

تاریخ دریافت مقاله: ۹۶/۱۰/۰۲
تاریخ پذیرش مقاله: ۹۷/۰۹/۰۴

روشی نوین برای شناسایی نویسنده متون با ترکیب الگوریتم‌های بهینه‌سازی توده ذرات و ماشین بردار پشتیبان

فرهاد سلیمانیان قره چیق*

استادیار، گروه مهندسی کامپیوتر، واحد ارومیه، دانشگاه آزاد اسلامی، ارومیه، ایران
پست الکترونیکی: farhad@iaurmia.ac.ir

محسن موتمن فر

کارشناس ارشد مهندسی کامپیوتر، گروه مهندسی کامپیوتر، واحد ارومیه، دانشگاه آزاد اسلامی، ارومیه، ایران
پست الکترونیکی: motaman.mohsen@yahoo.com

مهدی وفادار

کارشناس ارشد مهندسی کامپیوتر، گروه مهندسی کامپیوتر، واحد ارومیه، دانشگاه آزاد اسلامی، ارومیه، ایران
پست الکترونیکی: v.mehdi1364@gmail.com

چکیده

داده Reuter_50_50 انجام شده است. در روش پیشنهادی پس از استخراج ویژگی‌ها، داده‌ها با نسبت ۸۰ به ۲۰ به داده‌های آموزشی و آزمایشی تقسیم می‌شوند. مجموعه داده‌های آموزشی به‌عنوان ورودی به الگوریتم ماشین بردار پشتیبان داده می‌شوند و پس از آموزش دیدن توسط ماشین بردار پشتیبان، مدل مناسب ساخته می‌شود و سپس مجموعه داده‌های آزمایشی براساس مدل ساخته شده، اعتبارسنجی می‌شوند. تقسیم‌بندی داده‌ها در حالت ۸۰ (آموزش) درصد و ۲۰ (آزمایشی) درصد بهترین نتیجه را داشته است. نتایج بر مبنای معیارهای دقت و فراخوانی نشان می‌دهد که دقت تشخیص مدل پیشنهادی در مقایسه با مدل ماشین بردار پشتیبان بیشتر است. واژه‌های کلیدی: الگوریتم بهینه‌سازی توده ذرات، الگوریتم ماشین بردار پشتیبان، تشخیص نویسنده متون، استخراج ویژگی

تشخیص هوشمند نویسنده متون در زمینه‌های جرم‌شناسی دارای کاربردهای فراوانی می‌باشد اما به دلیل این‌که عوامل بسیاری در تشخیص نویسنده متون دخیل می‌باشند، نمی‌توان به‌طور دقیق و مطمئن نویسندگان متون را تشخیص داد. از این رو تحقیقات بسیاری در این زمینه از قرن نوزدهم صورت گرفته است اما تاکنون روشی با دقت ۱۰۰ درصد برای تمامی متون ارائه نشده است و محققان هر روزه روش‌های جدیدی را ارایه می‌دهند تا به دقت ۱۰۰ درصد نزدیک بشوند. در این مقاله برای حل این مسئله از ترکیب الگوریتم بهینه‌سازی توده ذرات و الگوریتم ماشین بردار پشتیبان استفاده شده است که از الگوریتم بهینه‌سازی توده ذرات به‌عنوان الگوریتم استخراج‌کننده ویژگی و از الگوریتم ماشین بردار پشتیبان به‌عنوان تشخیص‌دهنده نویسندگان استفاده شده است. ارزیابی بر روی مجموعه

* نویسنده مسئول

باگسترش روزافزون متون و داده‌های الکترونیکی در فضای مجازی، محققان برای جلوگیری از برخی سوء استفاده‌ها در محیط‌های مجازی، به دنبال ایجاد ابزارهای امنیتی بیشتری برای کاربران در جوامع مجازی هستند. بر همین اساس به‌منظور جلوگیری از دزدی‌های ادبی و علمی، محققان بسیاری به شناسایی افراد بر اساس سبک نگارش آن‌ها پرداخته‌اند. ایده شناسایی نویسنده، از مبحث طبقه‌بندی متون که خود زیرشاخه‌ای از علم پردازش زبان‌های طبیعی است گرفته شده است [۱] و در آن سعی می‌شود تا با تجزیه و تحلیل واژگان، دستور زبان، مفهوم یک جمله و با کمک گرفتن از دانش مربوط به واژگان، معنای آن جمله برای ماشین قابل درک شود [۲]. بر اساس فرضیه‌های مطرح شده توسط محققان، هر شخص الگوی خاصی در نحوه نگارش متون دارد که این موضوع به‌گونه‌ای همانند اثر انگشت هر نویسنده‌ای عمل می‌کند. شناسایی نویسنده از روی نثر، سبک و شیوه نوشتاری، یا به‌عبارت دیگر، ویژگی‌های نهفته در متون نوشته شده توسط نویسندگان، یکی از مباحث جدید در زمینه هوش مصنوعی و پردازش زبان طبیعی به‌شمار می‌رود.

برای حل مسئله شناسایی نویسنده متون باید ابتدا ویژگی‌های نهفته داخل متن استخراج گردند و سپس براساس ویژگی‌های استخراج شده و بااستفاده از روش‌های موجود در حوزه یادگیری ماشین به شناسایی نویسنده متون پرداخته می‌شود. میزان دقت در مرحله استخراج ویژگی‌ها تاثیر بسزایی در نحوه تشخیص نویسنده دارد به‌طوری که انتخاب درست ویژگی‌های استخراج شده باعث کاهش مدت زمان اجرای الگوریتم و افزایش میزان بهینگی الگوریتم می‌گردد. به‌طوری که اگر ویژگی‌های اضافی‌تر در تشخیص نویسنده استفاده گردد به‌جای بهبود دقت تشخیص، باعث کاهش دقت در تشخیص می‌گردد [۳].

در این مقاله برای تشخیص نویسنده متون از یک مدل

ترکیبی جدید بر مبنای الگوریتم بهینه‌سازی توده ذرات [۱۸] و ماشین بردار پشتیبان [۱۹] استفاده شده است. در روش پیشنهادی از الگوریتم بهینه‌سازی توده ذرات برای انتخاب ویژگی و از ماشین بردار پشتیبان برای طبقه‌بندی استفاده می‌شود. انتخاب ویژگی موجب کاهش سربار فضا و زمان برای پیاده‌سازی الگوریتم‌های طبقه‌بندی مانند ماشین بردار پشتیبان خواهد شد. برای ساخت یک روش طبقه‌بندی مناسب نیاز به داده‌های آموزشی با کیفیت است. این کیفیت با تعداد داده‌ها و ویژگی‌ها در مجموعه آموزش ارتباط دارد. انتخاب ویژگی‌ها به عنوان یکی از روش‌های پیش‌پردازش داده می‌تواند باعث افزایش کیفیت مجموعه داده آموزش برای ساخت مدل طبقه‌بندی گردد. انتخاب ویژگی‌ها دارای مزایای متعددی است. نخست زمان یادگیری در الگوریتم‌های طبقه‌بندی کاهش می‌یابد. و دیگر این که با انتخاب ویژگی‌ها نیاز کمتری به اندازه‌گیری و ذخیره‌سازی مقادیر ویژگی‌ها است.

هدف از انتخاب ویژگی، جستجوی یک مجموعه از ویژگی‌ها است که عملکرد طبقه‌بندی با روش‌های بانظارت مانند ماشین بردار پشتیبان را افزایش می‌دهد. در واقع باید ابعاد ویژگی کاهش یابد که هدف از این کار به‌دست آوردن بهترین و کوچکترین زیرمجموعه از ویژگی‌ها است. در روش انتخاب ویژگی‌ها، ویژگی‌هایی که در طبقه‌بندی مفید هستند، از یک مجموعه کامل ویژگی انتخاب می‌شوند. انتخاب ویژگی در تشخیص نویسنده متون اهمیت بسزایی دارد، زیرا تعداد زیادی ویژگی وجود دارند که بسیاری از آن‌ها یا بی‌استفاده هستند و یا این‌که بار اطلاعاتی چندانی ندارند و فقط به منظور انسجام و شکل‌گیری متون بیان شده‌اند. حذف نکردن این ویژگی‌ها مشکلی از لحاظ اطلاعاتی ایجاد نمی‌کند ولی بار محاسباتی را برای کاربرد مورد نظر بالا می‌برد.

ادامه مقاله به صورت زیر سازماندهی شده است. در بخش دوم، کارهای قبلی در مورد تشخیص نویسنده متون بررسی شده است. در بخش سوم ویژگی‌های استفاده

شده در تشخیص نویسنده متون، تحلیل شده است، در بخش چهارم روند طراحی و پیاده‌سازی روش پیشنهادی برای تشخیص نویسنده متون بیان شده است و در بخش پنجم به بررسی نتایج حاصله از روش پیشنهادی خواهیم پرداخت. و نهایتاً در بخش ششم به نتیجه‌گیری و کارهای آتی خواهیم پرداخت.

۲- روش‌های پیشین

مسائل موجود در حوزه پردازش زبان‌های طبیعی، از مباحثی می‌باشند که قابلیت حل شدن با روش‌های متن‌کاوی و یادگیری ماشین را دارا می‌باشند. تشخیص نویسنده متون نیز به‌عنوان یکی از مسائل مطرح در این حوزه می‌باشد که تاکنون برای حل این مسئله روش‌های فراوانی ارائه شده است که تعدادی از این روش‌ها را به این صورت می‌توان بیان نمود:

محققان [۴] نخستین بار برای تشخیص نویسندگان متون از نمایشنامه‌های شکسپیر استفاده نمودند که این تحقیق به‌عنوان اولین پژوهش در زمینه تشخیص نویسنده متون در قرن ۱۹ میلادی صورت گرفت و سپس در اواسط قرن بیستم مطالعات آماری در این زمینه توسط زیف^۱ در سال ۱۹۳۲ [۵] و یول^۲ در [۶] و [۷] انجام گردید. مطالعات دقیق‌تر توسط موستلر و همکارانش در سال ۱۹۶۴ بر روی پایگاه‌داده The Federalist Paper انجام شد که پس از این تحقیق محققان سعی در تشخیص نویسنده توسط روش‌های غیرسنتی نمودند [۸]. محققان [۹] برای شناسایی نویسنده متون موجود در اینترنت از الگوریتم‌های: بیزین ساده، ماشین بردار پشتیبان، نزدیک‌ترین همسایه و طبقه‌بندی حداقل مربعات منظم استفاده نمودند. اینان از وب‌نوشت‌های ۱۰۰ هزار نویسنده مختلف در سطح اینترنت به‌عنوان مجموعه‌داده استفاده نمودند.

باتوجه به تاثیر ویژگی‌های استخراج شده در شناسایی نویسندگان متون، محققان [۱۰] برای تشخیص نویسنده

متون از ویژگی‌های موجود در «تابع کلمه» استفاده نمودند که این ویژگی‌ها شامل حروف تعریف، حروف اضافه، حروف ربط، فعل و اسم می‌باشند. این ویژگی‌ها باتوجه به این‌که معنای متن را تغییر نمی‌دهند با تکرار زیادی در متن استفاده می‌گردند. محققان دیگری برای افزایش دقت نویسندگان متون از ویژگی‌های معنایی استفاده نمودند. ویژگی‌های معنایی به تحلیل دقیق‌تر متن و ویژگی‌های موجود در سطوح بالاتر می‌پردازند. این ویژگی‌ها شامل: (۱) گراف وابستگی معنایی که خود شامل ویژگی‌های دودویی معنایی و روابط اصلاح معنایی است، (۲) زمان و وجه فعل‌های به کار برده شده توسط نویسنده و شباهت‌های معنایی بین کلمات متن و همچنین شامل کلمات و عباراتی که جمله‌واره‌ها را به هم مربوط می‌کنند می‌باشد که این‌ها به‌عنوان حروف یا افزوده‌های ربطی شناخته می‌شوند [۱۱].

برکار دو و همکارانش برای شناسایی نویسندگان پیام‌های کوتاه برخط از الگوریتم n-گرام و روش‌های یادگیری بانظارت استفاده نمودند [۱۲]. برای شناسایی نویسنده متون عربی الگوریتم‌های مختلفی مبتنی بر الگوریتم بیزین ساده ارائه شده است [۱۳]. این الگوریتم‌ها شامل: الگوریتم چندجمله‌ای، الگوریتم چندجمله‌ای برنولی و الگوریتم چندجمله‌ای پواسون می‌باشند که از این الگوریتم‌ها برای شناسایی ۱۰ نویسنده مختلف استفاده شده است و براساس نتایج حاصله، الگوریتم چندجمله‌ای برنولی دارای بهترین عملکرد و دقت ۹۷٫۴۳ درصدی می‌باشد. محققان دیگری از ۳۳ ویژگی استخراج شده از معیارهای شباهت جمله‌متنی، شباهت مبتنی بر دانش، شباهت مبتنی بر پیکره، شباهت تعدادی، شباهت وابستگی نحوی و معیارهای ترجمه ماشینی به‌عنوان ویژگی‌های موثر در تشخیص نویسنده و از الگوریتم رگرسیون بردار پشتیبان برای محاسبه میزان شباهت در میان متون استفاده نمودند [۱۴].

هویدی و همکارانش برای تولید ویژگی‌های موردنیاز

1- Zipf
2- Yule

در شناسایی نویسندگان متون از سطح N گرام کلمات و سطح N گرام نویسه‌ها و برای شناسایی نویسندگان متون از الگوریتم بی‌زین ساده و ماشین بردار پشتیبان استفاده نمودند [۱۵]. محققان [۱۶] از الگوریتم‌های ماشین بردار پشتیبان و دسته‌بندی فازی برای شناسایی نویسندگان متون استفاده نموده‌اند. این محققان ابتدا ویژگی‌های مربوط به هر نویسنده را استخراج نموده و در مرحله بعدی براساس ویژگی‌های استخراج شده نویسندگان متون را شناسایی نموده‌اند. براساس نتایج ارائه شده توسط این محققان دقت پیش‌بینی الگوریتم ماشین بردار پشتیبان بیشتر از الگوریتم دسته‌بندی فازی می‌باشد و در نهایت دقت الگوریتم روش ترکیبی بیشتر از الگوریتم ماشین بردار پشتیبان می‌باشد که دارای دقت ۷۶ درصدی می‌باشد. برای تشخیص نویسندگان متون ترکیبی پژوهشگران [۱۷] روش نوینی را ارائه داده‌اند. این پژوهشگران ابتدا از میان ۳۵ ویژگی کلی، ۲۲ ویژگی را توسط الگوریتم جستجوی رتبه‌بندی انتخاب نموده و سپس ویژگی‌های انتخاب شده را به الگوریتم‌های یادگیری ماشین ارائه نموده‌اند که براساس نتایج حاصله الگوریتم چندجمله‌ای دارای دقت ۸۰ درصدی می‌باشد.

۳- ویژگی‌های استخراج شده برای شناسایی متون

استخراج ویژگی‌های موجود در درون متون به‌عنوان اصلی‌ترین مرحله در مسئله تشخیص نویسندگان متون می‌باشد. ویژگی‌های انتخاب شده جهت استخراج، تاثیر بسزایی در نحوه تشخیص نویسندگان متون دارد به‌طوری که اگر ویژگی‌های استخراج شده بهینه نباشد باعث کاهش دقت روش پیشنهادی خواهد شد. در این تحقیق نیز تعدادی از ویژگی‌های موجود در متون براساس سعی و خطا انتخاب و به‌عنوان ورودی به روش پیشنهادی ارائه شده است که این ویژگی‌ها در جدول (۱) بیان شده است [۱۲]. ویژگی‌های موجود در جدول (۱)، بر مبنای آزمایش بر روی مجموعه داده‌های تشخیص نویسندگان متون گردآوری و استخراج شده‌اند [۱۲]. این ویژگی‌ها در تشخیص نویسندگان

اسناد بسیار مهم هستند. چون بیشتر نویسندگان از دسته‌بندی‌های جدول (۱) در متون خود استفاده می‌کنند، لذا دسته‌بندی متون بر مبنای آن‌ها به الگوریتم‌های طبقه‌بندی کمک می‌کنند که به دنبال ویژگی‌های اصلی باشند و بر مبنای آن‌ها نوع رده متون را تشخیص دهند.

بر اساس جدول (۱) ویژگی‌های استخراج شده شامل ۵ دسته از ویژگی‌ها می‌باشد که هر کدام از ویژگی‌ها نیز دارای تعدادی زیرمجموعه می‌باشد. ۵ دسته اصلی شامل: ویژگی‌های لغوی، ویژگی‌های ساختاری، ویژگی‌های نحوی، کلمات دستوری، ویژگی‌های روانشناختی و زبانی می‌باشد [۱۲].

۴- روش پیشنهادی

با افزایش اهمیت مسئله تشخیص نویسندگان متون در حوزه‌های مختلف، این امر به‌عنوان یک امر ضروری و پراهمیت در حوزه پردازش زبان‌های طبیعی مطرح می‌گردد که نیازمند ارائه روش‌های جدیدی می‌باشد. لذا در این مقاله روش جدیدی مبتنی بر الگوریتم بهینه‌سازی توده ذرات و الگوریتم ماشین بردار پشتیبان ارائه شده است. با عنایت بر این‌که ویژگی‌های استفاده شده در مجموعه داده تاثیر بسزایی در میزان دقت تشخیص دارند لذا در روش پیشنهادی از الگوریتم بهینه‌سازی توده ذرات به‌عنوان الگوریتم استخراج کننده ویژگی‌ها استفاده شده است.

الگوریتم بهینه‌سازی توده ذرات [۱۸] یک الگوریتم مبتنی بر جمعیت است که در آن ذرات یک ازدحام (جمعیت) را تشکیل می‌دهند. جمعیت موجود در فضای مسئله حرکت می‌کند و براساس تجربیات فردی خود و تجربیات جمعی سعی می‌کنند تا راه حل بهینه در فضای جستجو را بیابند. الگوریتم بهینه‌سازی توده ذرات به‌عنوان یک الگوریتم بهینه‌سازی، یک جستجوی مبتنی بر جمعیت را فراهم می‌کند که در آن هر ذره با گذشت زمان موقعیت خود را تغییر می‌دهد. در الگوریتم بهینه‌سازی توده ذرات، ذرات در یک فضای جستجوی چند بعدی از راه‌حل‌های ممکن مسئله،

جدول ۱: ویژگی‌های استفاده شده برای تعیین نویسنده متون [۱۲]

دسته بندی	نوع ویژگی	توضیحات
	ویژگی‌های مبتنی بر نویسه	تعداد کل نویسه‌ها تعداد کل حروف / تعداد کل حروف کل تعداد نویسه‌های بزرگ / تعداد کل نویسه‌ها تعداد کل نویسه‌های عددی / تعداد کل نویسه‌ها تعداد کل فضای خالی / تعداد کل نویسه‌ها تعداد کل فضاها Tab / تعداد کل نویسه‌ها همه حروف نویسه‌های ویژه
	ویژگی‌های لغوی	تعداد کل کلمات تعداد کل کلمات مختصر / کل تعداد کلمات تعداد کل کلمات طویل / تعداد کل کلمات تعداد کل نویسه‌ها در کلمه / تعداد کل نویسه‌ها طول متوسط کلمات طول متوسط جملات به شکل نویسه طول متوسط جملات در قالب کلمات تعداد کل کلمات مختلف / کل تعداد کلمات تعداد کلمات که فقط یک بار در متن ظاهر میشود / تعداد کل کلمات تعداد کلمات که فقط دو بار در متن ظاهر میشود / تعداد کل کلمات
	ویژگی‌های ساختاری	مجموع خطوط تعداد کل جمله‌ها تعداد پاراگراف‌ها تعداد جملات در پاراگراف تعداد نویسه‌ها در پاراگراف جداکننده بین پاراگراف‌ها تعداد جملات شروع شده با حروف بزرگ / تعداد عبارات تعداد جملات شروع شده با حروف کوچک / تعداد جملات میانگین تعداد کلمات در هر جمله تعداد خطوط خالی / تعداد خطوط
	ویژگی‌های نحوی	فراوانی علائم نشانه‌گذاری مانند و . : ؛ ' ()
	کلمات دستوری	تعداد اسناد / تعداد کل کلمات تعداد ضمایر / تعداد کل کلمات تعداد افعال کمکی / تعداد کل کلمات تعداد پیوندها تعداد قیده‌ها / تعداد کل کلمات تعداد کلمات برای پاسخ کوتاه / تعداد کل کلمات
	ویژگی‌های روانشناختی و زبانی	منفی احساسات مثبت احساسات منفی اضطراب خشم غمگینی بینش شک یقین محدودیت توافق

تأثیر تجربه‌های خود و یا دانش همسایگانش بوده و رفتار جستجوی یک ذره در گروه تحت تأثیر ذرات دیگر است. همین رفتار ساده باعث پیدا شدن ناحیه‌های بهینه از فضای

حرکت می‌کنند. در این فضا یک معیار ارزیابی تعریف می‌شود و سنجش کیفیت راه‌حل‌های مسئله از طریق آن صورت می‌پذیرد. تغییر حالت هر ذره در یک گروه تحت

جستجو می‌گردد. بنابراین در الگوریتم بهینه‌سازی توده ذرات هر ذره به محض پیدا کردن موقعیت بهینه آن را به نحو مناسبی به اطلاع سایر ذرات می‌رساند و هر ذره بر اساس مقادیر به دست آمده برای تابع هزینه با یک احتمال معین تصمیم می‌گیرد تا از ذرات دیگر پیروی نماید و جستجو در فضای مسئله با استفاده از دانش قبلی ذرات انجام گیرد. این عمل باعث می‌شود تمام ذرات بیش از حد به یکدیگر نزدیک نشوند و به طور موثری از عهده حل مسائل بهینه‌سازی پیوسته برآیند.

در هر تکرار الگوریتم، هر ذره موقعیت بعدی خود را با توجه به دو مقدار تغییر می‌دهد، نخست بهترین موقعیتی که خود آن ذره تاکنون داشته است (pbest) و دیگری بهترین موقعیتی که تاکنون توسط کل ذرات جمعیت به وجود آمده است و در واقع بهترین pbest در کل جمعیت می‌باشد (gbest). بر مبنای مفهوم ریاضی، pbest برای هر ذره در واقع حافظه بیولوژیکی آن ذره محسوب می‌شود. gbest همان دانش عمومی جمعیت است و وقتی که افراد موقعیت خود را بر اساس gbest تغییر می‌دهند در واقع تلاش می‌کنند که سطح دانش خود را به سطح دانش جمعیت برسانند. از نظر مفهومی، بهترین ذره گروه همه ذرات گروه را به یکدیگر مرتبط می‌نماید. الگوریتم بهینه‌سازی توده ذرات الهام گرفته شده از حرکت دسته جمعی ذرات می‌باشد که در سال ۱۹۹۵ پیشنهاد شده است [۱۸]. در این الگوریتم برای حل مسئله از دو معیار سرعت و موقعیت ذرات استفاده می‌شود به طوری که در هر مرحله، ذرات بر اساس سرعت خاصی به سمت موقعیت بهینه حرکت می‌کنند.

در این الگوریتم برای بیان سرعت و موقعیت ذرات از فرمول‌های (۱) و (۲) استفاده می‌شود.

$$v_{ij}(t+1) = wv_{ij}(t) + c_1 r_1(t)(pbest_{ij}(t) - x_{ij}(t)) + c_2 r_2(t)(gbest(t) - x_{ij}(t)) \quad (1)$$

در فرمول (۱)، w به عنوان ضریب اینرسی، ضرایب c_1 و c_2 به عنوان ضرایب یادگیری و ضرایب r_1 و r_2 به عنوان اعداد تصادفی استفاده شده است. متغیر pbest بهترین جوابی

است که تا به حال توسط ذره i پیدا شده است و متغیر gbest نیز بهترین جوابی است که توسط کل ذرات موجود در جامعه، پیشنهاد شده است.

$$x_i(t+1) = x_i(t) + v_i(t+1) \quad (2)$$

در فرمول (۲) متغیر x_i بیانگر موقعیت فعلی ذره جاری و v_i بیانگر سرعت ذره جاری می‌باشد که توسط فرمول (۱) تولید می‌گردد.

در روش پیشنهادی پس از مشخص شده ویژگی‌ها توسط الگوریتم بهینه‌سازی توده ذرات، از الگوریتم ماشین بردار پشتیبان برای تشخیص نویسنده متون استفاده شده است. الگوریتم ماشین بردار پشتیبان در سال ۱۹۹۲ توسط واپنیک معرفی شده و بر پایه نظریه یادگیری آماری بنا گردیده است [۱۹]. در این الگوریتم داده‌ها بر اساس بردارهای پشتیبان به دسته‌های مختلفی تقسیم‌بندی می‌شوند به طوری که با تغییر بردارهای پشتیبان نتایج الگوریتم تغییر پیدا می‌کند. ماشین بردار پشتیبان [۱۹] یک رده‌بند دودویی است که دو رده را با استفاده از یک خط مرزی از هم جدا می‌کند. ماشین بردار پشتیبان قادر است به طور همزمان خطای رده‌بندی تجربی داده‌ها را کاهش و تفکیک‌پذیری رده را با استفاده از تغییر شکل‌های مختلف افزایش دهد. این مزیت ماشین بردار پشتیبان را قادر می‌سازد تا با داده با ابعاد زیاد و یا رده‌هایی با فضای توزیع چندبعدی بهتر عمل کند. در تقسیم خطی داده‌ها، هدف دستیابی به تابعی است که تعیین‌کننده ابرصفحه‌ای با بیشترین حاشیه می‌باشد. با حداکثر شدن حاشیه این ابرصفحه، تفکیک بین طبقات حداکثر می‌گردد. فرض کنید که $S = \{x_i, y_i\}$ یک نمونه آموزشی است که از دو رده $y_i = \pm 1$ و هر رده از $x_i, i=1, 2, \dots, m$ ویژگی تشکیل شده است.

در روش پیشنهادی به منظور نگاشت داده‌ها به صورت خطی از تابع پایه شعاعی استفاده شده است. تابع پایه شعاعی برای ایجاد ماشین‌هایی با انواع مختلف از سطوح غیرخطی در فضای داده‌ها، ضرب‌های داخلی تولید می‌کند [۱۹]. معمولاً تابع پایه شعاعی برای طبقه‌بندی و پیش‌بینی

عملکرد بهتری دارد. برای پیاده‌سازی این تابع از فرمول (۳) استفاده شده است.

$$k(x_i, x_j) = \exp(-\alpha |x_i - x_j|^2), \alpha > 0 \quad (3)$$

در فرمول (۳) پارامتر x_i و x_j بیانگر بردار α م و λ م باشد و پارامتر α بیانگر گاما م باشد که پارامتر گاما مشخص کننده شکل مرز تصمیم می‌باشد. ماشین بردار پشتیبان در مرحله آموزش، مرز تصمیم‌گیری را به گونه‌ای انتخاب می‌نماید که حداقل فاصله آن با هر یک از دسته‌های مورد نظر را بیشینه کند. این نوع انتخاب باعث می‌شود که تصمیم‌گیری در شرایط نوفه‌دار بهتر انجام شود و نمونه‌های بین دو رده با دقت بیشتری تشخیص داده شوند. برای بیان بهتر روش پیشنهادی نحوه عملکرد این الگوریتم در شکل (۱) نمایش داده شده است.

همان‌گونه که شکل (۱) نشان می‌دهد، پس از آن‌که بهینه‌ترین ویژگی‌ها از میان ویژگی‌های استخراج شده توسط الگوریتم بهینه‌سازی توده ذرات انتخاب گردید، مجموعه داده تولید شده توسط ویژگی‌های استخراج شده به دو مجموعه داده آزمایش و آموزش تقسیم‌بندی می‌شوند که مجموعه داده آموزش به‌عنوان ورودی به الگوریتم ماشین بردار پشتیبان داده می‌شود و پس از آموزش دیدن ماشین بردار پشتیبان، نویسنده متون موجود در مجموعه داده آزمایشی براساس ماشین بردار پشتیبان آموزش دیده، تعیین می‌گردند و در نهایت براساس دو معیار ارزیابی دقت و دوباره فراخوانی مورد ارزیابی قرار می‌گیرند. انتخاب ویژگی‌ها توسط الگوریتم بهینه‌سازی توده ذرات بر مبنای ویژگی‌های مطرح شده در متن اسناد مانند شمارش تعداد کلمات کلیدی انجام می‌گیرد. به هر ویژگی یک وزن اختصاص داده می‌شود و بر مبنای وزن داده شده به ویژگی‌ها، ذرات جستجو در فضای جستجو را شروع می‌کنند و همسایه‌های نزدیک به نقطه انتخاب شده را جستجو می‌کنند. لذا ویژگی‌هایی انتخاب می‌شوند که نزدیکی و تشابه بیشتری به یکدیگر دارند.

۵- ارزیابی و نتایج

در این مقاله برای ارزیابی روش پیشنهادی از چهار معیار ارزیابی: تعداد متونی که نویسنده آن‌ها درست تشخیص داده شده است، تعداد متونی که نویسنده آن‌ها درست تشخیص داده نشده است، درصد دقت و درصد فراخوانی از مجموعه داده Reuter_50_50 به‌عنوان مجموعه داده آموزشی و آزمایشی با نسبت ۸۰ به ۲۰ درصد استفاده شده است. در بیشتر موارد برای ارزیابی طبقه‌بندی و پیش‌بینی نمونه‌ها از ۸۰ درصد داده‌ها برای آموزش و از ۲۰ درصد داده‌ها برای آزمایش استفاده می‌شود. در مرحله آموزش، مدل ساخته می‌شود؛ که مدل بر مبنای قوانین و قواعد الگوریتم طبقه‌بندی و استفاده از شرط‌ها ایجاد می‌شود و سپس بر مبنای مدل ساخته شده در مرحله آموزش از ۲۰ درصد نمونه برای آزمایش استفاده می‌شود. نتایج حاصله از روش پیشنهادی براساس دو معیار دقت و فراخوانی [۲۰] مورد ارزیابی قرار گرفته است.

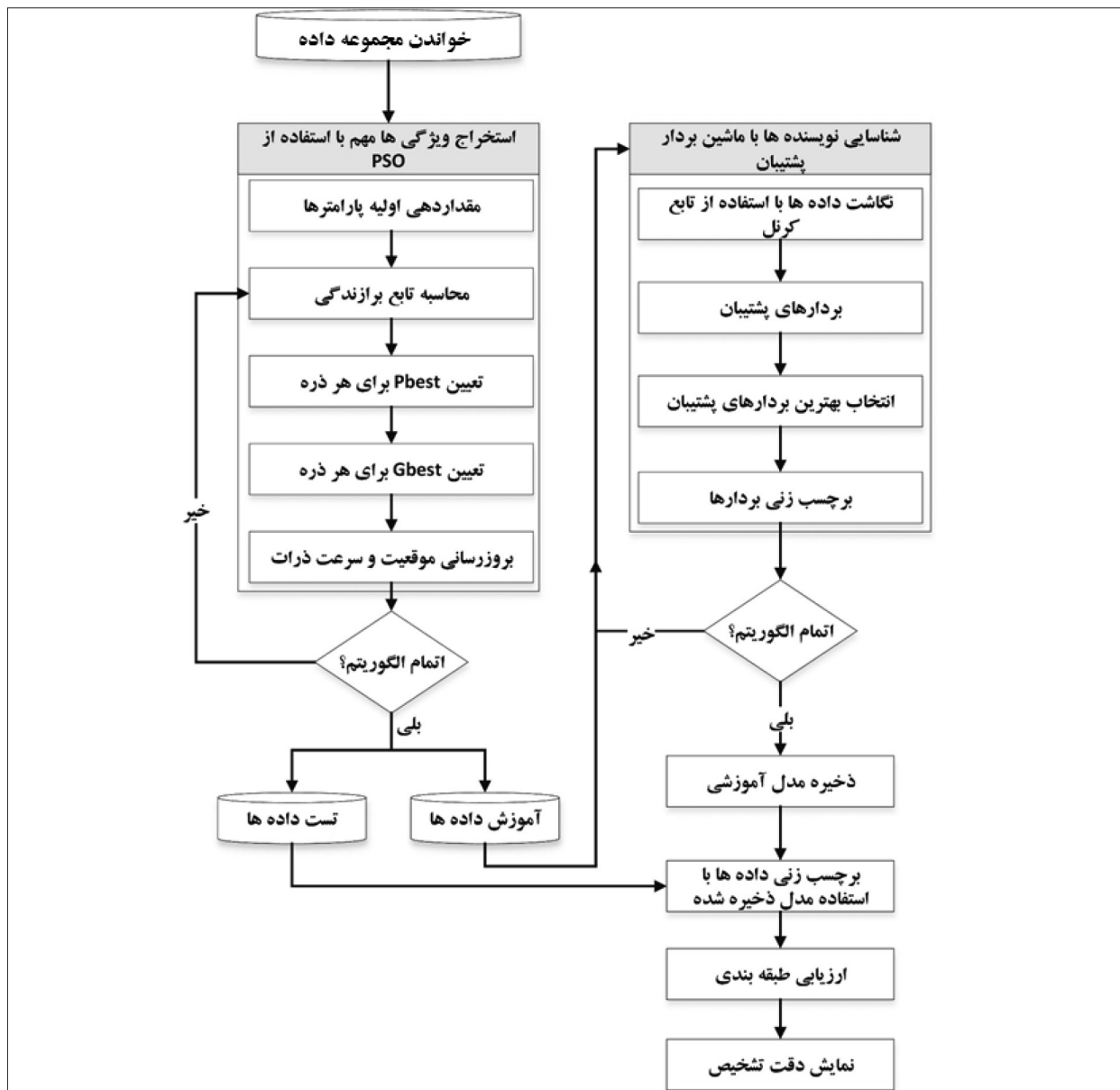
$$Precision = \frac{True\ positive}{True\ positive + False\ Positive} \quad (4)$$

از فرمول (۴) برای به دست آوردن دقت روش پیشنهادی و از فرمول (۵) برای به دست آوردن فراخوانی استفاده شده است.

$$Recall = \frac{True\ positive}{True\ positive + False\ Negative} \quad (5)$$

معیار دقت شامل پارامترهای درست مثبت (TP) و کاذب مثبت (FP) است. در این معیار نمونه‌های کاذب مثبت هم به عنوان مثبت درست در دقت تشخیص تاثیرگذار هستند. معیار فراخوانی شامل پارامترهای درست مثبت و کاذب منفی (FN) است. در این معیار نمونه‌های کاذب منفی هم در دقت تاثیرگذار هستند. نتایج حاصله از مقایسه روش پیشنهادی در اشکال (۲)، (۳) و جدول (۲) نمایش داده شده است.

براساس شکل (۲) که روش پیشنهادی براساس مجموعه داده آزمایشی با الگوریتم ماشین بردار پشتیبان

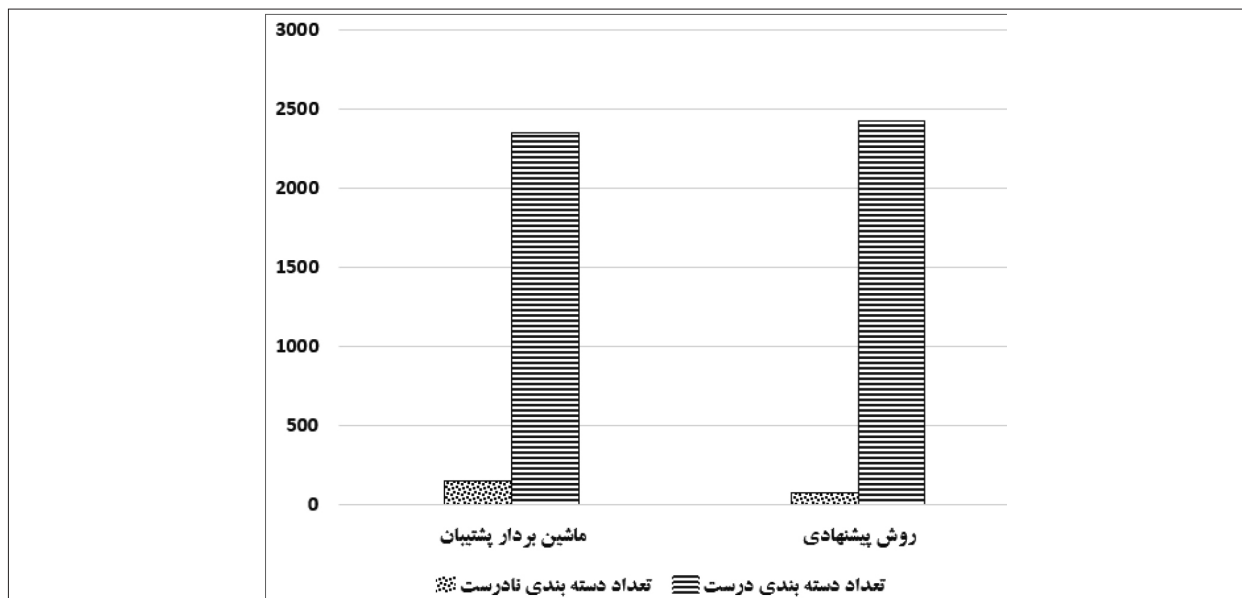


شکل ۱: نحوه عملکرد روش پیشنهادی

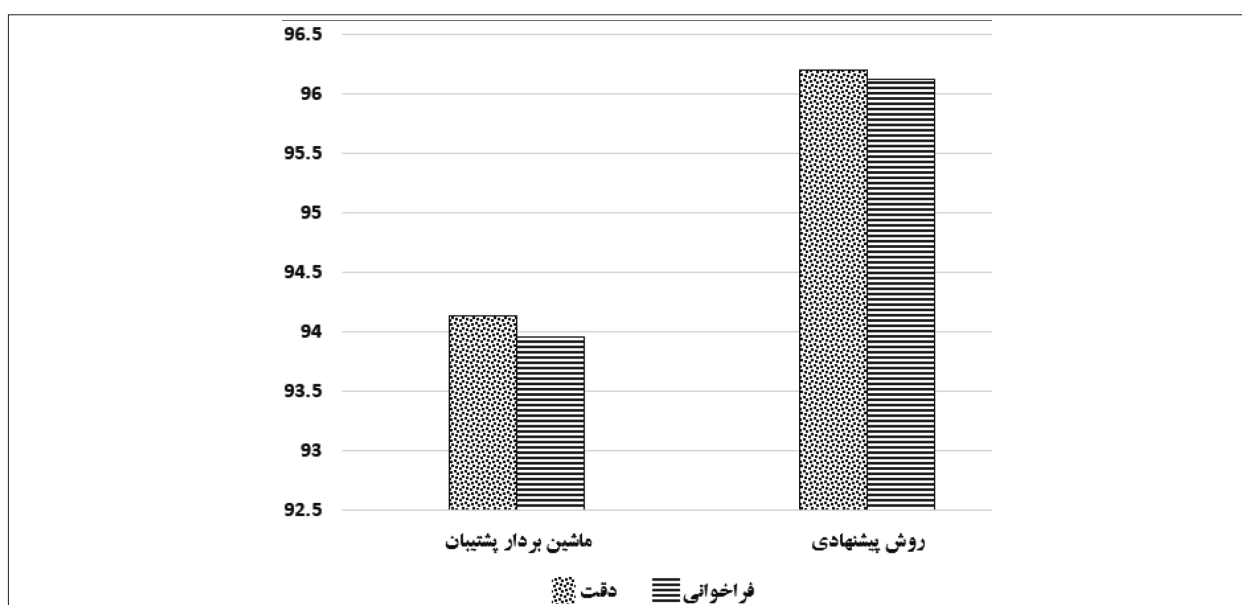
فراخوانی و دقت با الگوریتم ماشین بردار پشتیبان مورد مقایسه قرار گرفته است. براساس این نتایج می‌توان گفت که روش پیشنهادی دارای عملکرد بهینه‌تری نسبت به ماشین بردار پشتیبان می‌باشد. به دلیل این‌که در مدل پیشنهادی از الگوریتم بهینه‌سازی توده ذرات برای انتخاب ویژگی استفاده شده است.

در جدول (۲) روش پیشنهادی براساس چهار معیار تعریف شده با الگوریتم ماشین بردار پشتیبان مورد ارزیابی قرار گرفته شده است که نتایج حاصله از این جدول

مورد ارزیابی قرار گرفته است، می‌توان گفت که روش پیشنهادی فقط نویسنده ۷۷ متن را به درستی تشخیص نداده است؛ این درحالی می‌باشد که الگوریتم ماشین بردار پشتیبان نویسنده ۱۵۱ متن را به درستی تشخیص نداده است. لذا می‌توان بیان نمود که عملکرد روش پیشنهادی بهینه‌تر از الگوریتم ماشین بردار پشتیبان می‌باشد. تعداد دسته‌بندی درست در روش پیشنهادی و ماشین بردار پشتیبان به ترتیب برابر با ۲۴۲۳ و ۲۳۴۹ می‌باشد. در شکل (۳)، روش پیشنهادی براساس دو معیار



شکل ۲: مقایسه روش پیشنهادی با الگوریتم ماشین بردار پشتیبان



شکل ۳: مقایسه روش پیشنهادی با الگوریتم ماشین بردار پشتیبان

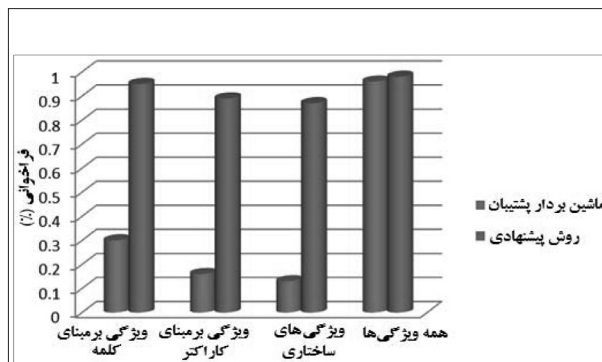
جدول ۲: مقایسه روش پیشنهادی با الگوریتم ماشین بردار پشتیبان

مدل‌ها	تعداد دسته‌بندی نادرست	تعداد دسته‌بندی درست	دقت	فراخوانی
ماشین بردار پشتیبان	151	2349	94,13	93,96
روش پیشنهادی	77	2423	96,20	96,12

۹۳,۹۶ درصد است. همچنین مقدار دقت در مدل پیشنهادی و ماشین بردار پشتیبان به ترتیب برابر با ۹۶,۲۰ درصد و ۹۴,۱۳ درصد است.

در این تحقیق از سه نوع ویژگی: مبتنی بر کلمات، مبتنی بر نویسه و ویژگی‌های ساختاری استفاده شده است. به منظور تعیین میزان تاثیرگذاری ویژگی‌های مطرح شده، هرکدام از ویژگی‌ها به صورت جداگانه به عنوان مجموعه داده در نظر گرفته شده‌اند که نتایج حاصله از این

بیانگر عملکرد مناسب روش پیشنهادی می‌باشد. جدول (۲)، نشان می‌دهد که مقدار فراخوانی در روش پیشنهادی برابر با ۹۶,۱۲ درصد است و در ماشین بردار پشتیبان برابر با



شکل ۵: مقایسه روش پیشنهادی با روش‌های مطرح شده در مجموعه داده آزمایشی با تفکیک ویژگی‌ها براساس فراخوانی

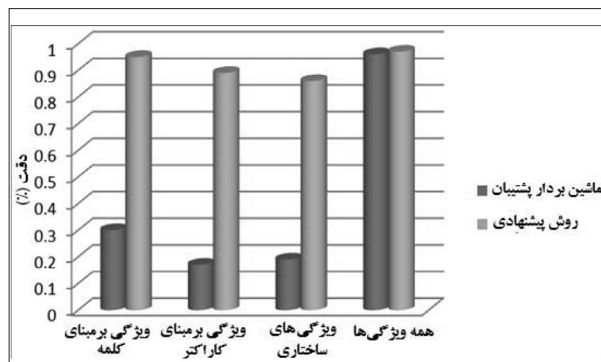
جدول ۳: مقایسه روش پیشنهادی با روش‌های مطرح شده در بخش پیشینه تحقیق

روش‌ها	دقت
الگوریتم نیوی بیز چندجمله‌ای برنولی [۱۳]	۹۷,۰۴
بیزین ساده [۱۵]	۷۱,۸۵
ماشین بردار پشتیبان [۱۵]	۶۲,۹۶
الگوریتم نیوی بیز چندجمله‌ای برمبنای جستجوی رتبه‌ای [۱۶]	۸۰,۰۰
ماشین بردار پشتیبان فازی [۱۷]	۷۶,۰۰
روش پیشنهادی	۹۶,۰۲

است دقت بیشتری داشته باشد و لذا الگوریتم مناسبی برای شناسایی نویسنده متون باشد. دلیل افزایش دقت روش پیشنهادی در مقایسه با چهار مدل دیگر استفاده از الگوریتم بهینه‌سازی توده ذرات برای انتخاب ویژگی است. به عبارتی الگوریتم بهینه‌سازی توده ذرات با انتخاب ویژگی، بهترین متون را شناسایی و انتخاب می‌کند و سپس آن‌ها را به ماشین بردار پشتیبان می‌دهد و ماشین بردار پشتیبان با داده‌های دقیق‌تر و با نوبه کمتر مواجه است.

۶- نتیجه گیری

با گسترش روزافزون متون الکترونیکی در سطح فضای مجازی، برای جلوگیری از وقوع تخلفات، افزایش امنیت در این فضا و کنترل بر محتوای تولید شده توسط کاربران باید روش‌های برای شناسایی نویسندگان متون نوشته شده ارائه گردد. از این رو در این مقاله سعی شده است تا با ارائه یک روش ترکیبی، میزان دقت را در



شکل ۴: مقایسه روش پیشنهادی با روش‌های مطرح شده در مجموعه داده آزمایشی با تفکیک ویژگی‌ها براساس دقت

ارزیابی در اشکال (۴) و (۵) نمایش داده شده است.

در شکل (۵)، درصد معیار فراخوانی برمبنای دسته‌های مختلف نشان داده شده است. روش پیشنهادی در مقایسه با ماشین بردار پشتیبان درصد فراخوانی بیشتری دارد. مقدار فراخوانی در روش پیشنهادی برمبنای همه ویژگی‌ها نزدیک به ۱۰۰ درصد است. روش پیشنهادی در همه دسته‌ها دارای درصد فراخوانی بیشتری است و لذا در روش پیشنهادی استفاده از انتخاب ویژگی تاثیرگذار بوده است.

براساس مقادیر نشان داده شده در اشکال (۴) و (۵) می‌توان بیان نمود که ویژگی‌های مبتنی بر کلمه دارای بیشترین تاثیر در شناسایی نویسنده می‌باشند. در جدول (۳) روش پیشنهادی با روش‌های مطرح شده در بخش پیشینه تحقیق مورد بررسی قرار گرفته است. نتایج این جدول بیانگر عملکرد بهینه روش پیشنهادی می‌باشد. روش پیشنهادی با الگوریتم چندجمله‌ای برنولی [۱۳]، بیزین ساده [۱۵]، ماشین بردار پشتیبان [۱۵]، الگوریتم چندجمله‌ای برمبنای جستجوی رتبه‌ای [۱۶] و ماشین بردار پشتیبان فازی [۱۷] مقایسه شده است.

جدول (۳)، نشان می‌دهد که درصد دقت روش پیشنهادی برابر ۹۶,۰۲ است که در مقایسه با مدل‌های [۱۵]، ماشین بردار پشتیبان [۱۵]، الگوریتم چندجمله‌ای برمبنای جستجوی رتبه‌ای [۱۶] و ماشین بردار پشتیبان فازی [۱۷] دقت بیشتری دارد. همچنین درصد دقت الگوریتم چندجمله‌ای برنولی [۱۳] در مقایسه با همه مدل‌ها بیشتر است. روش پیشنهادی با غلبه بر چهار مدل توانسته

On the Feasibility of Internet-Scale Author Identification, Draft. A Version of This Paper Will Appear At IEEE S&P 2012.

10-M.R.Davarpanah, M. SanjiAramideh, Farsi lexical analysis and stop word list, Library Hi Tech, Vol. 27 Issue 3, pp.435-449, 2009.

11-S.Argamon, C.Whitelaw, P.Chase, S.R.Hota, N.Garg, S.Levitan, Stylistic text classification using functional lexical features, Journal of the American Society for Information Science and Technology, 58(6), 802-822, 2007.

12-M.L.Brocardo, I.Traore, Sh.Saad, I.Woungang, Authorship Verification for Short Messages Using Styliometry, Conference on Computer, Information and Telecommunication Systems (CITS 2013), Piraeus-Athens, Greece, 2013.

13-A.S.Altheneyan and M.B. Menai, Naive Bayes Classifiers for Authorship Attribution of Arabic Texts, Journal of King Saud University – Computer and Information Sciences 26, 473-484, 2014.

14-T. T. ZHU And M. LAN, ECNUCS: Measuring Short Text Semantic Equivalence Using Multiple Similarity Measurements, Second Joint Conference On Lexical And Computational Semantics (*SEM), pp. 124-131, 2013.

15-F. Howedi, M.Mohd, Text Classification for Authorship Attribution Using Naive Bayes Classifier With Limited Training Data, Computer Engineering And Intelligent Systems, Vol.5, No.4, 2014.

16-M. Sudheep Elayidom, Ch.Jose, A.Puthussery, N.K Sasi, TEXT CLASSIFICATION FOR AUTHORSHIP ATTRIBUTION ANALYSIS, Advanced Computing: An International Journal (ACIJ), Vol.4, No.5, September 2013.

17-T.TAŞ, A.K. GORUR, Author Identification for Turkish Texts, Journal of Arts and Sciences Say: 7, May 2007.

18-J.Kennedy, RC.Eberhart. Particle swarm optimization, In: Proceedings of the IEEE conference on neural networks, Perth, Australia; pp. 1942-8, 1995.

19-C.Cortes and V.Vapnik, "Support-Vector Networks", Machine Learning, Vol. 20, Issue 3, pp. 273-297, 1995.

20-V. RAGHAVAN and G. JUNG, A Critical Investigation of Recall and Precision as Measures of Retrieval System Performance, ACM Transactions on Information Systems, Vol. 7, No. 3, Pages 205-229, July 1989.

تشخیص نویسندگان متون افزایش داد. در روش ترکیبی از الگوریتم بهینه‌سازی توده ذرات به عنوان الگوریتم استخراج‌کننده ویژگی‌های متون موجود در مجموعه داده، از الگوریتم ماشین بردار پشتیبان به عنوان شناسایی‌کننده نویسنده متون و از متون موجود در Reuter_50_50 به عنوان مجموعه داده استفاده شده است. براساس نتایج حاصله از ارزیابی روش پیشنهادی، الگوریتم پیشنهادی دارای درصد دقت ۹۶٫۲ و درصد فراخوانی ۹۶٫۱۲ می‌باشد این درحالی است که میزان درصد دقت و درصد دوباره فراخوانی برای الگوریتم ماشین بردار پشتیبان برابر با ۹۴٫۱۳ و ۹۳٫۹۶ می‌باشد. براساس این نتایج می‌توان گفت که روش پیشنهادی دارای عملکرد بهینه‌تری می‌باشد از این رو می‌توان از روش پیشنهادی در سایر مسائل مطرح در حوزه هوش مصنوعی استفاده نمود.

مراجع

- 1- D. D. Lewis, A new Benchmark Collection for Text Categorization Research, The Journal of Machine Learning Research, 5: 361-397, 2004.
- 2- J.Allen, Natural language Understanding, Benjamin/ Cummings Publishing Company, 1987.
- 3- L.S.Oliveira, A Methodology for Feature Selection Using Multi-objective Genetic Algorithms for Handwritten Digit String, International Journal of Pattern Recognition and Artificial intelligence, 17 (6), pp: 903-929, 2003.
- 4- T.C.Mendenhall, the Characteristic Curves of Composition, Science, IX, 237-249, 1887.
- 5- G.K. Zipf, Selected Studies of the Principle of Relative Frequency in Language, Cambridge, MA: Harvard University Press, 1932.
- 6- G.U. Yule, On Sentence-Length as a Statistical Characteristic of Style in Prose, With Application to Two Cases of Disputed Authorship, Biometrika, 363-390, 1938.
- 7- G.U. Yule, The Statistical Study of Literary Vocabulary, Cambridge University Press, 1944.
- 8-F. Mosteller and D.L.Wallace, Inference and Disputed Authorship: The Federalist, Reading, MA: Addison-Wesley, 1964.
- 9-A.Narayanan, H.Paskov, N.Z.Gong, J.Bethencourt,