., Ao, Q., 2012. Study a Text Classification Method Based on Neural Network Model. Proceedings of the MSEC International Conference on Multimedia, Software Engineering and Computing, Springer Berlin Heidelberg, 128, 471-475.
 23. Gharavi, E., Veisi, H., 2014. A Hidden Markov Model for Persian Text Classification. 3rd National Computational Linguistics Conference.
 ۲۴. احسان ریحانی آرانی، سیدمحمد رضا لاجوردی. ۱۳۹۵. بررسی روش‌های طبقه بندی خودکار اسناد متنی. سومین کنفرانس ملی مهندسی برق و کامپیوتر سیستم‌های توزیع شده و شبکه‌های هوشمند.
 25. Farhoodi, M., Yari, A., 2010. Applying Machine Learning Algorithms for Automatic Persian Text Classification. 6th International Conference on Advanced Information Management and Service (IMS), 318-323.
 26. Elahimanesh, M.H., Minaei-Bidgoli, B., Malekinezhad, H., 2012. Improving K-Nearest Neighbor Efficacy for Farsi Text Classification. The International Conference on Language Resources and Evaluation (LREC), 1618-1621.
 27. Maghsoodi, N. and Homayounpoor, M., 2011. Using Thesaurus to Improve Multiclass Text Classification. Part II, LNCS 6609, 244-253.
 28. Parchami, M., Akhtar, B., Dezfoulian, M.H., 2012. Persian Text Classification Based on K-NN Using Wordnet. IEA/AIE, LNAI, Springer-Verlag Berlin Heidelberg, 7345, 283-291.
 29. Agirre, E., Alfonseca, E., Hall, K., Kravalova, J., Pasca, M., Soroa, A., 2009. A Study on Similarity and Relatedness Using Distributional and WordNet-based Approaches. NAACL Proceedings of Human Language Technologies, Annual Conference of the North American Chapter of the Association for Computational Linguistics, 19-27.
 30. Fellbaum, C., 1998. WordNet: An Electronic Lexical Database and Some of its Applications. MIT Press, Cambridge.
 31. Farhoodi, M., Yari, A., Sayah, A., 2011. N-gram Based Text Classification for Persian Newspaper Corpus. Digital Content, Multimedia Technology and its Applications (IDCTA), IEEE, 55-59.
 32. Peng, F., Huang, X., 2007. Machine Learning for Asian Language Text Classification. Journal of Documentation, 63(3), 378-397.
 33. Ramasundaram, S., Victor, S.P, 2010, Text Categorization by Backpropagation Network. Intrenational Journal of Computer Applications, 8(6), 1-5.
 34. Chen, S.F., Goodman, J., 1999. An Emperical Study of Smoothing Techniques for Language Modeling. 13(4), 359-394.
 35. Kohonen, T., 1990. Improved versions of Learning Vector Quantization. In proceeding of the National Joint Conference of Neural Networks, San Diego, 545-550.
 36. Kohonen, T., 1995. Self-Organization and Associative Memory. Berlin: Springer-Verlag.
 37. Pilevar, M.T., Feili, H., Soltani, M., 2009. Classification