

بنام خداوند جان و خرد

کز این برتر اندیشه برنگذرد



دانشکده مهندسی کامپیوتر و فناوری اطلاعات
گروه هوش مصنوعی

دسته‌بندی اخبار فارسی با استفاده از تکنیک‌های هوش مصنوعی

زهرا رباطی

استاد راهنما:

دکتر مرتضی زاهدی

استاد مشاور:

مهندس مرضیه رحیمی

پایان‌نامه جهت اخذ درجه کارشناسی ارشد

شهریور ۱۳۹۳

دانشگاه صنعتی شاهرود

دانشکده: مهندسی کامپیوتر و فناوری اطلاعات

گروه: هوش مصنوعی

پایان نامه کارشناسی ارشد خانم زهرا رباطی

تحت عنوان: دسته‌بندی اخبار فارسی با استفاده از تکنیک‌های هوش مصنوعی

در تاریخ توسط کمیته تخصصی زیر جهت اخذ مدرک کارشناسی ارشد مورد ارزیابی و با درجه مورد پذیرش قرار گرفت.

امضاء	استاد مشاور	امضاء	استاد راهنما
	سرکار خانم مهندس مرضیه رحیمی		جناب آقای دکتر مرتضی زاهدی

امضاء	نماینده تحصیلات تکمیلی	امضاء	اساتید داور
	نام و نام خانوادگی:		نام و نام خانوادگی:
			نام و نام خانوادگی:

تقدیم بہ پدر و مادر عزیز و صبورم

بہ خانواده دوست داشتنی ام

و ہمہ خوبان زندگی ام

شکر و قدرانی

سپاس خدای را که سخوران، در ستودن او بماند و شمارندگان، شمردن نعمت های او ندانند و کوشندگان، حق او را کزاردن نتوانند.

از استاد راهنمای خود، جناب آقای دکتر زاهدی که علاوه بر راهنمایی های بی درنیشان در تدوین این پایان نامه، مشاوره و دلسوز دوستی

گرانبهادر زندگی اینجانب هستند، کمال شکر را دارم.

از سرکار خانم مهندس رحیمی که استاد مشاورم در این پایان نامه بودند و در مراحل مختلف پایان نامه، راهنمایی هایشان را از اینجانب

دریغ ننمودند نیز کمال شکر را دارم.

همچنین قدردان زحمات سایر اساتیدم، جناب آقای دکتر حسن پور و جناب آقای دکتر پویان که در دوره ی تحصیلی کارشناسی ارشد،

علاوه بر حمایت های علمی ایشان، درس زندگی از آنان آموختم، هستم.

در پایان از راهنمایی های استاد عزیزم جناب آقای مهندس نامور فرهاد صمیمانه شکر و قدر دانی می کنم.

تعهد نامه

اینجانب زهرا رباطی دانشجوی دوره کارشناسی ارشد رشته مهندسی کامپیوتر دانشکده کامپیوتر و فناوری اطلاعات دانشگاه صنعتی شاهرود نویسنده پایان نامه دسته بندی اخبار فارسی با استفاده از تکنیک-

های هوش مصنوعی تحت راهنمایی دکتر مرتضی زاهدی متعهد می شوم.

- تحقیقات در این پایان نامه توسط اینجانب انجام شده است و از صحت و اصالت برخوردار است .
- در استفاده از نتایج پژوهشهای محققان دیگر به مرجع مورد استفاده استناد شده است .
- مطالب مندرج در پایان نامه تاکنون توسط خود یا فرد دیگری برای دریافت هیچ نوع مدرک یا امتیازی در هیچ جا ارائه نشده است .
- کلیه حقوق معنوی این اثر متعلق به دانشگاه صنعتی شاهرود می باشد و مقالات مستخرج با نام « دانشگاه صنعتی شاهرود » و یا « Shahrood University of Technology » به چاپ خواهد رسید .
- حقوق معنوی تمام افرادی که در به دست آمدن نتایج اصلی پایان نامه تأثیرگذار بوده اند در مقالات مستخرج از پایان نامه رعایت می گردد.
- در کلیه مراحل انجام این پایان نامه ، در مواردی که از موجود زنده (یا بافتهای آنها) استفاده شده است ضوابط و اصول اخلاقی رعایت شده است .
- در کلیه مراحل انجام این پایان نامه، در مواردی که به حوزه اطلاعات شخصی افراد دسترسی یافته یا استفاده شده است اصل رازداری ، ضوابط و اصول اخلاق انسانی رعایت شده است .

تاریخ

امضای دانشجو

مالکیت نتایج و حق نشر

- کلیه حقوق معنوی این اثر و محصولات آن (مقالات مستخرج ، کتاب ، برنامه های رایانه ای ، نرم افزار ها و تجهیزات ساخته شده است) متعلق به دانشگاه صنعتی شاهرود می باشد . این مطلب باید به نحو مقتضی در تولیدات علمی مربوطه ذکر شود .
- استفاده از اطلاعات و نتایج موجود در پایان نامه بدون ذکر مرجع مجاز نمی باشد.

چکیده

با توجه به رشد روزافزون متون الکترونیکی و از جمله متون خبری، معرفی یک دسته‌بند کارآمد در بسیاری از برنامه‌های کاربردی مرتبط با متون از جمله وبسایت‌های خبری، اهمیت بیشتری یافته است. از مهم‌ترین مسائل دامنه‌ی دسته‌بندی متون، استخراج ویژگی‌های کارا برای دسته‌بندی است. در کارهای گذشته معیارهای ارزش‌گذاری ویژگی متافوتی ارائه شده و مورد استفاده قرار گرفته‌اند. در این پایان‌نامه نیز یک معیار ارزش‌گذاری ویژگی بنام E-Dominance ارائه شده است که با استفاده از آن، کاهش چشمگیری در تعداد ویژگی‌های انتخابی رخ می‌دهد. ویژگی‌های به کار رفته در این دسته‌بندی، ویژگی‌های هم‌رخداد هستند که تا به حال در دسته‌بندی متون فارسی مورد استفاده قرار نگرفته‌اند.

در کارهایی که در زبان انگلیسی با استفاده از ویژگی‌های هم‌رخداد انجام شده‌اند، روش وزن‌دهی باینری به کار گرفته شده‌اند. در این پایان‌نامه با توسعه‌ی روش وزن‌دهی TFIDF، یک روش وزن‌دهی برای ویژگی‌های هم‌رخداد به نام Co-occur TFIDF معرفی و به کار گرفته شده است که دقت بهتری را در مقایسه با وزن‌دهی باینری نشان می‌دهد. ویژگی‌های توسعه‌گر نیز برای کم کردن اثر همپوشانی کلاس‌ها و توسعه‌ی متون، در این تحقیق به کار گرفته شده‌اند. نتایج آزمایشات، نشان‌دهنده‌ی بهبود قابل توجهی در کارایی و دقت الگوریتم دسته‌بندی با بکارگیری معیار E-Dominance و استفاده از ویژگی‌های هم‌رخداد با روش وزن‌دهی Co-occur TFIDF است.

واژه‌های کلیدی: دسته‌بندی متون فارسی، انتخاب ویژگی، استخراج ویژگی، معیار E-

Dominance، ویژگی‌های هم‌رخداد، ویژگی‌های توسعه‌گر، وزن‌دهی Co-occur TFIDF.

لیست مقالات مستخرج از پایان نامه

1. Zahra Robati, Morteza Zahedi, Najme Fayazi Far, "*Feature Selection and Reduction for Text Classification*", International Journal of Computer Applications, IJCA (Accepted on September 2014).

۲. زهرا رباطی، مرتضی زاهدی، نجمه فیاضی فر، "استخراج ویژگی‌های هم‌خداد با وزن‌دهی TFIDF هم‌خدادی در دسته‌بندی متون فارسی"، سومین همایش ملی زبان‌شناسی رایانشی، دانشگاه صنعتی شریف (پذیرفته شده مهر ۱۳۹۳).

فهرست مطالب

عنوان	شماره صفحه
فصل اول: مقدمه.....	۱
۱-۱- پیش‌زمینه‌ای از دسته بندی متون.....	۲
۲-۱- چالش ها.....	۵
۱-۲-۱- همپوشانی متون در کلاس‌های متفاوت.....	۵
۲-۲-۱- ریشه‌یابی.....	۶
۳-۲-۱- وزندهی.....	۶
۴-۲-۱- برخورد با چالش‌ها.....	۶
۳-۱- دستاوردهای تحقیق.....	۷
۴-۱- ساختار پایان‌نامه.....	۷
فصل دوم: دسته‌بندی خودکار متون و مروری بر کارهای پیشین.....	۹
۱-۲- مقدمه.....	۱۰
۲-۲- دسته بندی متون.....	۱۰
۱-۲-۲- کاربردهای دسته‌بندی متون.....	۱۱
۲-۲-۲- اهمیت دسته‌بندی متون.....	۱۳
۳-۲-۲- انواع روش‌های دسته‌بندی.....	۱۳
۴-۲-۲- مشکلات دسته‌بندی متون.....	۱۴
۳-۲- مراحل مختلف دسته‌بندی متون به صورت کلی.....	۱۷
۱-۳-۲- پیش پردازش.....	۱۸
۲-۳-۲- انتخاب ویژگی از متن یا شاخص‌بندی متون.....	۲۶
۳-۳-۲- کاهش ابعاد داده.....	۳۸
۴-۳-۲- الگوریتم‌های دسته‌بندی.....	۳۹
۴-۲- مقالات موجود در دسته‌بندی متون.....	۴۱
۱-۴-۲- مقالات مبتنی بر زبان‌های غیر فارسی.....	۴۱
۲-۴-۲- مقالات مبتنی بر زبان فارسی.....	۴۹
۵-۲- نتیجه‌گیری.....	۵۴

۵۵	فصل سوم: معرفی روش پیشنهادی.....
۵۶	۱-۳- مقدمه
۵۶	۲-۳- معماری رویکرد پیشنهادی
۵۸	۱-۲-۳- اجزای تشکیل دهنده‌ی معماری پیشنهادی
۵۹	۲-۲-۳- انتخاب متون از پایگاه داده
۶۰	۳-۲-۳- پیش‌پردازش‌ها
۶۲	۴-۲-۳- وزندهی به ویژگیها و حذف لغات کم اهمیت با توجه به وزن آنها
۶۳	۵-۲-۳- انتخاب ویژگی‌های کلاسی با معیار E-Dominance
۶۷	۶-۲-۳- تشکیل زوج ویژگی‌های هم‌رخداد به عنوان ویژگیهای متن
۶۸	۷-۲-۳- کاهش ویژگی‌های هم‌رخداد با استفاده از معیار پیشنهادی E-Dominance
۶۸	۸-۲-۳- انتخاب ویژگی‌های توسعه‌گر و توسعه داده‌ها
۶۹	۹-۲-۳- وزندهی ویژگی‌های هم‌رخداد با روش وزندهی Co-Occur TFIDF
۷۱	۱۰-۲-۳- دسته‌بندی و ارزیابی نتایج
۷۲	۳-۳- نتیجه‌گیری
۷۳	فصل چهارم: نتایج تجربی، تحلیل و ارزیابی.....
۷۴	۱-۴- مقدمه
۷۴	۲-۴- معرفی معیارهای ارزیابی
۷۷	۳-۴- نتایج دسته‌بندی متون
۷۷	۱-۳-۴- نتایج بکارگیری معیار E-Dominance در دسته‌بندی متون
۸۱	۲-۳-۴- نتایج بکارگیری ویژگی‌های هم‌رخداد در دسته‌بندی متون
۹۰	۳-۳-۴- مقایسه نتایج روش پیشنهادی در ۵ حالت مختلف در ۵ کلاس
۹۰	۴-۴- مقایسه روش پیشنهادی با چند روش موجود
۹۱	۵-۴- نتیجه‌گیری
۹۳	فصل پنجم: نتیجه‌گیری و کارهای آینده.....
۹۴	۱-۵- خلاصه تحقیق
۹۵	۲-۵- پیشنهاد برای کارهای آینده.....
۹۶	مراجع.....

فهرست شکل‌ها

عنوان شکل.....شماره صفحه

شکل (۱-۲) مدل فضای برداری.....	۱۵
شکل (۲-۲) نمودار کلی فرآیند دسته‌بندی متون.....	۱۸
شکل (۱-۳) نمودار کلی معماری پیشنهادی.....	۵۷
شکل (۲-۳) نمونه‌ای از یک متن در گروه اقتصادی از پایگاه داده همشهری.....	۶۰
شکل (۳-۳) نمونه‌ای از متن با اولین پیش‌پردازش.....	۶۰
شکل (۴-۳) نمونه‌ای از متن با Stop Word های حذف شده.....	۶۲
شکل (۱-۴) تاثیر مقدار معیار E-Dominance در تعداد ویژگی‌های انتخابی.....	۷۹
شکل (۲-۴) تاثیر تعداد ویژگی‌های انتخابی در دقت دسته‌بندی.....	۷۹
شکل (۳-۴) تاثیر مقدار معیار E-Dominance در دقت دسته‌بندی.....	۸۰
شکل (۴-۴) تاثیر حد آستانه معیار E-Dominance در دقت دسته‌بندی با ویژگی‌های هم‌خداد با دسته‌بند SVM.....	۸۵
شکل (۵-۴) تاثیر تعداد ویژگی‌های انتخابی هم‌خداد در دقت دسته‌بندی با دسته‌بند SVM.....	۸۶
شکل (۶-۴) تاثیر تعداد ویژگی‌های توسعه‌گر در دقت دسته‌بندی.....	۸۸
شکل (۷-۴) تعداد ویژگی‌های توسعه‌گر نسبت به تعداد کل ویژگی‌ها.....	۸۹
شکل (۸-۴) تاثیر نسبت تعداد ویژگی‌های توسعه‌گر به تعداد کل ویژگی‌ها در دقت دسته‌بندی.....	۸۹
شکل (۹-۴) مقایسه‌ی معیار F1 برای روش‌های مختلف در ۵ کلاس.....	۹۰

فهرست جداول

عنوان جدول.....شماره صفحه

جدول (۱-۲)	نمونه‌ای از واژه‌های غیر ضروری (شامل اسم، صفت وحروف و وندها و افعال).....	۲۱
جدول (۱-۳)	داده‌های انتخابی از پایگاه داده همشهری ۲.....	۵۹
جدول (۲-۳)	نمونه‌هایی از کلمات ایست استخراجی.....	۶۲
جدول (۳-۳)	نمونه‌هایی از ویژگی‌های استخراج شده با معیار E-Dominance برای کلاس‌های مختلف.....	۶۷
جدول (۴-۳)	نمونه‌هایی از ویژگی‌های هم‌رخداد تشکیل شده.....	۶۸
جدول (۵-۳)	نمونه‌ای از ویژگی‌های توسعه‌گر انتخابی از کلاس.....	۶۹
جدول (۱-۴)	پارامترهای مورد نیاز اندازه‌گیری معیارهای کارایی.....	۷۴
جدول (۲-۴)	تعداد ویژگی‌های انتخابی از هر دسته.....	۷۷
جدول (۳-۴)	نتایج دسته‌بندی داده‌ها با معیارهای ارزیابی.....	۷۸
جدول (۴-۴)	مقایسه میزان دقت دسته‌بندی در کاهش ویژگی‌ها با معیار E-dominance و سایر معیارها.....	۸۱
جدول (۵-۴)	تعداد ویژگی‌های هم‌رخداد انتخابی از هر دسته.....	۸۲
جدول (۶-۴)	نتایج دسته‌بندی داده‌ها با ویژگی‌های هم‌رخداد با وزن‌دهی باینری KNN.....	۸۳
جدول (۷-۴)	نتایج دسته‌بندی داده‌ها با ویژگی‌های هم‌رخداد با وزن‌دهی باینری SVM.....	۸۳
جدول (۸-۴)	نتایج دسته‌بندی داده‌ها با ویژگی‌های هم‌رخداد با وزن‌دهی Co-occur TFIDF و دسته‌بند KNN.....	۸۴
جدول (۹-۴)	نتایج دسته‌بندی داده‌ها با ویژگی‌های هم‌رخداد با وزن‌دهی Co-occur TFIDF و دسته‌بند SVM.....	۸۴
جدول (۱۰-۴)	ویژگی‌های توسعه‌گر انتخابی از هر دسته.....	۸۷
جدول (۱۱-۴)	مقایسه روش پیشنهادی با چند روش دیگر.....	۹۲

فصل اول

مقدمه

۱-۱- پیش زمینه‌ای از دسته‌بندی متون

امروزه با گسترش روز افزون متون و نوشتار الکترونیکی، بخصوص رشد گسترده در وبسایت‌های خبری، وجود ابزار کارآمد دسته‌بندی خودکار متون و به طور ویژه متون خبری، اجتناب ناپذیر گردیده است. در حال حاضر دسته‌بندی متون در بسیاری از زمینه‌های تحقیقاتی و کاربردی، از شاخص‌گذاری متون بر مبنای یک لغت‌نامه کنترل شده تا فیلتر کردن متون، تولید خودکار فراداده، ابهام‌زدایی از کلمه، تولید کاتالوگ‌های سلسله‌مراتبی از منابع وبی و به طور کلی در هر کاربردی که به سازماندهی متون یا توزیع انتخابی و تطبیقی خاصی از متون نیاز دارد، کاربرد دارد [۱]. اما با اشراف بر این موضوع که دسته‌بندی متون امری لازم و ضروری است، باز هم حدود ۸۰٪ تا ۹۰٪ متون، بدون سازماندهی مشخص در ادارات و نهادها نگهداری می‌شوند [۲].

یکی از کاربردهای رو به رشد این حوزه، دسته‌بندی متون خبری است. در عصر حاضر، کار اصلی نرم-افزارهای خبری، تنها رساندن و نشر خبر نیست. یک وبسایت خبری هوشمند باید پاسخگوی تمامی نیازهای کاربر و بازدیدکننده، مثل دسته‌بندی و سازماندهی درست اخبار، امکان جستجوی سریع خبرها با واژه‌های کلیدی کارآمد، پیشنهادات مناسب برای کاربر بر اساس سوابق بازدید از صفحات مختلف، جذاب-سازی خبر جهت افزایش آمار بازدیدکنندگان و به طور کلی جوابگوی نیازهای روز افزون کاربر باشد، تا بتواند در عرصه‌ی رقابت باقی بماند. روش کلی کار در دسته‌بندی خبر و سایر متون یکسان است و در بیشتر کارهای انجام شده از متون خبری در آزمایشات استفاده شده است اما تأکید بر نوع متون نشده است و روش کار به صورت کلی بر روی همه‌ی متون بیان شده است.

دو دلیل عمده ما را بر آن داشته است که تمرکز بر متون خبری داشته باشیم، اولاً دادگان مورد استفاده در این تحقیق، متون خبری بوده و ثانیاً هدف ما ارتقاء یک نرم افزار خبری تحت وب ساده به یک

نرم‌افزار خبری هوشمند بوده است. در این پایان‌نامه هر جا از واژه‌ی "متون" نام برده شده، منظور متون خبری است.

در حال حاضر تکنیک‌های یادگیری ماشین، داده‌کاوی و پردازش زبان طبیعی به کمک هم آمده و سعی در دسته‌بندی دقیق متون الکترونیکی دارند [۲، ۳]. روش‌های متداول موجود در دسته‌بندی متون، شامل الگوریتم‌های احتمالی (تئوری بیز و مدل‌های n-gram)، الگوریتم‌های خطی مثل SVM و نیز الگوریتم‌های مبتنی بر استخراج قوانین است [۴]. بردار ویژگی مورد استفاده در این الگوریتم‌ها، معمولاً از روش مجموعه‌ی کیسه‌ی لغات^۱ یا مدل فضای بردار^۲ تولید می‌شود؛ یعنی متن اصلی با یک بردار از کلمات انتخابی و وزن این کلمات نمایش داده می‌شود [۵، ۶]. دلیل اصلی در گوناگونی روش‌های دسته‌بندی متون، شیوه‌ی ایجاد بردار ویژگی و نیز انتخاب ویژگی‌های مناسب، به دلیل اهمیت فوق‌العاده‌ی آن در نتایج دسته‌بندی می‌باشد. علاوه بر کیسه‌ی لغات، عبارات معنایی^۳ و N-gram های حروف و کلمات نیز از ویژگی‌های مطرح و پرکاربرد در حوزه‌ی دسته‌بندی متن هستند [۲]. از روش‌های دیگری که برای استخراج ویژگی وجود دارند، استخراج ویژگی‌های هم‌رخداد^۴ است [۷-۱۱]. این ویژگی، تا بحال به عنوان یک ویژگی برای دسته‌بندی متون فارسی استفاده نشده است. ویژگی‌های هم‌رخداد به دلیل ماهیتی که دارند، ویژگی‌های مطمئن‌تری از ویژگی‌های منفرد مثل کیسه‌ی لغات هستند و علاوه بر در نظر گرفتن ویژگی آماری، از دانش معنایی استخراج شده از داده‌های آموزشی نیز بهره می‌گیرند. این ویژگی‌ها نیز با معضل حجم بالای تعداد ویژگی‌ها روبرو هستند [۷]. در این نوع ویژگی‌ها، انتخاب ویژگی‌های اولیه و منفرد برای تشکیل زوج ویژگی‌ها، یکی از موثرترین عوامل در تعداد ویژگی‌های نهایی و دقت دسته‌بندی است [۷].

¹ Bag of Word

² Vector Space Model

³ Syntactic Phrases

⁴ CO-occurrence Features

مسئله‌ی وزن‌دهی و معیارهای ارزش‌گذاری به ویژگی‌ها نیز از مهم‌ترین مسائل دامنه‌ی دسته‌بندی متون هستند که هم در انتخاب و استخراج ویژگی‌ها نقش مهمی ایفا کرده و هم در دقت دسته‌بندی موثرند. از روش‌های متداول وزن‌دهی، وزن‌دهی باینری، بسامد عبارت^۱ (TF)، بسامد متن^۲ (DF)، TFIDF و ... هستند که در کارهای مختلف استفاده شده‌اند. از معیارهای ارزش‌گذاری ویژگی نیز بهره‌ی اطلاعات^۳، آنترپی^۴، مربع چی^۵ (χ^2) و Dominance و ... هستند که در مرحله‌ی انتخاب و استخراج ویژگی‌ها بسیار پرکاربرد هستند [۲، ۱۲]. گرچه این معیارها و متدهای وزن‌دهی همچنان در تحقیقات مورد استفاده قرار می‌گیرند، اما هیچ‌یک از آنها کارایی قابل ملاحظه‌ای در مقابل سایر روش‌ها ندارد و بسیاری از ویژگی‌های انتخابی بر اساس این معیار و شیوه‌های وزن‌دهی، نامناسب و ناکارآمد هستند.

صرف نظر از نحوه‌ی تشکیل بردار ویژگی، در میان کارهای ارائه شده با انواع مختلف دسته‌بندها مثل k-همسایه نزدیکتر^۶، بیزین^۷، مدل بیشترین آنترپی^۸، SVM^۹ و ...، به طور میانگین روش‌های مبتنی بر SVM از دقت بیشتری نسبت با سایر روش‌ها برخوردارند [۲، ۳، ۶، ۷، ۱۲]. در واقع این دسته‌بندها سعی می‌کنند ابرصفحه‌ای که بیشترین حاشیه‌ی مجاز بین ابرصفحه و نمونه‌های مثبت و منفی را در برداشته باشند، بیابند [۲، ۱۳]. اگرچه این دسته‌بندها به سبب سادگی ذاتی دارای پشتوانه‌ی نظری بسیار خوبی هستند، اما هزینه‌های محاسباتی، زمانی و حافظه‌ای آن به یک معضل تمام عیار تبدیل گشته است؛ بخصوص در این کاربرد که با داده‌های متنی و حجم عظیم کلمات و واژه‌ها و پردازش آنها سر و کار داریم، این معضل بیشتر دیده می‌شود. بعد از SVM نیز در کارهای انجام شده، دسته بند K-همسایه نزدیکتر دقت بالاتری را نسبت به سایر دسته‌بندها ارائه کرده است.

¹ Term Frequency

² Document Frequency

³ Information Gain

⁴ Entropy

⁵ Chi Square

⁶ K-Nearest Neighbors

⁷ Bayesian

⁸ Maximum Entropy Models

⁹ Support Vector Machine

ویژگی که متون خبری نسبت به سایر متون دارند این است که معمولا از قالب مشخص تری در هر دسته برای ارائه متن استفاده می کنند و ویژگی های مشخصی از هر دسته در قالب هر کلاس وجود دارند؛ مثلا به طور معمول متون خبری ورزشی از کلمات متداولی مثل "رقابت"، "برنده" و "امتیاز" استفاده می کنند. در فرآیند آموزش به دنبال یادگیری این قالب در مجموعه متون هر کلاس هستیم. این ویژگی در سایر متون مثلا مقالات علمی کمتر دیده می شوند؛ مثلا دو مقاله در شاخه فیزیک که یکی درباره عدسی ها و دیگری درباره ذرات بحث می کنند، ویژگی های مشترک کمتری دارند و در زیرگروه های مربوط به موضوع اصلی، ویژگی های مشترک بیشتری خواهند داشت. در این پروژه با ارائه یک معیار اهمیت داده کارآمد، ویژگی های همرخداد انتخاب شده و از میان آنها ویژگی های پرتکرارتر به عنوان ویژگی های توسعه گر برای توسعه متون به کار رفته اند. در انتخاب ویژگی های توسعه گر، ویژگی نامبرده متون خبری کمک شایانی به انتخاب ویژگی های توسعه گر می کند و قالب فراگیری شده از متون آموزشی به کل متون آموزشی کلاس و متون آزمایشی مشابه القا می شود. سپس با استفاده از طرح توسعه ی داده ها و نیز به کارگیری روش وزن دهی TFIDF همرخدادی¹ برای ویژگی های همرخداد، علاوه بر کاهش تعداد ویژگی ها، دقت نسبتا خوبی در نتایج دسته بندی بدست آمده است.

۲-۱- چالش ها

گرچه تحقیقات گسترده ای در زمینه ی دسته بندی متون در زبان های مختلف صورت گرفته است، اما برخی از مشکلات هنوز لاینحل باقی مانده است و انگیزه برای ادامه ی کار در این زمینه وجود دارد. از مهم ترین این مشکلات می توان به موارد زیر اشاره کرد:

۱-۲-۱- همپوشانی متون در کلاس های متفاوت: متون غیر منفک داده هایی هستند که

نمی توان به طور قطع نشان داد به کدامین کلاس تعلق دارند. برخی از متون به طور ذاتی به بیش از یک

¹ Co-Occur TFIDF

دسته تعلق دارند؛ یعنی ممکن است یک متن خبری با میزان تعلق متفاوت هم سیاسی باشد و هم اجتماعی. این امر باعث می‌شود که کلاس‌ها با هم همپوشانی داشته باشند که خود منجر به افزایش خطا در دسته‌بندی می‌گردد [۴].

۱-۲-۲- ریشه‌یابی: اکثر روش‌های تشکیل بردار ویژگی به منظور کاهش خطا، از ریشه‌یابی کلمات استفاده می‌کنند. برای ریشه‌یابی، از کلمات مورد پردازش، پیشوندها و پسوندها حذف شده و تبدیل به ریشه می‌شوند. از مشکلات ریشه‌یابی می‌توان به بیش ریشه‌یابی و کم ریشه‌یابی اشاره نمود. در بیش ریشه‌یابی، ریشه‌هایی تولید می‌شوند که هیچ معنایی ندارند و در کم ریشه‌یابی امکان تولید ریشه برای حالت‌ها استثنا وجود ندارد [۱۴].

۱-۲-۳- وزن‌دهی: اینکه هر یک از کلمات موجود در متن چه وزن و ارزشی خواهند داشت، یکی از اصلی‌ترین چالش‌هایی است که با آن روبرو هستیم. ثابت شده است که استفاده از عبارات بیشترین ارزش معنایی و کمتری ارزش آماری و از طرف دیگر استفاده از کلمات، کمترین ارزش معنایی و بیشترین ارزش آماری را دارد [۴]. برای به کار بردن هر کدام از این دو نوع ویژگی، یک روش وزن‌دهی موثر که میزان تاثیر و کارایی آن ویژگی را بیان کند، بسیار قابل توجه است. طرح‌های وزن‌دهی متفاوتی وجود دارند که در کارهای مختلف مورد استفاده قرار می‌گیرند، اما هنوز فقدان یک روش قدرتمند که منجر به یک دسته‌بندی با دقت بالا شود، احساس می‌شود.

۱-۲-۴- برخورد با چالش‌ها

در مواجهه با چالش داده‌های غیر منفک و همپوشانی کلاس‌ها در این پایان‌نامه، سعی شده است با طرح توسعه‌ی داده‌ها با استفاده از ویژگی‌های توسعه‌گر متون را غنی‌تر کرده و اثر همپوشانی کلاس‌ها به حداقل رسانده شود. این کار باعث بهبود نتیجه در دقت دسته‌بندی می‌گردد.

در رویارویی با چالش وزن‌دهی نیز در این پایان‌نامه، عبارات هم‌رخداد به عنوان ویژگی استفاده شده-
اند که ارزش معنایی را در بر دارند و از طرفی با توسعه‌ی روش TFIDF به Co-Occur TFIDF برای
ویژگی‌های هم‌رخداد از ارزش آماری ویژگی‌ها نیز استفاده شده است.

۳-۱- دستاوردهای تحقیق

- ارائه‌ی یک مجموعه‌ی ۱۰۸۰ کلمه‌ای از Stop Words برای زبان فارسی
- ارائه‌ی معیار^۱ E-Dominance به عنوان بهبودی از معیار Dominance برای اهمیت ویژگی‌ها
- کاربرد ویژگی‌های هم‌رخداد برای اولین بار در دسته‌بندی متون فارسی
- کاربرد ویژگی‌های توسعه‌گر برای توسعه‌ی متون برای اولین بار در دسته‌بندی متون فارسی
- توسعه‌ی طرح وزن‌دهی TFIDF پایه به Co-Occur TFIDF برای ویژگی‌های هم‌رخداد

۴-۱- ساختار پایان‌نامه

این پایان‌نامه در ۵ فصل تدوین شده است. پس از مقدمه، در فصل دوم به معرفی مفاهیم پایه در
حوزه‌ی دسته‌بندی متون پرداخته شده است. همچنین در این فصل یک مرور کلی بر روی کارهای انجام
شده و روش‌های موجود انجام شده است. فصل سوم به معرفی روش پیشنهادی اختصاص دارد و به صورت
گام به گام آن را تشریح می‌کند. در فصل چهارم نیز نتایج حاصل از آزمایشات انجام شده و مقایسه روش-
پیشنهادی با روش‌های دیگر بیان می‌گردد. در فصل پنجم یک نتیجه‌گیری کلی از تحقیق و نیز ارائه‌ی
پیشنهادات برای ادامه‌ی کار ارائه می‌گردد.

^۱ Extended Dominance

فصل دوم

دستبندی خودکار متون و

مروری بر کارهای پیشین

۲-۱- مقدمه

در این فصل ابتدا به تعریف مفهوم دسته‌بندی متون پرداخته و سپس اهمیت و لزوم دسته‌بندی و نیز کاربرد آن در زمینه‌ی یادگیری ماشین بیان می‌شود. دسته‌بندی متون دارای مشکلات و محدودیت‌هایی است که در این فصل به طور خلاصه در حوزه‌ی عمومی و نیز در حوزه‌ی نگارش زبان فارسی به آنها پرداخته می‌شود. انواع روش‌های دسته‌بندی نیز در این فصل بیان می‌شود و در ادامه مراحل مختلف سیستم دسته‌بندی خودکار متون مورد بررسی قرار می‌گیرد. در پایان نیز مروری بر کارهای گذشته در زبان فارسی و سایر زبان‌ها در حوزه‌ی دسته‌بندی متون انجام می‌شود.

۲-۲- دسته بندی متون

در صورتیکه مجموعه‌ای از متون $D = \{(d_1, y_2), \dots, (d_i, y_i), \dots, (d_n, y_n)\}$ داشته باشیم به طوریکه n تعداد متون و $d_i = [w_{i,1}, \dots, w_{i,k}, \dots, w_{i,|d_i|}]$ متن i ام این مجموعه باشد، آنگاه $w_{i,k}$ به وزن کلمه‌ی k ام در متن i ام و y_i به کلاسی که متن به آن متعلق است (یعنی $y_i \in C$ به طوریکه $C = \{c_1, c_2, \dots, c_{|C|}\}$ مجموعه کلاس‌های از پیش تعریف شده در سیستم باشد) اشاره می‌کند. هدف در دسته‌بندی متون، استنتاج یک تابع f است به نحوی که $y_i = f(d_i)$ باشد. دسته‌بندی متون تعیین یک مقدار بولی برای هر جفت $\langle d_j, c_i \rangle \in D * C$ ، در جایی است که D مجموعه‌ای از متون و C مجموعه کلاس‌های از پیش تعیین شده باشد. مقدار T تعیین می‌کند که متن d_j به کلاس c_i تعلق دارد و مقدار F نیز عدم تعلق را نشان می‌دهد. هدف در اینجا به دست آوردن تابع تخمین $\{T, F\} \rightarrow D * C$ می‌باشد [۴].

۲-۲-۱- کاربردهای دسته‌بندی متون

کاربردهای دسته‌بندی متون به کار آقای ماروت^۱ بر روی دسته‌بندی متون به صورت احتمالی باز می‌گردد. از آن پس دسته‌بندی متون برای کاربردهای متنوعی استفاده شده است. لازم به ذکر است که این گروه‌بندی مرز مشخصی نداشته و برخی از این گروه‌ها با یکدیگر هم‌پوشانی دارند که به طور خلاصه عبارتند از [۴]:

- **شاخص‌بندی^۲ برای سیستم‌های بازیابی اطلاعات:** در سیستم‌های بازیابی اطلاعات، برای هر متن یک یا چند کلید واژه (عبارت کلیدی) که محتوای آن متن را توصیف می‌کند، تعیین می‌شود. این کلید واژه‌ها متعلق به مجموعه لغت‌نامه‌ای کنترل‌شده با عنوان گنج‌واژه هستند. اغلب این گنج‌واژه‌ها ساختار سلسله‌مراتبی موضوعی^۳ دارند (مانند گنج‌واژه ناسا^۴ برای مرتب‌سازی مسائل فضایی و یا گنج‌واژه مش^۵ برای مسائل پزشکی).

- **سازمان‌دهی متون:** سازماندهی متون و بایگانی برای مقاصد شخصی، سازمانی و... توسط روش‌های دسته‌بندی خودکار متون امکان‌پذیر می‌باشد. برای مثال در ورودی تحریریه‌ی یک روزنامه، تبلیغات و آگهی‌ها برای انتشار باید تحت گروه‌های خدمات عمومی، فروش، خرید، درخواست کار و... دسته‌بندی شده باشند. [۱۵].

- **فیلترکردن متون:** فیلترکردن متون بر روی رشته‌ای از متن‌های ورودی که توسط یک تولیدکننده‌ی اطلاعات برای یک مصرف‌کننده‌ی اطلاعات بصورت غیرهمزمان توزیع می‌شود، اعمال می‌گردد [۱۶]. به عنوان مثال در تلکس‌های خبری، تولیدکننده‌ی خبر، یک خبرگزاری است و اخبار را برای

¹ marot

² Indexing

³ Thematic Hierarchical Thesaurus

⁴ NASA Thesaurus

⁵ MESH Thesaurus

یک مصرف‌کننده‌ی خبر که یک روزنامه است، می‌فرستد [۱۷]. در این حالت سیستم فیلترگذاری باید از دریافت خبرهایی که برای گیرنده جذاب نیستند، جلوگیری کند.

علاوه بر موارد ذکر شده، یک سیستم فیلتر می‌تواند متون را به دسته‌های موضوعی مورد نظر مصرف‌کننده نیز دسته‌بندی کند. اگر سیستم فیلتر در سمت تولیدکننده باشد، متون را به نحوی که مورد نظر مصرف‌کنندگان است، توزیع می‌کند. بدین منظور سیستم برای هر مصرف‌کننده یک نمایه^۱ منحصر به فرد تولید کرده و مرتباً آن را بروز می‌نماید. اگر سیستم فیلتر در سمت مصرف‌کننده باشد، از دریافت مطالب بی‌ربط به مصرف‌کننده جلوگیری می‌کند. در این حالت تنها یک نمایه مورد نیاز است. نمایه توسط کاربر مقداره‌ی اولیه شده و با استفاده از اطلاعات بازخوردی از کاربر (صریح و یا غیرصریح)، بروز می‌شود [۱۸]. این مسأله را فیلتر تطبیق‌پذیر می‌نامند.

- **رفع ابهام:** عملیات رفع ابهام در بسیاری از کاربردها، همچون پردازش زبان طبیعی و شاخص‌بندی مستندات با استفاده از ترجیح نقش کلمه بر خود کلمه در حوزه بازیابی اطلاعات مورد توجه است. در دسته‌بندی، تکرار کلمه در زمینه، به عنوان متن و نقش کلمه به عنوان دسته دیده می‌شود [۱۹، ۲۰].
- **دسته‌بندی سلسله مراتبی صفحات وب:** دسته‌بندی سلسله مراتبی صفحات وب منجر به افزایش کیفیت پرس‌وجو برای موتورهای جستجو می‌گردد. در این کاربرد، روش‌های دسته‌بندی مبتنی بر دسته^۲ انتخاب می‌شوند؛ چراکه اجازه اضافه‌کردن دسته‌های جدید و یا حذف دسته‌های بی‌مصرف را امکان‌پذیر می‌نمایند [۲۱، ۲۲].

از دیگر کاربردهای دسته‌بندی متون می‌توان به دسته‌بندی گفتاری که در واقع ترکیبی از دسته‌بندی متون و تشخیص گفتار است، دسته‌بندی متون چندرسانه‌ای^۳ از طریق عنوان‌های متنی، تشخیص نویسنده

¹ Profile

² Category-based Text Classification

³ Multimedia Document Categorization

برای متون ادبیاتی نامشخص، تشخیص زبان برای متونی که زبان آنها مشخص نیست، تشخیص خودکار جنس متن و رتبه‌بندی خودکار کیفیت نوشتار اشاره نمود [۱].

۲-۲-۲- اهمیت دسته‌بندی متون

دسته‌بندی متون بصورت دستی، فرآیندی زمانبر و نیازمند دانش افراد متخصص در زمینه‌های مختلف است و دسترسی به این افراد به طور مداوم امکان‌پذیر نخواهد بود. بدلیل کاربرد دانش انسانی در این فرآیند، دسته‌بندی متون عاری از خطا و اشتباه نخواهد بود. همچنین در صورت به‌کار گرفتن چند متخصص انسانی در حین فرآیند دسته‌بندی، سیستم بدلیل تفاوت تصمیم‌های اتخاذ شده افراد مختلف، دچار ناسازگاری می‌گردد. بخصوص در زمینه‌ی متون خبری، سرعت دسته‌بندی و سازمان‌دهی خبر برای نشر آن، امری ضروری است که به صورت دستی سرعت کار بسیار کند خواهد بود. بنابراین به دلیل حجم زیاد متون، دسته‌بندی بصورت دستی تقریباً غیرممکن است و به منظور صرفه‌جویی در زمان، هزینه، نیروی انسانی و از طرفی افزایش دقت و کارایی، در سال‌های اخیر به روش‌های دسته‌بندی خودکار روی آورده شده‌است.

۲-۲-۳- انواع روش‌های دسته‌بندی

در ادامه روش‌های مختلف دسته‌بندی متون مورد مطالعه قرار می‌گیرد [۴]. قابل توجه است که در اینجا فرض این است که دسته‌بندی متون بر مبنای دانش درون‌زا انجام می‌شود. یعنی یک متن فقط با استفاده از اطلاعات درونی اش و بدون استفاده از اطلاعات خارجی دسته‌بندی می‌شود.

۲-۲-۳-۱- دسته بندی تک‌برچسبی در مقابل چندبرچسبی

اگر مساله دسته‌بندی در حالتی باشد که هر متن فقط به یک دسته تعلق داشته باشد و کلاس‌ها با هم همپوشانی نداشته باشند، اصطلاحاً دسته‌بندی تک‌برچسبی گفته می‌شود. اما زمانی که یک متن به بیش از

یک کلاس تعلق داشته باشد، یعنی بین کلاس‌ها همپوشانی داشته باشیم، دسته‌بندی از نوع چندبرچسبی خواهد بود.

۲-۲-۳-۲ دسته‌بندی متون مبتنی بر کلاس^۱ در مقابل مبتنی بر متن^۲

در روش مبتنی بر متن، DPC، برای یک متن، باید تمام کلاس‌ها را که باید انتخاب شوند یافت. از طرف دیگر در روش مبتنی بر کلاس، CPC، برای هر کلاس، باید متونی که عضو آن کلاس هستند را یافت. این اختلاف بیشتر فلسفی است تا مفهومی، اما در جایی اهمیت پیدا می‌کند که مجموعه‌های C و D از ابتدا موجود نباشند. به طور عمده شرایط DPC از شرایط CPC متداول‌تر است و بیشتر مورد استفاده قرار می‌گیرد.

۲-۳-۳-۲ دسته‌بندی قطعی در مقابل دسته‌بندی رتبه‌ای

فرض کنید یک فرد خبره می‌خواهد مجموعه‌ای از متون را دسته‌بندی کند، سیستمی وجود دارد که متون را بدون هیچ‌گونه دسته‌بندی قطعی در کلاس‌های موجود رتبه‌بندی می‌کند. حالا فرد خبره، انتخاب کلاس برای هر متن را به بالای لیست به جای کل لیست محدود می‌کند و سریعتر و بهتر تصمیم می‌گیرد. از طرفی برای هر کلاس از مجموعه کلاس‌ها، سیستم توانایی رتبه‌دهی ساده متون را بر مبنای تناسب با آن کلاس خاص را دارد، یعنی برای دسته‌بندی تحت آن کلاس، کافی است فرد خبره متون با رتبه بالا را به جای کل متون آزمایش کند. [۴].

۲-۲-۴-۲ مشکلات دسته‌بندی متون

دسته‌بندی خودکار متون، همانند سایر داده‌ها، دارای مشکلاتی است. برخی از این مشکلات مربوط به کلیات و شیوه‌های دسته‌بندی است و برخی دیگر مربوط به داده‌ها. علاوه بر مشکلات عمده‌ی دسته‌بندی

¹ Category-Pivoted Categorization (CPC)

² Document-Pivoted Categorization (DPC)

متون، دسته‌بندی متون به زبان فارسی نیز مشکلات خاص خود را دارد که در بخش‌های بعدی به آن اشاره خواهیم کرد.

۲-۲-۴-۱- مشکلات و محدودیت‌های دسته‌بندی متون

شیوه نمایش متون از اولین و مهم‌ترین محدودیت‌های دسته‌بندی متون است. هر سند متنی باید از نسخه‌ی متنی به یک بردار ویژگی، تبدیل گردد. این فرآیند منجر به کاهش پیچیدگی متون و سهولت در مدیریت آنها می‌گردد. رایج‌ترین روش در شیوه نمایش متون "مدل فضای برداری"^۱ است، که بصورت گسترده‌ای در دسته‌بندی متون مورد استفاده قرار می‌گیرد [۵، ۶]. مدل فضای برداری، هر متن را به شکل یک بردار ویژگی با استفاده از عبارات درون متن که شامل وزن‌های آن عبارات است، نشان می‌دهد. این روش را از آنجائی که هر کلمه را به عنوان یک ویژگی در نظر می‌گیرد، روش کیسه کلمات^۲ می‌نامند. این روش نمایش متن دارای عیوبی مانند نادیده گرفتن وابستگی میان عبارات و عدم توجه به ساختار عبارات درون متون است. در شکل ۱-۲ یک نمونه از این نمایش متن بصورت مدل فضای برداری مشاهده می‌شود. در این مدل هر سطر نماینده یک متن و هر ستون نماینده یک عبارت کلیدی است. مقدار هر سلول از این ماتریس میزان اهمیت و وزن آن عبارت در متن مربوطه را نشان می‌دهد.

	F_1	F_2	F_m
d_1	$W_{1,1}$	$W_{1,2}$	$W_{1,m}$
d_2	$W_{2,1}$	$W_{2,2}$	$W_{2,m}$
.			
.			
d_n	$W_{n,1}$	$W_{n,1}$	$W_{n,m}$

شکل (۱-۲) مدل فضای برداری

¹ Vector Space Model (VSM)

² Bag Of Words

به طور کلی طراحی سیستم دسته‌بندی متون، با مشکلات زیر مواجه است:

- **بردار ویژگی با ابعاد زیاد**: بالا بودن بعد بردارهای نمایش ویژگی، یکی از مهم‌ترین مشکلات یک داده‌ی متنی است. این مشکل علاوه بر زمان‌بر بودن می‌تواند منجر به "بیش یادگیری"^۱ در سیستم گردد [۲۳].

- **نیاز به دانش پیش‌زمینه و قالب معنایی متون**: آنچه مسلم است مفاهیم و معانی که بیان‌کننده مفهوم و موضوع کلی یک متن هستند، اهمیت زیادی در دسته‌بندی ایفا می‌کنند، بنابراین نمی‌توان صرفاً به معیارهای آماری برای دسته‌بندی اکتفا نمود [۲۳].

- **اهمیت ارتباط میان کلمات**: برخی کلمات در کاربردهای متفاوت و در همراهی با سایر کلمات، معانی مختلفی خواهند داشت که این مسئله فرآیند دسته‌بندی را دچار مشکل خواهد کرد. همچنین ممکن است وابستگی کلمات به کلمات قبل و بعد از خود و همراهی این کلمات با سایر کلمات، معانی آنها را تغییر داده و کارایی دسته‌بندی را مختل کند.

به منظور حل مشکل بعد بالای مدل فضای برداری در دسته‌بندی متون، روشهای بسیاری به منظور کاهش بعد فضای مشخصه مورد بررسی قرار گرفته‌اند [۲۴، ۲۵].

۲-۲-۴-۲- مشکلات زبان فارسی با رویکرد دسته‌بندی

هدف ما در این تحقیق، پیاده‌سازی یک سیستم دسته‌بند برای متون و اخبار فارسی است. بدین جهت در این بخش به طور خلاصه به بررسی مشکلات مطرح در دسته‌بندی متون زبان فارسی می‌پردازیم.

^۱ Overfitting

زبان فارسی از دسته‌ی زبانهای هندی-اروپایی است. گرچه زبان فارسی همانند زبان انگلیسی دارای ریخت‌شناسی^۱ و ندادن‌افزای است و وجود پس‌وندها و پیش‌وندها در معانی لغات تغییر ایجاد می‌کند؛ اما پیچیدگی‌های ذاتی بیشتری نسبت به زبان انگلیسی دارد. در اینجا به دلیل حجم مطالب، از توضیح مشکلات زبان فارسی به طور کامل اجتناب شده است. برای مطالعه بیشتر به مرجع [۱] مراجعه شود. اما در کل اهم مشکلات پردازش متون فارسی در چند دسته زیر خلاصه می‌شود:

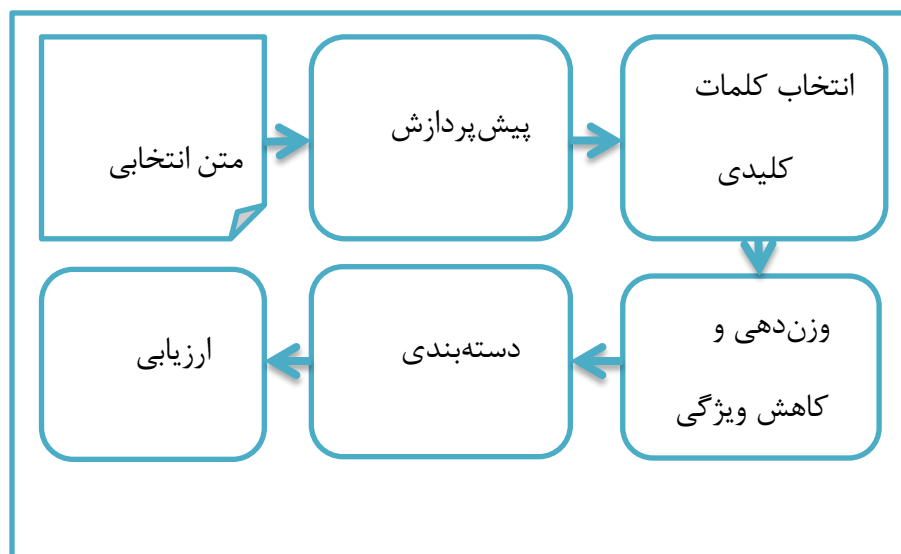
۱. مشکل تشخیص مرز کلمات (مسأله‌ی شیوه‌های نگارش متفاوت)
 ۲. مشکل تشخیص مرز گروه‌های اسمی (مسأله کسره‌ی اضافه نامرئی)
 ۳. از دست‌دادن اطلاعات گویشی
 ۴. مسأله ابهام
 ۵. افعال مرکب و اصطلاحات
 ۶. معناشناسی و مشکلات تحلیل معنایی
- موارد ذکر شده از ساختار زبان فارسی تنها بخشی از پیچیدگی این زبان را نشان می‌دهد. توجه به تمام این حالات به منظور قطعه‌بندی، ریشه‌یابی و نیز دسته‌بندی متون فارسی منجر به افزایش بار محاسباتی و نیز زمان مصرفی می‌گردد.

۲-۳- مراحل مختلف دسته‌بندی متون به صورت کلی

در شکل ۲-۲ مراحل مختلف و اصلی فرآیند دسته‌بندی متون آورده شده است. در حالت‌های خاص و کاربردهای متفاوت ممکن است مرحله‌ی جدیدی به آن اضافه گردد، اما کلیت کار دسته‌بندی به صورت

^۱ Morphology

زیر است. در بخش‌های بعدی این مراحل را به تفصیل بیان می‌کنیم. پس از شرح هر مرحله و یا به صورت همگام، مروری بر کارهای انجام شده در آن حوزه خواهیم داشت.



شکل (۲-۲) نمودار کلی فرآیند دسته‌بندی متون

۲-۳-۱- پیش‌پردازش

همانطور که می‌دانیم هر نوع داده‌ای که قرار است در سیستم دسته‌بندی قرار گرفته و برچسب مناسب کلاس خود را دریافت کند، نیاز به یک پیش‌پردازش اولیه دارد. این پیش‌پردازش اولیه باعث می‌شود که داده‌ی خام اولیه به یک داده‌ی مناسب برای فرآیند دسته‌بندی آماده گردد. به طور مثال در داده‌های مورد مطالعه ما وجود علائم نگارشی مانند نقطه و ویرگول و سایر علائم در متن مورد آزمایش، تاثیری در انتخاب برچسب کلاس آن متن ندارد. پس در گام نخست می‌توان این علائم را حذف کرد. اما به طور کلی پیش‌پردازش‌های عمده و معمولی وجود دارند که در کارهای مختلف انجام شده و متن را برای مرحله استخراج ویژگی‌ها و کلمات کلیدی آماده می‌سازد. قابل ذکر است که ممکن است در برخی مقالات بعضی

از مراحل ذکر شده جزء مرحله پیش‌پردازش نبوده و یکی از مراحل اصلی فرآیند دسته‌بندی مطرح گردند. صرف نظر از این نکته مراحل مختلف پیش‌پردازش به صورت زیر است:

۲-۳-۱-۱- پردازش لغوی و تعیین مرز کلمات

تعیین مرز کلمات و واژه‌ها در یک متن، از ابتدایی‌ترین مراحل در فرآیند دسته‌بندی اسناد است. این عملیات با توجه به پیچیدگی‌هایی که زبان فارسی دارد، یکی از مراحل دشوار در دسته‌بندی متون فارسی محسوب می‌گردد [۲۶-۲۸].

در بیشتر مواقع، تعیین مرز جمله‌ها از طریق بررسی علائم جداکننده انجام می‌شود. در این پایان‌نامه، شناسایی واژه‌ها و جمله‌ها با بررسی علائم نامبرده انجام گرفته است. علائمی که برای تعیین مرز جمله از آن‌ها استفاده می‌شود، عبارتند از: “.”، “؟”، “:”، “!”، “،”؛ باید توجه داشت که جستجو برای یافتن این علائم به تنهایی کافی نیست و در بسیاری از موارد مشکلاتی را به همراه دارد. به عنوان مثال، در زبان فارسی بعضی از واژه‌های اختصاری به صورت چند حرف که با نقطه از هم جدا شده‌اند، ظاهر می‌شوند (مانند «ه.ق.» که مخفف «هجری قمری» است). در صورتی که تعداد این حالات مشکل‌ساز در پیکره زیاد باشد، باید از روش‌های پیچیده‌تری استفاده شود، که قادر به شناسایی این موارد باشند. برای حل این مشکل از شمارش حروف و کلمه استفاده شده است؛ از آنجایی که کوچکترین جمله شامل دو کلمه است و کوچکترین کلمه از سه حرف تشکیل می‌شود (به غیر حروف اضافه مانند از، به) پس می‌توان واژه‌های اختصاری را مشخص کرد و به عنوان جمله در نظر نگرفت؛ از این رو فقط علائم در متن بررسی نشده‌اند بلکه اگر کلمات بین این علائم بیشتر از یک حد آستانه بود (که در اینجا دو در نظر گرفته شده) به عنوان جمله در آرایه قرار می‌گیرند.

شناسایی واژه‌ها نیز از طریق بررسی علائم قابل انجام است. این علائم عبارتند از: فضای خالی، علامت

خط جدید، “،” “،” “>” “<” “،” “[” “]” “_” “/” “-” “؛”

در [۲۹]، نویسنده علاوه بر بررسی روش‌های مختلف تعیین مرز واژه‌ها و جملات، مشکلات هر یک از آنها را مورد بررسی قرار داده است. [۳۰] نیز روشی برای شناسایی انتهای واژه‌ها معرفی کرده است که امکان جداسازی واژه‌ها را فراهم می‌آورد. در [۳۱] نیز یک روش آماری برای تعیین کسره اضافه معرفی شده است که با استفاده از آن می‌توان مرز گروه‌های اسمی را به طور دقیق‌تری شناسایی نمود.

۲-۳-۱-۲- حذف لغات غیر ضروری

در حوزه‌ی پردازش متن، کلماتی را که هیچ بار معنایی خاصی در بر ندارند به "لغات ایست"^۱ می‌شناسند. لغات ایست شامل حروف اضافه، بسیاری از قیود و صفات، حروف ربط و برخی از افعال هستند. این واژه‌ها در همه‌ی متون و با تکرار زیاد حضور داشته، اما در بردارنده اطلاعات معنایی نبوده و به همین دلیل دقت دسته‌بندی را کاهش می‌دهند. به طور مثال در زبان انگلیسی، لغات ایست شامل , that, the, an a و... هستند. در زبان فارسی نیز کلماتی مثل "از"، "برای"، "هر" و ... جزو کلمات ایست محسوب می‌گردند. اگر این کلمات حذف نشوند، بخش زیادی از شاخص‌ها را شامل می‌شوند. این کلمات که در متن پرتکرار هستند، در فرآیند وزن‌دهی که موثر از تعداد تکرار کلمات است، وزن بالایی به خود گرفته و اثر وزن کلمات دیگر را که در دسته‌بندی مهم بوده ولی پرتکرار نیستند، کم می‌کنند. حذف این دسته از کلمات نه تنها منجر به کاهش بعد فضای مشخصه می‌گردد، بلکه اثر نویزی این کلمات در دسته‌بندی را رفع می‌نماید. در حین استخراج کلمات اضافی باید به دامنه کلمات در اسناد و نیز پیکره مورد استفاده، توجه گردد [۳۲]. چرا که ممکن است برخی کلمات که در یک دامنه، کلمات اضافی هستند در برخی اسناد با دامنه‌ی دیگر دارای بار معنایی بوده و در دسته‌بندی اثربخش باشند. مرجع [۳۳] لغات غیر ضروری زبان فارسی را در دو جدول فعل و غیرفعل در دسترس قرار داده است. البته این مجموعه کامل نبوده و قابل گسترش است. جدول ۱-۲ نمونه‌ای از کلمات ایست را نشان می‌دهد.

^۱ Stop Words

جدول (۱-۲) نمونه‌ای از واژه‌های غیر ضروری (شامل اسم، صفت و حروف و وندها و افعال)

در	نیز	برای	یا	را	بلکه
به	تا	ها	دو	های	شاید
از	ما	آن	آنها	و	وقتی
که	باید	وی	اما	ممکن	اندرک
این	اند	یک	دیگر	هر	کمی
با	هم	خود	اگر	ای	بسیار
شما	همچنین	بر	چه	اندکی	زیاد
خواسته	داشتند	گرفتید	گیرند	یافتیم	خواسته

۲-۳-۱-۳- یکسان سازی نگارشی

در برخی زبان‌ها مثل فارسی وجود نویسه‌های متفاوت از یک حرف با کد اسکی متفاوت، باعث می‌شود که هنگام ساخت فرهنگ لغات و شمردن آن‌ها، دو واژه‌ی یکسان با املاهای متفاوت به عنوان دو واژه‌ی مختلف در نظر گرفته شوند. به طور مثال به کد اسکی حرف "ک" و "ك" که دو فرم فارسی و عربی از یک حرف هستند دقت کنید، ممکن است در دو متن مختلف یک کلمه با دو فرم متفاوت از این حرف ظاهر شده باشد. گرچه این دو کلمه یکسان هستند اما سیستم آن‌ها را دو کلمه‌ی مختلف تشخیص می‌دهد. برای جلوگیری از این معضل باید در زمان پردازش لغوی، حروف یک دسته را تبدیل به یک حرف یکسان نمود. به طور مثال در همه‌ی متون حرف "ك" به "ک" تبدیل گردد. علاوه بر حالت مذکور یک کلمه ممکن است به خاطر وجود پیشوند یا پسوند، املاهای متفاوتی داشته باشد. به عنوان مثال، پیشوند «می» و پسوند «ها» در ابتدا و انتهای واژه‌ها، ممکن است به سه صورت مختلف دیده شوند؛ کتابها

(چسبان)، کتاب‌ها (با نیم فاصله)، کتاب‌ها (با فاصله). به عنوان یک مثال از یکسان‌سازی، در مورد پسوند «ها» می‌توان شکل‌های «جدا با فاصله» و «جدا بدون فاصله» را نادیده‌گرفت و با مشاهده‌ی هر یک از این دو، آن را به شکل چسبان در آورد. در [۳۴]، یک مجموعه تبدیل‌های نسبتاً جامع و موثر برای یکسان‌سازی متون فارسی ارائه شده است. تجربه نشان می‌دهد که با استفاده از تبدیل‌های زیر، می‌توان پیکره‌ی متنی را به خوبی یکسان‌سازی کرد:

- تبدیل نویسه‌های «ی» و «ک» عربی به نوع فارسی آن.
- تبدیل نویسه‌های «ؤ» به «و»، «ئ» به «ی» و «أ» به «ا».
- تبدیل نویسه‌های «ة» به «ه» در آخر واژه‌ها.
- حذف «ی» از آخر واژه‌هایی مانند «خانه‌ی».
- حذف پسوندهای «تر» و «ترین» به آخر واژه‌ها.
- حذف فاصله بعد از پیشوند «بر» در واژه‌هایی مانند «برمی‌گردد».
- حذف شناسه‌ی «ء» در آخر بعضی واژه‌ها مانند «شهداء».
- حذف پیشوندهای «می»، «نمی»، «بی»، «درمی»، «برمی» از ابتدای واژه‌ها.
- حذف پسوندهای «ها»، «های»، «هایی»، «هایم»، «هایت»، «هایش»، «هایمان»، «هایشان»، «هایتان» از انتهای واژه‌ها.
- چسباندن پیشوند «به» به واژه‌هایی مانند «به‌ندرت» به ابتدای واژه‌ها.
- چسباندن پیشوند «هم» به واژه‌هایی مانند «هم‌چنین» به ابتدای واژه‌ها.

۲-۳-۱-۴- ریشه‌یابی

هدف از ریشه‌یابی حذف اضافات از کلمه و رسیدن به ریشه اصلی کلمه است به طوری‌که تمام کلمات هم‌ریشه در تمام متون، دارای یک نمایش منحصر‌بفرد به صورت ریشه‌ی اصلی کلمه باشند. این گام یکی

از مهم‌ترین و موثرترین مراحل پردازش لغوی است که در مراحل بعدی مثل انتخاب واژه‌های کلیدی، به عنوان ویژگی‌های اصلی متن، نقش بسیار تاثیرگذاری را ایفا می‌کند. روش‌های ریشه‌یابی را می‌توان به دو نوع عمده تقسیم نمود [۳۵]:

- **شیوه‌های مبتنی بر شبکه‌های واژه:** در این شیوه‌ها، از ارتباط واژه‌ها در یک شبکه معنایی استفاده می‌شود. واژه‌های هم‌ریشه در یک خوشه قرار می‌گیرند. این شبکه‌ها ساختاری گراف مانند دارند. این دسته از روش‌ها نیازمند بهره‌گیری گسترده از دانش انسانی هستند، به همین دلیل ایجاد و نگهداری آن‌ها بسیار پرهزینه است. هنوز الگوریتم‌هایی که بتواند به صورت خودکار این شبکه‌ها را به گونه‌ای قابل اطمینان ایجاد کنند، توسعه نیافته است. به دلیل کارایی بالای این دسته از روش‌ها، به طور گسترده در سامانه‌های داده‌کاوی و متن‌کاوی استفاده می‌شود.

- **شیوه‌های مبتنی بر ریخت‌شناسی:** رویکرد دوم استفاده از ریخت‌شناسی برای ریشه‌یابی است. این دسته از روش‌ها با بررسی ساخت واژه‌ها، ریشه‌ی کلمه را پیدا می‌کنند [۳۶]. این روش نیازی به معلوم بودن ریشه‌ی واژه‌ها ندارد و تنها نیازمند آگاهی از قواعد ساخت‌واژی است. از آنجا که در خلاصه-سازی بیشتر ریشه‌یابی تصریفی مورد نظر است و قواعد تصریفی اندک هستند، این دسته از روش‌ها به سادگی قابل پیاده‌سازی هستند.

به طور کلی از مشکلات اصلی ریشه‌یابی می‌توان به بیش‌ریشه‌یابی (تولید ریشه‌هایی که هیچ معنایی در زبان ندارند) و کم‌ریشه‌یابی (عدم امکان تولید ریشه برای حالات استثناء) اشاره کرد [۱۴]. اما سه روش معمول ریشه‌یابی که در بسیاری کاربردها مورد استفاده قرار می‌گیرند به شرح زیر است:

۱. ریشه‌یاب جدولی

این روش ریشه‌یابی از لحاظ پیچیدگی پیاده‌سازی، ساده‌ترین روش ممکن، اما بسیار زمانبر است. در این جداول هر کلمه به همراه ریشه‌اش نگهداری شده و برای یافتن ریشه‌ی یک کلمه باید جدول را برای

یافتن کلمه‌ی مورد نظر جستجو کرده سپس ریشه‌ی متناظر با آن را به دست آورد. در واقع این روش یک عمل جستجو است که منجر به ریشه‌یابی می‌گردد. این روش بیشترین دقت را نسبت به دیگر روشها دارا است؛ اما همانگونه که گفتیم سربرار زمانی بالایی دارد.

۲. ریشه‌یاب آماری

پایه اصلی این روش یک مجموعه از کلمات با ساخت‌های متفاوت است و عملکرد خوب این روش وابسته به کامل بودن این مجموعه است. در این دسته از ریشه‌یابها، با روش‌های آماری، وندهایی که در کلمه‌ها تکرار شده‌اند، شناسایی می‌گردد، از مزیت‌های این روش عدم وابستگی آنها به زبانی خاص است. باید توجه نمود که وجود کلمات نادرست در مجموعه مرجع این روشها، بر کارایی این ریشه‌یابها اثر منفی می‌گذارد [۳۷].

۳. ریشه‌یابی به کمک روش Porter

ریشه‌یاب پُرتِر^۱ [۳۸]، از معمول‌ترین شیوه‌های ریشه‌یابی در زبان انگلیسی است که بر مبنای زبان-شناسی و دسته‌بندی کلمه‌ها به کمک واجها و هجاها عمل می‌نماید. پایه و اساس روش‌های ریشه‌یابی براساس قاعده‌های زبان نیز همین روش پرتِر است. در طراحی این ریشه‌یابها، روال اصلی کار این است که یک ماشین پذیرنده متناهی تولید می‌گردد که معمولاً برخی قواعد و حالت‌های خاص در زمان طراحی از قلم افتاده و در انتها باید به آن اضافه گردد، یعنی به مرور زمان و در زمان‌های مقتضی کد برنامه باید تغییر کند. از این رو روند نگهداری و گسترش ریشه‌یاب بسیار پرهزینه می‌گردد.

¹ Porter

۲-۳-۱-۵- خلاصه سازی

منظور از خلاصه‌سازی، دریافت یک متن و تولید یا استخراج یک متن دیگر از آن متن است، به گونه‌ای که متن به دست‌آمده از متن اصلی کوتاه‌تر باشد، نکات اصلی و مهم آن را دربرداشته باشد و بین جملات آن پیوستگی وجود داشته باشد. خلاصه‌سازی باعث کاهش فضای ویژگی‌ها می‌شود.

به طور کلی، دو نوع اصلی برای خلاصه وجود دارد: خلاصه‌ی گزینشی و چکیده [۳۹].

- **خلاصه‌ی گزینشی**، با توجه به معیارهای آرمای، شهودی^۱ و یا ترکیبی از این دو تهیه می‌شود. از آنجا که در تولید این دسته از خلاصه‌ها، جملات متن تغییرات نحوی و معنایی ندارند، می‌توان آن را نوعی گزینش جملات قلمداد کرد. در واقع اگر متن خلاصه با انتخاب جملاتی از متن اصلی به دست آید، نوع خلاصه‌سازی، «استخراجی» یا «گزینشی» است. در حال حاضر خلاصه‌های استخراجی به دلیل آسان و ارزان بودن عمومیت بیش‌تری دارند.

- **خلاصه‌ی چکیده**، تفسیری از متن اولیه است. در تولید خلاصه‌ی چکیده، مفاهیم جملات متن اصلی به شکل کوتاه‌تر بازنویسی می‌شود. به عنوان مثال، جمله «اوسیب، انگور و گیلاس‌ها را خورد» را می‌توان به صورت «او میوه‌ها را خورد» نوشت. به عبارت دیگر اگر خلاصه متن پس از فهم مطالب موجود در متن اصلی تولید شود، خلاصه‌سازی از نوع «چکیده» است [۴۰].

۲-۳-۱-۶- ساده‌سازی

ساده‌سازی از طرق مختلف و در کاربردهای متفاوتی همچون مراجع [۴۱-۴۳]، بر روی متون انجام می‌گیرد. حذف کلمات هم‌معنی و تعمیم کلمات از جمله رویکردهای متداول در زمینه‌ی ساده‌سازی هستند که منجر به همگون‌سازی مستندات شده و سبب می‌شود که مستنداتی که دارای مضامین مشابهی هستند، دارای کلمات کلیدی تقریباً یکسانی باشند و در نتیجه کارایی و دقت سیستم دسته‌بندی

¹ Heuristic

خودکار متون را افزایش دهند. به عنوان مثال در یک مجموعه از اسناد، در دسته‌ی «ورزش»، اسنادی در زمینه‌ی شنا و اسنادی نیز در زمینه بسکتبال وجود دارند. تعمیم کلمات موجود در این دو دسته از اسناد، آنها را به سمت یک مجموعه از کلمات ریشه یعنی ورزش سوق می‌دهد. باید توجه نمود که برای اعمال روش‌های ساده‌سازی به طرق ذکرشده، به یک دانش پیش‌زمینه از لغات در زبان فارسی نیاز است.

۲-۳-۲- انتخاب ویژگی از متن یا شاخص‌بندی متون

یک متن نمی‌تواند به صورت مستقیم توسط یک دسته‌بند یا یک الگوریتم دسته‌بندی تفسیر شود، بلکه باید با استفاده از یک فرآیند شاخص‌بندی، متن را با استفاده از کلمات کلیدی یا ویژگی‌های استخراج شده از کلمات کلیدی نشان داد. به کلمات مهم هر متن که بیان‌کننده‌ی مفهوم اصلی متن هستند، **کلمه کلیدی** گفته می‌شود. به عبارت دیگر کلمات کلیدی خلاصه‌ای کوتاه از یک متن هستند. استخراج کلمات کلیدی یک گام مهم و تاثیرگذار در بازیابی اطلاعات، صفحات وب، دسته‌بندی متون، خلاصه‌سازی متون و غیره می‌باشد. این مهم کمک می‌کند تا یکنواختی و یک شکلی لازم برای متون مجموعه آموزشی، آزمایشی و اعتبارسنجی فراهم شود. انتخاب شاخص یا ویژگی برای متن، به مسائل مختلفی بستگی دارد، همانند:

- واحدهای معنایی (مساله واژگان معنایی^۱)

- قوانین معنایی طبیعی برای ترکیب این واحدها (مساله ترکیب معنایی^۲)

معمولا در دسته‌بندی متون مشکل دوم در نظر گرفته نمی‌شود (مشابه آنچه در بازیابی اطلاعات وجود دارد) و یک متن d_j ، با برداری از وزن عبارت‌های نشان داده می‌شود. به عبارت دیگر $d_j = \langle W_{1j}, W_{2j}, \dots \rangle$ دارد) و به طوریکه T مجموعه عبارت‌هایی است که دست کم یک بار در سرتاسر مجموعه آموزشی آمده

¹ Lexical Semantics

² Compositional Semantics

باشند (به آن ویژگی نیز گفته می‌شود) و $0 \leq W_{kj} \leq 1$ باشد. معمولاً تفاوت رهیافت‌ها در این زمینه به یکی از دلایل ذیل می‌باشد:

- تفاوت در تعریف چیزی که "عبارت" نامیده می‌شود.
 - تفاوت در زمینه محاسبه وزن ویژگی‌ها.
- یک انتخاب متداول برای مسئله اول، تعیین هر کلمه به عنوان یک عبارت می‌باشد. به این انتخاب، رهیافت مجموعه کلمات یا کیسه کلمات^۱ برای نمایش متون (مستقل از اینکه وزن‌ها دودویی باشند یا نباشند) گفته می‌شود [۴]. محققان از دو رهیافت مختلف با این مساله برخورد کرده‌اند:
- به طور صرفی و نحوی، تعریف عبارت با استفاده از قوانین گرامری زبان
 - به طور آماری، در نظر گرفتن یک ترکیب خطی نه به لحاظ قواعد گرامری بلکه بنا به همنشینی آن یا مجموعه یا دنباله‌ای از کلمات که الگوی خاصی را دنبال می‌نماید.
- لوئیز (Lewis) در [۴۴] استدلال کرد که استفاده از عبارت‌ها بیشترین ارزش معنایی و کمترین ارزش آماری را دارد. از طرف دیگر استفاده از کلمات کمترین ارزش معنایی و بیشترین ارزش آماری را دارد.
- به طور کلی بسته به کاربرد، کل یا بخشی از یک متن می‌تواند به عنوان شاخص‌بندی و ویژگی‌های یک متن مد نظر قرار گیرد. به عنوان مثال هنگامی که عنوان یک متن مشخص باشد، می‌توان در متن اصلی اهمیت بیشتری به کلماتی داد که در عنوان نیز حضور دارند. اما هنگامی که یک متن قالب خاصی ندارد و به صورت متن تنه‌است، تشخیص قسمتی که بیشترین ارتباط را با کل موضوع دارد، یک مساله بفرنج و نا مشخص است.

¹ Bag of Words

۲-۳-۱- روش‌های متداول در انتخاب کلمات کلیدی

در حالت کلی سه روش بسیار معمول برای انتخاب کلمات کلیدی وجود دارد که در بیشتر تحقیقات از آنها استفاده می‌شود [۴]:

- **روش‌های آماری**: در این روش‌ها میزان تکرار یک کلمه در یک متن و یا تعداد متون حاوی یک کلمه اهمیت خاصی دارند و بر اساس یک معیار آماری مثل TFIDF و یا معیارهای ارزش‌گذاری کلمه، کلمات کلیدی متن انتخاب می‌شوند.

- **روش یادگیری ماشینی**: در این شیوه با یک مجموعه متون آموزشی و کلمات کلیدی مشخص برای آن‌ها، فرآیند انتخاب کلمات کلیدی به عنوان یک مساله دسته‌بندی مدلسازی می‌شود. کلمات بر اساس مشخصه‌هایی که دارا هستند، به کلمات کلیدی و غیر کلیدی دسته‌بندی می‌شوند. احتمالات دسته‌بندی به صورت آماری از مجموعه‌ی آموزشی یاد گرفته می‌شود. این روش‌ها دارای انعطاف‌پذیری بالایی هستند.

- **ترکیب روش‌های تحلیل آماری و زبانشناختی**: علت ترکیب این دو روش این است که اطلاعات آماری بدون پالایش زبان‌شناختی نتایج مفیدی را در اختیار نمی‌دهند. بدون اطلاعات زبانشناختی، کلمات نامطلوبی مانند "از"، "او" و غیره هم به علت تعداد تکرار زیاد انتخاب خواهند شد. همانطور که گفته شد در روش‌های آماری با استفاده از معیارهای ارزش‌گذاری کلمات و معیارهای وزن‌دهی، به انتخاب کلمات کلیدی پرداخته می‌شود. در بخش بعدی این معیارها معرفی می‌شوند.

۲-۳-۲- طرح‌های وزن‌دهی به کلمات کلیدی

برای انتخاب کلمات کلیدی مناسب برای هر متن به عنوان ویژگی آن متن، آنچه بسیار مهم است، میزان اهمیت آن ویژگی در آن متن است. میزان اهمیت و تاثیرگذاری یک کلمه در یک متن را به عنوان

وزن آن کلمه تعریف می‌کنیم. چگونگی در نظر گرفتن میزان تاثیرگذاری و اهمیت کلمات و به اصطلاح وزن آنها، به قدری اهمیت دارد که حتی می‌تواند موثرتر از دسته‌بند انتخابی قلمداد گردد [۴۵].

در این بخش به معرفی چند شیوهی مختلف وزن‌دهی که در بازیابی اطلاعات و دسته‌بندی متون بسیار رایج هستند می‌پردازیم:

۲-۳-۲-۲-۱ روش‌های مبتنی بر TF

در این روش‌ها، وزن‌دهی ویژگی‌ها از توزیع آنها در هر یک از متون تبعیت می‌کند.

• روش TF

این روش برای اولین بار در مرجع [۴۶] ارائه شد و بدین شکل وزن یک ویژگی را تعیین می‌نمود که اگر ویژگی f_k در متن d_i حضور داشته‌باشد، وزن آن برابر تعداد تکرار آن خواهد بود.

$$W_{ik} = tf(f_k, d_i) = \begin{cases} \#(f_k, d_i) & \text{if } f_k \in d_i \\ 0 & \text{if } f_k \notin d_i \end{cases} \quad (1-2)$$

در رابطه‌ی ۱-۲ $\#(f_k, d_i)$ تعداد تکرار ویژگی f_k در متن d_i است.

• روش normTF

از آنجائیکه طول متون در مجموعه اسناد یکسان نیست، روش normTF به منظور رفع اثر طول مستند بر روی طرح وزن‌دهی و نگاشت آن در بازه‌ی (۰ و ۱) از نرمال‌سازی استفاده می‌کند (رابطه ۲-۲).

$$W_{ik} = norm_tf(f_k, d_i) = \frac{tf(f_k, d_i)}{\sqrt{\sum_k tf(f_k, d_i)^2}} \quad (2-2)$$

• روش logTF

در این روش به منظور رفع تاثیر منفی متفاوت بودن ماهیت ویژگی‌ها و مقادیرشان که منجر به کاهش دقت و کارایی دسته‌بندی می‌گردد، از عملگر لگاریتم استفاده شده‌است. این عملگر اثر ذکرشده را حذف و محدوده مقادیر تخصیص یافته به ویژگی‌ها را یکسان می‌نماید (رابطه ۳-۲).

$$W_{ik} = \text{Logtf}(f_k, d_i) = \log(\text{tf}(f_k, d_i)) \quad (3-2)$$

• روش ITF

این روش برای اولین بار در مرجع [47] معرفی شد که براساس آن وزن هر ویژگی از رابطه 4-2 بدست می آید.

$$W_{ik} = \text{ITF}(f_k, d_i) = 1 - \frac{r}{r + \text{tf}(f_k, d_i)} \quad (4-2)$$

در معادله ی 4-2 معمولاً r برابر 1 در نظر گرفته می شود.

• روش Sparck

روش Sparck از تئوری های آماری برای وزن دهی به ویژگی ها استفاده کرده و برای اولین بار در مرجع [48] ارائه شد.

$$W_{ik} = \text{Sparck}(f_k, d_i) = \text{tf}(f_k, d_i) * (k - \log(p_k)) \quad (5-2)$$

در رابطه 5-2، k تعداد کل ویژگی ها در مجموعه متون و $p_k = \sum_k \text{tf}(f_k, d_i)$ است.

2-2-2-3-2 روش های مبتنی بر IDF

در این دسته از روش ها، وزن دهی به مشخصه ها از توزیع ویژگی f_k در مجموعه متون تبعیت می نماید. مبنای اصلی این دسته از روش ها بر این اصل استوار است که هر چه تعداد متونی که دارای ویژگی f_k باشند کمتر باشد، این ویژگی برای تمایز بخشیدن به متون، مناسب تر خواهد بود.

• روش IDF

این روش برای اولین بار در مرجع [49] حوزه ی بازیابی اطلاعات مطرح شده است و در آن وزن ویژگی f_k از معادله 6-2 بدست می آید. این طرح وزن دهی از متغیر i که نشان دهنده شماره ی متن است، مستقل است.

$$W_{ik} = idf(f_k, d_i) = \log\left(\frac{|D|}{|D(f_k)|}\right) \quad (6-2)$$

در رابطه فوق $|D|$ تعداد کل متون مجموعه و $|D(f_k)|$ تعداد متون از مجموعه اسناد است که ویژگی f_k در آنها حضور دارد. همانطور که مشخص است وزن ویژگی f_k با افزایش $|D(f_k)|$ کاهش می‌یابد و این بدین معنی است که هر چه یک کلمه در متون بیشتری حضور داشته‌باشد، تمایزبخشی و اهمیت آن در دسته‌بندی متون کمتر است.

۲-۳-۲-۲-۳- روش‌های مبتنی بر TFIDF

این روش‌ها ابتدا در حوزه‌ی بازیابی اطلاعات مطرح شده و سپس در دسته‌بندی متون به منظور وزن‌دهی به ویژگی‌ها مورد استفاده قرار گرفت.

• روش TFIDF

روش TFIDF متداول‌ترین طرح وزن‌دهی به ویژگی‌ها در زمینه دسته‌بندی و خوشه‌بندی متون است. این روش در واقع حاصل ترکیب روش‌های مبتنی بر TF و IDF است. وزن ویژگی f_k در این طرح از رابطه (۷-۲) بدست می‌آید.

$$W_{ik} = TFIDF(f_k, d_i) = tf(f_k, d_i) * idf(f_k, d_i) \quad (7-2)$$

• روش normTFIDF

برای اطمینان از برابری شانس تمام متون با طول‌های متفاوت، بازیابی روش TFIDF در مرجع [۵۰] بصورت نرمال‌شده همانند معادله ۲-۸ ارائه شده‌است.

$$W_{ik} = normTFIDF(f_k, d_i) = \frac{tfidf(f_k, d_i)}{\sqrt{\sum_k tfidf(f_k, d_i)^2}} \quad (8-2)$$

در این بخش به معرفی چند روش وزن دهی پرداخته شد. استفاده از یک حد آستانه می تواند روش خوبی برای انتخاب و یا کاهش ویژگی ها باشد.

۲-۳-۲-۳- معیارهای ارزش گذاری ویژگی

در این بخش چند معیار پر کاربرد به عنوان معیاری برای سنجش میزان ارزش یک داده در یک کلاس و یا پایگاه داده معرفی می کنیم. در بسیاری از تحقیقات از این معیارها برای انتخاب ویژگی و فیلتر کردن ویژگی ها و در نهایت کاهش تعداد ویژگی ها استفاده شده است.

۲-۳-۲-۱- آنتروپی^۱

معیاری است که اهمیت داده را در کل کلاس ها نشان می دهد. اگر مقدار آنتروپی عبارت t_i صفر باشد، یعنی این عبارت فقط در یک کلاس رخ داده است [۲، ۱۲] و احتمالاً این عبارت قدرت تفکیک کنندگی خوبی برای دسته بندی خواهد داشت. اما اگر این عبارت در کلاس های مختلف و با تعداد رخداد مشابه ظاهر شود، مقدار آنتروپی برای آن زیاد خواهد بود و این ویژگی قدرت تفکیک کنندگی خوبی نداشته و برای کاهش ویژگی ها یک کاندید مناسب است. اگر n_i^k ، تعداد رخداد عبارت t_i در کلاس C_k باشد، فرکانس حضور t_i در کلاس C_k ، یعنی tf_i^k به صورت رابطه ۲-۹ تعریف می شود.

$$tf_i^k = \frac{n_i^k}{\sum_{k=1}^m n_i^k} \quad (۲-۹)$$

مقدار آنتروپی هم برای عبارت t_i با استفاده از ویژگی مقدار فرکانس حضور عبارت در کلاس به صورت رابطه ۲-۱۰ محاسبه می شود.

$$Entropy(t_i) = - \sum_{k=1}^r (tf_i^k) * \log_2(tf_i^k) \quad (۲-۱۰)$$

^۱ Entropy

۲-۳-۲-۳-۲- بهره اطلاعات^۱

بهره‌ی اطلاعات که به اختصار Info-Gain یا IG نامیده می‌شود، یک معیار برای ارزش‌گذاری ویژگی-هاست. این معیار برای هر ویژگی نسبت به کلاس‌های مختلف مقداردهی می‌شود، نه نسبت به متون. در کارهای [۶، ۷، ۱۲] از این معیار استفاده شده است. برای کلمه‌ی t_i به عنوان یک ویژگی و کلاس c_k ، بهره اطلاعات $IG(t_i, c_k)$ ، با فرمول ۱۱-۲ تعریف می‌شود.

$$IG(t_i, c_k) = d(t_i, c_k) \log\left(\frac{d(t_i, c_k)}{d(t_i)d(c_k)}\right) + d(t'_i, c_k) \log\left(\frac{d(t'_i, c_k)}{d(t'_i)d(c_k)}\right) \quad (۱۱-۲)$$

$d(t_i, c_k)$: تعداد متون از کلاس c_k که کلمه‌ی t_i در آنها رخ داده است.

$d(t'_i, c_k)$: تعداد متون از کلاس c_k که کلمه‌ی t_i در آنها رخ نداده است.

$d(t_i, c'_k)$: تعداد متون از تمام کلاس‌ها به جز کلاس c_k که کلمه‌ی t_i در آنها رخ داده است.

$d(t'_i, c'_k)$: تعداد متون از تمام کلاس‌ها به جز کلاس c_k که کلمه‌ی t_i در آنها رخ نداده است.

$d(t_i)$: تعداد متون از کل پایگاه داده که کلمه‌ی t_i در آنها رخ داده است.

$d(t'_i)$: تعداد متون از کل پایگاه داده که کلمه‌ی t_i در آنها رخ نداده است.

$d(c_k)$: تعداد متون کلاس c_k است.

$d(c'_k)$: تعداد متون همه‌ی کلاس‌ها به غیر از کلاس c_k است.

۲-۳-۲-۳-۳- معیار آماری مربع چی^۲ یا χ^2

χ^2 یا میزان اثربخشی، درجه وابستگی و مشارکت میان یک عبارت و دسته را تعیین می‌نماید. این کاربرد را می‌توان برپایه این فرضیه دانست که عبارتی که تعداد تکرار بیشتری در یک دسته دارد، می‌تواند ویژگی متمایزکننده‌ی مناسبی برای آن دسته باشد. بنابراین ویژگی‌هایی که دارای مقدار χ^2 کمتری هستند، از اهمیت پایین‌تری برخوردار بوده و در نتیجه حذف آنان تأثیری بر دسته‌بندی نخواهد داشت و با

^۱ Information Gain

^۲ Chi Square

حذف آنان، فضای ویژگی کاهش می‌یابد، زمانی که χ^2 مقدار صفر دارد یعنی این ویژگی و کلاس کاملاً مستقل از هم هستند [۶، ۱۲، ۵۱]. رابطه‌ی ۲-۱۲ میزان اثربخشی عبارت t_i را برای کلاس c_k ، محاسبه می‌کند.

$$\chi^2(t_i, c_k) = \frac{N \cdot [d(t_i, c_k) \cdot d(t'_i, c'_k) - d(t'_i, c_k) \cdot d(t_i, c'_k)]}{d(t_i) \cdot d(t'_i) - d(c_k) \cdot d(c'_k)} \quad (12-2)$$

متد χ^2 یک متد نظارت‌شده چندمتغیره است که انتخاب عبارات موثر در تفکیک‌پذیری را نه تنها براساس تعداد تکرار آنها در هر دسته، بلکه با توجه به وابستگی عبارات با یکدیگر و نیز وابستگی عبارات با دسته‌ها انجام می‌دهد.

۲-۳-۲-۳-۴ معیار Dominance

این معیار نیز اهمیت یک عبارت را در یک کلاس نشان می‌دهد و نسبت تعداد متون کلاس K که کلمه t_i در آن رخ داده است را به کل متونی از همه‌ی کلاس‌ها که کلمه t_i در آنها رخ داده است، بیان می‌کند و مقدار آن بین صفر و یک است. اگر مقدار Dominance برای یک عبارت نسبت به یک کلاس "۱" باشد، یعنی آن عبارت فقط در آن کلاس رخ داده و در کلاس دیگری رخ نداده است [۷]. فرمول ۲-۱۳ مقدار این معیار را برای عبارت t_i را نسبت به کلاس c_k ، محاسبه می‌کند.

$$Dominance(t_i, c_k) = \frac{d(t_i, c_k)}{d(t_i, c_k) + d(t_i, c'_k)} \quad (13-2)$$

در بخش قبل در مورد انتخاب کلمات کلیدی بحث شد. در بخش بعدی به معرفی روش‌های استخراج ویژگی می‌پردازیم. استخراج ویژگی، فرآیند تشکیل ویژگی‌های جدید از ویژگی‌های انتخابی اولیه است.

۲-۳-۲-۴ روش‌های موجود در استخراج ویژگی

تحقیقات در زمینه‌ی استخراج خودکار ویژگی‌ها جهت دسته‌بندی متون به قبل از سال ۱۹۹۰ برمی-

گردد [۷]. در این بخش به متداول‌ترین روش‌ها اشاره می‌کنیم:

۲-۳-۲-۴-۱- عبارات نحوی^۱

اولین روش‌های استخراج ویژگی‌ها به آغاز دهه‌ی ۹۰ میلادی بازمی‌گردد، زمانی که برخی محققان از عبارات نحوی به عنوان ویژگی‌ها، برای دسته‌بندی متون استفاده کردند [۵۲-۵۴]. عبارات نحوی به گروهی از کلمات در معنای دستوری، از قبیل عبارات اسمی، عبارات فعلی و عبارات مصدری و غیره گفته می‌شود که با استفاده از فرآیندهای تجزیه، مثل تکنیک‌های پردازش زبان طبیعی به دست می‌آیند. بیشتر کارهای انجام شده، نشان داده‌اند که استفاده از عبارات گرامری بهتر از کیسه‌ی کلمات در دسته‌بندیها عمل می‌کند. لوئیز در [۵۳] دلایل موفقیت دسته‌بندی را بر اساس چهار ویژگی مطلوب نشان داد. این چهار ویژگی عبارتند از:

- فقدان افزونگی در شاخص‌بندی عبارات

- توزیع نسبتاً یکنواخت ارزش‌ها در شاخص‌بندی

- فقدان ابهام ناشی از شاخص‌بندی

- نویز پایین در ارزش کلمات شاخص‌بندی

لوئیز آزمایشاتی انجام داد که نشان داد عبارات نحوی، سعی در ایجاد چهارمین ویژگی دارند و نه سایر ویژگی‌ها. به دلیل این که اولاً، چون فرکانس عبارات نحوی در میان کلاس‌های مختلف، متنوع است و ثانیاً، ابهام ناشی از ایجاد مدل که نتیجه‌ی کاربرد عبارات مترادف به عنوان معیارهای هم‌ارزی معنایی در میان عبارات نحوی، که از لحاظ لغوی متفاوتند، بسیار قابل توجه است. همچنین با در نظر گرفتن هر کلمه‌ی تنها از یک عبارت T کلمه‌ای، که S مفهوم مترادف دارد، S^T عبارت نحوی با معنای مترادف دارد.

¹ Syntactic Phrases

لوئیز نتیجه گرفت که فرکانس پایین عبارات متمایز معنایی و ابعاد بالای فضای مترادفها، ظرفیت بالقوه‌ی عبارات نحوی را به عنوان یک ویژگی تفکیک کننده کاهش می‌دهد. در [۵۵] نیز، نتیجه‌ی قابل قبولی در دسته‌بندی با استفاده از عبارات نحوی بدست نیامد. علیرغم این نتایج، نویسندگان [۵۶] یک تحقیق مشابه با لوئیز با تغییراتی در جزئیات پیاده سازی، برای انتخاب ویژگی‌های بهتر انجام دادند. برای این کار از الگوریتم‌های Rule-Based استفاده کردند. فرض آنان این بود که برخی قوانین، عبارت نحوی متمایز کننده‌تری در نظر می‌گیرند و در نتیجه دسته‌بندی‌های بهتری تولید می‌کنند. با وجود این، با انجام آزمایشات نه تنها با استفاده از بهترین عبارات نحوی متمایز کننده، نتایج بهتری نسبت به کاربرد کیسه‌ی لغات بدست نیامد، بلکه کیفیت دسته‌بندی هم کاهش پیدا کرد. نویسندگان مقاله با توجه به نتایج به دست آمده، نتیجه گرفتند که ادامه‌ی کار بر روی عبارات نحوی بی نتیجه خواهد بود.

۲-۳-۲-۲-۴-۲ n-grams کلمات^۱

نویسندگان در [۵۷]، کاربرد n-grams کلمات (دنباله‌ای از n کلمه‌ی مجاور)، با اندازه‌های متفاوت را به عنوان ویژگی در دسته‌بندی متون بررسی کردند. آنها نشان دادند که n-grams های سه کلمه‌ای (و نه بیشتر از سه کلمه) نتیجه‌ای بهتر از کیسه‌ی لغات دارند. کار [۵۸] نشان داد که n-grams کلمات، نه تنها بهتر از کیسه‌ی لغات عمل نمی‌کند، بلکه در مواردی بدتر از آن نتیجه می‌دهد. مقاله [۵۹] نتایج خوبی را در کاربرد n-grams های دو کلمه‌ای گزارش کرد. اما تمرکز بروی کلاس‌هایی بود که تنها ۲,۸٪ از کل متون را تشکیل می‌داد. حتی خود نویسنده هم در پایان اظهار می‌دارد که نتایج حاصله به دلیل وسعت کم داده‌ها، خیلی قابل نتیجه‌گیری نیست. در [۶۰] از یک ریشه‌یاب برای در نظر گرفتن n-gram های متفاوت مورفولوژیکی، اما با معنی مشابه استفاده شد. در کنار این n-grams ها از ویژگی‌های منفرد تفکیک کننده نیز استفاده شد. نتایج نشان داد که این روش کارایی لازم را ندارد. [۶۱] نیز از ترکیب n-gram ها و

^۱ Word n-grams

ویژگی‌های منفرد استفاده کرد و فقط ۲٪ از n-gram‌های تفکیک کننده برای این کار استفاده شدند نتایجی بهتر از کاربرد کیسه‌ی لغات ارائه شد. در [۶۲] نیز n-gram‌ها برای دسته‌بندی ایمیل‌ها به کار گرفته شد.

به طور خلاصه مطالعه‌ی کارهای مرتبط با استفاده از n-gram‌ها، نشان داد که موفقیت‌ها در شرایط و یا داده‌های خاص بوده و یک موفقیت کلی بدست نیامده است.

۲-۳-۲-۴-۳-ویژگی‌های هم‌رخداد غیر مجاور^۱ (C-Features)

از ایده‌ی n-gram‌ها استفاده شده و با حذف محدودیت مجاورت کلمات و فاصله‌ی آنها از هم، ویژگی موثری بنام ویژگی‌های هم‌رخداد غیر مجاور معرفی شد. منظور از غیرمجاور الزاما غیر مجاور بودن نیست بلکه تاکید این است که مانند n-gram‌ها اجباری در نزدیکی و مجاورت کلمات تشکیل دهنده‌ی ویژگی وجود ندارد. ترتیب و فاصله‌ی کلمات یک ویژگی هم‌رخداد محدودیتی نخواهد داشت. این ویژگی‌ها بیشتر در حوزه‌ی بازیابی اطلاعات متنی به کار رفته‌اند و در حوزه‌ی دسته‌بندی کمتر به کار گرفته شده‌اند. در [۸] رابطه‌ی میان ویژگی‌های منفرد، با در نظر گرفتن کلمات مجزا در متون مشابه بررسی شد. فرض مسئله این بود که نویسندگان تمایل دارند که ایده‌های مرتبط را با استفاده از چندین ویژگی منفرد در عبارات‌های مختلف در متون مشابه تکرار کنند. در این کار از چکیده‌ی متون علمی که معمولا عبارات تکراری و مشخصی دارند استفاده شده است.

نویسنده در [۱۰]، در دسته‌بندی چند کلاسه‌ی متون، از ویژگی‌های هم‌رخداد پرتکرار برای توسعه‌ی متون استفاده کرده و به حذف ویژگی‌هایی که فرکانس حضور آنها در چندین کلاس بالاست، پرداخته است. در این کار نویسنده از تکنیک‌های هرس برای کاهش کلمات کاندید در ساخت ویژگی‌های هم‌رخداد

^۱ Non-adjacent Co-occurrences Features

در مقابله با تاثیر حجم نمایی فضای جستجوی کلمات استفاده کرده است، اما جزئیاتی در مورد پیاده سازی آن ارائه نکرده است.

در [۱۱] نیز کار مشابهی انجام شده و وزن مرتبط با فرکانس حضور برای ویژگی‌ها در نظر گرفته شده است. در این مقاله از ویژگی‌های هم‌رخداد در استخراج قواعد انجمنی برای دسته‌بندی‌های انجمنی^۱ استفاده شده است. در [۹]، نویسندگان نشان داد که ویژگی‌های با فرکانس متوسط بهترین و تفکیک کننده‌ترین ویژگی‌ها هستند.

مقاله [۷] نیز نشان داد ویژگی‌های کم تکرار نیز می‌توانند مهم و تفکیک کننده باشند. کارهای قبلی از ویژگی‌های هم‌رخداد با هر تعداد کلمات مختلف استفاده کرده بودند، اما مقاله‌ی مذکور از ویژگی‌های دو کلمه‌ای استفاده کرده است. در ابتدا با استفاده از معیار Info-Gain برخی ویژگی‌های منفرد را حذف کرده، سپس با معرفی معیار Dominance و در نظر گرفتن یک حد آستانه برای آن، به حذف ویژگی‌های بیشتری می‌پردازد. پس از آن با استفاده از ویژگی‌های باقیمانده به تشکیل ویژگی‌های هم‌رخداد پرداخته است. برای توسعه‌ی متون نیز از ویژگی‌های هم‌رخداد انتخابی بر اساس معیار Dominance استفاده کرده است. سپس از دسته‌بندی‌های KNN، Naïve Bayesian و SVM برای دسته‌بندی متون چهار پایگاه داده Reuters، 20 Newsgroups، OHSUMED و ACM استفاده کرده است.

۲-۳-۳- کاهش ابعاد داده

ابعاد بزرگ فضای ویژگی‌ها در دسته‌بندی متون (یعنی مقدار بزرگ $|T|$) معمولاً دردسرساز می‌باشد. در حقیقت، بزرگ شدن فضای ویژگی‌ها باعث پیچیدگی بیشتر (صرف هزینه زمانی و حافظه‌ی بیشتر) شده و از طرفی عدم وابستگی بین اطلاعات کمتر می‌گردد. داده‌هایی که به هم وابسته‌اند، عموماً ارزش دسته‌بندی ندارند. بیشتر دسته‌بندی‌ها در این موارد به بیش یادگیری برمی‌خورند. بنابراین عموماً قبل از دسته-

^۱ Associative Classifier

بندی، داده‌ها از فرآیند کاهش ابعاد داده می‌گذرند. روش‌های کاهش بعد داده‌ها در دو حوزه عمومی و محلی (داخل کلاسی) قابل بحث است.

بنا بر آنچه گفته شد، مساله کاهش بعد داده‌ها یکی از زمینه‌های اصلی تحقیق در حوزه‌ی متون است. عمده این تحقیقات بر دو محور اصلی استوار است [۴]:

- **انتخاب ویژگی:** در این روش یک مجموعه از عبارات بر مبنای تئوری اطلاعات یا ویژگی‌های آماری از متون انتخاب می‌شود. این رهیافت را "فیلترکردن" می‌نامند چراکه عبارات بی‌ربط از عبارات استخراج شده فیلتر می‌شوند [۶۳].

- **استخراج ویژگی:** در این رویکرد عبارات از طریق تابع تبدیل خاصی به یک فضای ویژگی جدید منتقل می‌شوند. این روشها با ایجاد عبارات ترکیبی جدید از مجموعه اصلی، یک فضای ویژگی جدید ساخته و سعی می‌کنند تا با انجام جایگزینی کلمات با مفاهیمشان، بعد فضای عبارت را کاهش دهند. به طوری که ممکن است ویژگی‌های نهایی کاملاً متفاوت از ترم‌های اولیه باشند.

از شیوه‌های متداول در کاهش ویژگی‌ها، روش‌های PCA^۱ [۶۴، ۶۵]، LDA^۲ [۶۶] و ICA^۳ [۶۷] را می‌توان نام برد که در دامنه‌های متفاوت کاری و بر روی داده‌های مختلف به خوبی عمل می‌کند. این روش‌ها در دسته‌بندی متون نیز به کار گرفته شده اند.

۲-۳-۴- الگوریتم‌های دسته‌بندی

الگوریتم‌هایی که بصورت خودکار ویژگی‌های دسته‌ها را از متون نمونه یاد می‌گیرند، الگوریتم‌های با ناظر نامیده می‌شوند. در فاز آموزش این الگوریتم‌ها، هر متن تعدادی ویژگی خواهد داشت که در حوزه

¹ Principal Component Analysis

² Linear Discriminant Analysis

³ Independent Component Analysis

بازیابی اطلاعات به آن عبارت^۱ می‌گویند و منجر به ساخت یک مدل دسته‌بند می‌گردد. در این کار منظور از ویژگی، همان کلمات کلیدی استخراج‌شده و یا ویژگی‌های هم‌رخداد از متون است.

الگوریتم‌های دسته‌بند باناظر را می‌توان به سه دسته کلی تقسیم نمود [۱]:

۱. **دسته‌بندهای مبتنی بر قاعده**^۲: این دسته‌بندها با استنباط یک مجموعه از قواعد از اسناد که از پیش دسته‌بندی شده‌اند، یادگیری را انجام می‌دهند. الگوریتم ریپر^۳ از جمله دسته‌بندهای مبتنی بر قاعده است. قواعد تصمیم ممکن است بصورت درخت تصمیم مانند C4.5 در مرجع [۶۸] باشند.

۲. **دسته‌بندهای خطی**^۴: در این دسته‌بندها، برای هر دسته یک نمایه محاسبه می‌شود که در واقع برداری از اوزان، براساس فرکانس و احتمالات حضور ویژگی خاص است. برای هر دسته و هر سند، امتیاز براساس ویژگی‌های دسته و سند محاسبه می‌شود. از دسته‌بندهای این دسته می‌توان به دسته‌بند بیز اشاره نمود که براساس تخمین شرایط احتمالی عمل می‌کند [۶۹]. ماشین‌های بردار پشتیبان [۱۳]، دسته‌بند خطی بهینه را با استفاده از انتقال به فضای ویژگی بدست می‌آورند. این گروه از دسته‌بندها همچنین شامل الگوریتم‌های یادگیری اکتشافی از هوش مصنوعی مانند پرسپترون می‌باشند که در آن اوزان از راهی پیچیده‌تر بدست می‌آیند [۷۰].

۳. **دسته‌بندهای مبتنی بر نمونه**^۵: این دسته‌بندها یک سند جدید را با پیدا کردن K تا از اسناد نزدیکتر به آن در مجموعه آموزشی، دسته‌بندی می‌کنند و با رأی‌گیری، آن را به دسته‌ی نزدیکتر تخصیص می‌دهند [۷۱] مثل KNN.

¹ Term

² Rule Based Classifiers

³ Ripper

⁴ Linear Classifiers

⁵ Example Based Classifier

۲-۴- مقالات موجود در دسته‌بندی متون

تحقیقات بسیاری در زمینه دسته‌بندی متون انجام شده است که برخی از آنها را به طور خلاصه می‌آوریم. برای مطالعه‌ی بهتر این بخش، مقالات را به دو دسته مقالات فارسی و غیر فارسی تقسیم می‌کنیم.

۲-۴-۱- مقالات مبتنی بر زبان‌های غیر فارسی

در سال ۲۰۱۰ جیانگ^۱ و همکاران در [۷۲]، ادعا کردند روش مبتنی بر گرافی که ارائه داده‌اند، از روش مجموعه کلمات در دسته‌بندی‌های مختلف کارایی بهتری دارد. در این روش یک الگوریتم کاوش گراف وزن‌دار، زیر گراف‌ها را به صورت مکرر استخراج می‌کند. سپس مهم‌ترین زیرگراف‌ها به عنوان ویژگی برای هر متن انتخاب می‌شوند.

کوئی^۲ و همکاران [۷۳]، پس از انتخاب ویژگی‌های هر کلاس به صورت محلی، با استفاده از معیارهای TF، DF و IDF، با اعمال روش CHI بر روی ویژگی‌های محلی، یک ویژگی کلی به دست آورده که نتایج خوبی را به بار آورده است.

همین‌طور هو^۳ و همکاران [۶۷]، از تحلیل اجزای مستقل (ICA) برای انتخاب ویژگی‌ها استفاده کردند. آنها این روش را با برخی روش‌های سنتی دیگر مثل CHI، DF، IG و MI ترکیب کردند و بهترین نتیجه را از ترکیب ICA با CHI به دست آوردند.

چنگ^۴ و همکاران در [۷۴]، روشی بنام TEERS، یک فرمول جدید برای اهمیت ویژگی‌ها ارائه کرده و بر مبنای قابلیت دسته‌بندی، ویژگی‌های شرطی ارائه کردند. پس از کاربرد فرمول فوق از تحلیل میزان همبستگی ویژگی‌ها برای کاهش موثر ویژگی‌ها استفاده کردند. آنها در گام اول نیز Wordnet را برای رفع مشکل کلمات هم‌معنی و کاهش ابعاد داده به کار بردند.

¹ Jiang

² Li-Quing Qui

³ M.Hu

⁴ Y.Cheng

پارک^۱ نیز در [۷۵]، نشان داد که در دسته‌بندی متون چند برچسبه، خروجی‌های روش‌های خطی چندگانه می‌تواند برای ساختن ویژگی‌های جدید استفاده شود و از آن برای نمایش متن در ابعاد کمتر استفاده نمود. برای این کار از خروجی‌های SVM و LSR^۲ در کاهش بعد استفاده کرده و در نهایت از ترکیب آنها ویژگی‌های جدیدی ارائه کرد.

ژو^۳ و همکاران [۷۶]، نشان دادند که TFIDF سنتی تنها به خصوصیت تکرار یا بسامد توجه کرده است و توجهی به مکان قرار گرفتن آن‌ها نکرده است. لذا با در نظر گرفتن خصوصیات همچون: بسامد، موقعیت و روابط متقابل، یک شیوه‌ی جدید وزن‌دهی بنام TF-IDF-Rel را بر مبنای توسعه کلمات موضوعی ارائه کردند.

در سال ۲۰۰۷، گوپال^۴ و همکاران [۷۷]، با ترکیب بیز ساده و شبکه عصبی پرسپترون، بر مشکل توزیع بد داده در مجموعه آموزشی فائق آمده و نتایج خوبی را گزارش کرده‌اند.

در همین سال، چن^۵ و همکاران [۷۸]، از پرسپترون افزایشی برای دسته‌بندی متن استفاده کردند. مدل ارائه شده دارای این قابلیت است که با افزایش تعداد کلاس‌ها مشکلی در آموزش شبکه ایجاد نشوند.

در سال ۲۰۰۸، وانگ^۶ و همکاران [۶۶]، برای دسته‌بندی متون، ابتدا بردار ویژگی هر متن را که دارای ابعاد زیادی است، با استفاده از الگوریتم LDA کاهش داده سپس با استفاده از الگوریتم SVM به طبق-بندی داده‌ها پرداخته است.

در همین سال فنگ و وارن^۱ [۷۹]، با ترکیب الگوریتم ژنتیک و شبیه‌سازی تبرید و نیز با استفاده از نوعی الگوریتم فیدبک بنام Rochio، سیستمی را ارائه کردند که با تولید قوانینی سعی در دسته‌بندی متون دارد.

¹ C.H.Park

² Least Square Regression

³ Zhou

⁴ Goyal

⁵ Chen

⁶ Zigiang Wang

سو^۲ و همکاران [۸۰]، طی مقاله‌ای یک روش جدید بنام AC-Closed معرفی می‌کردند که پس از استخراج یک سری قوانین از دادگان آموزشی، به طبقه‌بندی متون می‌پردازد.

شی^۳ و همکاران [۸۱]، ابتدا از استفاده از شبکه عصبی مصنوعی ISOMAP، به کاهش ابعاد ویژگی-های دادگان آموزشی پرداختند و سپس با استفاده از SVM تک کلاسه به طبقه‌بندی چند کلاسه‌ی داده‌ها پرداختند. آن‌ها ادعا کرده‌اند که کارایی الگوریتم، بدون تغییر در میزان درستی آن افزایش یافته است.

ژو^۴ [۸۲]، روشی ارائه کرد که از مجموعه‌های سخت برای استخراج و کاهش ابعاد استفاده نمود. حاصل این استخراج ویژگی به یک شبکه عصبی FCM داده شد که در آن روابط غیر مهم به روابط همبستگی تغییر داده شده‌اند.

وانگ^۵ و همکاران [۸۳]، برای دسته‌بندی متن، ابتدا از مجموعه‌های سخت برای کاهش ابعاد فضای آموزشی استفاده کردند. سپس با استفاده از الگوریتم ژنتیک به انتخاب ویژگی پرداختند. عمل دسته‌بندی را یک SVM انجام می‌دهد که پارامترهای آن توسط الگوریتم ژنتیک بهینه انتخاب شده‌اند.

یانگ^۶ و همکاران [۸۴]، یک چارچوب کلی برای الگوریتم‌های دسته‌بندی کننده‌ی متن برای استفاده از آنتولوژی پیاده‌سازی کرده و ادعا کردند کارایی خوبی هم دارد.

دینگ^۷ و همکاران [۸۵]، با تغییراتی در الگوریتم بیز ساده، شرط عدم همبستگی داده‌ها با یکدیگر را در آن کمرنگ کرده‌اند. با این تغییر، الگوریتم در مورد داده‌های دارای همبستگی نیز خوب عمل می‌کند.

¹ ZHU Zhen-fang Lu Run

² Zhitong Su

³ Lukui Shi

⁴ Yu

⁵ Wang

⁶ Yang

⁷ Wang Ding

ژو^۱ و همکار [۸۶]، طی مقاله‌ای از مدل تخصیص پنهان دیریکله برای کاهش ابعاد و نیز برای دسته‌بندی استفاده کردند. آزمایش‌ها نشان داده که روش ارائه شده از بیز ساده و مشابه SVM عمل کرده است.

چن^۲ و همکاران [۸۷]، برای دسته‌بندی متون، از مجموعی از دسته‌بندها استفاده کردند. این مجمع شامل سه دسته‌بند SVM، KKN و Centroid است. در این روش داده‌ها به دو کلاس تقسیم می‌شوند. داده‌هایی که به آسانی کلاسشان مشخص می‌شود. داده‌هایی که تعیین کلاسشان سخت‌تر است. برای گروه اول یک دسته‌بند تصمیم می‌گیرد و برچسب آن را مشخص می‌کند اما برای گروه دوم از رای اکثریت استفاده می‌شود.

یوان^۳ و همکاران [۸۸]، برای دسته‌بندی متن از ترکیب SVM و KNN استفاده کردند. به این صورت که ابتدا با SVM مرزها را مشخص کرده و سپس در داخل مرزها با استفاده از KNN دسته‌بندی را انجام می‌دهند.

لیو^۴ و همکاران [۸۹]، برای دسته‌بندی متون از ترکیب درخت تصمیم دودویی با SVM بهره بردند. همچنین برای کاهش فضای حالت و بالابردن دقت و سرعت از خوشه‌بندی استفاده کرده‌اند.

چن^۵ و همکاران [۹۰]، با ارائه‌ی دو معیار ارزش‌گذاری ویژگی بنام‌های MOR^۶ و CDM^۷ به استخراج و کاهش ویژگی‌ها پرداخته و از دسته‌بند بیزین ساده برای دسته‌بندی استفاده کردند. آنها نشان دادند که این معیارهای انتخاب ویژگی، از سایر روش‌ها بهتر عمل می‌کند.

¹ Zhou

² Qingxuan Chen

³ Pingpeng Yuan

⁴ Lizhen Liu

⁵ Jingnian Chen

⁶ Multi-class Odds Ratio

⁷ Class Discriminating Measure

وانگ^۱ و همکاران [۹۱]، با استفاده از مفاهیم و ارتباط کلمات و مترادفات و روابط انجمنی در ویکی‌پدیا به توسعه‌ی ویژگی‌ها پرداخته و نشان داده که این روش در مقایسه با سایر روش‌های بازنمایی متن قوی‌تر است.

سان^۲ و همکاران [۹۲]، یک مطالعه‌ی مقایسه‌ای بر روی روش‌هایی که با داده‌های نامتوازن در کلاس‌ها کار می‌کنند مثل *resampling*، *instance weighting* و ...، انجام دادند و با استفاده از دسته‌بند SVM به بررسی ده روش موجود در زمینه‌ی دسته‌بندی نامتوازن^۳ متون پرداختند.

اسیرام^۴ و همکاران در [۹۳]، بدلیل کوتاه بودن متون در سرویس‌هایی همچون توئیتر، با استفاده از یک مجموعه از کلمات از پیش تعریف شده در دامنه‌های خاص، به دسته‌بندی متون کوتاه در توئیتر پرداخته، تا مشکل خام بودن متون مورد استفاده کاربران را در این سرویس‌ها از بین ببرد.

جیانگ^۵ و همکاران [۹۴]، از گراف‌ها در دسته‌بندی متون استفاده کرده‌اند. آنها با استفاده از یک الگوریتم کاوش گراف وزن‌دار^۶، به جستجوی زیرگراف‌ها به عنوان ویژگی‌های متون پرداخته‌اند. این روش با استفاده از چند دسته‌بند مختلف نتایج خوبی را نشان داده است و نشان داده شده است که در مقابل افزایش داده‌ها و بزرگ شدن گراف‌ها نیز خوب عمل می‌کند.

شینگ هوا^۷ و همکاران [۹۵]، با ترکیبی از دو دسته‌بند بیزین و دسته‌بندهای انجمنی ضمن رفع نقصان هریک از آنها به صورت منفرد، دقت بالاتری را در دسته‌بندی ارائه می‌دهند. آنها متون آموزشی را با دسته‌بند بیزین دسته‌بندی کرده و حد آستانه‌های اطمینان را برای کلاس‌های متفاوت با استفاده از داده‌های دسته‌بندی شده با بیزین، برای انتخاب قوانین انجمنی تنظیم کرده‌اند. سپس نتایج دسته‌بندی را با

¹ Pu Wang

² Sun

³ Imbalanced Text Classification

⁴ Sriram

⁵ Jiang

⁶ Weighted Graph Mining Algorithm

⁷ Shing Hwa

استفاده از این قوانین انجمنی کلاس‌ها بهینه کرده‌اند. نتایج نشان داده است که ترکیب دو دسته‌بندی، بهتر از هر یک از آنها نتیجه می‌دهد.

لی^۱ و همکاران [۹۶]، نشان می‌دهند که داده‌های آموزشی منفی اثر بسیار بدی در دسته‌بندی دارند. ابتدا داده‌های آموزشی مثبت و منفی با فرض طور توزیع یکنواخت، به دو دسته مثبت و منفی دسته‌بندی شده‌اند. در نهایت نشان داده شده است که داده‌های مثبت توزیع یکنواخت دارند، اما این یکنواختی توزیع در داده‌های منفی وجود ندارد. همچنین نشان می‌دهند که اثر منفی و مخرب داده‌های منفی، در تنظیم مقادیر و حدود آستانه در معیارهای مختلف است.

خان^۲ و همکاران در [۲]، یک مطالعه‌ی مروری بر روش‌های موجود در دسته‌بندی متون و استخراج ویژگی‌ها انجام داده‌اند. در پایان نشان داده‌اند که پرکاربردترین روش‌های استخراج ویژگی استفاده از Chi Square و InfoGain هستند، گرچه روش‌های دیگر استخراج ویژگی که نتایج بهتری را نیز نشان داده‌اند معرفی شده و نیاز به توسعه‌ی بیشتر دارند. همچنین پرکاربردترین دسته‌بندها، KNN، NB و SVM معرفی شده‌اند که روش‌های ترکیبی معمولاً نتایج بهتری را نشان داده‌اند.

ژانگ^۳ و همکاران [۹۷]، نیز یک مطالعه‌ی مروری در بازنمایی متن بر روی روش‌های LSI، TFIDF و multi-word داشته‌اند. نتایج نشان داده است که LSI در دسته‌بندی متون انگلیسی و چینی با دقت بیشتری عمل می‌کند. همچنین در بازیابی متون انگلیسی نیز بهتر عمل کرده است.

جیانگ^۴ و همکاران [۹۸]، برای کاهش ویژگی‌ها، با استفاده از یک الگوریتم خودسازی بر مبنای شباهت فازی^۱، به خوشه‌بندی ویژگی‌ها پرداخته‌اند. کلمه‌های شبیه به هم در خوشه‌های یکسان قرار

¹ Li
² Khan
³ Zhang
⁴ Jiang

گرفته‌اند. هر خوشه با یک تابع عضویت و یک میانگین و انحراف مشخص شده است. در نهایت ویژگی استخراج شده متناظر با هر خوشه، ترکیب وزنی کلمات هر خوشه می‌باشد. نتایج نشان داده است که این الگوریتم سریعتر و با کارایی بیشتر نسبت به سایر روش‌ها عمل می‌کند.

هارون^۲ و همکاران [۹۹]، با استفاده از روش شاخص‌گذاری معنایی پنهان، یک روش دومرحله‌ای برای استخراج ویژگی‌ها ارائه دادند. در مرحله‌ی اول با ارائه‌ی یک روش کاهش ویژگی، به انتخاب ویژگی‌ها پرداخته، سپس یک فضای معنایی جدید بین کلمات با استفاده از روش شاخص‌گذاری معنایی پنهان ایجاد کرده است. این روش بر روی کاربردهایی که با فیلتر کردن اسپم‌ها سر و کار دارند، بهتر از سایر روش‌ها نتیجه داده است.

وان^۳ و همکاران [۱۰۰]، ترکیبی از دسته‌بند KNN و SVM را برای دسته‌بندی متن با عنوان SVM-NN ارائه می‌دهد. این دسته‌بند ترکیبی برای به حداقل رساندن تاثیر پارامتر K در دسته‌بند KNN ارائه شده است. در مرحله‌ی آموزش، SVM به کار گرفته شده است تا نمونه‌های آموزشی هر دسته به بردارهای پشتیبان هر کلاس کاهش پیدا کنند. تصمیم‌گیری دسته‌بندی بر اساس دسته‌ای است که بردارهای پشتیبان آن کمترین فاصله را به داده‌ی تست دارند.

کورسات^۴ و همکاران [۱۰۱]، برای کاهش ویژگی‌ها، از یک روش جدید انتخاب ویژگی احتمالاتی، بر اساس فیلتر کردن، بنام DFS^۵ برای دسته‌بندی متون ارائه می‌دهند. این روش با روش‌های معروفی همچون Information Gain و Chi Square و غیره مقایسه شده است.

¹ Fuzzy Similarity-Based Self-Constructing Algorithm

² Harun

³ Wan

⁴ Kursat

⁵ Distinguishing Feature Selector

یانگ^۱ و همکاران [۱۰۲]، الگوریتم CMFS را به عنوان یک روش نوین در انتخاب و کاهش ویژگی‌ها ارائه کردند. این الگوریتم به طور جامع میزان تاثیر یک کلمه را در داخل و خارج یک کلاس اندازه‌گیری می‌کند. این الگوریتم بر روی داده‌های چند مجموعه با استفاده از SVM و بیزین ساده آزمایش شده است. نتایج دسته‌بندی با شش روش معروف و شناخته شده مثل IG، CHI، DF و غیره مقایسه شده است و با استفاده از هر دو دسته‌بند، از همه‌ی روش‌ها نتیجه‌ی بهتری داشته است.

جیانگ^۲ و همکاران در [۱۰۳]، برای افزایش کارایی الگوریتم دسته‌بندی KNN، یک نسخه‌ی بهبود یافته از این الگوریتم را ارائه داده است. این الگوریتم ترکیبی از یک الگوریتم خوشه‌بندی یک مرحله‌ای محدود و دسته‌بند KNN است. در آزمایشات بر روی داده‌های مختلف، این دسته‌بند از سه دسته‌بند معروف KNN، SVM و NB نتایج بهتری داشته است.

صدیق^۳ و همکاران در [۱۰۴]، پس از پیش‌پردازش و انتخاب ویژگی‌های تفکیک‌کننده‌ی مناسب با استفاده از تئوری Rough Set به دسته‌بندی متون می‌پردازد. بنابراین این مدل از یک زوج مفاهیم دقیق بر اساس تئوری موجود، بنام تقریب‌های بالا و پایین^۴ برای دسته‌بندی متون آزمایشی به یک دسته یا بیشتر و حتی زیر دسته‌ها، با دقت بسیار خوب استفاده می‌کند.

سلامت^۵ و همکاران در [۱۰۵]، روش WPCM^۶ را برای دسته‌بندی متون خبری وب ارائه می‌دهد. در این کار از ویژگی‌های CPBF^۷ و دسته‌بند شبکه‌ی عصبی برای دسته‌بندی استفاده می‌کند و برای کاهش ویژگی‌ها نیز از روش PCA استفاده می‌کند. آزمایشات انجام شده بر روی متون خبری ورزشی وب، با دقت خوبی گزارش شده است.

¹ Yang

² Jiang

³ Sadiq

⁴ Lower and Upper Approximations

⁵ Selamat

⁶ Web Page Classification Method

⁷ Class Profile-Based Features

بیلوسوس^۱ و همکاران در [۱۰۶]، یک عامل هوشمند طراحی کرده است که اخبار روزانه را برای کاربران گردآوری می‌کند. این سیستم بر مبنای بازخورد کاربر، اخبار مورد علاقه‌ی کاربر را جمع‌آوری می‌کند. تمرکز این کار بروی اجزای مدل کردن کاربر سیستم است. این سیستم علاقه‌های کوتاه مدت و بلند مدت کاربر را نیز در نظر می‌گیرد. برای این کار از KNN برای دسته‌بندی علاقه‌های کوتاه مدت و از Naïve Bayesian برای دسته‌بندی علاقه‌مندی بلند مدت کاربر استفاده می‌کند. متد وزن‌دهی بکار رفته نیز TFIDF برای ویژگی‌های استخراجی است. نتایج ارائه شده در این کار رضایت‌بخش نبوده است. ژائو لیو^۲ و همکار در [۱۰۷]، توانسته است با بکارگیری ترکیبی از کلمات منفرد و کلمات هم‌رخداد و نیز ترکیب دسته‌بند SOM و KNN به نتایج بسیار خوبی در دسته‌بندی متون دست پیدا کند. نویسنده ادعا کرده است که این روش با کاهش بعد بالا راه حل مناسبی برای مجموعه داده‌های بزرگ است. پکار^۳ و همکاران در [۱۰۸]، با استفاده از ویژگی‌های هم‌رخداد به دسته‌بندی کلمات بر اساس معنای آن‌ها پرداختند. در این تحقیق از دسته‌بندهای بیزین و KNN چندین روش وزن‌دهی مثل Information Gain, Odds Ratio, Gain Ratio و Mutual Information باهم مقایسه شده‌اند. نشان داده شده است که ویژگی‌های کلاسی بسیار موثرتر از ویژگی‌های منفرد عمل می‌کنند.

۲-۴-۲- مقالات مبتنی بر زبان فارسی

تحقیقات انجام شده در دسته‌بندی متون فارسی به گستردگی سایر زبان‌ها نیست. در سال ۱۳۸۶ خانم عربی و همکاران [۱۰۹]، در مقاله خود به استخراج کلمات کلیدی و استفاده از آن برای دسته‌بندی متون فارسی پرداختند. آن‌ها برای این کار از ترکیب روش‌های Porter و Wordnet

¹ Billsus

² Xiao Luo

³ Pekar

استفاده کرده‌اند. روش آن‌ها برای استخراج کلمات کلیدی استفاده از ساختار درهم ریزی Trie است و دادگان مورد استفاده در این مقاله، دادگان همشهری بوده است.

در مقاله دیگری که در همین سال و توسط همین گروه صورت گرفت [۱۱۰]، از الگوریتم کاوشگر مورچه^۱ برای استخراج قوانین دسته‌بندی از میان داده‌ها استفاده شده است.

کار دیگری که در همین سال انجام شده است مربوط به احسان بصیری و همکاران [۱۱۱] است. آن‌ها در مقاله‌ای که ارائه کردند، دو الگوریتم KNN و FKNN را با دو روش استخراج ویژگی IG و DF آزمودند، آن‌ها نشان دادند که بهترین ترکیب استفاده از روش FKNN و IG است. با این ترکیب آن‌ها به میانگین دقت ۸۰٪ دست پیدا کرده‌اند.

در سال ۱۳۸۷، مسلم محمدی و همکاران در [۱۱۲]، از شبکه عصبی CC4 بهره جسته‌اند. CC4، یک نوع از الگوریتم‌های دسته‌بند کناره‌ای برای شبکه عصبی سه لایه رو به جلو می‌باشد. این شبکه برای دسته‌بندی اسناد حجیم فارسی مناسب است. هم‌چنین اگر مقیاس بین فایل‌های متنی نزدیک به هم باشد، این شبکه نتایج بهتری ارائه می‌دهد. در این روش در مرحله پیش‌پردازش کلمات عمومی حذف و سایر کلمات ریشه‌یاب می‌شوند. سپس ویژگی‌های آماری برای کلمات مختلف محاسبه شده، با استفاده از یک روش مبتنی بر فراوانی کلمه، اسناد با اندازه‌های مختلف به یک فضای K بعدی با اندازه‌ی ثابت نگاشت می‌شوند. دقت این دسته‌بند روی دادگان ایسنا حدود ۹۰٪ گزارش شده است.

در همین سال، ایوب باقری و همکاران در [۱۱۳]، با استفاده از بیز ساده و کاهش فضای حالت با روش TFV^۲ دسته‌بندی قابل قبولی ارائه دادند.

سید مهدی نقیبی و همکاران [۱۱۴]، نیز در همین سال برای دسته‌بندی متون از شبکه عصبی استفاده کرده‌اند. به این صورت که به تعداد کلاس‌های از پیش تعیین شده شبکه MLP باینری وجود

^۱ Ant_miming

^۲ Term Frequency Variance

دارد. متن ورودی به هر یک از این شبکه‌ها وارد می‌شود. خروجی این شبکه‌ها به یک MLP ترکیب کننده‌ی نتایج رفته و به این ترتیب برچسب متن ورودی مشخص می‌شود. آنها همچنین با پیشنهاد روش بسامد سند موضوعی ادعا کرده‌اند که تا ۹۹٪ منجر به کاهش فضای حالت شده‌اند.

همچنین در مقاله دیگری که بینا و همکاران [۱۱۵]، ارائه دادند، روش‌های tri-gram و quad-gram با معیارهای اندازه‌گیری فاصله منهتن^۱، اندازه‌ی دایس^۲ و ضرب نقطه‌ای، به همراه روش یادگیری KNN مورد بررسی قرار گرفته و بهترین نتیجه از ترکیب quad-gram و ضرب نقطه‌ای به دست آمد.

مقصودی و همکار [۱۱۶]، برای افزایش دقت دسته‌بندی از گنجواژه‌ی معنایی بهره گرفتند. آنها با استفاده از معیارهای ICF^۳ و Uni^۴، ویژگی‌ها را انتخاب کرده و با گسترش آنها با استفاده از گنجواژه و وزن‌دهی TFIDF و نیز کاربرد دسته‌بند SVM به دسته‌بندی متون پرداختند. آنها آزمایشات خود را در تعداد کلاس‌های مختلف و با تعداد متون مختلف انجام دادند و به نتایج پایداری در این آزمایشات دست یافتند.

سرخی و همکار [۱۱۷]، از روش وزن دهی ویژگی مبتنی بر اطلاعات کلاس درحوزه طبقه بندی مستندات TFCRF برای انتخاب ویژگیها استفاده کردند. برای استخراج ویژگیهای متون فارسی نیز از روش تجزیه و تحلیل مولفه های اصلی PCA استفاده کردند. نتایج به دست آمده نشان داده است که با در نظر گرفتن روشهای دسته‌بندی KNN و Bayesian در روش پیشنهادی بهبود قابل توجهی در طبقه بندی متون فارسی و کاهش مدت زمان تست با ویژگیهای استخراج شده به دست خواهد آمد.

حسن پور و همکاران در [۱۱۸]، با استفاده از روش PCA و الگوریتم ژنتیک، به استخراج و کاهش ویژگی‌ها پرداخته‌اند. آنها روش وزن‌دهی TFCRF را با استفاده از فاکتورهای ارتباط مثبت و منفی، به کار

¹ Manhattan Distance

² Dice Measure

³ Inverted Conformity Frequency

⁴ Uniformity

برده‌اند. فاکتور ارتباط مثبت، نسبت تعداد متون یک دسته که حاوی ویژگی مورد نظر هستند را به کل متون آن دسته نشان می‌دهد. فاکتور ارتباط منفی نیز، نسبت تعداد متون سایر دسته‌ها که حاوی ویژگی مورد نظر هستند را به کل متون سایر دسته‌ها نشان می‌دهد. آنها با دو دسته‌بند KNN و بیزین به دسته-بندی متون پرداختند و نسبت به کارهای قبلی نتایج بهتری گرفتند.

طاهری‌نیا [۱۱۹]، به بررسی یکی از روش‌های یادگیری هوشمند به نام یادگیری نیمه نظارت شده در دسته بندی متون فارسی پرداخت. بسیاری از روش‌های یادگیری هوشمندانه مانند یادگیری نظارت شده، فقط بر روی داده‌های آموزشی برجسب دار تکیه می‌کنند، در شرایطی که بدست آوردن این داده‌های آموزشی دارای برجسب بسیار پرهزینه است. حال آنکه حجم زیادی از داده‌های بدون برجسب به سرعت زیاد و با هزینه‌ی کم در دسترس هستند. در مقابل روش‌هایی مانند روش یادگیری بدون نظارت فقط بر روی داده‌های بدون برجسب تکیه می‌کنند. وی به بررسی روش یادگیری نیمه نظارت شده که ما بین روش‌های یادگیری نظارت شده و یادگیری بدون نظارت قرار دارد و از ترکیبی از مثال‌های آموزشی برجسب دار و بدون برجسب برای یادگیری استفاده می‌کند پرداخته و از این تکنیک برای دسته بندی متون فارسی استفاده کرده و نتایج نسبتاً خوبی گرفت.

زمانی و همکاران [۱۲۰]، از روش آنالیز معنایی پنهان احتمالاتی و از مبانی احتمالات و قوانین بیزین، برای عملیات دسته بندی متون فارسی پایگاه داده همشهری و با استفاده از روش فاصله یابی اقلیدسی در فضای ماتریس‌های کاهش بعد یافته استفاده کرده است. همچنین روش TFIDF را جهت انتخاب و وزن-دهی کلمات کلیدی موجود در مجموعه متون مورد استفاده، به کار گرفته است. در مرحله هرس، علاوه بر کلمات دستوری مانند حروف اضافه و ربط، کلمات کلیدی زاید نیز به صورت دستی حذف شده‌اند. نتایج حاصل شده در این مقاله نشان می‌دهد که روش بکار رفته در بهبود عملکرد سیستم نقش موثری دارد.

برفامی و همکار [۱۲۱]، از ترکیب دوشبکه عصبی چندلایه پرسپترون MLP در دسته بندی مستندات نیمه ساختیافته XML بر روی پایگاه داده روزنامه همشهری استفاده کرده اند. روش TFCRF نیز برای وزن دهی به مستندات انتخاب شده سپس به کمک روش تجزیه و تحلیل مولفه های اصلی PCA به ارزش دهی ویژگیها پرداخته و با الگوریتم KNN تعدادی از آنها انتخاب و به عنوان ورودی یکی از شبکه های عصبی استفاده شده اند. در مرحله بعد از ویژگیهای ارزش دهی شده به کمک الگوریتم GA تعدادی ویژگی انتخاب شده و به عنوان ورودی یکی دیگر از شبکه های عصبی استفاده میشود سپس نتایج خروجی این دوشبکه به کمک جمع جبری بایکدیگر ترکیب شده در نهایت مورد تست و ارزیابی قرار گرفته است.

عابدی نیا و همکاران در [۱۲۲]، دسته بندی لایه ای اطلاعات به عنوان روش اصلی مطرح کرده است و نیز از مدل فضای برداری برای دسته بندی متون فارسی در وب استفاده کرده است. سپس روش جدیدی برای کاوش در متون فارسی بر اساس تحلیل معنایی کلمات فارسی مطرح می کند که شامل سه بخش استخراج کلمات تحلیل ترکیبات جملات و تحلیل های معنایی جملات است روش پیشنهادی این امکان را به کاربران فارسی زبان می دهد که بادقت و سرعت بیشتری متن مورد نظر خود را دروب بیابند.

صفیان و همکاران در [۱۲۳]، با بیان این مسئله که روش سنتی KNN پیچیدگی محاسباتی بالایی دارد و این نقیصه از کاربردی بودن KNN می کاهد، خوشه بندی را یک راه حل برای این مشکل می داند. به این روش که تعداد نمونه های یادگیری در دسته بندی را کاهش داده و به این طریق از پیچیدگی زمانی محاسبات می کاهد. در این مقاله، یکی از روشهای ساده مبتنی برای خوشه بندی الگوریتم KNN که قبلا پیشنهاد شده است، انتخاب شده است. این روش بوسیله خوشه بندی پیچیدگی زمانی را تا حد زیادی کاهش می دهد ولی مقداری از کارایی دسته بندی می کاهد. لذا یک مقدار وزنی به مرکز خوشه ها نسبت داده شده و همچنین مقدار K در الگوریتم KNN کاهش داده شده است. آزمایشها بر روی مجموعه داده Reuter-21578 نشان داده که این روش می تواند کارایی دسته بندی را افزایش دهد.

حسن‌پور و همکاران [۵۱]، از مشخصه معنایی و توجه به بافت معنایی متون، استفاده کرده‌اند. در رویکرد پیشنهادی، ابتدا مفاهیم توسط هستان شناسی لغوی فارس نت استخراج شده و پس از اعمال فرایند رفع ابهام، با توجه به بافت متن، اوزان مفهوم مرتبط در حین فرایند وزن دهی به وزن عبارات کلیدی اضافه می‌شود. این مقاله، رویکردی جدید در حین استفاده از روش اولین مفهوم به منظور رفع ابهام از مفاهیم استخراج شده، معرفی می‌کند که کارایی این روش را بهبود می‌دهد. معیار χ^2 نیز برای انتخاب ویژگی و روش وزن دهی ویژگی نرمال شده TFIDF در وزن دهی به کار گرفته شده است. نتایج نشان داده است که اعمال دانش معنایی حاصل از هستان شناسی لغوی در فرایند طبقه بندی متون یک گام موثر در بهبود کارایی دسته بند است.

۲-۵ نتیجه‌گیری

در این فصل ضمن تعریف مفاهیم پایه در دسته‌بندی متون، مراحل مختلف سیستم دسته‌بندی متون معرفی شد. انواع روش‌های استخراج و کاهش ویژگی و نیز روش‌های وزن‌دهی به ویژگی‌ها که در تحقیقات مختلف به کار رفته‌اند مطرح گردید. متداول‌ترین روش‌های دسته‌بندی متون نیز معرفی بیان شد. در پایان نیز مروری بر مقالات فارسی و سایر زبان‌ها در حوزه‌ی دسته‌بندی متون انجام شد.

فصل سوم

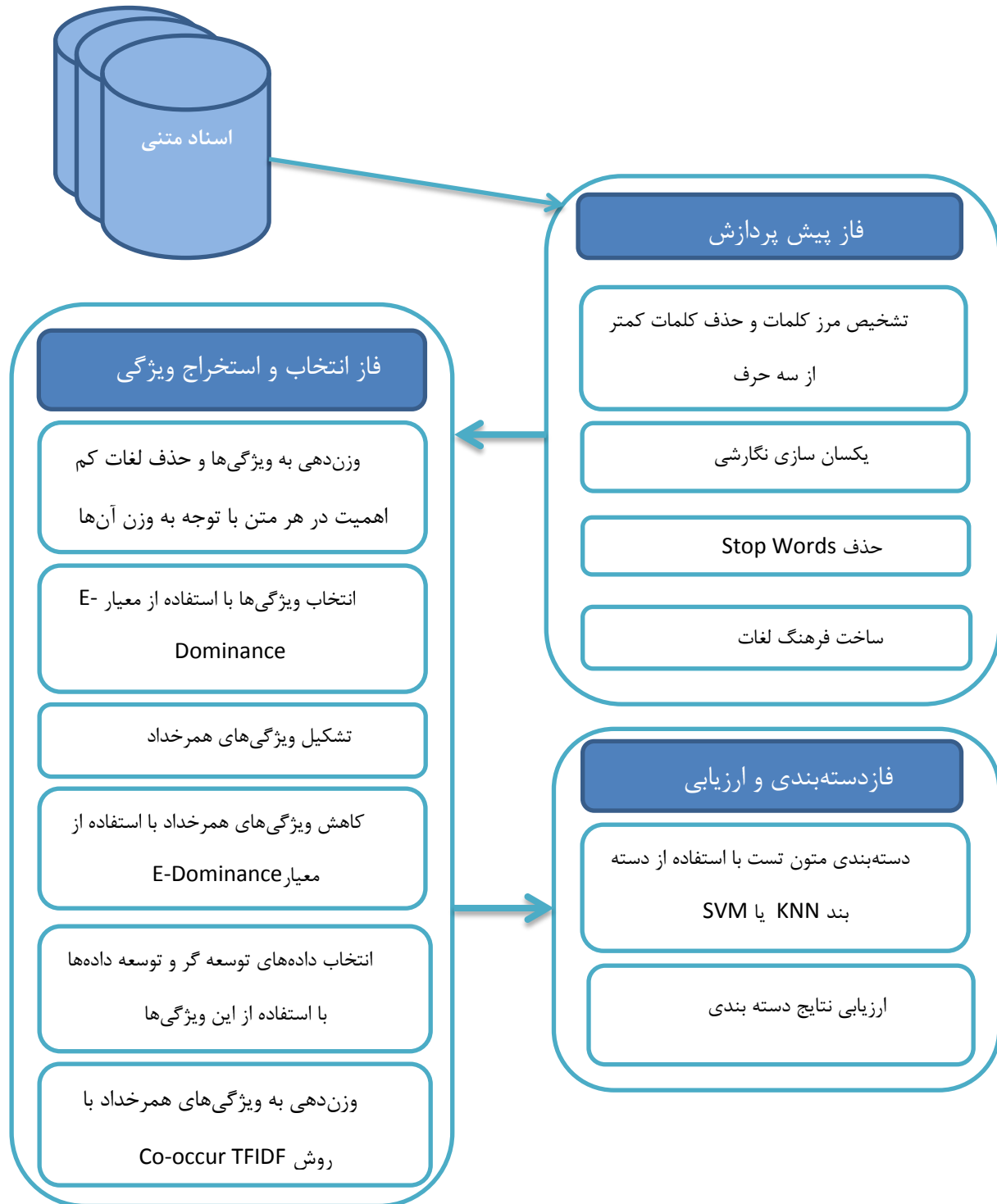
معرفی روش پیشنهادی

۳-۱- مقدمه

در فصل قبل، معماری کلی سیستم دسته‌بندی متون و مراحل آن، تشریح شد و به طور ضمنی با بررسی کارهای گذشته به محدودیت‌ها و نکات ضعف روش‌های موجود، اشاره گردید. در این پایان‌نامه، با ارائه یک معیار اهمیت ویژگی و نیز ارائه یک طرح توسعه یافته از طرح‌های موجود در وزن‌دهی برای ویژگی‌های هم‌رخداد، سعی شده است با مشکلات روش‌های قبلی، یعنی تعداد ویژگی‌های زیاد و فقدان یک روش وزن‌دهی کارآمد برای ویژگی‌های هم‌رخداد مقابله شود. نتایج حاصل از اعمال این روش‌ها، نزدیک شدن به هدف مورد نظر را نشان داده است. در این فصل، سیستم دسته‌بندی متون فارسی ارائه و تشریح می‌شود، سپس در فصل بعدی نتایج حاصل از ارزیابی آزمایشات انجام شده ارائه می‌گردد. از این‌رو در این فصل به طراحی چارچوب و پیاده‌سازی این سیستم پرداخته می‌شود.

۳-۲- معماری رویکرد پیشنهادی

در این بخش ابتدا یک بلوک دیاگرام از معماری پیشنهادی برای سیستم دسته‌بندی متون ارائه داده، سپس به طور مفصل مراحل مختلف آن تشریح می‌شود. شکل ۳-۱، معماری پیشنهادی را تا حدودی با جزئیات لازم نشان می‌دهد، اما جزئیات بیشتر در قسمت‌های بعدی به صورت مفصل‌تر بیان شده است. معماری پیشنهادی شامل چهار فاز اصلی و هر فاز شامل بخش‌های مختلفی هستند. در اکثر روش‌های دسته‌بندی متون، این مراحل به صورت کلی مشترک بوده و آنچه باعث متفاوت و مجزا بودن کارهای مختلف می‌گردد، جزئیات کار و شیوهی متفاوت در هر مرحله است. بخش‌های مختلف معماری پیشنهادی در ذیل آمده است.



شکل (۳-۱) نمودار کلی معماری پیشنهادی

۳-۲-۱- اجزای تشکیل دهنده معماری روش پیشنهادی

- فاز انتخاب متون از پایگاه داده
 - فاز پیش پردازش
 - i. تشخیص مرز کلمات و حذف کلمات کمتر از سه حرف
 - ii. یکسان سازی نگارشی
 - iii. حذف Stop Words
 - iv. ساخت فرهنگ لغات
 - فاز انتخاب و استخراج ویژگی
 - i. وزن دهی به ویژگی‌ها و حذف لغات کم اهمیت در هر متن با توجه به وزن آنها
 - ii. انتخاب ویژگی‌ها با استفاده از معیار پیشنهادی E-Dominance
 - iii. تشکیل ویژگی‌های هم‌رخداد
 - iv. کاهش ویژگی‌های هم‌رخداد با استفاده از معیار پیشنهادی E-Dominance
 - v. انتخاب داده‌های توسعه‌گر و توسعه متون با استفاده از این ویژگی‌ها
 - vi. وزن دهی به ویژگی‌های هم‌رخداد با استفاده از روش Co-occur TFIDF
 - فاز دسته‌بندی و ارزیابی
 - i. دسته‌بندی متون تست با استفاده از دسته‌بندهای KNN یا SVM
 - ii. ارزیابی نتایج دسته بندی
- حال پس از ارائه‌ی معماری روش پیشنهادی ، توضیح مفصل تری از مراحل مختلف آن بیان می‌گردد.

۳-۲-۲- انتخاب متون از پایگاه داده

پایگاه داده همشهری [۱۲۴]، به دلیل مطابقت با استانداردهای TREC [۱۲۵]، در سال‌های اخیر پرکاربردترین پایگاه داده در حوزه‌ی تحقیقات زبان فارسی بوده است. این پایگاه داده توسط گروه تحقیقاتی دانشگاه تهران ارائه شده است. نسخه‌ی یک این مجموعه نمونه‌ای است که در همایش‌های CLEF در سال‌های ۲۰۰۸ و ۲۰۰۹ برای ارزیابی سامانه‌های بازیابی اطلاعات تک‌منظوره^۱ مورد استفاده قرار گرفته است. نسخه‌ی دو این پیکره، آخرین نسخه‌ی مجموعه بوده و نسبت به نسخه‌ی یک بزرگتر و جامع‌تر است.

در این تحقیق برای انجام آزمایشات، از پایگاه داده‌ی همشهری ۲ استفاده شده است. این نسخه شامل ۳۱۸ هزار سند در قالب فایل‌های XML می‌باشد. داده‌های انتخابی در این پایان‌نامه مطابق جدول ۳-۱ است.

جدول (۳-۱) داده‌های انتخابی از پایگاه داده همشهری ۲

پایگاه داده	دسته‌ها	متون آموزش	متون تست
همشهری ۲	۵ دسته (ورزشی، اقتصادی، سیاسی، علمی، فرهنگ و هنر)	۱۰۰ متن از هر دسته	۵۰ متن از هر دسته

شکل ۳-۲ یک نمونه از متون این پایگاه داده را که به دسته‌ی "اقتصادی" تعلق دارد، نشان می‌دهد.

^۱ Ad Hoc

۳-۲-۳- پیش پردازش‌ها

همانطور که گفتیم متون انتخابی باید بعد از پیش‌پردازش‌هایی وارد مرحله انتخاب ویژگی شوند. پیش-

پردازش‌های انجام شده بر روی متون، جهت آماده‌سازی داده‌ها برای دسته‌بندی به صورت ذیل است:

۳-۲-۳-۱- تشخیص مرز کلمات و حذف کلمات کمتر از سه حرف

در این مرحله، با استفاده از کاراکترهای “.”، “؟”، “:”، “؛”، “،”، “” و ... مرز شروع و پایان یک

کلمه را شناسایی کرده و با شمارش کاراکترهای موجود در این دو مرز، کلمات کمتر از سه حرف حذف

می‌شوند. شکل ۳-۳، نمونه‌ای از یک متن انتخابی را پس از اولین مرحله از پیش‌پردازش نشان می‌دهد.

تازه‌ترین بررسی‌های آماری نشان می‌دهد، در روزهای اول اردیبهشت ماه، ارزش سهام بانک‌ها و لیزینگ‌ها تحت تأثیر خبر احتمال تثبیت نرخ سود تسهیلات با صعودی جالب توجه روبه‌رو شده است. در روزهای اخیر ارزش سهام ۵ بانک و لیزینگ فعال در بورس تهران بین ۱/۱ تا ۷/۱ درصد رشد کرد، در این مدت بالاترین سود ناشی از رشد قیمت نصیب سهامداران لیزینگ ایران شده است. ارزش سهام این شرکت لیزینگ در پنج روز اول اردیبهشت ۷/۱ درصد افزایش یافته است. در این میان سودآورترین سهام در بین بانک‌ها، سهام بانک کارآفرین بوده، قیمت سهام این بانک در روزهای گذشته ۴ درصد رشد کرده است. براین اساس رایان سایپا با رشد قیمتی معادل ۳/۹ درصد، بانک پارسیان با صعودی برابر ۱/۹ درصد و بانک اقتصاد نوین با افزایش قیمتی برابر ۱/۱ درصد ارکستر صعود سهام بانک‌ها را تکمیل کرده‌اند.

شکل (۳-۲) نمونه‌ای از یک متن در گروه اقتصادی از پایگاه داده همشهری

تازه‌ترین بررسی‌های آماری نشان دهد روزهای اول اردیبهشت ماه ارزش سهام بانک لیزینگ تحت تأثیر خبر احتمال تثبیت نرخ سود تسهیلات صعودی جالب توجه روبه‌رو شده است روزهای اخیر ارزش سهام بانک لیزینگ فعال بورس تهران بین درصد رشد کرد این مدت بالاترین سود ناشی رشد قیمت نصیب سهامداران لیزینگ ایران شده است ارزش سهام این شرکت لیزینگ پنج روز اول اردیبهشت درصد افزایش یافته است این میان سودآورترین سهام بین بانک سهام بانک کارآفرین بوده قیمت سهام این بانک روزهای گذشته درصد رشد کرده است این اساس رایان سایپا رشد قیمتی معادل درصد بانک پارسیان صعودی برابر درصد بانک اقتصاد نوین افزایش قیمتی برابر درصد ارکستر صعود سهام بانک تکمیل کرده

شکل (۳-۳) نمونه‌ای از متن با اولین پیش‌پردازش

۲-۳-۲-۳- یکسان سازی نگارشی

در این مرحله حروف مشابه را دسته بندی کرده و با استفاده از کد اسکی آنها، یک حرف را به جای حروف مشابه در تمامی کلمات جایگزین کرده تا یک کلمه با دو حرف مشابه اما غیر یکسان نداشته باشیم؛ مثل «ؤ» به «و»، «ی» به «ی» و «أ» به «ا» و تبدیل نویسه‌های «ی» و «ک» عربی به نوع فارسی آن. همچنین برخی کلمات که در برخی متون با املاي فاصله دار و در برخی متون با املاي بدون فاصله حضور داشتند، همگی به شکل بدون فاصله درآمده‌اند.

۳-۳-۲-۳ حذف Stop Words

پس از تشخیص یک کلمه، باید بررسی کنیم که آیا کلمه‌ی مورد نظر جزو کلمات بی اهمیت یا اصطلاحاً Stop Words هست یا نه؟ این مرحله کمک شایانی به تشکیل فرهنگ لغت اولیه می‌کند. کلمات ایست نه تنها کمکی به دسته بندی نکرده و نتایج دسته‌بندی را تضعیف می‌کنند، بلکه به سبب بزرگ شدن فرهنگ لغت یا همان بردار ویژگی اولیه، حجم پردازش‌ها و نیز جستجوها را افزایش داده و باعث افزایش قابل توجه زمان فرآیند دسته‌بندی می‌شوند. برای زبان فارسی، هیچ منبع جامعی از Stop Words ارائه نشده است. در کارهای انجام شده در زبان فارسی فقط اشاره شده است که کلمات ایست از فرهنگ لغت تشکیل شده به عنوان بردار ویژگی، حذف شده‌اند و تنها نمونه‌هایی از آنها ارائه شده است. ما در این پایان نامه، با استفاده از معیار DF مجموعه‌ای ۱۰۸۰ کلمه‌ای از کلمات ایست ارائه داده‌ایم که عدم وجود این کلمات در متون، نه تنها تاثیر منفی در نتایج دسته‌بندی ندارد، بلکه زمان دسته‌بندی را نیز به طور محسوسی کاهش می‌دهد. نمونه‌ای از کلمات ایست استخراجی از متون در جدول ۲-۳ آمده است. شکل ۳-۴ نیز متن انتخابی را پس از حذف کلمات ایست نشان می‌دهد. کلماتی که پس از حذف کلمات ایست از لیست لغات حذف می‌شوند به عنوان فرهنگ لغات و یا بردار ویژگی استفاده می‌شوند.

جدول (۲-۳) نمونه‌هایی از کلمات ایست استخراجی

اگرچه	گرچه	زیرا	بود	کمی
این	اند	یک	دیگر	هر
ما	آن	آنها	و	وقتی
نیز	برای	یا	را	بلکه
باید	وی	اما	ممکن	اندک

تازه بررسی آماری نشان روزهای اردیبهشت ماه ارزش سهام بانک لیزینگ تأثیر خبر احتمال تثبیت نرخ سود تسهیلات صعودی جالب توجه روزهای ارزش سهام بانک لیزینگ فعال بورس تهران درصد رشد مدت بالاترین سود ناشی رشد قیمت نصیب سهامداران لیزینگ ایران ارزش سهام شرکت لیزینگ پنج روز اردیبهشت درصد افزایش یافته سودآورترین سهام بانک سهام بانک کارآفرین قیمت سهام بانک روزهای گذشته درصد رشد اساس رایان سایپا رشد قیمتی معادل درصد بانک پارسیان صعودی برابر درصد بانک اقتصاد نوین افزایش قیمتی برابر درصد ارکستر صعود سهام بانک تکمیل

شکل (۴-۳) نمونه‌ای از متن با Stop Word های حذف شده

۳-۲-۴- وزن دهی به ویژگی‌ها و حذف لغات کم اهمیت با توجه به وزن آنها

پس از خواندن کلمات متون و حذف کلمات ایست و کلمات کمتر از سه حرف، حدود ۱۶۸۹۳ کلمه منحصر بفرد در بردار اولیه ویژگی‌ها قرار گرفتند. این تعداد کلمه به عنوان ویژگی، در مراحل بعدی پردازش و دسته‌بندی، بسیار حجیم و زمانبر هستند. بنابراین در مرحله اول از کاهش ویژگی‌ها، با روش TFIDF ویژگی‌ها را وزن دهی کرده و با استفاده از سعی و خطا و در نظر گرفتن مقدار آستانه ۵ برای این معیار، در هر متن، لغاتی که دارای وزن کمتری از آستانه‌ی موردنظر هستند، از لیست ویژگی‌های هر کلاس حذف می‌شوند. نهایتاً بردار ویژگی اولیه به بردار ویژگی با ۱۶۴۳۰ کلمه کاهش پیدا کرده است.

۳-۲-۵- انتخاب ویژگی‌های کلاسی با معیار E-Dominance

در مرجع [۷] معیار Dominance معرفی شده است. این معیار، میزان اهمیت یک کلمه در هر کلاس را نسبت به کل داده‌ها نشان می‌دهد؛ به طوری که اگر یک کلمه فقط در متون یک کلاس اتفاق بیفتد، آنگاه مقدار Dominance برای آن کلمه در آن کلاس یک خواهد بود. در رابطه‌ی ۱-۳، $df(f_i, c_j)$ تعداد متون کلاس c_j است که ویژگی f_i در آنها رخ داده است. یکی از ویژگی‌های Dominance این است که مجموع آن برای یک عبارت f_i در کلاس‌های متفاوت یک است (رابطه‌ی ۲-۳). در رابطه‌ی ۲-۳، k تعداد کلاس‌هاست.

$$\text{Dominance}(f_i, c_j) = \frac{df(f_i, c_j)}{\sum_{j=1}^k df(f_i, c_j)} \quad 1-3$$

$$\sum_{j=1}^k \text{Dominance}(f_i, c_j) = 1 \quad 2-3$$

مرجع مذکور، از این معیار برای کاهش ویژگی‌ها استفاده کرده است؛ یعنی انتخاب یک حد آستانه برای معیار فوق، کلماتی که مقادیری کمتر از آن آستانه دارند، از لیست ویژگی‌ها حذف می‌شوند. بر اساس این معیار، اگر یک کلمه، فقط و فقط در یک متن خبری حضور داشته باشد و در هیچ متن دیگری از متون پایگاه داده رخ نداده باشد، مقدار Dominance برای این کلمه، در کلاس شامل متن "۱" خواهد بود و به عنوان یک ویژگی با اهمیت از لیست ویژگی‌ها حذف نخواهد شد. فرض کنیم این کلمه فقط در دو یا سه متن از همان کلاس رخ داده باشد، باز هم شرایط مشابهی داریم و مقدار Dominance برای این کلمه در آن کلاس "۱" و در سایر کلاس‌ها صفر خواهد بود. سوالی که در اینجا مطرح است این است که آیا کلمه‌ای که تنها در یک متن از متون آموزشی رخ داده است، می‌تواند به عنوان یک ویژگی مهم از آن متن

یا آن کلاس قلمداد گردد؟ آیا این ویژگی خاصیت تفکیک کنندگی خوبی دارد؟ یا احتمال اینکه این کلمه در متون آزمایشی رخ داده و این ویژگی باعث شباهت متن مذکور به متن آزمایشی جدید شود چقدر خواهد بود؟ با بررسی فرهنگ لغت تولید شده از متون آموزشی و اعمال ویژگی Dominance برای کاهش تعداد کلمه‌ها در این فرهنگ لغت یا همان بردار ویژگی‌ها، مشاهده می‌شود بسیاری از کلماتی که در برخی کلاس‌ها، مقدار Dominance مساوی با "۱" و یا نزدیک به "۱" و بیشتر از آستانه‌ی مورد نظر دارند، نه تنها نقش مثبتی در دسته‌بندی ایفا نکرده، بلکه امکان کاهش دقت نیز وجود دارد؛ چرا که احتمال رخداد این کلمات در متون مشابه و هم‌دسته از متون آزمایشی، بسیار کم بوده و در محاسبه میزان شباهت آنها مثلا در دسته‌بند KNN که معمولا از فاصله‌ها ی اقلیدسی و کسینوسی استفاده می‌کند، تاثیر منفی می‌گذارد. علاوه بر تاثیر منفی در محاسبه‌ی میزان شباهت و کاهش کارایی، ضعف دیگر این معیار در کاهش ویژگی‌ها این است که پس از اعمال این معیار، باز هم کلمات بی‌اهمیتی در بردار ویژگی حضور داشته و حجم و زمان پردازش‌ها را تحت‌الشعاع قرار می‌دهند. با وجود این ضعف‌ها، مزیتی که این ویژگی دارد این است که نشان می‌دهد اگر این ویژگی در یک کلاس مقداری نزدیک به "۱" داشته باشد، می‌توان تا حدود زیادی مطمئن بود، که این ویژگی متعلق به این کلاس است و در سایر کلاس‌ها کمتر اتفاق می‌افتد در نتیجه این ویژگی قدرت تفکیک‌کنندگی خوبی خواهد داشت.

۳-۲-۵-۱- بهبود معیار Dominance

در بخش قبل گفتیم یکی از نقاط ضعف معیار Dominance، انتخاب برخی از ویژگی‌هایی است که به دلیل رخداد کم در متون یک کلاس، با وجود داشتن قدرت تفکیک‌کنندگی بالا، فقط اندازه بردار ویژگی را افزایش می‌دهند، چرا که این ویژگی‌ها به ندرت اتفاق افتاده و کاربرد موثری در دسته‌بندی ندارند. اما حالت ایده‌آل برای یک ویژگی قدرتمند و تفکیک کننده از یک کلاس این است که اولاً آن ویژگی فقط در متون همان کلاس رخ دهد و ثانياً اینکه در همه و یا تعداد موثری از متون آن کلاس رخ دهد، نه تعداد

متون اندک. یعنی اگر یک ویژگی در متون زیادی از یک کلاس رخ دهد، می‌توان با اطمینان گفت آن ویژگی متعلق به این کلاس است. در حالتی که در یک متن و یا تعداد متون اندک رخ دهد، آن ویژگی، ویژگی کلاس نبوده و ویژگی آن متن محسوب می‌گردد. با توجه به مطالب مذکور چگونه یک ویژگی قدرتمند از یک کلاس پیدا کنیم که در آن کلاس جامعیت داشته باشد؟ با توجه به رابطه ۱-۳، مقدار Dominance برای ویژگی f_i در کلاس c_j برابر است با تعداد متونی از کلاس c_j که ویژگی در آنها رخ داده باشد به تعداد متونی از همه‌ی کلاس‌ها (k تعداد کلاس‌ها) که ویژگی f_i در آنها رخ داده باشد. برای مثال اگر یک کلمه فقط در یک متن از یک کلاس رخ داده و در هیچ متن دیگری رخ ندهد، مقدار صورت و مخرج کسر در رابطه‌ی ۱-۳ برابر بوده و مقدار Dominance برابر با یک خواهد بود. اما همانطور که گفته شد احتمال حضور این ویژگی تفکیک کننده در متون تست نیز خیلی کم خواهد بود.

در این پایان‌نامه از ویژگی مثبت معیار Dominance، که غالب بودن یک کلاس بر سایر کلاس‌ها در ویژگی است، استفاده شده است. از طرف دیگر سعی شده است نقطه‌ی ضعف این معیار - که برای حضور کم تعداد یک ویژگی در متون یک کلاس، مقدار بالا تخصیص می‌دهد - بر طرف شود. این کار را با اضافه کردن یک ضریب به معیار Dominance انجام شده است.

$$\alpha_{ij} = \frac{df(f_i, c_j)}{d(c_j)} \quad ۳-۳$$

در رابطه ۳-۳، صورت کسر نشان‌دهنده‌ی تعداد اسناد کلاس c_j است که ویژگی f_i در آنها رخ داده است، و $d(c_j)$ ، تعداد کل اسناد کلاس c_j است. در نظر بگیرید، اگر $d(c_j) = df(f_i, c_j)$ ، ضریب α_{ij} مساوی با یک خواهد شد زیرا تعداد متون کلاس c_j که ویژگی f_i در آنها رخ دهد، برابر با تعداد کل اسناد کلاس c_j می‌باشد. حالا اگر ضریب α_{ij} را که از این به بعد ضریب جامعیت Dominance^۱ می‌نامیم در

^۱ Dominance Comprehensiveness Factor

عبارت اصلی Dominance ضرب کنیم رابطه (۳-۴)، مقدار E-Dominance^۱ را خواهیم داشت (رابطه‌ی ۳-۵). عبارت کلی یا E-Dominance، زمانی برابر با ۱ خواهد بود که هم Dominance مقدار یک داشته باشد و هم ضریب α .

$$E - \text{Dominance}(f_i, c_j) = \frac{df(f_i, c_j)}{d(c_j)} * \frac{df(f_i, c_j)}{\sum_{j=1}^k df(f_i, c_j)} \quad ۴-۳$$

$$E - \text{Dominance}(f_i, c_j) = \alpha_{ij} * \text{Dominance}(f_i, c_j) \quad ۵-۳$$

حال اگر یک ویژگی مقدار E-Dominance، در یک کلاس خاص، یک باشد، معنایش این است که اولاً این کلمه فقط در همان کلاس رخ داده است و در کلاس‌های دیگر حتی یک بار هم در هیچ متنی ظاهر نشده است، ثانیاً این کلمه در تمامی متون همان کلاس حضور داشته است. با این تفسیر از این معیار، می‌توان مطمئن شد که این ویژگی برای این کلاس یک ویژگی قدرتمند محسوب می‌شود. در فصل بعدی که آزمایشات و نتایج آن‌ها بیان می‌شود، این ادعا با ارزیابی نتایج نشان داده خواهد شد. همان‌طور که گفته شد این معیار زمانی که برای یک کلاس مقدار یک داشته باشد، یک ویژگی منحصر بفرد از یک کلاس خواهد بود، اما مقدار یک، یک حالت ایده آل است که خیلی کم اتفاق می‌افتد. برای بدست آوردن مقدار مناسب برای این معیار از روش سعی و خطا استفاده شده و یک حد آستانه برای این ویژگی تعیین شده است. اگر در کلاسی مقدار این ویژگی بیشتر از آن آستانه‌ی انتخابی بود، آن ویژگی به عنوان یکی از ویژگی‌ها که در بردار نهایی ویژگی قرار می‌گیرد، انتخاب خواهد شد. در این تحقیق در مرحله انتخاب ویژگی از معیار E-Dominance برای انتخاب ویژگی‌های مناسب برای مراحل بعدی استفاده شده است. حد آستانه برای این معیار با استفاده از آزمون و خطا، ۰,۰۴ در نظر گرفته شد. چند

^۱ Extended Dominance

نمونه از ویژگی‌های استخراجی برای هر کلاس با استفاده از معیار E-Dominance در جدول ۳-۳ آمده است.

جدول (۳-۳) نمونه‌هایی از ویژگی‌های استخراج شده با معیار E-Dominance برای کلاس‌های مختلف

اقتصادی	سیاسی	ورزشی	فرهنگ و هنر	علمی
بورس	نمایندگان	رقابت	سینما	آموزش
رقابت	انتخابات	امتیاز	هنرمندان	کتاب
سهام	خارجی	ورزش	آموزش	یادگیری

۳-۲-۶- تشکیل زوج ویژگی‌های هم‌خداد به عنوان ویژگی‌های متن

در مرجع [۷]، پس از حذف ویژگی‌های بی‌اهمیت و انتخاب واژه‌های مهم با استفاده از معیار Dominance، واژه‌های هم‌خداد به عنوان ویژگی‌های نهایی داده‌ها، تشکیل شده و در دسته‌بندی مورد استفاده قرار گرفته‌اند. بردار ویژگی حاصل از تشکیل ویژگی‌های هم‌خداد بسیار بزرگ و قابل ملاحظه است. این حجم بالا از ویژگی‌ها، در این حقیق به عنوان چالشی بزرگ مورد مطالعه قرار گرفته و سعی شده است، با بهبود ویژگی Dominance - که در بخش ۳-۲-۵-۱ توضیح داده شد - و انتخاب ویژگی‌های مناسب و کارا تر جهت تشکیل ویژگی‌های هم‌خداد، کاهش چشمگیری در اندازه‌ی این بردار ویژگی ایجاد شود. پس از انتخاب ویژگی‌های مناسب و تفکیک کننده با استفاده از معیار پیشنهادی، به بررسی هم‌خدادی این ویژگی‌ها در متون مورد مطالعه پرداخته شده است. این ویژگی‌های هم‌خداد، با تشکیل زوج کلمات، که هر کلمه یکی از کلمات انتخابی در مرحله‌ی قبل است، ایجاد می‌شوند. در بررسی اولیه تعداد این ویژگی‌های هم‌خداد بسیار زیاد خواهد بود. اگر N ، تعداد ویژگی‌های منفرد باشد، اما $\frac{N(N-1)}{2}$ ویژگی هم‌خداد تشکیل می‌شود. این تعداد ویژگی برای کل متون، بسیار زیاد خواهد بود. اما در زمان تشکیل این ویژگی‌ها، دیده می‌شود که بسیاری از این زوج ویژگی‌ها در هیچ‌یک از متون رخ نداده

و حذف خواهند شد. اما با این وجود، باز هم تعداد این ویژگی‌ها زیاد بوده و در مراحل بعدی باید فرآیند کاهش ویژگی دنبال شود. در جدول ۳-۴، نمونه‌هایی از ویژگی‌های هم‌رخداد تشکیل شده آمده است.

جدول (۳-۴) نمونه‌هایی از ویژگی‌های هم‌رخداد تشکیل شده

بورس-سهام	نمایندگان-قانون	رقابت-امتیاز	سینما-فیلم	آموزشی-علمی
سرمایه-سود	سیاست-خارجی	فوتبال-گل	هنرمند-موسیقی	الکترونیکی-اطلاعات
رقابت-سود	وزیر-وزارتخانه	ورزشگاه-تماشاگر	آموزشی-فرهنگی	اینترنت-سایت
وزیر-اینترنت	سایت-فوتبال	فرهنگی-وزیر	خارجی-اطلاعات	هنرمند-برگزاری
آموزشی-بورس	برگزاری-رقابت	گل-رقابت	سود-سینما	سیاست-سرمایه

۳-۲-۷- کاهش ویژگی‌های هم‌رخداد با استفاده از معیار پیشنهادی E-Dominance

پس از تشکیل ویژگی‌های هم‌رخداد و توسعه‌ی آنها با استفاده از ویژگی‌های هم‌رخداد توسعه‌گر در مرحله‌ی بعدی، کاهش ویژگی‌ها انجام می‌شود. برای این کار، مثل مرحله‌ی اول کاهش ویژگی‌ها، که برای پیدا کردن عبارات کاندید برای ساخت ویژگی‌های هم‌رخداد انجام شد، در این مرحله نیز از معیار E-Dominance استفاده شده و تعداد ویژگی‌ها کاهش داده می‌شود. در این مرحله حد آستانه برای معیار E-Dominance مقدار ۰,۰۳۵ با استفاده از آزمون و خطا به دست آمده است.

۳-۲-۸- انتخاب ویژگی‌های توسعه‌گر^۱ و توسعه داده‌ها

یکی از نقاط قوت مرجع [۷] توسعه و غنی‌سازی داده‌های یک کلاس با ویژگی‌های مهم آن کلاس است. غنی‌سازی داده‌ها به این صورت انجام شده است که در فرآیند آموزش پس از یافتن بهترین واژه‌های هم‌رخداد به عنوان ویژگی‌های اصلی، در این میان، ویژگی‌های هم‌رخدادی که در تعداد متون بیشتری از یک حد آستانه در یک کلاس خاص تکرار شده‌اند و یا بر اساس یک معیار ویژه، مهم تلقی می‌شوند، به عنوان داده‌ی توسعه‌گر یا همان داده‌ای که برای غنی‌سازی داده‌ها به کار خواهد رفت، انتخاب

^۱ Developer Features

می‌شوند. به طور مثال اگر دو کلمه‌ی "بورس" و "سرمایه" به عنوان ویژگی هم‌رخداد توسعه‌گر در نظر گرفته شوند و جزء ویژگی‌های با اهمیت کلاس "اقتصادی" محسوب شوند، در فرآیند توسعه و غنی‌سازی داده‌ها، این زوج ویژگی یعنی "بورس و سرمایه" را به تمام داده‌های کلاس اقتصادی در داده‌های آموزشی اضافه می‌کنیم. با این کار این ویژگی در این کلاس نسبت به سایر کلاس‌ها بارزتر و در میان داده‌های این کلاس عمومی‌تر و پرتکرارتر خواهد شد. در سایر کلاس‌ها هم در شرایطی این ویژگی به داده اضافه خواهد شد، که یکی از ویژگی‌های این زوج هم‌رخداد، در داده‌ی مورد نظر رخ داده باشد، در این صورت این ویژگی تنها به همین داده از کلاس غیر اقتصادی اضافه می‌گردد. در فرآیند غنی‌سازی داده‌ها در داده‌های آزمایشی، بدلیل عدم آگاهی از کلاس داده‌ی مورد نظر، از شیوه‌ی دوم استفاده کرده و در صورت وجود هر تک واژه از زوج واژه‌های هم‌رخداد، این ویژگی به داده‌ی مورد نظر اضافه می‌شود. نمونه‌ایی از ویژگی‌های توسعه‌گر در جدول ۳-۵ آمده است. در این پایان‌نامه نیز از تکنیک توسعه متون استفاده شده است.

جدول (۳-۵) نمونه‌ای از ویژگی‌های توسعه‌گر انتخابی از هر کلاس

اقتصادی	سیاسی	ورزشی	فرهنگ و هنر	علمی
بورس - سهام	نمایندگان - قانون	رقابت - امتیاز	سینما - فیلم	آموزشی - علمی

۳-۲-۹- وزن‌دهی ویژگی‌های هم‌رخداد با روش وزن‌دهی Co-Occur TFIDF

علاوه بر کاهش حجم بردار ویژگی‌ها، هدف دیگری که در بهبود کاربرد واژه‌های هم‌رخداد دنبال شده است، ارائه‌ی یک طرح وزن‌دهی برای ویژگی‌های هم‌رخداد است. در کارهای مختلفی از ویژگی‌های هم‌رخداد به عنوان ویژگی‌های تفکیک‌کننده و اصلی برای دسته‌بندی استفاده شده است، اما هیچ‌گاه یک طرح وزن‌دهی برای این ویژگی‌ها معرفی نشده است. در مرجع [۷] نیز از وزن‌دهی باینری یعنی "صفر" برای عدم حضور ویژگی و "یک" برای حضور ویژگی استفاده شده است. این یک شیوه‌ی پایه برای وزن-

دهی است. اما همانطور که می‌دانیم، دلیل اصلی در ارائه‌ی طرح‌های وزن‌دهی که برای ویژگی‌های ساده معرفی شد، این است که میزان تکرار یا فرکانس حضور یک ویژگی خیلی با اهمیت تر از صرفاً حضور یا عدم حضور یک ویژگی قلمداد می‌شود، یعنی داده‌ی با حضور بیشتر باید سهم بیشتری در دسته‌بندی ایفا کند. لذا در ویژگی‌های هم‌رخداد نیز پا را از اهمیت حضور یا عدم حضور فراتر گذاشته و به ارائه‌ی یک طرح وزن‌دهی برای دخالت دادن میزان حضور ویژگی‌های هم‌رخداد در دسته‌بندی می‌پردازیم. در کارهای سابق که ویژگی‌های مورد استفاده برای دسته‌بندی کلمات یا عبارات متن انتخاب شده‌اند، یکی از رایج‌ترین شیوه‌های استاندارد وزن‌دهی بنام TFIDF جهت وزن‌دهی و تعیین میزان اهمیت کلمات در هر متن به کار گرفته شده است. در این معیار (رابطه ۳-۶)، آنچه وزن یک کلمه را به عنوان ویژگی آن متن تعیین می‌کند، اولاً تعداد رخداد کلمه یا ویژگی f_k در متن d_i ، یعنی مقدار $tf(f_k, d_i)$ و ثانیاً تعداد متونی که ویژگی f_k در آنها حضور دارد یعنی $|D(f_k)|$ ، می‌باشد (رابطه ۳-۷). در این رابطه $|D|$ ، نیز تعداد کل متون است.

$$TFIDF(f_k, d_i) = tf(f_k, d_i) * idf(f_k, d_i) \quad ۶-۳$$

$$idf(f_k, d_i) = \log\left(\frac{|D|}{|D(f_k)|}\right) \quad ۷-۳$$

از این معیار استفاده کرده و آن را برای ویژگی‌های هم‌رخداد بسط می‌دهیم. اگر ویژگی $f_{k1, k2}$ شامل دو مولفه یعنی f_{k1} و f_{k2} باشد - که این دو مولفه، مولفه‌های ویژگی هم‌رخداد هستند - و $|D(f_{k1}, f_{k2})|$ تعداد متونی باشند که پس از عملیات غنی‌سازی داده‌ها، شامل هر دو ویژگی f_{k1} و f_{k2} باشند آنگاه وزن ویژگی هم‌رخداد $f_{k1, k2}$ با استفاده از رابطه‌ی ۳-۸ محاسبه می‌شود:

$$\text{Co - Occur TFIDF}(f_{k1,k2}, d_i) = \text{Max}(tf(f_{k1}, d_i), tf(f_{k2}, d_i)) * idf(f_{k1,k2}, d_i) \quad 8-3$$

$$idf(f_{k1,k2}, d_i) = \log\left(\frac{|D|}{|D(f_{k1}, f_{k2})|}\right) \quad 9-3$$

$\text{Max}(tf(f_{k1}, d_i), tf(f_{k2}, d_i))$: مقدار ماکسیمم برای تعداد حضور مولفه اول و دوم ویژگی است. دلیل انتخاب مقدار ماکسیمم این است که در مرحله‌ی توسعه و غنی سازی داده‌ها، به ازای وجود هر مولفه از ویژگی، مولفه‌ی بعدی نیز اضافه شده است. به دلیل هم‌رخداد بودن این ویژگی‌ها، منطقی است که باید به ازای هر بار حضور مولفه‌ی پرتکرارتر، مولفه‌ی دیگر نیز باید حضور داشته باشد، بنابراین این مولفه اضافه می‌شود. تعداد اضافه شدن مولفه دوم، به تعداد اختلاف تعداد مولفه‌های اول و دوم است تا جایی که تعداد آنها برابر شود. بنابراین در نهایت تعداد رخداد هر مولفه، به تعداد مولفه‌ی بیشتر خواهد بود.

$|D(f_{k1}, f_{k2})|$: تعداد متونی است که شامل هر دو مولفه‌ی ویژگی مورد نظر، یعنی f_{k1} و f_{k2}

هستند.

۳-۲-۱۰- دسته‌بندی و ارزیابی نتایج

در بخش ۲-۳-۴ به انواع دسته‌بندیها در حوزه‌ی دسته‌بندی متون اشاره شد. در این پایان‌نامه از دسته بند SVM و KNN برای دسته‌بندی استفاده شده است. مقدار K در دسته‌بند KNN، با سعی و خطا، ۷ در نظر گرفته شده و فاصله‌ی کسینوسی نیز به عنوان معیار شباهت در این دسته‌بند به کار گرفته شده است. پس از دسته‌بندی داده‌ها، آنچه که می‌تواند درستی و کارایی روش پیشنهادی را مورد تایید قرار دهد، یک ارزیابی درست از نتایج به دست آمده است. در حوزه‌ی بازیابی و دسته‌بندی متون، معیارهای دقت، صحت، فراخوانی و F_1 متداول‌ترین معیارهایی هستند که به کار برده می‌شوند. این معیارها در فصل بعدی به طور کامل توضیح داده می‌شوند.

۳-۳- نتیجه‌گیری

در این فصل، یک معماری کلی از رویکرد پیشنهادی ارائه شد. سپس به توضیح کاملی از هر گام از این رویکرد پرداخته شد. معیار E-Dominance به عنوان معیاری برای انتخاب و کاهش ویژگی‌ها معرفی گردید. پس از تشکیل ویژگی‌های هم‌رخداد و کاهش آنها با معیار مذکور، نحوه انتخاب و بکارگیری ویژگی‌های توسعه‌گر نیز بیان شد. سپس با توسعه روش وزن‌دهی TFIDF پایه برای ویژگی‌های هم‌رخداد، روش وزن-دهی Co-occur TFIDF نیز معرفی شد. نتایج حاصل از اعمال این الگوریتم‌ها بر روی داده‌ها، در فصل بعدی ارائه می‌شود.

فصل چہارم

نتیجہ تجربی،

تخلیل و ارزیابی

۴-۱- مقدمه

در فصل سوم، سیستم دسته‌بندی متون پیشنهادی معرفی شد. در این فصل ابتدا به معرفی معیارهای ارزیابی مورد استفاده در این پایان‌نامه پرداخته می‌شود. سپس نتایج حاصل از اعمال روش پیشنهادی بر روی داده‌ها، ارائه می‌گردد. روش پیشنهادی در ۵ حالت مختلف آزمایش شده است. در پایان نیز یک ارزیابی کامل از نتایج و مقایسه آن با روش‌های موجود ارائه شده و یک نتیجه‌گیری کلی از روش پیشنهادی ارائه می‌شود.

۴-۲- معرفی معیارهای ارزیابی

معیارهای ارزیابی برای نشان دادن کارایی الگوریتم و نیز مقایسه روش پیشنهادی با سایر روش‌ها، به کار می‌روند. معیارهای Accuracy، Precision، Recall و F-measure پرکاربردترین معیارها در دسته‌بندی و به ویژه در زمینه‌ی دسته‌بندی متن هستند [۱۲۶]. برای معرفی معیارهای ارزیابی، نیاز به معرفی مفاهیم و پارامترهایی داریم که در این معیارها به کار گرفته شده‌اند. این مفاهیم در جدول ۴-۱ آمده‌اند:

جدول (۴-۱) پارامترهای مورد نیاز اندازه‌گیری معیارهای کارایی

نسبت داده شده به دسته C_i	نسبت داده نشده به دسته C_i	
TP_i	FN_i	متعلق به دسته C_i
FP_i	TN_i	متعلق به دسته‌ای غیر از C_i

پس از تعریف مفاهیم اصلی، معیارهای ارزیابی الگوریتم به شرح ذیل می‌باشد :

Accuracy: نسبت نمونه‌های درست دسته‌بندی شده را به کل نمونه‌ها نشان می‌دهد (رابطه ۴-۱):

$$Accuracy : A_i = \frac{TP_i + TN_i}{TP_i + TN_i + FP_i + FN_i} \quad (\text{دقت}) \quad ۴-۱$$

Precision: این معیار، نسبت تعداد متون درست نسبت داده شده را به کل متون نسبت داده شده به همان دسته نشان می‌دهد (رابطه ۲-۴). مقدار معیار صحت، درجه پایداری^۱ الگوریتم دسته‌بند را نشان می‌دهند.

$$Precision : P_i = \frac{TP_i}{TP_i + FP_i} \quad (\text{صحت}) \quad 2-4$$

Recall: معیار بازخوانی نسبت تعداد متونی را که به درستی در یک دسته قرار گرفته‌اند به تعداد کل متونی که بایستی در آن دسته قرار می‌گرفتند، نشان می‌دهد (رابطه ۳-۴). این معیار، درجه‌ی تمامیت^۲ الگوریتم دسته‌بند را نمایش می‌دهد.

$$Recall : R_i = \frac{TP_i}{TP_i + FN_i} \quad (\text{بازخوانی}) \quad 3-4$$

$F_{\beta} Measure$: معیاری است که از ترکیب دو معیار $Precision$ و $Recall$ بدست می‌آید، چرا که معیارهای صحت و فراخوانی به تنهایی برای اندازه‌گیری کارایی الگوریتم دسته‌بندی مؤثر نبوده و استفاده‌ی مجزای آنها ممکن است نتایج نادرستی از سیستم را ارائه دهد. از این‌رو لازم است که این معیارها به گونه‌ای با یکدیگر ترکیب شوند. این معیار میزان تاثیر معیار میزان صحت و معیار میزان فراخوانی را در ارزیابی دسته‌بند در نظر می‌گیرد و از رابطه‌ی ۴-۴ بدست می‌آید.

$$F_{\beta} Measure = \frac{(\beta^2 + 1) * Recall * Precision}{(\beta^2)(Recall + Precision)} \quad 4-4$$

¹ Soundness

² Completeness

معمولا مقدار β در حوزه بازیابی و دسته‌بندی متون ۱ در نظر گرفته می‌شود، مقدار ۱ برای β به معنی تاثیری برابر از *Precision* و *Recall* خواهد بود. در رابطه‌ی ۴-۵ معیار F_1 تعریف شده‌است.

$$F_{1Measure} = \frac{2 * Recall * Precision}{(Recall + Precision)} \quad ۵-۴$$

همانطور که گفته شد، بر اساس پارامترهای تعریف شده در فرمول این معیارها، معیارهای معرفی شده، به طور مجزا برای هر دسته یا کلاس محاسبه می‌شوند. برای ارزیابی نتایج رویکرد پیشنهادی به صورت کلی باید از مقادیر میانگین هر معیار در تمام دسته‌ها استفاده شود. فرمول میانگین این معیارها در روابط ۴-۶ تا ۴-۱۰ آمده است.

$$Mean_{F1} = \frac{\sum_{i=1}^m F_1(i)}{m} \quad ۶-۴$$

$$Mean_{Precision} = \frac{\sum_{i=1}^m Precision(i)}{m} \quad ۷-۴$$

$$Mean_{Precision} = \frac{\sum_{i=1}^m Precision(i)}{m} \quad ۸-۴$$

$$Mean_{Recall} = \frac{\sum_{i=1}^m Recall(i)}{m} \quad ۹-۴$$

$$Mean_{Accuracy} = \frac{\sum_{i=1}^m Accuracy(i)}{m} \quad ۱۰-۴$$

معیارهای فوق، معیارهایی هستند که در بیشتر تحقیقات در حوزه متون به کار رفته‌اند، لذا معیار مناسبی برای مقایسه‌ی رویکرد پیشنهادی با سایر رویکردها نیز هستند. در این پایان‌نامه نیز معیارهای نامبرده برای ارزیابی کارایی روش پیشنهادی به کار رفته‌اند.

۳-۴- نتایج دسته‌بندی متون

۳-۴-۱- نتایج بکارگیری معیار E-Dominance در دسته‌بندی متون

در بخش ۲-۴، چهار معیار ارزیابی برای بررسی و ارزیابی نتایج حاصل از اعمال روش پیشنهادی بر روی داده‌های انتخابی از پایگاه داده، معرفی شد. جدول ۳-۴ مقادیر این معیارها را به تفکیک کلاس‌ها و نیز به صورت میانگین، برای حالتی از دسته‌بندی نشان می‌دهد که معیار E_Dominance بر روی ویژگی‌های اولیه که همان کلمات منفرد در متون هستند، اعمال شده است. با اعمال این معیار بسیاری از ویژگی‌ها فیلتر شده و با استفاده از ویژگی‌های باقیمانده از این مرحله، بردار ویژگی تشکیل داده شده است. تعداد ویژگی‌ها در این دسته‌بندی، ۱۰۴۱ عبارت انتخابی با معیار E_Dominance بوده و دسته‌بند KNN نیز در این دسته‌بندی استفاده شده است. جدول ۲-۴ تعداد ویژگی‌های انتخابی را به تفکیک هر دسته و نسبت به کل ویژگی‌ها نشان می‌دهد. تعداد مجموع ویژگی‌ها در تمام کلاس‌ها از تعداد ویژگی‌های کل، به دلیل همپوشانی کلاس‌ها بیشتر است.

نرخ تشخیص سیستم که برای مقایسه‌ی روش پیشنهادی با سایر روش‌ها استفاده می‌شود، مقدار میانگین برای معیار Recall و F1 است، که در جدول ۳-۴ مقدار ۰,۹۱۲ را دارا هستند.

جدول (۲-۴) تعداد ویژگی‌های انتخابی از هر دسته

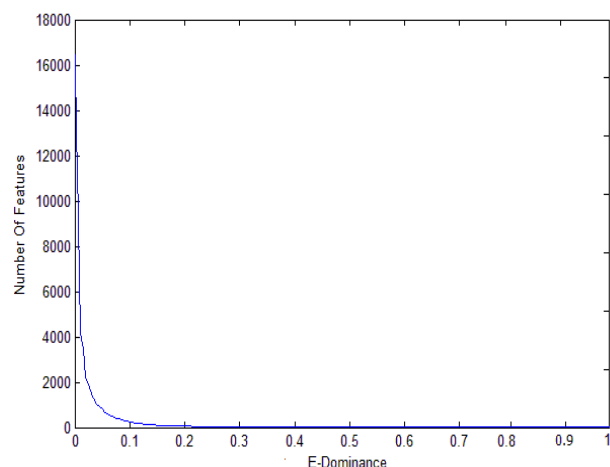
کلاس	علمی	فرهنگ و هنر	ورزشی	سیاسی	اقتصادی	
تعداد ویژگی‌ها	۲۱۶	۱۴۷	۲۱۳	۲۲۳	۲۸۰	۱۰۷۹
نسبت تعداد ویژگی انتخابی به کل ویژگی‌ها	٪۲۱	٪۱۴	٪۲۰	٪۲۱	٪۲۷	-

جدول (۳-۴) نتایج دسته‌بندی داده‌ها با معیارهای ارزیابی

دسته	<i>Recall</i>	<i>Precision</i>	F_1	<i>Accuracy</i>
اقتصادی	۰,۹۴	۰,۸۷۰	۰,۹۰۴	۰,۹۶
سیاسی	۰,۹۴	۰,۸۸۷	۰,۹۱۳	۰,۹۶۴
ورزشی	۰,۹۰	۰,۸۸۲	۰,۸۹۱	۰,۹۵۶
فرهنگ و هنر	۰,۹۲	۰,۹۳۹	۰,۹۲۹	۰,۹۷۲
علمی	۰,۸۶	۱	۰,۹۲۵	۰,۹۷۲
میانگین	۰,۹۱۲	۰,۹۱۶	۰,۹۱۲	۰,۹۶۵

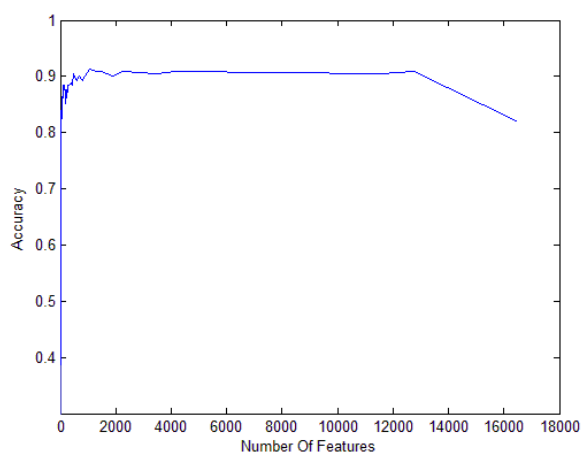
در روش پیشنهادی، ویژگی E-Dominance بعنوان یک معیار اهمیت ویژگی، برای انتخاب ویژگی‌های مناسب معرفی شد. با محاسبه‌ی مقدار این معیار برای ویژگی‌های اولیه و با استفاده از یک حد آستانه برای آن، ویژگی‌ها انتخاب شده‌اند. شکل ۴-۱، بیان کننده تاثیر مقدار حد آستانه این معیار در تعداد ویژگی‌های انتخاب شده است. هر چه مقدار یا حد آستانه انتخابی کمتر باشد، تعداد ویژگی‌های بیشتری انتخاب می‌شوند. مقدار ۰,۰۴ برای این معیار، تعداد ۱۰۴۱ ویژگی را با بیشترین دقت به دست آمده، انتخاب می‌کند.

همان‌طور که گفته شد، انتخاب مقدار آستانه برای E-Dominance، در تعداد ویژگی‌های انتخابی موثر است. مسلماً تعداد ویژگی‌های انتخاب شده نیز بر روی دقت سیستم، تاثیرگذار هستند.

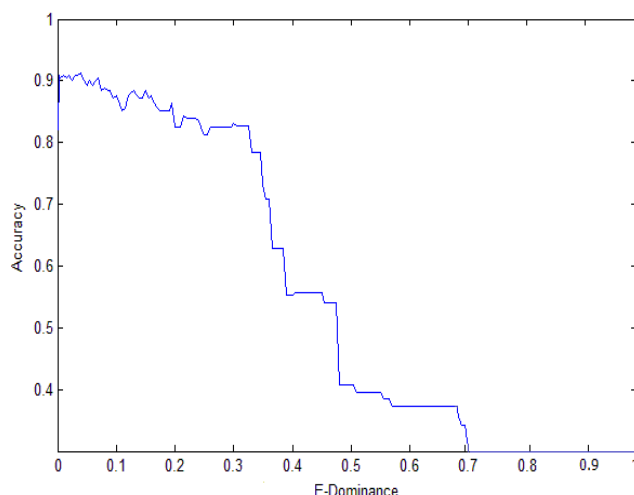


شکل (۱-۴) تاثیر مقدار معیار E-Dominance در تعداد ویژگی‌های انتخابی

شکل ۲-۴ تاثیر تعداد ویژگی‌های انتخابی را بر روی دقت دسته‌بندی نشان می‌دهد. با انتخاب ۱۰۴۱ ویژگی از ویژگی‌های موجود، دقت ۹۱,۲٪ بدست آمده است. این تعداد ویژگی و این دقت دسته‌بندی با انتخاب حد‌آستانه ۰,۰۴ برای E-Dominance بدست آمده است. بدون استفاده از آستانه و با کل ویژگی‌ها، دقت ۸۲٪ بدست آمد. در هر مرحله از انتخاب آستانه E-Dominance، مقدار ۰,۰۰۵ به مقدار قبلی اضافه شده است. در آستانه ۰,۰۰۵ دقت سیستم از ۸۲٪ اولیه به ۸۹,۲٪ رسیده و تعداد ویژگی‌ها نیز به ۱۲۷۳۸ عدد رسید. شکل ۳-۴ نیز تاثیر مقدار آستانه‌ی انتخابی برای E-Dominance را در دقت دسته‌بندی نشان می‌دهد.



شکل (۲-۴) تاثیر تعداد ویژگی‌های انتخابی در دقت دسته‌بندی



شکل (۳-۴) تاثیر مقدار معیار E-Dominance در دقت دسته‌بندی

۴-۳-۱-۱- مقایسه معیار E-Dominance با سایر معیارهای اهمیت ویژگی

در جدول ۴-۴، معیار E-Dominance با چند معیار مختلف که پرکاربردترین معیارهای انتخاب ویژگی مورد استفاده در حوزه‌ی متن‌کاوی هستند و در بخش فصل دوم معرفی شده‌اند، مقایسه شده‌است. این مقایسه تاثیر هر یک از معیارها را بر تعداد ویژگی‌های انتخابی و دقت دسته‌بندی نشان می‌دهد. تعداد ویژگی‌های هر روش، با توجه با بهترین دقت دسته‌بندی آورده شده است.

همانطور که در جدول ۴-۴ می‌بینید، بهترین میزان دقت را معیار Dominance و سپس روش پیشنهادی یعنی معیار E-Dominance با اختلاف ناچیز نشان می‌دهند. البته اختلاف دقت دسته‌بندی در میان این روش‌ها قابل ملاحظه نیست. هدف اصلی نیز در اینجا مقایسه‌ی میزان دقت دسته‌بندی نیست. آنچه در این مقایسه مهم است، تعداد ویژگی‌های انتخابی در هر روش است. در مورد تعداد ویژگی‌های انتخابی، معیار E-Dominance که معیار پیشنهادی است، اختلاف قابل ملاحظه‌ای با سایر معیارها دارد. این معیار توانسته است تعداد ویژگی‌های اولیه که ۱۶۴۳۰ ویژگی بوده‌اند را به کمتر از ۶٪ کاهش دهد و میزان کاهش ۹۳٫۶۶٪ را داشته باشد. پس از آن نیز معیار DF، با میزان کاهش ۸۷٫۰۱٪ و معیار Chi-

Square با ۸۱٫۸۴٪، در کاهش ویژگی‌ها بهترین عملکرد را داشته‌اند. معیار Dominance نیز که بهترین دقت را داراست، ویژگی‌ها را به حدود ۶۰٪ ویژگی‌های اولیه رسانده است. دو معیار TFIDF و Info-Gain نیز در کاهش ویژگی‌ها ضعیف تر عمل کرده‌اند و نتایج نزدیک به هم دارند. معیار Entropy نیز بدترین عملکرد را در کاهش ویژگی‌ها داشته و توانسته است، تنها حدود ۱٪ از داده‌ها را کاهش دهد.

جدول (۴-۴) مقایسه میزان دقت دسته‌بندی در کاهش ویژگی‌ها با معیار E-dominance و سایر معیارها

کاهش در ویژگی‌ها	Recall	تعداد ویژگی‌ها	معیار
بدون کاهش	۰٫۸۲	۱۶۴۳۰	-
۸۷٫۰۱٪	۰٫۹۰	۲۱۳۳	DF
۴۳٫۸۸٪	۰٫۸۸۸	۹۲۲۱	TFIDF
۰٫۹۹٪	۰٫۸۸۸	۱۴۵۸۶	Entropy
۴۹٫۶۸٪	۰٫۹۰	۸۲۶۷	Info-Gain
۸۱٫۸۴٪	۰٫۹۰	۲۹۴۸	Chi-Square
۳۹٫۹۵٪	۰٫۹۱۶	۹۸۶۶	Dominance
۹۳٫۶۶٪	۰٫۹۱۲	۱۰۴۴	E-Dominance (معیار پیشنهادی)

۴-۳-۲- نتایج بکارگیری ویژگی‌های هم‌خداد در دسته‌بندی متون

۴-۳-۲-۱- بکارگیری ویژگی‌های هم‌خداد با دو روش وزن‌دهی و دو دسته‌بند

در بخش ۴-۳-۱، نتایج دسته‌بندی با استفاده از ویژگی‌های ساده که همان کلمات منفرد حاصل از فیلتر شدن با معیار E-Dominance هستند، نشان داده شد. ویژگی‌های هم‌خداد به کار رفته در دسته‌بندی، زوج کلماتی هستند که در مرحله‌ی قبل با استفاده از معیار E-Dominance کاندید شده‌اند. پس از

تشکیل عبارات هم‌رخداد از ویژگی‌های کاندید شده و شکل‌گیری ۹۹۴۳۹ زوج عبارت، معیار E-Dominance برای ویژگی‌های هم‌رخداد نیز محاسبه شده و با استفاده از حد آستانه‌ی ۰,۰۳۵ برای این معیار، ویژگی‌ها به ۲۱۵۰ ویژگی هم‌رخداد کاهش داده شده، سپس داده‌ها وارد مرحله‌ی دسته‌بندی شده‌اند. جدول ۴-۵ تعداد ویژگی‌های انتخابی از هر دسته را نشان می‌دهد. جمع تعداد ویژگی‌های همه‌ی دسته‌ها به دلیل همپوشانی برخی کلاس‌ها و انتخاب برخی از ویژگی‌ها در بیش از یک کلاس، از کل ویژگی‌های انتخابی بیشتراست. سهم دسته‌ی اقتصادی در انتخاب ویژگی‌های هم‌رخداد، بیشتر از سایر دسته‌هاست.

جدول (۴-۵) تعداد ویژگی‌های هم‌رخداد انتخابی از هر دسته

علمی	فرهنگ و هنر	ورزشی	سیاسی	اقتصادی	
۴۲۰	۳۹۸	۴۷۰	۴۲۳	۵۰۸	تعداد ویژگی انتخابی
٪۲۰	٪۱۹	٪۲۱	٪۲۰	٪۲۴	نسبت تعداد ویژگی‌های هر دسته به کل ویژگی‌ها

جدول ۴-۶، نتایج دسته‌بندی متون را با استفاده از ویژگی‌های هم‌رخداد، با وزن‌دهی باینری و با استفاده از دسته‌بند KNN نشان می‌دهد. جدول ۴-۷، نتایج دسته‌بندی متون را با استفاده از ویژگی‌های هم‌رخداد، با وزن‌دهی باینری و با استفاده از دسته‌بند SVM نشان می‌دهد. با مقایسه‌ی نتایج در جداول ۴-۶ و ۴-۷، دسته‌بند SVM عملکرد بهتری از دسته‌بند KNN داشته است. میانگین Recall به دست آمده برای KNN، ۰,۹۴ است و میانگین Recall به دست آمده برای SVM، ۰,۹۵۲۰ است. از میان ۲۵۰ متن آزمایشی، دسته‌بند SVM، ۳ متن آزمایشی را بیشتر از دسته‌بند KNN، درست دسته‌بندی کرده است.

جدول (۴-۶) نتایج دسته‌بندی داده‌ها با ویژگی‌های هم‌رخداد با وزن‌دهی باینری و دسته‌بند KNN

<i>Accuracy</i>	<i>F₁</i>	<i>Precision</i>	<i>Recall</i>	دسته
۰,۹۸	۰,۹۵۱	۰,۹۲۵	۰,۹۸	اقتصادی
۰,۹۷۲	۰,۹۳۵	۰,۹۲۲	۰,۹۴	سیاسی
۰,۹۷۶	۰,۹۴	۰,۹۴	۰,۹۴	ورزشی
۰,۹۸	۰,۹۵۰	۰,۹۴۱۱	۰,۹۶	فرهنگ و هنر
۰,۹۷۲	۰,۹۲۶	۰,۹۷۸	۰,۸۸	علمی
۰,۹۷۶	۰,۹۴	۰,۹۴۱	۰,۹۴	میانگین

جدول (۴-۷) نتایج دسته‌بندی داده‌ها با ویژگی‌های هم‌رخداد با وزن‌دهی باینری و دسته‌بند SVM

<i>Accuracy</i>	<i>F₁</i>	<i>Precision</i>	<i>Recall</i>	دسته
۰,۹۸	۰,۹۵۲۳	۰,۹۰۹۰	۱	اقتصادی
۰,۹۸۴	۰,۹۶	۰,۹۶	۰,۹۶	سیاسی
۰,۹۸	۰,۹۴۹	۰,۹۵۹	۰,۹۴	ورزشی
۰,۹۸	۰,۹۵۰	۰,۹۴۱	۰,۹۶	فرهنگ و هنر
۰,۹۸	۰,۹۴۷	۱	۰,۹۰	علمی
۰,۹۸۱	۰,۹۵۲	۰,۹۵۴	۰,۹۵۰	میانگین

جداول ۴-۸ و ۴-۹ نتایج دسته‌بندی متون را با استفاده از ویژگی‌های هم‌رخداد و روش وزن‌دهی Co-

TFIDF occur با دسته‌بندهای KNN و SVM نشان می‌دهد. در اینجا نیز دسته‌بند SVM بهتر از KNN

عمل کرده است.

جدول (۴-۸) نتایج دسته‌بندی داده‌ها با ویژگی‌های هم‌رخداد با وزن‌دهی Co-occur TFIDF و دسته‌بند KNN

دسته	Recall	Precision	F ₁	Accuracy
اقتصادی	۱	۰,۹۶۲	۰,۹۸۰	۰,۹۹۲
سیاسی	۰,۹۸۰	۰,۹۴۲	۰,۹۶۱	۰,۹۸۴
ورزشی	۰,۹۶۰	۰,۹۶	۰,۹۶۰	۰,۹۸۴
فرهنگ و هنر	۱	۰,۹۶۲	۰,۹۸۰	۰,۹۹۲
علمی	۰,۸۸	۱	۰,۹۳۶	۰,۹۷۶
	۰,۹۶,۴	۰,۹۶۵	۰,۹۶۳	۰,۹۸۶

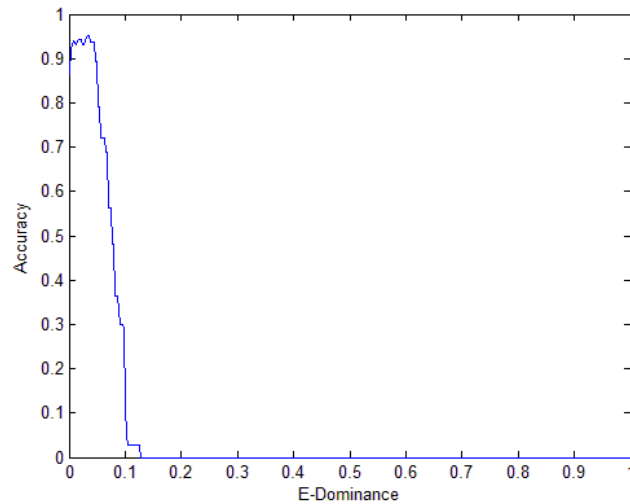
جدول (۴-۹) نتایج دسته‌بندی داده‌ها با ویژگی‌های هم‌رخداد با وزن‌دهی Co-occur TFIDF و دسته‌بند SVM

دسته	Recall	Precision	F ₁	Accuracy
اقتصادی	۱	۰,۹۴۳	۰,۹۷۱	۰,۹۸
سیاسی	۰,۹۸۰	۰,۹۶۱	۰,۹۷۰	۰,۹۹۶
ورزشی	۰,۹۸۰	۰,۹۸۰	۰,۹۸۰	۰,۹۹۲
فرهنگ و هنر	۱	۰,۹۸۰	۰,۹۹۰	۰,۹۸۸
علمی	۰,۹۰	۱	۰,۹۴۷	۰,۹۸۸
میانگین	۰,۹۷۲	۰,۹۷۳	۰,۹۷۲	۰,۹۸۹

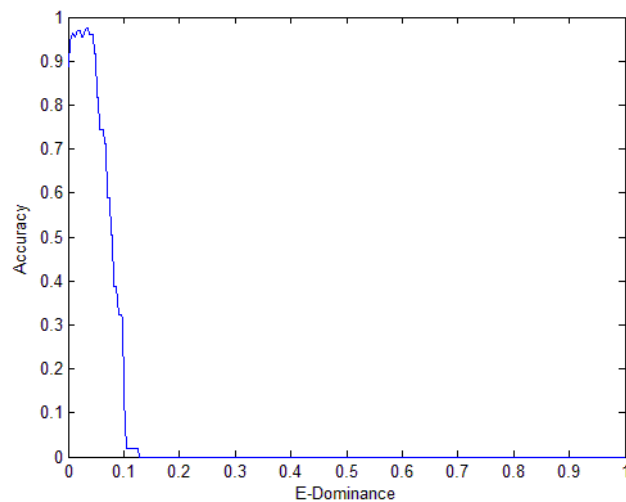
شکل ۴-۴ تاثیر مقدار معیار E-Dominance را در دقت دسته‌بندی سیستم با استفاده از ویژگی‌های

هم‌رخداد و دسته‌بند SVM، نشان می‌دهد. در شکل ۴-۴-الف ویژگی‌ها به صورت باینری وزن‌دهی شده-

اند و در شکل ۴-۴ ب ویژگی‌ها با روش Co-occur TFIDF وزن‌دهی شده‌اند. حد آستانه انتخابی برای این معیار، مقدار ۰,۰۳۵ است که در این حد آستانه، سیستم، بیشترین دقت را داشته است.



الف



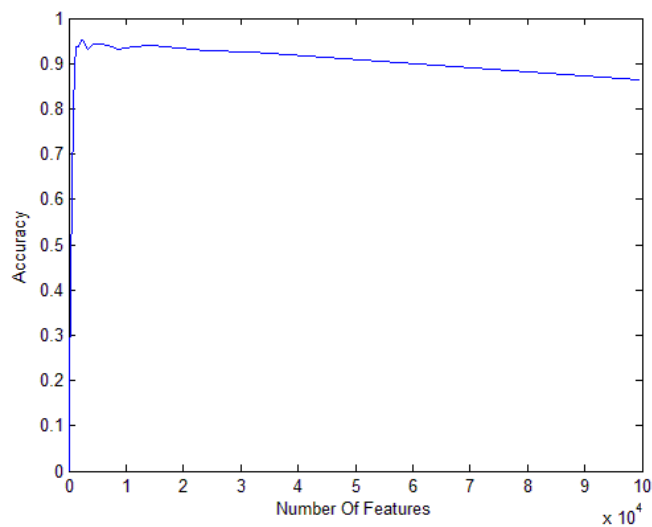
ب

شکل (۴-۴) تاثیر حد آستانه معیار E-Dominance در دقت دسته‌بندی با ویژگی‌های هم‌رخداد با دسته‌بند SVM (شکل الف مربوط به وزن‌دهی باینری و شکل ب مربوط به وزن‌دهی Co-occur TFIDF)

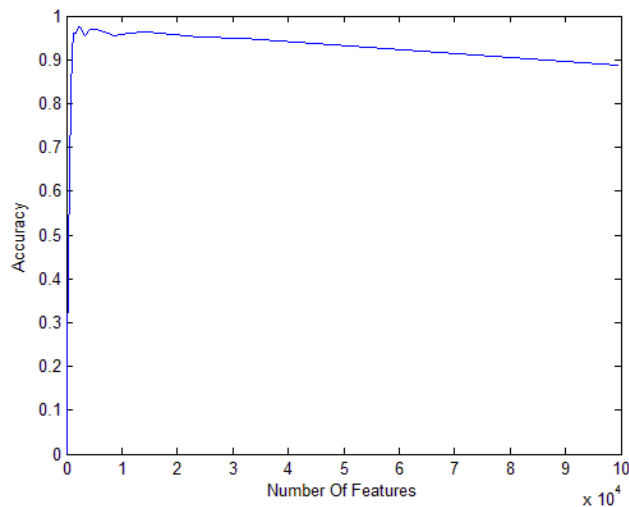
شکل ۴-۵، نیز تاثیر تعداد ویژگی‌های انتخابی هم‌رخداد را در دقت دسته‌بندی سیستم با استفاده از

دسته‌بند SVM نشان می‌دهد. در شکل ۴-۵ الف ویژگی‌ها به صورت باینری وزن‌دهی شده‌اند و در شکل

۴-۵-ب ویژگی‌ها با روش Co-occur TFIDF وزن‌دهی شده‌اند. سیستم با تعداد ۲۱۵۰ ویژگی هم‌خداد و وزن‌دهی باینری، توانسته است بهترین دقت را داشته باشد.



الف



ب

شکل (۴-۵) تاثیر تعداد ویژگی‌های انتخابی هم‌خداد در دقت دسته‌بندی با دسته‌بند SVM (شکل الف مربوط به وزن‌دهی باینری و شکل ب مربوط به وزن‌دهی Co-occur TFIDF)

۴-۳-۲- بررسی تاثیر تعداد ویژگی‌های توسعه‌گر در دقت دسته‌بندی

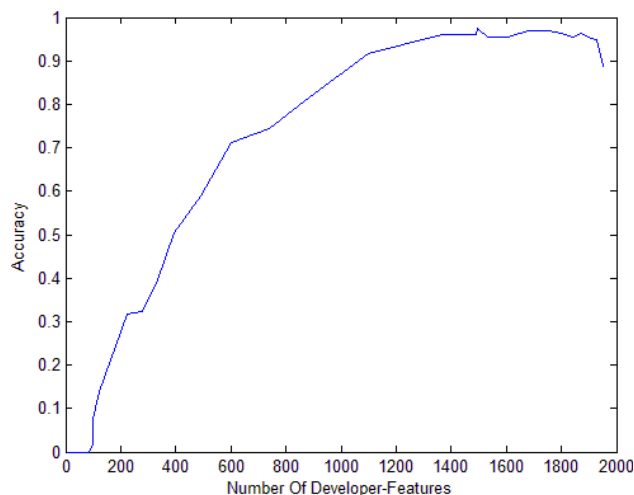
در حالی که تعداد ویژگی‌های انتخابی بهترین دقت را در سیستم نشان می‌دهد، ۲۱۵۰ ویژگی هم‌رخداد وجود دارد. تعداد ویژگی‌های توسعه‌گر نیز، ۱۴۹۵ ویژگی هستند جدول ۴-۱۰ تعداد ویژگی‌های توسعه‌گر انتخابی به تفکیک هر دسته را نشان می‌دهد. جمع تعداد ویژگی‌های همه دسته‌ها از کل ویژگی‌های توسعه‌گر بیشتر است. این امر ناشی از همپوشانی برخی کلاس‌ها و انتخاب برخی از ویژگی‌های توسعه‌گر در بیش از یک کلاس است. همچنین این جدول نسبت تعداد ویژگی‌های توسعه‌گر انتخابی از هر دسته را نسبت به کل ویژگی‌های انتخابی آن دسته و نسبت به کل ویژگی‌های توسعه‌گر نشان می‌دهد. نسبت‌ها نشان می‌دهد که دسته اقتصادی بیشترین و دسته فرهنگی و هنر کمترین نقش را در تعداد ویژگی‌های توسعه‌گر دارند.

جدول (۴-۱۰) ویژگی‌های توسعه‌گر انتخابی از هر دسته

علمی	فرهنگ و هنر	ورزشی	سیاسی	اقتصادی	
۲۰۵	۱۲۷	۳۰۲	۳۲۱	۴۰۳	تعداد ویژگی توسعه‌گر انتخابی
٪۴۸	٪۳۱	٪۶۴	٪۷۵	٪۷۹	نسبت تعداد ویژگی‌های توسعه‌گر دسته به کل ویژگی‌های دسته
٪۱۴	٪۸	٪۲۰	٪۲۸	٪۳۴	نسبت تعداد ویژگی‌های توسعه‌گر دسته به کل ویژگی‌های توسعه‌گر

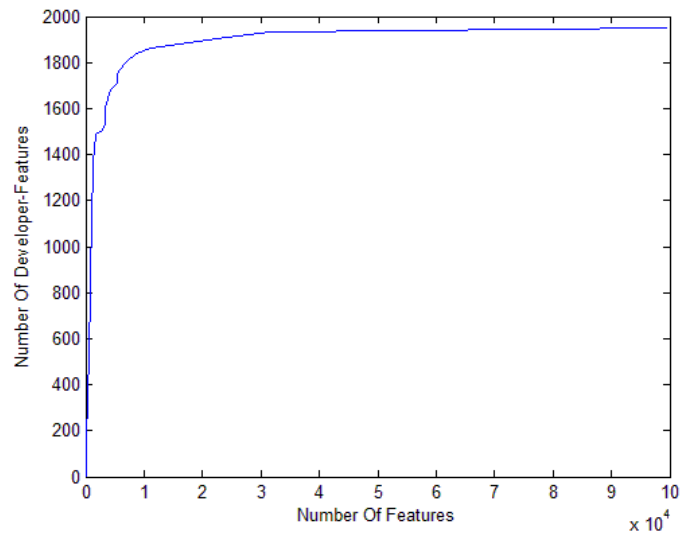
شکل ۴-۶ تاثیر تعداد ویژگی‌های توسعه‌گر انتخابی در دقت سیستم دسته‌بندی در حالی که ویژگی‌ها با استفاده از معیار E-Dominance انتخاب شده و ویژگی‌های هم‌رخداد را تشکیل داده‌اند و با استفاده از ویژگی‌های توسعه‌گر توسعه داده شده‌اند و پس از توسعه، ویژگی‌ها با استفاده از روش وزن‌دهی Co-

TFIDF occur وزن‌دهی شده و در نهایت با SVM، دسته‌بندی شده‌اند، را نشان می‌دهد. همانطور که می‌بینید در جاییکه تعداد ویژگی‌های توسعه‌گر ۱۴۹۵ بوده است، سیستم بهترین دقت را داشته است.



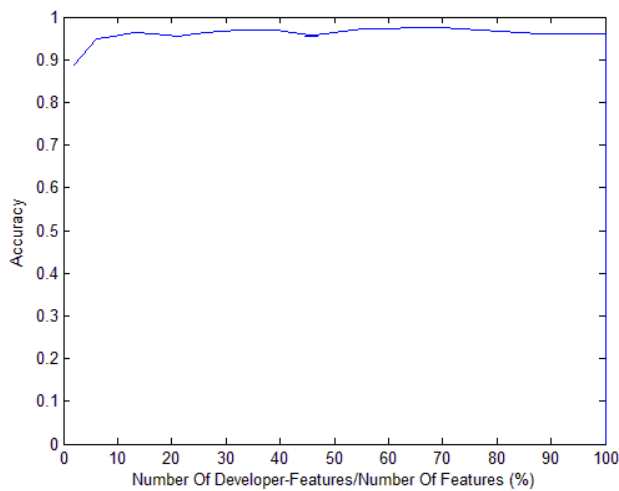
شکل (۴-۶) تاثیر تعداد ویژگی‌های توسعه‌گر در دقت دسته‌بندی

شکل ۴-۷ تعداد ویژگی‌های توسعه‌گر انتخابی را نسبت به کل ویژگی‌های هم‌رخداد انتخابی، در روند کاهش ویژگی‌ها نشان می‌دهد. در حالی که تعداد کل ویژگی‌ها ۹۹۴۳۸ بوده است، تعداد ویژگی‌های توسعه‌گر، ۱۹۵۰ عدد بوده است. در بهترین حالت سیستم که تعداد ویژگی‌های کل به ۲۱۵۰ ویژگی کاهش یافته است، تعداد ویژگی‌های توسعه‌گر ۱۴۹۵ عدد بوده و ۷۰٪ از ویژگی‌های کل بوده است. در روند کاهش ویژگی‌ها با معیار E-Dominance، از زمانی که تعداد ویژگی‌های کل به ۱۳۶۵ عدد رسید، تعداد ویژگی‌های توسعه‌گر نیز با تعداد ویژگی‌های انتخابی برابر بوده است. در شکل ۴-۷ بخش اول نمودار، یک خط راست بوده و بعد از آن به صورت منحنی درآمدن است.



شکل (۴-۷) تعداد ویژگی‌های توسعه‌گر نسبت به تعداد کل ویژگی‌ها

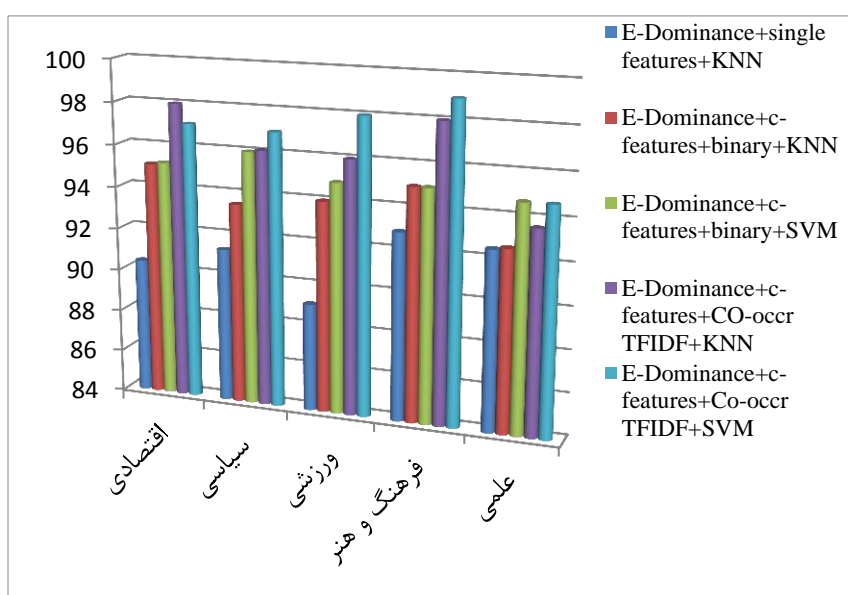
شکل ۴-۸ نیز تاثیر نسبت تعداد ویژگی‌های توسعه‌گر را به تعداد کل ویژگی‌ها در روند کاهش ویژگی‌ها، برای دقت دسته‌بندی نشان می‌دهد. در شکل می‌بینید که وقتی نسبت تعداد ویژگی‌های توسعه‌گر را به تعداد کل ویژگی‌ها، ۷۰٪ بوده است، سیستم بهترین دقت را ارائه داده است.



شکل (۴-۸) تاثیر نسبت تعداد ویژگی‌های توسعه‌گر به تعداد کل ویژگی‌ها در دقت دسته‌بندی

۳-۳-۴- مقایسه نتایج روش پیشنهادی در ۵ حالت مختلف در ۵ کلاس

روش پیشنهادی در ۵ حالت مختلف، پیاده سازی و نتایج آن در بخش‌های قبل ارائه گردید. شکل ۴-۹، معیار F1 را در ۵ حالت مختلف و در ۵ کلاس مختلف نشان می‌دهد. تقریباً در همه‌ی کلاس‌ها، معیار F1، با به‌کارگیری ویژگی‌های هم‌رخداد انتخابی با معیار E-Dominance و وزن‌دهی به روش Co-occur TFIDF و دسته‌بند SVM، بیشترین مقدار را داشته است.



شکل (۹-۴) مقایسه‌ی معیار F1 برای روش‌های مختلف در ۵ کلاس

۴-۴- مقایسه روش پیشنهادی با چند روش موجود

در جدول ۴-۱۱، روش پیشنهادی در حالت‌های مختلف با استفاده از دو روش وزن‌دهی باینری و Co-occur TFIDF و دو دسته‌بند KNN و SVM و با استفاده از ویژگی‌های منفرد و هم‌رخداد آورده شده است. در این جدول نتایج مربوط به مراجع [۵۱، ۱۱۶، ۱۱۸] نیز که برای مقایسه با روش پیشنهادی انتخاب شده‌اند، آورده شده است. مراجع نامبرده به این دلیل انتخاب شده‌اند که در زبان فارسی آزمایشات خود را انجام داده و از پایگاه داده‌ی همشهری استفاده کرده‌اند. همچنین تعداد متون آموزشی و آزمایش و

تعداد کلاس‌ها و نیز موضوعات کلاس‌ها، در این مقالات مشابه این پایان‌نامه است. مرجع [۷] نیز به عنوان مرجع اصلی این پایان‌نامه نیز برای مقایسه انتخاب شده است. برای مقایسه روش پیشنهادی با مرجع [۷]، روش مورد نظر بر روی متون انتخابی این پایان‌نامه پیاده‌سازی و تست شده است. مرجع اصلی نتایج خود را بر روی متون زبان انگلیسی ارائه داده است. در فصل دوم روش‌های استفاده شده در این مراجع آورده شده است. همان‌طور که نتایج در جدول ۴-۱۱ نشان می‌دهد، روش پیشنهادی در حالتی که ویژگی‌ها با استفاده از معیار E-Dominance انتخاب شده و ویژگی‌های هم‌رخداد را تشکیل داده‌اند و نیز با استفاده از روش Co-occur TFIDF وزن‌دهی شده و در نهایت با دسته‌بند SVM، دسته‌بندی شده‌اند، بهترین نتیجه را در میان نتایج ارائه داده است.

نکته قابل توجه در نتایج روش پیشنهادی این است که مقادیر معیارهای F1 و Recall در بیشتر حالات یکسان و یا بسیار نزدیک به هم بوده است. دلیل نزدیکی این مقادیر، نزدیک بودن مقادیر Recall و Precision است که تشکیل دهنده‌ی معیار F1 هستند.

۴-۵- نتیجه‌گیری

در این فصل پس از معرفی معیارهای ارزیابی، نتایج آزمایشات ارائه گردید. معیار E-Dominance با چند معیار ارزش‌گذاری داده دیگر مقایسه شده و کارایی بیشتر را نسبت به دیگر معیارها نشان داد. آزمایشات انجام شده در ۵ حالت مختلف و با استفاده از دسته‌بندهای KNN و SVM انجام شد. دو روش وزن‌دهی باینری و Co-occur TFIDF نیز برای ویژگی‌ها هم‌رخداد به کار رفته‌اند. در پایان، در حالتی که ویژگی‌ها با استفاده از معیار E-Dominance انتخاب شده و ویژگی‌های هم‌رخداد را تشکیل داده‌اند و ۷۰٪ از ویژگی‌ها به عنوان ویژگی‌های توسعه‌گر، به توسعه و غنی‌سازی داده‌ها پرداخته‌اند و نیز با استفاده از روش Co-occur TFIDF وزن‌دهی شده و در نهایت با دسته‌بند SVM، دسته‌بندی شده‌اند، بهترین نتیجه را در میان نتایج ارائه داده است.

جدول (۴-۱۱) مقایسه روش پیشنهادی با چند روش دیگر

توضیح	F1	Recall	تعداد ویژگی‌ها	وزن‌دهی	معیار کاهش ویژگی	ویژگی‌ها	دسته‌بند	روش
کاربرد فارسی‌نت	۹۳,۳۰	۹۴,۳۹	۳۱۹۰	TFIDF	Chi-square	کیسه لغات	NN	مرجع [۵۱]
کاربرد گنجوازه معنایی	۸۶	۸۳	-	TFIDF	ICF+UNI	کیسه لغات	SVM	مرجع [۱۱۶]
-	۸۹	۸۷	-	TFCRF	PCA + Genetic Algorithm	کیسه لغات	1-NN	مرجع [۱۱۸]
کاربرد ویژگی‌های توسعه‌گر	۹۲,۶۱	۹۳,۲	۷۰۰۷	باینری	Dominance	زوج کلمات هم‌خداد	KNN	مرجع [۷]
کاربرد ویژگی‌های توسعه‌گر	۹۴	۹۴,۸	۷۰۰۸	باینری	Dominance	زوج کلمات هم‌خداد	SVM	مرجع [۷]
-	۹۱,۲۳	۹۱,۲	۱۰۴۱	TFIDF	E-Dominance	کیسه لغات	KNN	روش پیشنهادی
کاربرد ویژگی‌های توسعه‌گر	۹۴	۹۴	۲۱۵۰	باینری	E-Dominance	زوج کلمات هم‌خداد	KNN	روش پیشنهادی
کاربرد ویژگی‌های توسعه‌گر	۹۵,۵	۹۵	۲۱۵۰	باینری	E-Dominance	زوج کلمات هم‌خداد	SVM	روش پیشنهادی
کاربرد ویژگی‌های توسعه‌گر	۹۶,۴	۹۶,۴۰	۲۱۵۰	Co-occur TFIDF	E-Dominance	زوج کلمات هم‌خداد	KNN	روش پیشنهادی
کاربرد ویژگی‌های توسعه‌گر	۹۷,۲	۹۷,۲	۲۱۵۰	Co-Occur-TFIDF	E-Dominance	زوج کلمات هم‌خداد	SVM	روش پیشنهادی

فصل پنجم

نتیجه‌گیری و کارهای آینده

۵-۱- خلاصه تحقیق

انتخاب ویژگی‌ها در متون و نحوه‌ی وزن‌دهی به ویژگی‌های انتخاب شده، از مهم‌ترین مسائل دسته‌بندی خودکار متون هستند. در این تحقیق یک معیار ارزش‌گذاری ویژگی، بنام E-Dominance معرفی شد. این معیار با سایر معیارهای ارزش‌گذاری و وزن‌دهی متداول در حوزه‌ی دسته‌بندی متون مقایسه شد و نشان داده شد که با استفاده از معیار فوق می‌توان به دقتی همسان و یا بیشتر از روش‌های متداول با تعداد بسیار کمتری ویژگی، دست پیدا کرد. کاهش ویژگی‌ها با استفاده از این معیار، حدود ۹۴٪ بوده است. هدف از معرفی این معیار، انتخاب و کاهش ویژگی‌ها برای تشکیل و استخراج ویژگی‌های هم‌رخداد است. با انتخاب کلمات با اهمیت در هر دسته با استفاده از معیار E-Dominance، ویژگی‌های هم‌رخداد تشکیل شد. ویژگی‌های هم‌رخداد هیچ‌گاه در زبان فارسی به عنوان ویژگی‌های متون برای دسته‌بندی استفاده نشده‌اند. در سایر زبان‌ها نیز بندرت به عنوان ویژگی در نظر گرفته شده‌اند. کاربرد اصلی ویژگی‌های هم‌رخداد در بازیابی اطلاعات است. از ویژگی‌های هم‌رخداد تشکیل شده، تعداد زیادی بر اساس حد آستانه معیار فوق، از بردار ویژگی‌ها حذف شدند سپس ویژگی‌هایی که تعداد رخداد بیشتری در یک دسته داشته‌اند، به عنوان ویژگی‌های توسعه‌گر، برای توسعه‌ی داده‌های آموزش و آزمایش انتخاب می‌شوند. توسعه‌ی داده‌ها باعث کم شدن شباهت‌های بین کلاسی و زیاد شدن شباهت داخل کلاسی می‌شود. توسعه‌ی داده‌ها در داده‌های آزمایشی نیز به گونه‌ای اعمال می‌شود که بدون توجه به کلاس داده، بتوان داده‌های غنی‌تری را برای دسته‌بندی به سیستم ارائه کرد. پس از توسعه‌ی داده‌ها از روش وزن‌دهی Co-occur TFIDF که یک روش توسعه یافته از TFIDF معمول است، برای ویژگی‌های هم‌رخداد به کار گرفته شده است. در کارهای گذشته فقط از روش وزن‌دهی باینری برای ویژگی‌های هم‌رخداد استفاده شده است. پس از استخراج ویژگی و وزن‌دهی به آنها، فاز بعدی دسته‌بندی متون است. در

این کار، از دسته‌بندهای KNN و SVM برای دسته‌بندی استفاده شده است. در بیشتر تحقیقات انجام شده، معمولا SVM از دقت بالاتری نسبت به سایر دسته‌بندها برخوردار بوده است. در روش پیشنهادی نیز، این دسته‌بند نتایج بهتری را ارائه داده است. نتایج آزمایشات در مقایسه با سایر کارها نشان داده است که روش پیشنهادی در حالتی که ویژگی‌ها با استفاده از معیار E-Dominance انتخاب شده و ویژگی‌های هم‌رخداد را تشکیل داده‌اند و نیز با استفاده از روش Co-occur TFIDF وزن‌دهی شده و در نهایت با دسته‌بند SVM، دسته‌بندی شده‌اند، بهترین نتیجه را ارائه داده است. برای نشان دادن کارایی روش پیشنهادی در مقایسه با سایر تحقیقات، از معیارهای Recall و F1 استفاده شده است. این دو معیار در همه‌ی حالت‌ها در روش پیشنهادی، مقادیر بیشتری نسبت به سایر روش‌ها داشته و بهبود دسته‌بندی را نشان دادند.

۵-۲- پیشنهاد برای کارهای آینده

گرچه این کار تحقیقاتی پایان یافته است، اما مسلما دارای نقاط ضعف و کمبودهایی است که در آینده برطرف خواهد شد. از جمله پیشنهادات برای اتمام‌ی کار حاضر عبارتند از:

- طراحی یک ریشه‌یاب موثر در زبان فارسی و کاربرد آن در انتخاب ویژگی‌ها
- توسعه‌ی کلمات ایست استخراجی در زبان فارسی به لیست کامل و استاندارد که مورد استفاده همگان باشد.
- ارائه‌ی یک معیار موثر در انتخاب ویژگی‌های توسعه‌گر برای توسعه‌ی داده‌ها، در این پایان نامه از حد آستانه DF برای این کار بر اساس آزمون و خطا استفاده شده است.
- توسعه داده‌ها با استفاده از ویژگی‌های توسعه‌گر در چند سطح، به نحوی که در هر مرحله ویژگی‌های توسعه‌گر جدید استخراج و مورد استفاده در توسعه داده‌ها قرار گیرند.

مراجع

- [۱] مدنی. ص. س. ، "دسته بندی اسناد فارسی به کمک هستان شناسی فارس نت" ، پایان نامه کارشناسی ارشد، دانشکده کامپیوتر و فن آوری اطلاعات، دانشگاه صنعتی شاهرود، بهمن ۱۳۹۱.
- [۴] اصدقی. ف. ، "طبقه بندی اسناد متنی با استفاده از سیستم ایمنی مصنوعی" ، پایان نامه کارشناسی ارشد، دانشکده کامپیوتر و فن آوری اطلاعات، دانشگاه علم و صنعت ایران ، اسفند ۱۳۸۸.
- [۲۸] دستجردی ا. ، زنجانی م. ، "روش جدید خوشه بندی مستندات متنی الکترونیکی فارسی به کمک واژه-هستان شناسی فارس نت" ، اولین کنفرانس ملی دانش پژوهان کامپیوتر و فناوری اطلاعات، ۱۳۹۰.
- [۲۹] ایران پور م. ، "بررسی مشکلات تعیین حدود جمله و کلمه" ، سمینار کارشناسی ارشد، دانشگاه علم و صنعت ایران، ۱۳۸۶.
- [۳۰] داداش میری پ. ، "تشخیص انتهای کلمات و ایجاد فاصله میان کلمات" ، پایان نامه کارشناسی، دانشگاه علم و صنعت ایران، ۱۳۸۰.
- [۳۱] بیجن خان م. ، "تشخیص کسره ی اضافه" ، طرح تحقیقاتی، پژوهشگاه فرهنگ، هنر و ارتباطات، تهران، ۱۳۸۴.
- [۳۴] حافظی م. ، ثامتی ح. ، منصورى ن. ، بحرانی، موثق ، "ارائه یک مدل دستوری برای بهبود دقت سیستم های بازشناسی گفتار پیوسته ی فارسی" ، دومین گارگاه پژوهشی زبان فارسی و رایانه، ص. ۸۰ تا ۹۱، تهران ، ۱۳۸۵.
- [۳۵] حسامی فرد ر. ، قاسم خانی غ.ر. ، "طراحی یک الگوریتم ریشه یابی برای زبان فارسی" ، یازدهمین کنفرانس بین المللی انجمن کامپیوتر ایران، ص. ۵۱۵ تا ۵۱۹، تهران، ۱۳۸۴.
- [۳۶] مشکی م. ، "خلاصه سازی گزینشی چند سندی متون فارسی" ، پایان نامه ارشد، دانشکده کامپیوتر، دانشگاه علم و صنعت ، ۱۳۸۸.
- [۳۷] ابوالحسنی ح. ، اسماعیلی س. ، محمدی م. ، "یک خلاصه ساز آماری برای زبان فارسی" ، یازدهمین کنفرانس بین المللی انجمن کامپیوتر ایران، ۱۳۸۴.
- [۴۰] آنالویی م. ، مشکی م. ، "خلاصه سازی چند سندی متون فارسی با استفاده از یک روش مبتنی بر خوشه بندی" ، اولین کنفرانس ملی مهندسی نرم افزار ایران، تهران، ۱۳۸۷.

- [۵۱] حسن پور ح. ، مدنی ص. س. ، "بهبود دقت سیستم دسته بندی خودکار اسناد فارسی به کمک هستان شناسی فارس نت،" مجله علمی پژوهشی رایانش نرم و فن آوری اطلاعات، جلد ۳، شماره ۱، ۱۳۹۳ .
- [۶۵] حسن پور ح. ، قنبری سرخی ع. ، پارسی ا. ، "استخراج بهترین ویژگی از متون فارسی با استفاده از تجزیه و تحلیل مولفه های اصلی با کمک میانگین یادآوری و الگوریتم ژنتیک" ، نخستین کنفرانس بین المللی پردازش خط و زبان فارسی، ۱۳۹۱.
- [۱۰۹] عربی نژی س. ، وحیدی اصل م. ، مینایی بیگدلی ب. ، "استخراج کلمات کلیدی جهت طبقه بندی متون فارسی"، اولین کنفرانس داده کاوی ایران، ۱۳۸۶.
- [۱۱۰] عربی نژی س. ، وحیدی اصل م. ، مینایی بیگدلی ب. ، "رویکردی جدید جهت طبقه بندی مستندات فارسی با استفاده از الگوریتم کاوشگر مورچه"، هفتمین همایش زبان شناسی ایران، ۱۳۸۶.
- [۱۱۱] احسان بصیری ا. ، نعمتی ش.، قاسم آقایی ن. ، "مقایسه طبقه بندی متون فارسی با استفاده از الگوریتم *KNN* و *FKNN* و انتخاب ویژگی ها بر اساس بهره ی اطلاعات و فرکانس سند"، سیزدهمین کنفرانس سالانه انجمن کامپیوتر ایران، ۱۳۸۶.
- [۱۱۲] محمدی م.، مینایی بیگدلی ب.، آنالویی م. ، "استفاده از شبکه های عصبی *CC4* برای طبقه بندی اسناد فارسی"، دومین کنفرانس داده کاوی ایران، ۱۳۸۷.
- [۱۱۳] احمدزاده، ا.ب.ح.ف.م.ح.س.م.، ، "طبقه بندی متون خبری فارسی با استفاده از الگوریتم *Naive Bayes*"، دومین کنفرانس داده کاوی ایران، ۱۳۸۷.
- [۱۱۴] محمدی م. ، مینایی بیگدلی ب.، آنالویی م. ، "طبقه بندی اسناد متنی با استفاده از شبکه های عصبی پرسپترون چند لایه"، دومین کنفرانس داده کاوی ایران، ۱۳۸۷.
- [۱۱۶] مقصودی ن. ، همایون پور م. ، "ارائه روشی جدید در طبقه بندی متون فارسی با استفاده از دانش معنایی"، پانزدهمین کنفرانس بین المللی سالانه انجمن کامپیوتر ایران، ۱۳۸۸.
- [۱۱۷] قنبری سرخی ع. ، ابراهیمی ف. ، "بهبود عملکرد طبقه بندی متون فارسی با استفاده از تجزیه و تحلیل مولفه های اصلی با کمک معیار میانگین یادآوری و دقت" ، چهاردهمین کنفرانس دانشجویی مهندسی برق ایران، ۱۳۹۰.
- [۱۱۸] حسن پور ح. ، قنبری سرخی ع. ، پارسی ا.، "استخراج بهترین ویژگی از متون فارسی با استفاده از تجزیه و تحلیل مولفه های اصلی با کمک میانگین یادآوری و الگوریتم ژنتیک" ، نخستین کنفرانس بین المللی پردازش خط و زبان فارسی، ۱۳۹۱.

- [۱۱۹] طاهری نیا م. ، "دسته بندی متون فارسی با استفاده از یادگیری نیمه نظارت شده"، چهارمین کنفرانس مهندسی برق و الکترونیک ایران، ۱۳۹۱.
- [۱۲۰] زمانی م. ، هدایت ر. ، صادق زاده م. ، "دسته بندی متون فارسی با استفاده از روش آنالیز معنایی پنهان احتمالاتی"، اولین همایش ملی کاربرد سیستم های هوشمند (محاسبات نرم) در علوم و صنایع، ۱۳۹۲.
- [۱۲۱] برفامی م. ، فاطری س. ، "استفاده از ترکیب شبکه های عصبی جهت دسته بندی متون فارسی مبتنی بر الگوریتم های GA, KNN, PCA جهت کاهش ویژگی"، اولین همایش ملی رویکردهای نوین در مهندسی کامپیوتر و بازیابی اطلاعات، ۱۳۹۲.
- [۱۲۲] عابدی نیا م. ، الله دادی ل. ، شیخی ف. ، "کاوش متون فارسی دروب با استفاده از تحلیل معنایی و روش طبقه بندی چندلایه ای اطلاعات"، اولین همایش ملی رویکردهای نوین در مهندسی کامپیوتر و بازیابی اطلاعات، ۱۳۹۲.
- [۱۲۳] صفیان ز. ، بلداجی ، نادرری دهکردی م. ، رستگار ح. ، "دسته بندی متون به کمک الگوریتم KNN بهبود یافته مبتنی بر خوشه های وزن دار"، اولین کنفرانس ملی نوآوری در مهندسی کامپیوتر و فناوری اطلاعات، ۱۳۹۲.
- [2] A.Khan, B.B., L.Hong Lee, K.khan, "A Review of Machine Learning Algorithms for Text-Documents Classification" , *Journal of Advances in Information Technology*, VOL. 1, NO. 1, February 2010.
- [3] B.Zhang, "Learning Features for Text Classification" , *Thesis of Doctor of Philosophy, University of Washington*, 2013.
- [5] F.Sebastiani, "Machine learning in automated text categorization", *ACM Computing Surveys*, Vol. 34, pp. 1-47 ,2002.
- [6] M.YUAN, Y.X.O., Z. XIONG, "A Text Categorization Method using Extended Vector Space Model by Frequent Term Sets" , *Journal of Advances in Information Science and Engineering*, 29, 99-114 ,2013.
- [7] F. Figueiredo, L.R., T. Couto, T. Salles, M. A. Goncalves, W. Meira Jr, "Word co-occurrence features for text classification", *Information Systems*, 36, pp.843-858, 2011.
- [8] J.Feng, H.L., "Sat-mod:moderate itemset fittest for text classification", *in:WWW'05:Special Interest Tracks and Posters of the 14th International Conference on World Wide Web, New York, USA*, pp.1054-1055, 2005.

- [9] J.H.H.Cheng, X.Y., C.Hsu, "Discriminative frequent pattern analysis for effective classification", in: *ICDE '07: Proceedings of International Conference on Data Engineering, Istanbul, Turkey, pp.716–725, 2007.*
- [10] O.R.Zaiane, M.L.A., "Classifying text documents by associating terms with text categories, in: *CRPITS '02 :Proceedings of The 13th Australasian Conference on Database Technologies, Australian Computer Society,Inc. , Darling hurst, Australia, pp.215–222, 2002.*
- [11] R.Rak, W.S., O.R.Zaiane, M.L.Antonie, "Considering re-occurring features in associative classifiers", in: *PAKDD '05: 9th Pacific Asia Conference on Knowledge Discovery and Data Mining, Hanoi, Vietnam, pp.240–24, 2005.*
- [12] C. Largeton, C.M., M. Gery, "Entropy based feature selection for text categorization", *ACM Symposium on Applied Computing, TaiChung : Taiwan, Province Of China, 2011.*
- [13] Joachims, T., "Text categorization with support vector machines: Learning with many relevant features", In *Proceedings of ECML '98, 10th European Conference on Machine Learning, No.1398, pp137-142, Chemnitz, DE, Springer Verlag, Heidelberg, DE, 1998.*
- [14] T.Gaustad, G.B., R. Groningen, "Accurate Stemming of Dutch for Text Classification," *Language and Computers, 2002.*
- [15] K.Tzeras, S.H., "Automatic indexing based on Bayesian inference networks," presented at the *Proceedings of the 16th annual international ACM SIGIR conference on Research and development in information retrieval, Pittsburgh, Pennsylvania, United States, 1993.*
- [16] N. J. Belkin , W.B.C., , "Information filtering and information retrieval: two sides of the same coin? ", *Commun. ACM, vol. 35, pp. 29-38, 1992.*
- [17] P. J. Hayes, P.M.A., I. B. Nirenburg, L. M. Schmandt, "TCS: a shell for content-based text categorization," presented at the *Proceedings of the sixth conference on Artificial intelligence applications, Santa Barbara, California, United States, 1990.*
- [18] E. D. Liddy, W.P., and E. S. Yu, "Text categorization for multiple users based on semantic features from a machine-readable dictionary," *ACM Trans. Inf. Syst., vol. 12, pp. 278-295, 1994.*
- [19] G.Escudero, L.M., G.Rigau, "Boosting Applied to Word Sense Disambiguation," in *Machine Learning: ECML 2000. vol. 1810, R. López de Mántaras and E. Plaza, Eds., ed: Springer Berlin Heidelberg, pp. 129-141, 2000.*
- [20] Gale, W.A., Kenneth W. Church, and David Yarowsky, "A method for disambiguating word senses in a large corpus", *Computers and the Humanities*26, no. 5 ,PP. 415-439, 1992.

- [21] Fürnkranz, J., *"Exploiting Structural Information for Text Classification on the WWW," presented at the Proceedings of the Third International Symposium on Advances in Intelligent Data Analysis, 1999.*
- [22] S.Chakrabarti, B.D., P .Indyk, *"Enhanced hypertext categorization using hyperlinks," presented at the Proceedings of the 1998 ACM SIGMOD international conference on Management of data, Seattle, Washington, United States, 1998.*
- [23] S.Eyheramendy, A.G., W.-H. Ju, D. D. Lewis, and D. Madigan, *"Sparse Bayesian Classifiers for Text Categorization," Joint Statistical Meeting in San Francisco, 2003*
- [24] A.Dasgupta, P.D., B.Harb, V.Josifovski, M.W.Mahoney,, *"Feature selection methods for text classification," presented at the Proceedings of the 13th ACM SIGKDD international conference on Knowledge discovery and data mining, San Jose, California, USA, 2007.*
- [25] H.Kim, P.H., H.Park, *"Dimension Reduction in Text Classification with Support Vector Machines," J. Mach. Learn. Res., vol. 6 ,pp. 37-53, 2005.*
- [26] M.Hassel, N.M., *"FarsiSum: a Persian text summarizer," presented at the Proceedings of the Workshop on Computational Approaches to Arabic Script-based Languages, Geneva, Switzerland, 2004.*
- [27] Megerdooonian, R.Z.K., *"Processing Persian Text: Tokenization in the Shiraz Project," Memoranda in Computer and Cognitive Science, vol. MCCS-00-322, April 2000.*
- [32] M.R.Davarpanah, M.S., M. Aramideh, *"Farsi lexical analysis and stop word list," Library Hi Tech, vol. 27 pp. 435-449, 2009.*
- [33] K.Taghva, R.B., M.Sadeh, " *A List of Farsi Stopwords", ISRI Technical Report, 2003.*
- [38] K. Taghva, R.B., and M. Sadeh, , *"A stemming algorithm for the Farsi language", in Information Technology: Coding and Computing. ITCC 2005. International Conference on, pp. 158-162 Vol. 1, 2005.*
- [39] Hovy, E., *"Parsimonious and Profligate Approaches to the Question of Discourse Structure Relations", in Workshop on NLG, 1990.*
- [41] C.Bouras, V.T., *"W-kmeans: clustering news articles using wordNet," presented at the Proceedings of the 14th international conference on Knowledge-based and intelligent information and engineering systems: Part III, Cardiff, UK, 2010.*

- [42] C.Yi-Hsing , H.H.-Y., "An Automatic Document Classifier System based on Naive Bayes Classifier and Ontology," *International Conference on in Machine Learning and Cybernetics*, pp. 3144-3149., 2008.
- [43] Z.Karimi , M.S., "Automatic Text Summarization Systems Farsi," in *12th International Conference of Computer Society of Iran*, 2006.
- [44] D.D.Lewi, "Evaluating and optimizing autonomous text classification systems", In *Proceeding of SIGIR '95 ,18th ACM International Conference on Research and Development in Information Retrieval*. pp. 246-254, 1995.
- [45] D.Chai, "A Trainable Web Page Document Summarizer," *Natural Language Processing, CS224N/Ling237*, 2000.
- [46] H.P.Luhn, "A statistical approach to mechanized encoding and searching of literary information," *IBM J. Res. Dev.*, vol. 1, pp. 309-317, 1957.
- [47] L.E. , K.J., "Text Categorization with Support Vector Machines. How to Represent Texts in Input Space," *Machine Learning*, vol. 46, pp. 423-444 , 2002.
- [48] Jones, K.S., "Indexing Term Weighting," *Information Storage and Retrieval*, vol. 9, pp. 619-633, 1973.
- [49] G.Salton, C.S.Y., "On the specification of term values in automatic indexing," *Journal of Documentation*, vol. 29, pp. 351-372, 1973.
- [50] G.Salton, J.A., A.Singhal,, "Automatic text decomposition and structuring," *Information Processing & Management*, vol. 32, pp. 127-138, 1996.
- [52] C.Apte , F.D., S.M.Weiss, "Automated learning of decision rules for text categorization", *ACM Transactions on Information Systems* 12 (3) , pp.233–251, 1994.
- [53] D.D.Lewis, "Representation and learning in information retrieval", *Ph.D. Thesis, Amherst, MA, USA*, 1992.
- [54] D.D.Lewis, "An evaluation of phrasal and clustered representations on a text categorization task", in: *SIGIR '92: Proceedings of the 15th Annual International ACM SIGIR Conference on Research and Development in Information Retrieval*, ACM, New York, NY, USA, pp. 37–50, 1992.
- [55] S.Dumais, J.P., D.Heckerman, M.Sahami, "Inductive learning algorithms and representations for text categorization", in: *CIKM '98: Proceedings of the 7th International Conference on Informa- tion and Knowledge Management*, ACM, New York, NY, USA, pp. 148–155, 1998.

- [56] S.Scott, S.M., "Feature engineering for text classification", in: *ICML '99: Proceedings of the 16th International Conference on Machine Learning*, Morgan Kaufmann Publishers, Inc., San Francisco, CA, USA, pp. 379–388, 1999.
- [57]] D.Mladenic , M.G., "Word sequences as features in text learning", in: *Proceedings of ERK-98, the 7th Electrotechnical and Computer Science Conference, Ljubljana,SL*,pp.145–148, 1998.
- [58] urnkran, J.F., " A study using n-gram features for text categorization", *Technical Report,OEFAI-TR-9830,Austrian Institute for Artificial Intelligence*, 1998.
- [59] G.Forman, "Feature selection for text classification",in:*Computational Methods of Feature Selection*",Chapman and Hall/CRC, pp. 257–276, 2007.
- [60]. M.F.Caropreso, S.M., F.Sebastiani,Alearner, "independent evaluation of the usefulness of statistical phrases for automated text categorization",pp.78–102,2001.
- [61] C.M.Tan, Y.F.W., C.D.Lee" ,*The use of bigrams to enhance text categorization*", *Information Processing and Management*38(4) 529–546, 2002.
- [62] E.Crawford, I.K., J.Patrick, "Phrases and feature selection in e-mail classification",in:*P.Bruza,A.Moffat,A.Turpin(Eds.),ADCS, University of Melbourne,Department of Computer Science*, pp. 59–62,2004.
- [63] " Y.Yang, J.O.Pedersen., "A Comparative Study on Feature Selection in Text Categorization," presented at the *Proceedings of the Fourteenth International Conference on Machine Learning*, 1997.
- [64] Selamat, A., Hidekazu Yanagimoto, and Sigeru Omatu, "Web news classification using neural networks based on PCA." In *SICE 2002. Proceedings of the 41st SICE Annual Conference*, vol. 4, pp. 2389-2394. IEEE, 2002.
- [66] W.Z. , X.Q., "Text categorization based on LDA and SVM." In *Computer Science and Software Engineering, 2008 International Conference on*, vol. 1, pp. 674-677. IEEE, 2008.
- [67] Hu, M., Shijun Wang, Anhui Wang, and Lei Wang, "Feature Extraction Based on the Independent Component Analysis for Text Classification." In *Fuzzy Systems and Knowledge Discovery, 2008. FSKD'08. Fifth International Conference on*, vol. 2, pp. 296-300. IEEE, 2008.

- [68] W.W.Cohen , Y.S., "A simple, fast, and effective rule learner," presented at the *Proceedings of the sixteenth national conference on Artificial intelligence and the eleventh Innovative applications of artificial intelligence conference innovative applications of artificial intelligence, Orlando, Florida, United States, 1999.*
- [69] D.D.Lewis, "Naive (Bayes) at forty: The independence assumption in information retrieval," *MACHINE LEARNING: ECML-98 Lecture Notes in Computer Science, 1998*
- [70] I.Dagan, Y.K., D.Roth, , "Mistake-driven learning in text categorization," *The Second Conference on Empirical Methods in Natural Language Processing, pp. 55-63, 1997*
- [71] Elkan, C" ,".Nearest neighbor classification." *University of California–San Diego, 2007.*
- [72] Jiang, C., Frans Coenen, Robert Sanderson, and Michele Zito, "Text classification using graph mining-based feature extraction." *Knowledge-Based Systems 23, no. 4 : 302-308 ,2010.*
- [73] Qiu, L.-Q., Ru-Yi Zhao, Gang Zhou, and Sheng-Wei Yi, "An extensive empirical study of feature selection for text categorization." *In Computer and Information Science, 2008. ICIS 08. Seventh IEEE/ACIS International Conference on, pp. 312-315. IEEE, 2008.*
- [74] Cheng, Y., Ruiling Zhang, Xiufeng Wang, and Qiushuang Chen, "Text Feature Extraction Based on Rough Set." *In Fuzzy Systems and Knowledge Discovery, 2008. FSKD'08. Fifth International Conference on, vol. 2, pp. 310-314. IEEE, 2008.*
- [75] Park ,C.H., "Dimension reduction using least squares regression in multi-labeled text categorization." *In Computer and Information Technology, 2008. CIT 2008. 8th IEEE International Conference on, pp. 71-76. IEEE, 2008.*
- [76] Sh.Zhou, T.W.L., J.Hu, A.Zhou, , "Fast Text Classification : A Training-Corpus Pruning Based Approach", *in Eighth International conference on Database Systems for Advanced Applications (DASFAA '03), IEEE, 2003.*
- [77] G.R.D., "Knowledge Based Neural Network for Text classification", *in IEEE International conference on Granular Computing, IEEE, pp. 542-547, 2003.*
- [78] Z.Chen, L.H.a.Y.L.M., "Incremental Learning for Text", *In International Joint conference on Neural Networks, IEEE, Orlando, USA, pp. 1-4, 2007. .*
- [79] Zhen-fang, Z., L. Pei-yu, Lu Ran, "Research of text classification technology based on genetic annealing algorithm." *In Computational Intelligence and Design, 2008. ISCID'08. International Symposium on, vol. 1, pp. 265-269. IEEE, 2008.*
- [80] Meng, Z.S., Wei Song Dan, Jinhong Li, "a new associative classifier for text categorization." *Intelligent System and Knowledge Engineering, 2008. ISKE 2008. 3rd International Conference on, 2008.*

- [81] L.Shi, J.Z., E.Liu, P.He, "Text classification Based on Nonlinear Dimensionality Reduction Techniques and Support Vector Machines", in *Third International conference on Natural Computation, IEEE*, pp. 1-4, 2007.
- [82] Zhou, X., Huaxiang Zhan, "An algorithm of text categorization based on similar rough set and fuzzy cognitive map." In *Fuzzy Systems and Knowledge Discovery, 2008. FSKD'08. Fifth International Conference on*, vol. 3, pp. 127-131. *IEEE*, 2008.
- [83] Wang, X., Zhen Hua, Rujiang Bai, "A hybrid text classification model based on rough sets and genetic algorithms." In *Software Engineering, Artificial Intelligence, Networking, and Parallel/Distributed Computing, 2008. SNPD'08. Ninth ACIS International Conference on*, pp. 971-977. *IEEE*, 2008.
- [84] Yang, Y., Seán Slattery, and Rayid Ghani, "A study of approaches to hypertext categorization." *Journal of Intelligent Information Systems* 18, no. 2-3 (2002): 219-241.
- [85] Ding, W., Songnian Yu, Qianfeng Wang, Jiaqi Yu, and Qiang Guo. , "A Novel Naive Bayesian Text Classifier." In *Information Processing (ISIP), 2008 International Symposiums on*, pp. 78-82. *IEEE* ,2008.
- [86] Zhou, S., Kan Li, and Yushu Liu, "Text Classification Based on a Novel Bayesian Hierarchical Model." In *Fuzzy Systems and Knowledge Discovery, 2008. FSKD'08. Fifth International Conference on*, vol. 2, pp. 218-221. *IEEE*, 2008.
- [87] Chen, Q., Dequan Zheng, Tiejun Zhao, and Sheng Li, "A Fusion of Multiple Classifiers Approach Based on Reliability function for Text Categorization." In *Fuzzy Systems and Knowledge Discovery, 2008. FSKD'08. Fifth International Conference on*, vol. 2, pp. 338-342. *IEEE* ,2008.
- [88] Yuan, P., Yuqin Chen, Hai Jin, and Li Huang, "MSVM-kNN: Combining SVM and k-NN for Multi-class Text Classification." In *Semantic Computing and Systems, 2008. WSCS'08. IEEE International Workshop on*, pp. 133-140. *IEEE*, 2008.
- [89] Wang, L.L.Z.L.A ,.and Haijun He, "the research of decision support vector machine in web information classification." *Computer Supported Cooperative Work* , 2008
- [90] Chen, J., Houkuan Huang, Shengfeng Tian, and Youli Qu, "Feature selection for text classification with Naïve Bayes." *Expert Systems with Applications* 36, no. 3: 5432-5435, 2009.
- [91] Wang, P., Jian Hu, Hua-Jun Zeng, and Zheng Chen, "Using Wikipedia knowledge to improve text classification." *Knowledge and Information Systems*19, no. 3: 265-281, 2009.
- [92] Sun, A ,.Ee-Peng Lim, and Ying Liu, "On strategies for imbalanced text classification using SVM: A comparative study." *Decision Support Systems* 48, no. 1: 191-201, 2009.

- [93] Sriram, B., Dave Fuhry, Engin Demir, Hakan Ferhatosmanoglu, and Murat Demirbas, . "Short text classification in twitter to improve information filtering." *In Proceedings of the 33rd international ACM SIGIR conference on Research and development in information retrieval*, pp. 841-842. ACM, 2010.
- [94] Jiang, C., Frans Coenen, Robert Sanderson, and Michele Zito, "Text classification using graph mining-based feature extraction." *Knowledge-Based Systems 23*, no. 4, PP. 302-308, 2010.
- [95] Lu, S.-H., Ding-An Chiang, Huan-Chao Keh, and Hui-Hua Huang, "Chinese text classification by the Naïve Bayes Classifier and the associative classifier with multiple confidence threshold values." *Knowledge-based systems 23*, no. 6 , PP. 598-604, .2010.
- [96] Li, X.-L., Bing Liu, and See-Kiong Ng, "Negative training data can be harmful to text classification." *In Proceedings of the 2010 conference on empirical methods in natural language processing*, pp. 218-228. Association for Computational Linguistics, 2010.
- [97] Zhang, W., Taketoshi Yoshida, and Xijin Tang, "A comparative study of TF* IDF, LSI and multi-words for text classification." *Expert Systems with Applications 38*, no. 3,pp 2758-2765, 2011.
- [98] Jiang, J.-Y., Ren-Jia Liou, and Shie-Jue Lee, "A fuzzy self-constructing feature clustering algorithm for text classification." *Knowledge and Data Engineering, IEEE Transactions on 23*, no. 3 , pp.: 335-349, 2011.
- [99] Uğuz, H., "A two-stage feature selection method for text categorization by using information gain, principal component analysis and genetic algorithm." *Knowledge-Based Systems 24*, no. 7, pp. 1024-1032, 2011.
- [100] Wan, C.H., Lam Hong Lee, Rajprasad Rajkumar, and Dino Isa, "A hybrid text classification approach with low dependency on parameter by integrating K-nearest neighbor and support vector machine." *Expert Systems with Applications 39*, no. 15 , pp. 11880-11, 2012.
- [101] Uysal, A.K., and Serkan Gunal, "A novel probabilistic feature selection method for text classification." *Knowledge-Based Systems 36*, pp. 226-235, 2012.
- [102] Yang, J., Yuanning Liu, Xiaodong Zhu, Zhen Liu, and Xiaoxu Zhang, "A new feature selection based on comprehensive measurement both in inter-category and intra-category for text categorization." *Information Processing & Management 48*, no. 4 , pp. 741-754, 2012.
- [103] Jiang, S., Guansong Pang, Meiling Wu, and Limin Kuang, "An improved K-nearest-neighbor algorithm for text categorization." *Expert Systems with Applications 39*, no. 1 , pp. 1503-1509, 2012.

- [104] Helvetic, D.S., T. Ahmed, and Sura Mahmood Abdullah., "Hybrid Intelligent Techniques for Text Categorization." *International Journal of Advanced Computer Science and Information Technology (IJACSIT) Vol 2* , pp. 23-40, 2014.
- [105] Selamat, A., Hidekazu Yanagimoto, and Sigeru Omatu, "Web news classification using neural networks based on PCA." *In SICE 2002. Proceedings of the 41st SICE Annual Conference, vol. 4*, pp. 2389-2394. IEEE, 2002.
- [106] Billsus, D., and Michael J. Pazzani, " A hybrid user model for news story classification". Springer Vienna, 1999.
- [107] [Xiao Luo Zincir-Heywood, A.N., "combining word based and word co-occurrence based sequence analysis for text categorization." *Machine Learning 2004*.
- [108] Pekar, V., Michael Krkoska, and Steffen Staab, "Feature weighting for co-occurrence-based classification of words." *In Proceedings of the 20th international conference on Computational Linguistics*, p. 799. Association for Computational Linguistics, 2004.
- [115] Bina, B., Mohamad Hasan Ahmadi, and Maseud Rahgozar, "Farsi Text Classification Using N-Grams and Knn Algorithm A Comparative Study." *InDMIN*, pp. 385-390, 2008.
- [124] A. AleAhmad, H.A., E. Darrudi, M. Rahgozar, and F. Oroumchian, "Hamshahri: A standard Persian text collection", [Online]. Available: <http://ece.ut.ac.ir/dbrq/hamshahri/index.html>.
- [125] "<http://trec.nist.gov>."
- [126] Y.Yang, "An Evaluation of Statistical Approaches to Text Categorization," *Inf. Retr.*, vol. 1, pp. 69-90, 1999.

Abstract

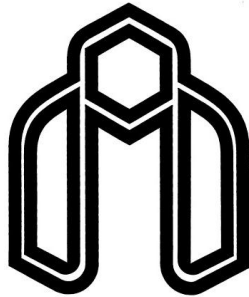
Due to increasing growth of electronic texts, such as news texts, the problem of developing an efficient classifier becomes an important issue in many text-related applications such as news websites. The most important issues in the domain of text classification are efficient feature selection and feature extraction. In previous works, different feature validation criteria have been proposed and used.

In this thesis, a feature validation criterion called E-Dominance is proposed. Using this criterion causes a significant reduction of number of selected features. Features used in this categorization, are called co-occurrence features which have not been used in classification of Persian texts yet.

In English text classification, researches which use co-occurrence features, usually use binary weighing method. In this thesis is proposed a weighing method for co-occurrence features which is called Co-occur TFIDF. IN this study, developer features are used to develop texts, Therefore the problem of class overlapping is solved partially.

Experiments show a significant improvement in the efficiency and accuracy of classification algorithms using E-Dominance criterion and co-occurrence features using Co-occur TFIDF weighting method.

Keywords: Persian text classification, feature selection, feature extraction, E-Dominance criterion, co-occurrence features, developer features, Co-occur TFIDF weighting method.



Shahrood University of Technology
Faculty of Computer Engineering

Persian News Classification Using Artificial Intelligence

Thesis Submitted in Partial Fulfillment of the Requirement for the Degree of
Master of Science (M.Sc.)

Zahra Robati

Supervisor
Dr. Morteza Zahedi

Associate Supervisor
Marzieh Rahimi

Date: September 2014