

ارائه یک روش دسته بندی متون با تکنیک کاهش ویژگی فیلتری و یادگیری ماشین

پیمان جلالی^{۱*}، سید محسن هاشمی^۲، عارف سیاحی^۳، عادل حیدری^۴

-۱- آموزشکده فنی و حرفه ای سما، دانشگاه آزاد اسلامی، واحد سوسنگرد، سوسنگرد، ایران

-۲- آموزشکده فنی و حرفه ای سما، دانشگاه آزاد اسلامی، واحد سوسنگرد، سوسنگرد، ایران

-۳- آموزشکده فنی و حرفه ای سما، دانشگاه آزاد اسلامی، واحد سوسنگرد، سوسنگرد، ایران

-۴- آموزشکده فنی و حرفه ای سما، دانشگاه آزاد اسلامی، واحد سوسنگرد، سوسنگرد، ایران

خلاصه

پایگاه داده های متونی با توجه به افزایش میزان اطلاعات موجود در فرم الکترونیکی، مانند نشریات الکترونیکی، انواع مختلف مدارک الکترونیکی، پست الکترونیکی و شبکه جهانی وب به سرعت در حال رشد می باشند. یکی از مهم ترین تکنیک های متون کاوی، دسته بندی متون است. متون کاوی را می توان به عنوان متدها و الگوریتم هایی از فیلدهای یادگیری ماشین و آماری برای متون ها با هدف پیدا کردن الگوهای مفید در نظر گرفت. در این تحقیق الگوریتم انتخاب ویژگی فیلتری بهره اطلاعات مورد استفاده قرار گرفته است. این مقاله روشی برای دسته بندی متون پیشنهاد شده است که در آن ابتدا یک دسته بند ماشین بردار پشتیبان ، بیز ساده و درخت تصمیم با تعداد کمی نمونه برچسب دار ساخته می شود، سپس با استفاده از یادگیری فعال و بکارگیری روش نمونه گیری بر اساس عدم اطمینان به همراه ایده جدید مشابهت و انتخاب گروهی نمونه ها ، به صورت هدفمند نمونه های مفید را برای برچسب گذاری به کاربر می دهد تا در آموزش دسته بند از آنها استفاده کند.

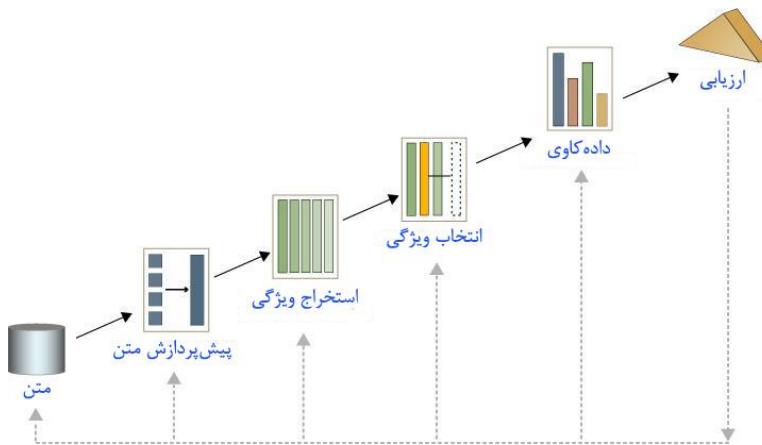
کلمات کلیدی: متون کاوی، داده کاوی، درخت تصمیم، بیز ساده، یادگیری ماشین

۱. مقدمه:

امروزه با توجه به رشد روزافزون دسترسی به استناد الکترونیکی، دسته بندی خودکار اهمیت ویژه ای یافته است. دسته بندی متون به عمل برچسب گذاری موضوعی متون زبان طبیعی بر مبنای یک مجموعه از پیش تعیین شده می باشد. در طول سالیان اخیر، طبقه بندی های مختلف و با نگرش های مختلفی برای این کار مطرح شده است. نکته قابل توجه، امکان خودکارسازی این طبقه بندی ها به ازای متون جدید است [۹]. با افزایش روزافزون حجم اطلاعات، وجود سیستمی برای دسته بندی خودکار متون ضروری به نظر می رسد. از این رو در این پژوهش سعی خواهد شد با استفاده از شیوه کار متون کاوی موجود و نیز تکنیک های یادگیری ماشین، یک مدل با دقت و کار آیی مناسب برای دسته بندی خودکار متون فارسی ارائه شود تا بتوان در مواردی همچون دسته بندی متون در بایگانی استناد، دسته بندی صفحات وب، یادگیری

¹ p.jalali05@yahoo.com

خودکار علایق مطالعاتی و پژوهشی کاربران، فیلتر کردن خودکار پست الکترونیکی بر اساس محتوا و غیره از آن استفاده کرد. شکل (۱) فرایند کلی استخراج دانش از متن را نشان می‌دهد [۱۱].



شکل (۱): فرایند استخراج دانش از متن [۱]

۲. مسئله دسته‌بندی متون

دسته‌بندی متن شامل اسناد به یکی از چند گروه از پیش تعیین شده است. برای نایل شدن به این هدف اسناد ورودی توسط یک مجموعه از مشخصات که معمولاً خصوصیات نامیده می‌شود توصیف می‌شوند. برخلاف خوشبندی که شامل آموزش بدون ناظارت است، در دسته‌بندی یک مجموعه آموزشی از داده‌ها با برچسب‌گذاری قبلی نیاز است. هدف دسته‌بندی تحلیل داده‌های ورودی و ایجاد یک مدل دقیق برای هر دسته با استفاده از این خصوصیات است. اسناد جدید در داخل یکی از این دسته‌ها دسته‌بندی می‌شوند. در مسئله دسته‌بندی متن، مشخصات کلماتی هستند که درون اسناد متنی قرار دارند. در بسیاری از موارد قبل از یادگیری ماشین انتخاب خصوصیت‌صورت می‌گیرد تا فضای خصوصیات را کاهش دهد [۹]. در دسته بندی هدف بدست آوردن یکتابع نگاشت بین اسناد و مجموعه کلاس‌ها یا گروه‌ها به وسیله یک سری سند که به آنها مجموعه آموزشی [۱۰-۱۱] می‌گویند، می‌باشد. که از این تابع نگاشت برای تعیین خودکار گروه سند جدید استفاده می‌شود. توجه شود که در هنگام تولید تابع نگاشت مجموعه اسناد برچسب خورده هستند یعنی گروه هر یک از اسناد مشخص شده است و به وسیله این اسناد برای هر گروه یک سری ویژگی و صفات منحصر به فرد استخراج می‌کنیم.

۳. کاربردهای دسته‌بندی متون

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

۴. مروری بر کارهای پیشین

در مرجع [۳] یک الگوریتم انتخاب ویژگی ناهموار مؤثر بادید دانه‌بندی برای مجموعه داده‌ها در مقیاس بزرگ ارائه دادند که در آن یک مجموعه داده در مقیاس بزرگ داده شده است، این الگوریتم ابتدا دانه‌های کوچک مختلف را انتخاب می‌کند و سپس بر روی هر دانه کوچک، میزان کاهش را از مجموعه داده اصلی برآورد می‌کند.

در مرجع [۴] یک نمونه جدید به طور همزمان تکاملی و الگوریتم انتخاب ویژگی ارائه داده‌اند که برای میلیون‌ها مورد از هزاران ویژگی مقیاس‌پذیر است. این پیشنهاد بر اساس اصل تقسیم و غلبه همراه با ساماندهی است، اصل تقسیم و غلبه موجب اجرای الگوریتم در زمان خطی می‌شود. همچنین این روش در محیط‌های موازی به راحتی قابل اجرا است و می‌تواند بدون بارگذاری کل مجموعه داده در حافظه کار کند.

در مرجع [۵] یکتابع معیار کلی در مورد اطلاعات متقابل در انتخاب ویژگی ارائه شده است، که می‌تواند بسیاری از اندازه‌گیری‌های اطلاعاتی در الگوریتم‌های پیشین را به یکدیگر نزدیک کند. در انتخاب‌گرهای سنتی، اطلاعات متقابل در کل فضای نمونه برآورد می‌شوند که نمی‌تواند دقیقاً نشان‌دهنده ارتباط میان ویژگی‌ها باشد. برای مقابله با این مشکل، هدف دوم این مقاله پیشنهاد یک الگوریتم انتخاب ویژگی جدید بر اساس اطلاعات متقابل پویا است، که تنها در موارد بدون برچسب تخمین زده شده است.

در مرجع [۶] یک روش نیمه نظارت بیزی برای انتخاب ویژگی طبقه‌بندی شده را ارائه داده‌اند. با توجه به اینکه در برنامه‌های کاربردی دنیای واقعی، دریافت کردن نمونه‌های بدون برچسب معمولاً آسان است، اما به دست آوردن برچسب‌های دقیق، مربوط به نمونه‌ها گران است، این امر منجر به زباله‌های بالقوه از اطلاعات طبقه‌بندی شده ارزشمند، در نمونه‌های بدون برچسب می‌شود که دور ریخته می‌شوند. درنتیجه این روش نمونه‌های فاقد برچسب را در مشکل انتخاب ویژگی‌های طبقه‌بندی شده، حل می‌کند.

در مرجع [۷] یک روش انتخاب ویژگی دومرحله‌ای برای طبقه‌بندی متن را ارائه دادند، که مشکل فضای موردنیاز برای استناد با ابعاد بالا را برای نشان دادن اسناد، حل می‌کند و کارش به این صورت است که برای بالا رفتن دقت و کارایی طبقه‌بندی از انتخاب ویژگی دومرحله‌ای استفاده می‌کند که در مرحله اول، یک روش انتخاب ویژگی جدید برای کاهش اصطلاحات کم‌اهمیت اعمال می‌کند و در مرحله دوم یک فضای معنایی جدید، بین واژگان، بر اساس روش نمایه‌سازی، معنایی نهان ایجاد می‌کند.

۳. روش تحقیق:

برای ارائه یک مدل برای دسته‌بندی خودکار متون لاتین ابتدا روش‌های پیش پردازش انجام گردیده است و ار روش انتخاب ویژگی بهره اطلاعات که در زیر توضیح داده خواهد شد، استفاده می‌شود. با تجزیه و تحلیل دقیق داده‌ها با ارزش، روش جدیدی ارائه می‌شود که بتواند نقاط ضعف روش‌های پیشین را پوشش دهد.

۳.۱. مجموعه داده:

استناد ورودی یا مجموعه داده انتخابی، مجموعه داده اخبار رویترز- ۲۱۵۷۸ می‌باشد [۱۲]. این مجموعه داده در سال ۱۹۸۷ گردآوری و توسط گروهی از کارکنان خبرگزاری رویترز تهیه و شاخص گذاری شد. با توجه به اینکه در این تحقیق، روش تک برچسبی - چند دسته مورد توجه است، از زیرمجموعه (8) R استفاده می‌شود. این زیرمجموعه شامل ۸ دسته اصلی و دارای ۷۶۷۴ سند متنی می‌باشد. هر سند متنی تنها به یکی از دسته‌های موجود تعلق دارد. با توجه به روش جداسازی استناد آموزشی از استناد آزمایشی ModApte Split، مجموعه استناد آموزشی شامل ۵۴۸۵ سند و مجموعه استناد آزمایشی شامل ۲۱۸۹ سند می‌باشد. جزئیات زیرمجموعه (8) R اخبار رویترز در جدول (۱) نشان داده شده است.

جدول (۱): جزئیات زیرمجموعه اخبار رویترز

تعداد کل استناد دسته	تعداد استناد آزمایشی	تعداد استناد آموزشی	دسته
2292	696	1596	Acq
374	121	253	Crude
3923	1083	2840	Earn
51	10	41	Grain
271	81	190	Interest
293	87	206	Money-fx
144	36	108	Ship
326	75	251	Trade

۳.۲. پیش پردازش:

برای دسته بندی متون پس از فراهم کردن مجموعه داده، نوبت به مرحله پیش پردازش می رسد. در این مرحله استناد متنی خام باید به فرمی که قابل استفاده در مرحله انتخاب ویژگی و الگوریتم یادگیری باشند، تبدیل شوند. در مرحله پیش پردازش روش پیشنهادی، عملیات Transform Case، Tokenization، Stemming، Filter StopWords و Stemming انجام می شود. عملکرد هر کدام در زیر بیان شده است.

Transform Case: در این مرحله تمامی کاراکترهای موجود در متن به فرم یکسان تبدیل می شوند. در این مرحله تمامی کاراکترها به حروف کوچک تبدیل می شوند.

Tokenization: در این مرحله کل متن به صورت کلمات متوالی جدا از هم تقسیم می شود.

Filter StopWords: در این مرحله کلمات زائد در زبان انگلیسی حذف خواهند شد.

Stemming: در این مرحله ریشه کلمات در زبان انگلیسی با استفاده از الگوریتم ریشه یابی پورتر برای حذف پسوندها و پیشوندهای کلمات در جهت کاهش طول کلمات و تا زمانی که به فرم ریشه شان تبدیل شوند، استفاده شده است.

Generate n-gram: در این مرحله برای شاخص بندی و کاهش ابعاد متن از n-gram استفاده شده است. با استفاده از n-gram می توان متن را به صورت یک سری از کلمات متوالی به طول n نشان داد. این مدل ابتدا برای مسائل پردازش گفتار مطرح شد. ولی اکنون نکارش های مختلفی از این مدل برای مسائل پردازش زبان طبیعی و دسته بندی متون مطرح شده است [۱۶ و ۱۷]. با توجه به آزمایش های انجام شده بر روی مقادیر مختلف n و برای جلوگیری از افزایش پیچیدگی، از n-gram با میزان n=2 استفاده شد. پس از مراحل فوق، وزن دهی ویژگی ها انجام می شود. در این تحقیق از روش وزن دهی TF-IDF استفاده شده است که عملکرد دسته بندها با استفاده از این روش وزن دهی مورد ارزیابی قرار می گیرند.

۳.۳. انتخاب ویژگی با روش بهره اطلاعات (IG):

این روش یکی از پر کاربرد ترین و محبوب ترین روش ها می باشد که در بسیاری از پژوهش ها مورد استفاده قرار گرفته و نتایج مطلوبی داشته است. در سال های اخیر، سودمندی اطلاعات به عنوان معیاری برای میزان سودمندی واژه در زمینه یادگیری ماشین مطرح شده است. این معیار تعدادی قلم های اطلاعاتی را برای پیشگویی دسته، با کار بر روی وجود و عدم وجود یک حمله در یک دیتا است به دست می آورد. به عبارتی در این روش میزان سودمندی یک حمله در دیتا است تعداد حملاتی است که برای پیشگویی دسته، با توجه به وجود یا عدم وجود حمله در دیتا است به دست می آید. رابطه (۱) مربوط به این روش می باشد [۱۳].

$$IG(t_k, c_i) = \sum_{c \in \{c_i, \bar{c}_i\}} \sum_{t \in \{t_k, \bar{t}_k\}} p(\bar{t}, c) \log_2 \frac{P(t, c)}{P(\bar{t})P(c)} \quad (1)$$

با داشتن یک مجموعه آموزشی، می توان بهره اطلاعاتی را برای هر واژه محاسبه نمود. حملاتی که بهره اطلاعاتی آنها از یک حد آستانه کمتر است از فضای ویرگی حذف می شوند. طبق تعریف محاسبه شامل تخمینی از احتمال های شرطی از دسته تعیین کننده واژه و محاسبات آنتروپی است. در مرحله یادگیری از الگوریتم های یادگیری ماشین زیر استفاده می نماییم.

۳.۴. الگوریتم Naïve Bayes

یکی دیگر از الگوریتم های مطرح شده در زمینه دسته بندی [۱]، بیز ساده است. این روش به دلایل متعددی اهمیت دارد. اینکه ساخت آن بسیار ساده است و نیازی به برنامه های تخمین پارامتر تکرار شونده پیچیده ندارد. بدین معنی که می توان از آن برای مجموعه داده های بسیار وسیع نیز استفاده نمود. این الگوریتم یکی از سریع ترین الگوریتم ها در زمینه دسته بندی می باشد [۸].

۳.۵. درخت تصمیم

درخت های تصمیم گیری ابزار استانداردی در داده کاوی هستند. این الگوریتم ها هم در متغیرها و هم سایز مجموعه آموزش سریع و همچنین مقیاس پذیر^۱ هستند. یکی از مشکلات درخت های تصمیم گیری برای متن کاوی این است که تنها به تعداد کمی از ترم ها وابسته است [۱۴]. دسته کننده درخت تصمیم گیری یکی از پر استفاده ترین روش های یادگیری با نظارت است که برای کاوش [۱۱] داده ها مورد استفاده قرار می گیرد. این دسته کننده برای تفسیر آسان است و می تواند بصورت قوانین IF-THEN-ELSE بازنمایی گردد. در این دسته کننده تابعی توسط مناطق ثابت تکه های تقریب زده می شود و به هیچ دانش قبلی از توزیع داده ها نیاز ندارد. این دسته کننده بر روی داده های پارازیت دار به خوبی کار می کند.

۳.۶. ماشین بردار پشتیبان (SVM)

یکی از روش های یادگیری با نظارت است که از آن برای طبقه بندی و رگرسیون استفاده می کنند. این روش از جمله روش های نسبتاً جدیدی است که در سال های اخیر کارایی خوبی نسبت به روش های قدیمی تر برای طبقه بندی از جمله شبکه های عصبی پرسپترون نشان داده است. مبنای کاری دسته بندی کننده SVM دسته بندی خطی داده ها است و در تقسیم خطی داده ها سعی می کنیم خطی را انتخاب کنیم که حاشیه اطمینان بیشتری داشته باشد. حل معادله پیدا کردن خط بهینه برای داده ها به وسیله روش های QP که روش های شناخته شده ای در حل مسائل محدودیت دار هستند صورت می گیرد [۱۵].

۴. نتایج:

۴.۱. ارزیابی روش بر روی دسته بند SVM

در این مرحله، جهت ارزیابی روش پیشنهادی در مرحله یادگیری از دسته بند SVM استفاده شده است. جدول ۲ نتایج ارزیابی روش پیشنهادی را با استفاده از دسته بند SVM نشان می دهد. نتایج به دست آمده نشان دهنده کارایی بالای روش پیشنهادی با استفاده از دسته بند ماشین بردار پشتیبان در میانگین دقت برابر ۹۹.۲۵٪، میانگین فراخوانی برابر ۹۹.۱۷٪، میانگین F1 برابر ۹۹.۲۰٪، معیار درستی برابر ۹۹.۶۳٪ و خطای دسته بندی با میزان ۰.۳۷٪ می باشد.

^۱ Scalable

جدول (۲): نتایج ارزیابی روش با استفاده از دسته بند SVM

دسته	Precision	Recall	F1
Acq	99.58	99.73	99.65
Trade	99.88	100	99.94
Ship	100	99.60	99.80
Interest	97.35	98.07	97.70
Grain	100	98.99	99.49
Crude	99.78	99.83	99.80
Earn	99.92	99.81	99.86
Money-fx	97.45	97.31	99.38
Average precision: 99.25			
Average recall: 99.17			
Average F1: 99.20			
Accuracy: 99.63			
Classification error: 0.37			

۴.۲. ارزیابی روش بر روی دسته بند Naive Bayes

در این مرحله، جهت ارزیابی روش پیشنهادی در مرحله یادگیری از دسته بند Naive Bayes استفاده شده است. جدول (۳) نتایج ارزیابی روش پیشنهادی را با استفاده از دسته بند Naive Bayes نشان می دهد. نتایج به دست آمده روش پیشنهادی با استفاده از دسته بند Naive Bayes در میانگین دقت برابر ۷۶.۶۰٪، میانگین فراخوانی برابر ۸۹.۱۰٪، میانگین F1 برابر ۸۲.۳۷٪، معیار درستی برابر ۹۰.۳۷٪ و خطای دسته بندی با میزان ۹.۶۳٪ می باشد که نسبت به ماشین بردار پشتیبان در کلیه معیارها دارای عملکرد ضعیفتری می باشد.

جدول (۳): نتایج ارزیابی روش با استفاده از دسته بند Naive Bayes

دسته	Precision	Recall	F1
Acq	87.23	91.25	89.194
Trade	92.15	85.20	85.54
Ship	47.17	86.32	61
Interest	80.74	90.88	85.51
Grain	43.04	99.25	60
Crude	80.52	83.84	82.14
Earn	99.70	91.32	95.32
Money-fx	81.31	84.70	82.97
Average precision: 76.60			
Average recall: 89.10			
Average F1: 82.37			
Accuracy: 90.37			
Classification error: 9.63			

۴.۳. ارزیابی روش بر روی دسته بند J48

در این مرحله، جهت ارزیابی روش پیشنهادی در مرحله یادگیری از دسته بند J48 استفاده شده است. جدول (۴) نتایج ارزیابی روش پیشنهادی را با استفاده از دسته بند J48 نشان می دهد. نتایج به دست آمده روش پیشنهادی با استفاده از دسته بند J48 در میانگین دقت برابر ۸۳.۶۵٪، میانگین فراخوانی برابر ۷۹.۹۶٪، میانگین F1 برابر ۸۱.۷۶٪، معیار درستی برابر ۹۲.۵۸٪ و خطای دسته بندی با میزان ۷.۴۲٪ می باشد که نسبت به ماشین بردار پشتیبان در کلیه معیارها دارای عملکرد ضعیفتری می باشد و لی در معیارهای میانگین دقت، درستی و خطای دسته بندی دارای عملکرد بهتری نسبت به دسته بند Naive Bayes می باشد.

جدول (۴): نتایج ارزیابی روش با استفاده از دسته بند J48

دسته	Precision	Recall	F1
Acq	89.97	92.96	91.44
Trade	90.08	90.80	90.43
Ship	55.05	55.56	55.30
Interest	88.48	90.57	89.51
Grain	61.11	46.81	53.01
Crude	92.36	90.07	91.20
Earn	97.10	96.28	96.68
Money-fx	82.51	75.50	78.84
Average precision: 83.65			
Average recall: 79.96			
Average F1: 81.76			
Accuracy: 92.58			
Classification error: 7.42			

۵. نتیجه گیری:

با توجه به اهمیت دسته بندی متون، هدف در این پژوهش طراحی و پیاده سازی یک مدل کارآمد جهت دسته بندی خودکار برای متون لاتین با استفاده از تکنیک های داده کاوی و یادگیری ماشین بوده است که با استفاده از آن بتوان متون را با کارایی بالا دسته بندی نمود. در این قسمت، نتایج به دست آمده از مدل پیشنهادی با استفاده از الگوریتم SVM که دارای بالاترین کارایی نسبت به الگوریتم های بیز ساده و درخت تصمیم بر روی روش پیشنهادی است، مورد مقایسه قرار گرفت. نتایج به دست آمده نشان دهنده برتری روش پیشنهادی با استفاده از الگوریتم یادگیری SVM در میانگین دقت با میزان ۹۹.۲۵٪، میانگین فراخوانی با میزان ۹۹.۱۷٪، میانگین F1 با میزان ۹۹.۲۰٪، درستی با میزان ۹۹.۶۳٪ و خطای دسته بندی با میزان ۳٪.۰ در مقایسه با دسته بندی با استفاده از الگوریتم های منفرد می باشد.

مراجع:

- Shahbaz, M., Ahsen, S. M., Shaheen, M., Shaheen, M., Masood, S. A., 2011. An Effective Preprocessing Methodology for Textual Data Classification, Journal of American Science, 7 (6), pp. 944-951.
- Ruta, D., Gabrys, B., 2000. An Overview of classifier fusion Methods, University of Paisley, Computing & Information Systems, 7 (1), pp. 1-10.
- Jiye Liang , FengWanga, Chuangyin Dangb, Yuhua Qian , 2012. An efficient rough feature selection algorithm with a multi-granulation view , In: International Journal of Approximate Reasoning 53 pp:912–926 .
- Nicolas Garcia-Pedrajas, Aida de Haro-Garcia, Javier Perez-Rodriguez , 2013. A scalable approach to simultaneous evolutionary instance and feature selection, Information Sciences 228 pp:150–174 .
- Huawen Liu, JiguiSun, LeiLiu, HuijieZhang , 2009. Feature selection with dynamic mutual information , Pattern Recognition 42 pp: 1330 – 1339 .
- Ruichu Cai , ZhenjieZhang , ZhifengHao , 2011. BASSUM: A Bayesian semi-supervised method for classification feature selection , Pattern Recognition 44 pp: 811–820 .
- Jiana Meng, Hongfei Lin , Yuhai Yu , 2011. A two-stage feature selection method for text categorization , Computers and Mathematics with Applications 62 pp: 2793–2800 .
- Ganiz, M. C., George, C., & Pottenger, W. M. (2011). Higher order Naive Bayes: A novel non-IID approach to text classification. Knowledge and Data Engineering, IEEE Transactions on, 23(7), 1022-1034.
- Aghdam, M. H., Aghaei, N. G., Basiri, M. E., 2009. Text feature selection using ant colony optimization, Expert Systems with Applications, 36 (3), pp. 6843-6853.

- 10) Al-Mubaid, H., Umair, S. A., 2006. A New Text Categorization Technique Using Distributional Clustering and Learning Logic, IEEE Transactions on Knowledge and Data Engineering, 18 (9), pp. 1-10.
- 11) Arturo, M. R., 2006. Automatic Text Categorization of documents in the High Energy Physics domain, Ph. D. Thesis, University of Granada.
- 12) Dataset, available in: <http://archive.ics.uci.edu/ml/datasets/Reuters-21578+Text+Categorization+Collection>, Retrieved on 28 July 2013.
- 13) Jing, L. P., Huang, H. K., & Shi, H. B. (2002). Improved feature selection approach TFIDF in text mining. In Machine Learning and Cybernetics, 2002. Proceedings. 2002 International Conference on (Vol. 2, pp. 944-946). IEEE.
- 14) Bi, Y., Bell, D., Wang, H., Guo, G., Dubitzky, W., 2004-a. Classification Decision Combination for Text Categorization: An Experimental Study, Database and Expert Systems Applications, 3180, pp. 222-231.
- 15) Kecman, V. (2001). Learning and soft computing: support vector machines, neural networks, and fuzzy logic models. MIT press.
- 16) Peng, F., Schuurmans, D., Wang, S., 2003. Language and Task Idependent Text Categorization with Simple Language Models, In Proceedings of the HLT-NAACL.
- 17) Wei, Z., Miao, D., Hugues, J., Zhao, R., Li, W., 2009. N-grams based feature selection and text representation for Chinese Text Classification, International Journal of Computational Intelligence Systems, 2 (4), pp. 365-374.