

بهبود دسته بندی متون فارسی با ترکیب روش دو مرحله ای انتخاب ویژگی و الگوریتم های یادگیری ماشین

سید محسن هاشمی^{۱,*}

۱- گروه کامپیوتر، موسسه آموزش عالی رهنما، اهواز، ایران

خلاصه

امروزه با افزایش حجم داه ها و استناد متنی نیاز به دسته بندی استناد احساس می شود. از جمله کاربردهای متن کاوی می توان به دسته بندی متون، شاخص گذاری متون بر اساس یک لغتنامه، فیلتر کردن متون، تولید خودکار فرا داده، تولید کاتالوگ های سلسله مراتبی از منابع وبی و دیگر کاربرد های مشابه این زمینه نام برد. در این پژوهش به دسته بندی متون فارسی با ۲۳ دسته می پردازیم. پژوهش های پیشین اغلب بر تعداد کلاس های کمتر تمرکز داشته اند و لذا با افزایش تعداد کلاس ها دقت و کارایی به نسبت بسیار زیادی کاهش خواهد یافت. لذا به ضرورت موضوع در این مقاله، یک روش ترکیبی دو مرحله ای در قسمت انتخاب ویژگی با الگوریتم های CFS و IG پیشنهاد نموده ایم. بنابراین در مرحله یادگیری ابتدا با الگوریتم های یادگیری ماشین بصورت منفرد آزمایشاتی انجام شده است و سپس با الگوریتم های ترکیبی یادگیری ماشین آزمایش نهایی را بررسی نمودیم. نتایج نشان دهنده برتری ایده و روش پیشنهادی با درصد کارایی ۹۵.۴۹٪ باشد.

کلمات کلیدی: دسته بندی، انتخاب ویژگی، دسته بندی متون، یادگیری ماشین، متن کاوی

-۱ مقدمه

با افزایش اطلاعات، میزان پایگاه داده های متنی در شبکه جهانی وب به سرعت در حال رشد می باشد. این اطلاعات در فرم الکترونیکی، مانند نشریات الکترونیکی، انواع مختلف مدارک الکترونیکی، پست الکترونیکی می باشند. سیستم های زیادی برای پردازش اطلاعات وجود دارد که بطور عمده می توان از سیستم های اطلاعات مدیریتی، سیستم های مدیریت پایگاه داده ها، سیستم های تصمیم یار، سیستم های پرسش و پاسخ و سیستم های بازیابی اطلاعات نام برد. کاربران نیاز به ابزارهایی برای مقایسه استناد مختلف، رتبه بندی اهمیت و ارتباط استناد و یا پیدا کردن الگوها و روندها در استناد متعدد هستند [۵]. بنابراین، متن کاوی به موضوعی محبوب و ضروری در داده کاوی تبدیل شده است.

* Email: Sayedmohsenhashemi@yahoo.com

۲- ضرورت دسته بندی

در سالهای اخیر به علت افزایش سریع اطلاعات و در دسترس قرار گرفتن این اطلاعات به اشكال مختلف دیجیتالی (متن، فایلهاي صوتی و تصویری و ...) نیازی مبرم به نوعی مدیریت جهت دسته بندی، پردازش و ارائه حس می گردد. دسته بندی متون به علت اهمیت بالای متن در ثبت اطلاعات این نیاز را دوچندان می کند. عمل برچسب گذاری موضوعی متون بر اساس یک مجموعه مشخص یکی از این موارد می باشد [۲]. دسته بندی متون به صورت دستی برای زمینه های تخصصی نیاز به دانش افراد خبره می باشد و از آنجا که برچسب گذاری دستی براساس تجربه و دانش فرد می باشد بسیار خطأ پذیر است به همین علت تصمیم دو فرد خبره در برچسب گذاری می تواند متفاوت و حتی ناسازگار باشد.

۳- مروری بر کارهای پیشین

تا کنون اغلب کارهای انجام شده بر روی اسنادی به زبان انگلیسی بوده و اخیرا در مورد زبانهایی مثل چینی، عربی و ... فعالیت هایی صورت گرفته است. روش های معمول در زمینه دسته بندی متون، روش های یادگیری ماشین هستند در چند سال اخیر الگوریتم های زیادی برای مسئله دسته بندی اسناد، پیشنهاد شده است. در زمینه پردازش زبان طبیعی و به خصوص پردازش متن، یکی از پایه ای ترین کارها رده بندی یا دسته بندی خودکار متون است [۲]. دسته بندی خودکار در کاهش زمان و هزینه لازم برای دسته بندی صدھا یا هزاران سندی که در یک روز وجود دارد، بدون نیاز به افراد متخصص، بسیار مفید می باشد. روش ها و الگوریتم های متفاوتی برای دسته بندی متن بکار برده شده است، که این روش ها در دقت و محاسبات متفاوت می باشند.

دسته بندی متن شامل مراحل مختلفی می باشد و پروسه دسته بندی به الگوریتم های استفاده شده بستگی دارد بر مبنای این روش، تعدادی از مستندات آموزشی طبقه بندی شده به صورت دستی، داده می شود [۶]. قوانین طبقه بندی متن باید به نحوی آموزش داده شوند که بتوانند مسئله یادگیری باظهارت را به راحتی حل کنند [۷]. در مرجع [۹] هاشمی و همکاران، پژوهشی به نام دسته بندی متون فارسی با استفاده از فیلتر کردن واژه ها بر اساس اندازه آنها انجام داده اند که در این مقاله با انجام مراحلی از پیش پردازش داده ها و استفاده از الگوریتم های یادگیری ماشین از خانواده بیز کارایی دسته بندی را با مدل اعتبار سنجی ۱۰ مرحله ای متقابل به ۸۴.۶۲ درصد با ۸ دسته رسانند.

در مرجع دیگری [۱۰] هاشمی و همکاران، پژوهشی به نام استفاده از تکنیک های متن کاوی برای دسته بندی متون فارسی با مجموعه داده همشهری انجام داده اند که با به کارگیری الگوریتم های یادگیری بیز کارایی دسته بندی را با مجموعه داده همشهری با ۷۰ درصد آموزش و ۳۰ درصد آزمایش مورد بررسی قرار داده ایم. نتایج کارایی به ۸۵.۴۸ درصد و خطای ۱۴.۵۲ درصد با ۸ دسته بوده است.

در مرجع دیگری مهدی پور و همکاران [۱۵] سیستم خلاصه ساز خودکار متن فارسی با استفاده از الگوریتم ترکیبی SA-GA را ارائه کرده اند که این سیستم پس از ریشه یابی کلمات با استفاده از ترکیب روش های مبتنی بر گراف و TF-IDF با ترکیب الگوریتم SA-GA عمل می کند. نتایج حاصل از ارزیابی متن خلاصه شده توسط سیستم ققنوس و مقایسه آن با سیستم پارسینا نشان می دهد که کیفیت متن خلاصه شده در سیستم ققنوس ۶۴.۳۵٪ بوده که بیشتر از سیستم پارسینا با ۵۹.۸۶٪ است.

۴- مجموعه داده:

یک مجموعه های متنی ابزارهای مهمی برای پیشبرد تحقیقات در تعدادی از شاخه های علوم کامپیوتر مانند بازیابی اطلاعات^۱، زبانشناسی پیکره ای^۲ و زبانشناسی محاسباتی^۳ هستند. این مجموعه آزمایش متعلق به پیکره همشهری در سال ۲۰۰۷ می باشد [۸]. مجموعه آزمایش همشهری یکی از معتبرترین منابع در زبان فارسی است. از این مجموعه در همایش های معتبر بین المللی استفاده شده است. تعداد هر کدام از استناد در جدول زیر قابل مشاهده می باشد. در این پژوهش از ۲۳ دسته که در جدول ۱ نشان داده ایم، برای فرآیند دسته بندی متون استفاده می نماییم.

جدول شماره ۱: مجموعه داده ۲۳ کلاسی همشهری

دسته ها	اسامی دسته ها	نام فارسی
۱	Economy	اقتصاد
۲	Economy. Bank and Bourse	اقتصاد.بانک و بورس
۳	Literature and Art	ادب و هنر
۴	Literature and Art. Art	ادب و هنر.هنر
۵	Literature and Art .Art. Cinema	ادب و هنر.هنر.سینما
۶	Literature and Art .Art. Music	ادب و هنر.هنر.موسیقی
۷	Literature and Art. Art. Theater	ادب و هنر. هنر.تئاتر
۸	Literature and Art. Literature	ادب و هنر.ادب
۹	Miscellaneous	گوناگون
۱۰	Miscellaneous. Happenings	گوناگون.حوادث
۱۱	Miscellaneous. Urban	گوناگون.شهری
۱۲	Miscellaneous. World News	گوناگون.خبر جهانی
۱۳	Politics	سیاسی
۱۴	Politics. Iran Politics	سیاسی. ایران سیاسی
۱۵	Science and Culture	علمی فرهنگی
۱۶	Science and Culture. Science	علمی فرهنگی.علمی
۱۷	Science and Culture. Science. Book	علمی فرهنگی.علمی.کتاب
۱۸	Science and Culture. Science. Information and Communication Technology	علمی فرهنگی.علمی.ارتباطات و فناوری اطلاعات
۱۹	Science and Culture. Science. Medicine and Remedy	علمی فرهنگی.علمی.پزشکی و درمان
۲۰	Social	اجتماعی
۲۱	Social. Religion	اجتماعی.جهان
۲۲	Sport	ورزشی
۲۳	Tourism	گردشگری

¹Information Retrieval²Corpus Linguistics³Computational Linguistics

۵- روش کار

۱-۵- مرحله پیش پردازش

در حقیقت پیش پردازش وظیفه نگاشت متن داده شده به یک نمای منطقی را بر عهده دارد. به عبارت دیگر استخراج ویژگی و وزن دهی و کاهش ابعاد در این قسمت انجام می‌گیرد. بسته به کاربرد استخراج ویژگی می‌تواند بسیار ساده و یا بسیار مفصل باشد. تحلیل واژگانی شامل عملیات مربوط به یکسان سازی متن، قواعد مربوط به شناسه گذاری ها و مرزبندی بین کلمات می‌باشد. در این مرحله داده های ورودی باید به صورتی تبدیل شوند که قابل پردازش توسط مراحل بعد باشند. در مرحله پیش پردازش داده های فارسی معمولاً بر روی داده های ورودی عملیات زیر صورت می‌گیرد:

۱. جداسازی کلمات: در این مرحله تمام متن به صورت کلمات متمایز تبدیل می‌شود و همچنین علائم نگارشی، تگها و ... حذف می‌شوند.

۲. وزن دهی کلمات (ویژگی ها): پس از انجام مراحل فوق متن به صورت برداری از کلمات که ویژگی های متن هستند می‌تواند نمایش داده شود. حال می‌تواند این ویژگی ها را با توجه به میزان اهمیت آنها نسبت به سند متنی و دسته آن وزن دهی کرد. یکی از روش های وزن دهی TF-IDF می‌باشد. این وزن دهی بر مبنای تکرار کل کلمه و معکوس تعداد متونی که این کلمه را در بر دارند، محاسبه می‌گردد . به عبارت دیگر، هرچه تعداد تکرار کلمه در متن بیشتر و تعداد متن هایی که این کلمه را در بر دارند کمتر باشد، وزن آن ویژگی بیشتر است.

۳. انتخاب ویژگی با روش بهره اطلاعات (IG): این روش یکی از پر کاربرد ترین و محبوب ترین روش ها می‌باشد 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 دیتابست به دست می آورد. به عبارتی در این روش میزان سودمندی یک حمله در دیتابست تعداد حملاتی است که برای پیشگویی دسته، با توجه به وجود یا عدم وجود حمله در دیتابست به دست می آید. رابطه (۱) مربوط به این روش می باشد [۱۴].

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

با داشتن یک مجموعه آموزشی، می‌توان بهره اطلاعاتی را برای هر واژه محاسبه نمود. حملاتی که بهره اطلاعاتی آنها از یک حد آستانه کمتر است از فضای ویژگی حذف می‌شوند. طبق تعریف محاسبه شامل تخمینی از احتمال های شرطی از دسته تعیین کننده واژه و محاسبات آنتروپی است.

۴. انتخاب ویژگی با روش CfsSubsetEval: یک روش فیلتر است که زیر مجموعه های ویژگی ها را بر اساس قابلیت پیش بینی به همراه درجه افزونگی بین آنها ارزیابی می کند. در نتیجه ویژگی هایی که همبستگی بالایی با متغیر هدف داشته و در عین حال همبستگی بین خود ویژگی ها پایین باشد ترجیح داده می شود [۱۴].

۲-۵- مرحله یادگیری

در مرحله یادگیری از الگوریتم های یادگیری ماشین زیر استفاده می نماییم.

۱-۲-۵- 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.

۲-۲-۵- الگوریتم Multinomial Naive Bayes

وظیفه کلاسه بندی متنی را می توان از یک دیدگاه یادگیری بیزین دانست که توزیع های واژه ای در اسناد از طریق مدل پارامتری خاصی بوجود آمده و پارامترها می توانند از طریق داده های یادگیری محاسبه گردند. با شناسایی (کاوش) مجموعه داده یادگیری ما می توانیم همه فراوانی های مورد نیاز را بدست آورده و سپس احتمال های شرطی متناظر را محاسبه نماییم. یکی از مزایای مدل MNB این است که پیش بینی ها را با راندمان بالا انجام می دهد. بنابراین برای انجام پیش بینی ها MNB تنها نیازمند جستجو میان واژه های دارای شمارگان صفر می باشد. بطور متداول در کنار نمایش تکرار کم واژه های متنی بکار گرفته می شود. بطور عکس دیگر کلاسه بندی شبکه ها با یستی نیازمند جستجو در میان همه واژگان در مجموعه یادگیری می باشند. از این رو که در کلاسه بندی متنی اندازه واژگان اغلب خیلی بزرگتر از طول سند می باشد. الگوریتم MNB می تواند پیش بینی ها را سریع تر از کلاسه بنده های دیگر به جز خانواده بیز انجام دهد [۳].

۳-۲-۵- الگوریتم Discriminative Multinomial Naive Bayes

بیز ساده چند جمله ای متمایز گر (DMNB) یکی از روش یادگیری ماشین است. در فرانس [۴] این الگوریتم را شرح داده شده است که از یک مدل بیز ساده چند جمله ای بهره می گیرد. هدف DMNB حفظ اطلاعات فراوانی در حالی که به ماهیت قابل تمايز کلاسه بندی توجه دارد و بنابراین ترکیبی از یادگیری عمومی و متمایز می باشد. الگوریتم DMNB از طریق اسناد یادگیری تولید می گردد. مشخص است که DMNB دارای همان خصوصیات محاسباتی MNB می باشد؛ این الگوریتم دقیقتر و نسبت به MNB سریعتر است و همچنین یادگیری آنلاین را پشتیبانی می نماید.

۴-۲-۵- الگوریتم Complement Naive Bayes

دسته بندی بیز ساده مکمل (CNB) ضعف طبقه بندی بیز ساده را با پارامترهایی از داده ها در تمام دسته های حساس بهبود می بخشد [۳].

الگوریتم ترکیبی **Bagging** ۲-۵-۵

الگوریتم Bagging تنها یک روش یادگیری دارد. با نمونه‌گیری‌های متفاوت همراه با جایگزینی از مجموعه آموزشی، چندین مجموعه آموزشی جدید با همان اندازه اولیه به دست می‌آید و بدین طریق از هر کدام می‌توان یک دسته‌بند جداگانه را آموزش داد. در عمل نشان داده شده است که Bagging نتایج بسیار خوبی به ویژه وقتی که اندازه مجموعه داده‌ها کوچک باشد به دست می‌دهد [۱۲].

الگوریتم ترکیبی **Boosting** ۲-۶-۵

در روش Boosting دسته‌بندی که به آن الگوریتم یادگیری ضعیف می‌گویند [۱۳]. چندین بار و هر بار با داده‌های آموزشی متفاوت که با توجه به اجرای قبلی انتخاب می‌شوند اجرا شده و در نهایت آن جوابی که بیشتر تکرار شده انتخاب می‌شود. اگر چه این روش وقت‌گیر است ولی جواب‌های آن مطمئن خواهد بود. در این روش به هر نمونه از مجموعه داده‌های آموزشی یک وزن انتساب می‌شود. در ابتدا وزن همه نمونه‌ها برابر است و در تکرارهای بعدی با توجه به کارایی دسته‌بند، وزن هر نمونه اصلاح می‌شود، به طوری که نمونه‌هایی که اشتباه دسته‌بندی شده‌اند وزنشان افزایش و نمونه‌هایی که درست دسته‌بندی شده‌اند وزنشان کاهش می‌یابد. AdaBoost* مشهورترین روش Boosting می‌باشد [۱۱].

۶ نتایج

جدول ۲: مقایسه کارایی و خطای دسته بندی با 1Gram و روش وزن دهنده TF-IDF

Algorithms	1Gram+ TF-IDF		
	Accuracy	Error classification	Time(s)
NB	۹۲.۹۲	۷.۰۸	۹۹.۶۹
MNB	۹۳.۲۶	۶.۷۴	۰.۵۶
DMNB	۹۳.۴۳	۶.۵۷	۱۰.۲۵
CNB	۹۳.۲۵	۶.۶۵	۰.۶۶

درجول ۲ با روش 1-gram تعداد ۸۰۰۷ ویژگی تولید شده است که با چهار الگوریتم آزمایش هایی را انجام دادیم. نتایج نشان می دهد که الگوریتم DMNB با کارایی ۹۳.۴۳٪ بالاترین کارایی و کمترین درصد خطا دسته بندی را دارد.

جدول ۳: مقایسه کارایی و خطای دسته بندی با 2Gram و روش وزن دهنده TF-IDF

Algorithms	2Gram+ TF-IDF		
	Accuracy	Error classification	Time(s)
NB	۹۰.۱۵	۹.۸۵	۲۵۴.۹۳
MNB	۹۱.۰۹	۸.۹۱	۰.۶۹
DMNB	۹۳.۲۵	۶.۶۵	۱۲.۶۶
CNB	۹۳.۰۹	۶.۹۱	۰.۸۱

* Adaptive Boosting

مهندسی برق و علوم کامپیوتر

The International Conference in New Research of Electrical Engineering and Computer Science

درجول ۳ با روش ۲-gram تعداد ۱۸۸۹۷ ویژگی تولید شده است که چهار الگوریتم بالا را آزمایش نمودیم. نتایج نشان می دهد که الگوریتم DMNB با کارایی ۹۳.۳۵٪ بالاترین کارایی و کمترین درصد خطا دسته بندی را دارد. ولی در مقایسه با ۱-gram کارایی کمتری را دارد.

با توجه به جدول های ۲ و ۳ نتایج نشان دهنده این مساله است که ۱-gram ویژگی های بهتری را در انتخاب و دسته بندی دارد و لذا در این مرحله از ۱-gram استفاده خواهد شد. در مرحله بعد از الگوریتم IG برای شناسایی ارزش ویژگی های مفید استفاده گردید که به شناسایی ۵۰۰ ویژگی مفید منجر شد. این ویژگی ها از ۸۰۰۷ مورد به ۵۰۰ ویژگی تقلیل داده شد ولی ایده پژوهش حاضر بر کاهش مجدد ویژگی می باشد. ویژگی های کاهش یافته از خروجی روش IG ذخیره و با روش انتخاب ویژگی CFS به تشخیص ویژگی های مفید پرداختیم. نتایج ویژگی های مفید به ۱۰۷ مورد رسیده است که کاهش بسیار مطلوبی در این روش داشته ایم.

جدول ۴: مقایسه کارایی روش انتخاب ویژگی پیشنهادی با ۱۰۷ ویژگی موثر

Algorithms	1Gram+ TF-IDF+(IG+CFS)		
	Accuracy	Error classification	Time(s)
NB	۸۳.۳۰	۱۶.۷۰	۱.۵۶
MNB	۹۱.۰۱	۸.۹۹	۰.۱۶
DMNB	۹۳.۳۱	۶.۶۹	۱.۲۵
CNB	۸۷.۸۲	۱۲.۱۸	۰.۳۱

در جدول ۴ بالا توانستم با ۱۰۷ ویژگی با انتخاب و کاهش آن ویژگی های به کارایی ۹۳.۳۱٪ برسیم که نتایج نشان دهنده کاهش در زمان تمامی الگوریتم ها می باشیم. مسئله کاهش ویژگی و انتخاب موثر آنها به زمان اجرای آزمایش الگوریتم های ما کمک بسزایی نموده است.

جدول ۵: مقایسه روش ترکیبی انتخاب ویژگی با روش ترکیبی یادگیری Bagging

Algorithms	Bagging		
	Accuracy	Error classification	Time(s)
NB	۹۲.۹۳	۷.۰۷	۸.۴۹
MNB	۹۲.۷۰	۷.۳۰	۰.۵۴
DMNB	۹۳.۳۲	۶.۶۸	۹.۱۴
CNB	۷۸.۷۸	۲۱.۲۲	۰.۲۳

در جدول ۵ روش ترکیبی انتخاب ویژگی با الگوریتم ترکیبی Bagging بهمراه ۴ الگوریتم بیز را آزمایش نموده و نتایج در این جدول بهبود بعضی روش ها را شامل شده است.

جدول ۶: مقایسه روش ترکیبی انتخاب ویژگی با روش ترکیبی یادگیری Adaboost

Algorithms	Adaboost		
	Accuracy	Error classification	Time(s)
NB	۹۳.۳۱	۶.۶۹	۶.۶۹
MNB	۹۳.۳۱	۶.۶۹	۲.۴۶
DMNB	۹۳.۴۷	۶.۵۳	۱.۳۴
CNB	۹۵.۴۹	۴.۵۱	۱.۵۵

در جدول ۶ روش ترکیبی انتخاب ویژگی با الگوریتم ترکیبی Adaboost به همراه ۴ الگوریتم بیز را آزمایش نموده و به بالاترین درصد با زمان نسبتاً مناسبی دست یافته ایم. کارایی نشان می دهد که الگوریتم بیز ساده مکمل با ۹۵.۴۹٪ خطای دسته بندی ۴.۵۱٪ بدست آمده است.

۷- نتیجه گیری

با گسترس اطلاعات و فراگیر شدن آنها در شبکه هایی مانند شبکه های اجتماعی، سایت ها و دیگر رسانه های ارتباطی جهت کاوش متون نیازی احساس گردید که اغلب حجم انبیه داده ها و گاهاً برخی از آنها تکراری و بی ارزش هستند که ممکن است که موتور های جستجو نتوانند به درستی یک متن را کاوش نمایند. حذف ویژگی های تکراری و انتخاب ویژگی های مناسب در این مقاله بیشترین اهمیت را داشته است. در این مقاله یک روش ترکیبی انتخاب ویژگی ارائه گردید که باعث کاهش زمان یادگیری الگوریتم و تسريع در امر پاسخگویی می باشد. در روش های ترکیبی همیشه زمان بیشتری را صرف یادگیری می کنند ولی با ارائه روش انتخاب ویژگی، زمان یادگیری زمان بر را کاهش داده ایم همچنین نتایج ارائه شده برای ۲۳ دسته با ۹۵.۴۹٪ بالاترین کارایی در این پژوهش و کمترین خطای دسته بندی را دارد.

مراجع:

- [1] 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.
- [2] Sebastiani, F. (2002). Machine learning in automated text categorization. *ACM computing surveys (CSUR)*, 34(1), 1-47.
- [3] Rennie, J. D., Shih, L., Teevan, J., & Karger, D. R. (2003, August). Tackling the poor assumptions of naive bayes text classifiers. In *ICML* (Vol. 3, pp. 616-623).
- [4] Su, J., Zhang, H., Ling, C. X., & Matwin, S. (2008, July). Discriminative parameter learning for Bayesian networks. In *Proceedings of the 25th international conference on Machine learning* (pp. 1016-1023). ACM.
- [5] Berry, M. W. (2004). Survey of text mining. *Computing Reviews*, 45(9), 548.
- [6] Dumais, S., Platt, J., Heckerman, D., Sahami, M. (1998). Inductive learning algorithms and representations for text categorization. In *Proceedings of the 7th International Conference on Information and Knowledge Management*.
- [7] Joachims, T. (1998). Text categorization with support vector machines: Learning with many relevant features. In *Proceedings of the European Conference on Machine Learning (ECML)*.
- [8] <http://ece.ut.ac.ir/dbrg/hamshahri/faindex.html>
- [9] Hashemi, S. M., & Dorfeshan, z., &Javadi, S., &Dezfouli, M.M. (2014). Using text mining techniques to categorize Persian texts dataset Hamshahri, *International Conference on Engineering, Arts and the Environment*.

- [10] Hashemi, S. M., & Dorfeshan, z., & Javadi, S., & Dezfuli, M.M. (2014). Categories of Persian texts Using filter words by size, International Conference on Engineering, Arts and the Environment.
- [11] Shipp, C. A., & Kuncheva, L. I. (2002, January). An investigation into how adaboost affects classifier diversity. In Proceedings of 9th International Conference on Information Processing and Management of Uncertainty in Knowledge-Based Systems (pp. 203-208).
- [12] Polikar, R. (2006). Ensemble based systems in decision making. Circuits and Systems Magazine, IEEE, 6(3), 21-45.
- [13] Freund, Y., & Schapire, R. E. (1996, July). Experiments with a new boosting algorithm. In ICML (Vol. 96, pp. 148-156).
- [14] 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.
- [۱۵] الهام مهدی پور، معصومه باقری قرقوک، افسانه رضایی، ۱۳۹۲ . سیستم خلاصه ساز خودکار متن فارسی با استفاده از الگوریتم ترکیبی SA-GA، هشتمین سمپوزیوم بین المللی پیشرفت‌های علوم و تکنولوژی مشهد، موسسه آموزش عالی خاوران.