

دسته‌بندی بیزین جمعی با استفاده از انتخاب ویژگی رپر مبتنی بر الگوریتم ژنتیک در تشخیص هرزنامه

مدیریت

اطلاعات

دوره ۶، شماره ۲

پاییز و زمستان ۱۳۹۹

وحید نصرتی^۱

دانشجوی دکتری، مهندسی کامپیوتر، دانشکده فنی و مهندسی، دانشگاه اراک، اراک، ایران

محسن رحمانی

دانشیار، مهندسی کامپیوتر، دانشکده فنی و مهندسی، دانشگاه اراک، اراک، ایران

چکیده: جایگاه ایمیل در ارتباطات، با ورود پدیده‌ای به نام هرزنامه با تهدید جدی مواجه شده است. تاکنون، به‌منظور مقابله با این پدیده، روش‌های فراوانی پیشنهاد شده که یکی از مهم‌ترین این روش‌ها، دسته‌بندی آنها بر اساس محتوا به دو دسته هرزنامه و غیرهرزنامه است. دسته‌بندی بر اساس محتوا با استفاده از کلمات بهعنوان ویژگی انجام می‌شود که به‌دلیل تعداد زیاد ویژگی‌ها، استفاده از یک سازوکار انتخاب ویژگی کارآمد موضوعی حیاتی به نظر می‌رسد. بر این اساس، تمرکز روش پیشنهادی در این مقاله روی انتخاب ویژگی‌های مفید بوده و یک فرایند انتخاب ویژگی رپر با بهره‌گیری از الگوریتم قدرتمند ژنتیک و با همکاری دسته‌بند بیزین که دارای کارایی بالایی در مسائل دسته‌بندی متون است، ارائه می‌شود. روش کار نیز به این صورت است که ابتدا یک بردار ویژگی اولیه ساخته شده، سپس با ضرب کردن آن در یک ماتریس با عنوان ماتریس انتقال، با استفاده از الگوریتم ژنتیک، روی آن عملیات بهینه‌سازی اعمال شده و در پایان، k بردار ویژگی نهایی ساخته می‌شوند. عملیات دسته‌بندی نیز بهصورت جمعی و با اعمال k دسته‌بند بیزین روی بردارهای ویژگی اعمال شده و از بین آنها رأی‌گیری انجام می‌شود. روش پیشنهادی روی دو پایگاه داده اجرا شده که بر اساس نتایج، روش پیشنهادی با مقدار $k = 7$ دارای نرخ صحت ۸۷/۹۱ و ۸۷/۷۶ در دو پایگاه داده PU۱ و PU۲ است. همچنین نتایج مقایسه روش پیشنهادی، حاکی از کارآمدی روش پیشنهادی در مقایسه با بیزین پایه و دو دسته‌بند SVM و KNN است.

کلیدواژه‌ها: ایمیل، هرزنامه، دسته‌بندی، الگوریتم ژنتیک، انتخاب ویژگی، ماتریس انتقال، یادگیری جمعی.

مقدمه

سرویس ایمیل امروزه به عنوان یک روش ارتباطی مهم، در بین افراد جامعه، دارای محبوبیت فراوان و روزافزونی است. همین محبوبیت و رایگان بودن ایمیل، برای تبلیغات با قیمت ارزان، زمینه مناسبی را فراهم می‌آورد. نامه‌های ناخواسته که معمولاً بعنوان هرزنامه شناخته می‌شوند، یکی از تهدیدها برای امنیت اینترنت به شمار می‌روند. هرزنامه‌ها، صندوق پستی افراد را با حجم وسیعی از پیام‌های ناخواسته اشغال می‌کنند که علاوه بر هدررفت پنهانی باند شبکه و فضای ذخیره‌سازی، به توزیع سریع اطلاعات نادرست و گسترش کدهای مخرب در بین کاربران اینترنت منجر شده و افراد خردسال را نیز در معرض محتوای نامناسب قرار می‌دهد (Faris et al, 2019). گسترش نامه‌های مخرب به حدی است که بر اساس Saidani et al, 2020: آنچه این مسئله را بدتر می‌کند این است که هرزنامه‌نویس‌ها به‌طور پیوسته، برای دور زدن فیلترها، روش‌های جدیدی ابداع می‌کنند. از طرف دیگر، حجم گسترشده جریان داده بین صدها میلیون نفر و تعداد زیاد صفات، این مشکل را بسیار پیچیده‌تر می‌کند، بنابراین، پیشنهاد مدل‌های کارآمد و سازگار با تشخیص هرزنامه به یک ضرورت تبدیل شده است (Faris et al, 2019: 67).

تاکنون، برای مقابله با هرزنامه روش‌های مختلفی ارائه شده است. دادا و همکاران^۱ (۲۰۱۹)، مسائل و مشکلات موجود در این زمینه و آخرین روش‌های موجود را بررسی کرده‌اند. به‌طور کلی، به‌منظور شناسایی پیام‌های مخرب، از دو نوع اطلاعات شامل اطلاعات محتوایی^۲ و غیرمحتوایی^۳ استفاده می‌شود (Hu et al, 2010). در روش‌های مبتنی بر اطلاعات غیرمحتوایی از سرصفحه‌های ایمیل و اطلاعات فرستنده مربوطه مانند دامنه آدرس ایمیل فرستنده (آدرس IP)، سابقه^۴، سبک نوشتن و زمان ارسال استفاده می‌شود و روش‌های مبتنی بر محتوا از اطلاعات متنی موجود در متن ایمیل‌ها و موضوعات ایمیل استفاده می‌کنند. در پژوهش حاضر، تمرکز اصلی بر روش‌های مبتنی بر محتوا برای شناسایی هرزنامه قرار دارد. در بحث تشخیص و شناسایی هرزنامه، با مجموعه‌ای از سندهای متنی با عنوان ایمیل سروکار داریم که می‌توان برای این منظور، از روش‌های پردازش متنی کمک گرفت. یکی از روش‌هایی که در این زمینه کاربرد فراوانی دارد، دسته‌بندی ایمیل‌ها بر اساس محتوا با عنوان عملیات دسته‌بندی متن^۵ است که در آن، تعدادی فایل متنی (ایمیل) وجود دارد که بایستی (با تجزیه و تحلیل آنها) دسته‌بندی شوند.

الگوریتم‌های دسته‌بند بر اساس یک سری ویژگی به دسته‌بندی پیام‌های رسیده اقدام می‌کنند. همین موضوع باعث می‌شود کارایی آنها تا حدود زیادی مشروط به انتخاب ویژگی مناسب باشد. هرچقدر این ویژگی‌ها به صورت بهینه‌تر انتخاب شوند، باعث افزایش کارایی الگوریتم دسته‌بندی می‌شوند. به همین

1. Kaspersky
2. Dada et. al
3. Content based information
4. Non-content
5. Header
6. Reputation
7. Text classification

دلیل تاکنون برای بهبود کارایی دسته‌بند با بهینه کردن روش انتخاب ویژگی تلاش‌های مختلفی شده است (Kołcz et al, 2004).

انتخاب ویژگی می‌تواند به صورت فیلتر^۱، رپر^۲ یا جاسازی شده^۳ انجام شود. روش رپر به دلیل بهره‌گیری از یک الگوریتم دسته‌بندی، در مقایسه با فیلتر معمولاً نتایج بهتری دارد. انتخاب ویژگی رپر مبتنی بر الگوریتم ژنتیک، برای انتخاب ویژگی روشی رایج به شمار می‌رود و کارایی آن در حیطه‌های مختلف همچون صنعت پزشکی، پردازش تصویر، پردازش متن، بیوانفورماتیک و کاربردهای صنعتی اثبات شده است (JadHAV et al, 2018)، بنابراین در پژوهش حاضر، یک روش انتخاب ویژگی مبتنی بر الگوریتم ژنتیک در مسئله تشخیص هرزنامه ارائه می‌شود. برخلاف روش‌های پیشین، در روش پیشنهادی، به عنوان جواب نهایی از یک کروموزوم استفاده نمی‌شود و برای انتخاب k بردار از ویژگی‌های برتر از مجموعه‌ای از k کروموزوم برتر استفاده شده و در مرحله بعد، k دسته‌بند بیزین روى آنها اعمال شده و دسته‌بندی نهایی با رأی‌گیری ساده از بین دسته‌بندها انجام می‌شود. با توجه به اینکه الگوریتم ژنتیک شناخته شده‌ترین روش بهینه‌سازی هوشمند و الگوریتم تکاملی به شمار می‌رود، می‌توان با قرار دادن آن در کنار روش آماری همچون بیزین که خود یکی از روش‌های آماری مهم محسوب می‌شود، نتایج بهتری به دست آورد.

ایده استفاده از k الگوریتم دسته‌بند از مفهوم یادگیری جمعی^۴ الهام گرفته شده است که در آنها برخلاف سازوکارهایی که فقط مبتنی بر یک دسته‌بند هستند، از مجموعه‌ای از دسته‌بندها استفاده می‌شود. به طور کلی، این نوع دسته‌بندها در تشخیص هرزنامه در مقایسه با روش‌های تکی کارایی بهتری دارند (Chinavle et al, 2009). بوستینگ^۵ و بکینگ^۶ از روش‌های مهمی هستند که به صورت گروهی عمل می‌کنند. روش بوستینگ با اندازه‌گیری خطای دسته‌بندی هر گروه^۷ وزن مثال‌های اشتباہ دسته‌بندی شده را برای گروه دسته‌بند بعدی افزایش می‌دهد. روش بکینگ نیز با آموزش دادن تعدادی دسته‌بند با مجموعه‌های مستقل از مثال‌های آموزشی سعی در کاهش واریانس دارد (DeBarr et al, 2012).

با توجه به اهمیت دو فرایند دسته‌بندی و انتخاب ویژگی در سیستم‌های تشخیص هرزنامه، باید برای یافتن روش کارآمد، به هر دوی این فرایندها توجه شود. در حالی که در اکثر پژوهش‌های حاضر تلاش شده است تا فقط یکی از این دو فرایند بهبود بخشیده شود، در پژوهش حاضر به دلیل وابستگی این دو بخش و اینکه یک الگوریتم انتخاب ویژگی کارآمد تأثیر سیار زیادی بر کارایی فرایند دسته‌بندی دارد، به هر دوی آنها توجه شده و برای بهبود آنها تلاش شده است. همچنین در پژوهش حاضر، به منظور ساخت یک مدل قدرتمند تشخیص هرزنامه، به طور همزمان از دو روش الگوریتم ژنتیک و یادگیری جمعی، به عنوان روش‌هایی موفق در عرصه یادگیری ماشین، بهره گرفته شده است. در بخش انتخاب ویژگی، یک

1. Filter

2. Wrapper

3. Embedded

4. Ensemble learning

5. Boosting

6. Bootstrap aggregation (bagging)

7. Ensemble

الگوریتم رپر مبتنی بر الگوریتم ژنتیک و در بخش دسته‌بندی، یک سازوکار جمعی مبتنی بر الگوریتم بیزین ارائه می‌شود.

در ادامه، مقاله به این صورت ساختاربندی شده است: ابتدا، برخی روش‌های موجود در زمینه موضوع پژوهش بیان می‌شوند. سپس در بخش بعد، مبحث فیلترینگ بررسی می‌شود که خود دارای دو قسمت است؛ ابتدا در قسمت نخست، خلاصه‌ای از الگوریتم بیزین مطرح می‌شود و سپس در قسمت دوم، به مبحث انتخاب ویژگی به عنوان یکی از مهم‌ترین بخش‌های الگوریتم دسته‌بند پرداخته می‌شود. در ادامه، شرح مختصراً از الگوریتم ژنتیک داده شده و پس از آن، الگوریتم پیشنهادی مطرح می‌شود. در بخش بعدی، نتایج بیان شده و در نهایت، بخش پایانی به نتیجه‌گیری اختصاص دارد.

پیشینه پژوهش

تاکنون پژوهش‌های مختلفی به مطالعه و گروه‌بندی روش‌های دسته‌بندی پرداخته‌اند. میچین و همکاران¹ (۱۹۹۴) در پژوهش خود روش‌های موجود را به سه دسته کلی روش‌های آماری، یادگیری ماشین و شبکه‌های عصبی تقسیم‌بندی کردند. از روش‌های آماری مهم می‌توان بیزین را نام برد که به دلیل سادگی و قدرت بالا، یکی از روش‌های موفق در این زمینه بوده و تاکنون نیز برای بهبود آن تلاش‌های سیاری شده است (Rajaram et al., 2011)، اگرچه در بعضی مقاله‌ها بیزین را به عنوان یک روش یادگیری ماشین نیز در نظر گرفته‌اند (Hua Li et al, 2012). یادگیری ماشین نیز از روش‌های موفق دیگر در زمینه دسته‌بندی متون است. تاکنون الگوریتم‌های یادگیری ماشینی متعددی ارائه شده است که از مهم‌ترین آنها می‌توان درخت تصمیم و k نزدیک‌ترین همسایه² (Crawford et al, 2004) را نام برد. اساس کار این گونه الگوریتم‌ها، یادگیری یک وظیفه خاص بر اساس تعدادی نمونه آموزشی است. شبکه‌های عصبی نیز شاخه‌ای از هوش مصنوعی هستند که از عملکرد مغز انسان الهام گرفته‌اند (Wang et al, 2006). تاکنون، برای مسئله دسته‌بندی، نسخه‌های مختلفی از این الگوریتم ارائه شده است (Xu et al, 2010 & Hua Li et al, 2012). علاوه بر این روش‌ها، روش‌هایی نیز وجود دارند که برای دسته‌بندی هرزنامه از ترکیب چندین روش استفاده می‌کنند (Mohammad et al, 2011& Su et al, 2010).

مبحث دیگری که به مبحث دسته‌بندی بسیار مرتبط است و حتی در بسیاری از پژوهش‌ها از آن به عنوان پیش‌پردازش نام بردۀ‌اند، انتخاب ویژگی است. رویکرد رپر، سازوکار انتخاب ویژگی را با بهره‌گیری از الگوریتم دسته‌بندی انجام داده و زیرمجموعه‌ای از ویژگی‌ها را انتخاب می‌کند. روش رپر، به خصوص برای حل مسائلی مفید است که در آنها نمی‌توان تابع ارزندگی را به راحتی با یک معادله ریاضی دقیق بیان کرد. این روش به دلیل اجرای ساده، توجه بسیاری از پژوهش‌ها را به خود جلب کرده است (Kohavi, 1998 & George, 1998). صحت الگوریتم دسته‌بند، به عنوان معیار ارزیابی برای انتخاب ویژگی‌ها استفاده می‌شود. مزیت روش‌های رپر این است که چون یک الگوریتم دسته‌بندی زیرمجموعهٔ نهایی انتخاب شده را

1. Michie et. al

2. K nearest neighbor

تعیین می‌کند، بر کل فرایند انتخاب ویژگی کنترل بیشتری وجود دارد. همچنین استفاده از ظرفیت یادگیری ارائه شده توسط الگوریتم دسته‌بندی می‌تواند دقت سیار بالایی ایجاد کند (Jadhav et al., 2018). فرایند انتخاب ویژگی هنگام محاسبه ارزندگی یک زیرمجموعه از ویژگی، همه زیرمجموعه‌ها در ویژگی‌های موجود را جستجو می‌کند. بهدلیل اینکه جستجوی جامع از نظر محاسباتی سیار گران است، برای انتخاب ویژگی، به روش‌های جستجوی فرالبتکاری، مانند الگوریتم ژنتیک (GA) علاقه نشان داده شده است.

تاکنون، به منظور انتخاب ویژگی در کاربردهای مختلف از جمله تشخیص هرزنامه، روش‌های مختلفی از الگوریتم ژنتیک ارائه شده است. امینی و یونزه^۱ (۲۰۲۱) یک روش انتخاب ویژگی دولایه با استفاده از الگوریتم ژنتیک و شبکه‌الاستیک را ارائه دادند. در لایه نخست، الگوریتم ژنتیک (GA) به عنوان یک رپر برای جستجوی زیرمجموعه بهینه استفاده شده که هدف آن، کاهش ابعاد و خطای پیش‌بینی است. شبکه‌الاستیک (EN) در لایه دوم به روش پیشنهادی اضافه شده است تا پیش‌بینی‌های اشتباہ احتمالی برای بهبود دقت پیش‌بینی از بین برود.

محمدزاده و همکاران^۲ (۲۰۲۱)، برای انتخاب ویژگی در زمینه شناسایی هرزنامه ایمیل بر اساس مفهوم یادگیری مبتنی بر مخالفت^۳ (OBL)، یک الگوریتم جدید از ترکیب الگوریتم‌های بهینه‌سازی وال (WOA) و گردهافشانی گل (FPA) ارائه دادند. روش کار بدین صورت است که ابتدا WOA اجرا می‌شود و هم‌زمان در حین اجرا، جمعیت OBL توسط WOA تغییر می‌یابد و برای افزایش دقت و سرعت هم‌گرایی، از آن به عنوان جمعیت اولیه FPA استفاده می‌شود.

فارسیس و همکاران^۴ (۲۰۱۹)، برای شناسایی هرزنامه و شناسایی مهم‌ترین ویژگی‌ها بر اساس الگوریتم ژنتیک (GA) و شبکه‌های وزنی تصادفی (RWN)^۵، یک سیستم هوشمند ارائه دادند. یک قابلیت شناسایی خودکار نیز در سیستم پیشنهادی تعییه شده است تا مهم‌ترین ویژگی‌ها را طی فرایند شناسایی انتخاب کند. هوانگو همکاران^۶ (۲۰۰۷)، برای انتخاب ویژگی رپ^۷ بر اساس اطلاعات متقابل، یک الگوریتم ژنتیک ترکیبی ارائه داده‌اند که دارای دو مرحله بهینه‌سازی بیرونی و درونی است. مرحله بیرونی به جستجوی سراسری برای بهترین زیرمجموعه ویژگی‌ها می‌پردازد که در آن اطلاعات متقابل بین کلاس‌های پیش‌بینی شده و واقعی به عنوانتابع ارزش الگوریتم ژنتیک عمل می‌کند و بهینه‌سازی داخلی جستجوی محلی را به صورت فیلتر انجام می‌دهد.

با توجه به کارایی سیار مناسب الگوریتم ژنتیک، در پژوهش‌های بسیاری، از آن به منظور بهینه‌سازی فرایندهای مختلفی بهره‌گیری شده است که در پژوهش‌های ذکر شده به برخی از آنها اشاره شد. با توجه به کاربرد فراوان الگوریتم ژنتیک در کاربردهای مختلف و به خصوص در انتخاب ویژگی، در این مقاله قصد

1. Amini & Yunze

2. Mohammadzadeh et. al

3. Opposition-based learning

4. Faris et. al

5. Random Weight Networks

6. Huang et. al

7. Wrapper

داریم برای تشخیص هرزنامه، سازوکار نوینی ارائه دهیم که الگوریتم ژنتیک را برای انتخاب ویژگی خود به کار می‌برد و با الهام از روش‌های یادگیری گروهی، چندین مدل دسته‌بندی مبتنی بر دسته‌بندی بیزین ایجاد کرده و دسته‌بندی نهایی بر اساس رأی گیری بین این مدل‌ها به دست می‌آید.

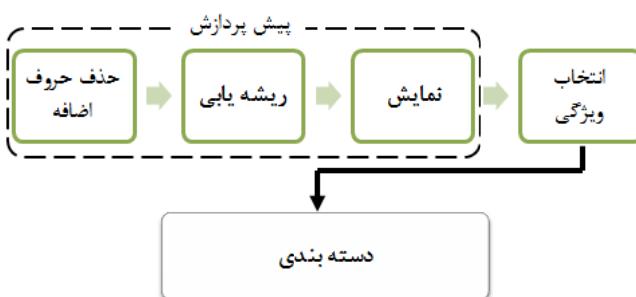
دسته‌بندی به عنوان روشی برای فیلترینگ هرزنامه‌ها

مراحل یک سیستم دسته‌بندی هرزنامه، در شکل ۱ مشاهده می‌شود. همان‌طور که در این شکل دیده می‌شود، عملیات فیلترینگ شامل سه پیش‌پردازش، انتخاب ویژگی و دسته‌بندی است که در ادامه هر یک از این سه فرایند بررسی می‌شوند.

پیش‌پردازش

همان‌طور که در شکل ۱ نشان داده شده است، قبل از انتخاب ویژگی بایستی روی نمونه‌های آموزشی، سه مرحله که به آنها پیش‌پردازش^۱ گفته می‌شود، انجام شود. هدف اصلی پیش‌پردازش، حذف داده‌هایی است که درباره کلاس یک پیام هیچ اطلاعات مفیدی ندارند. همچنین یکی دیگر از اهداف پیش‌پردازش حذف افزونگی^۲ است. این سه مرحله عبارت‌اند از:

- **حذف حروف اضافه بی تأثیر**^۳: شامل لیستی از کلمات است که در فرایند انتخاب ویژگی هیچ نقشی ندارند مانند 'for', 'the', 'as', 'a' و ... و حذف می‌شوند.
- **ریشه‌یابی کلمات**^۴: این مرحله عبارت است از برگرداندن کلمات به ریشه اصلی خود، زیرا عموماً بهتر است که فقط ریشه کلمات استفاده شوند.
- **نمایش داده**: در این مرحله هر پیام ایمیل به صورت برداری از ویژگی‌ها (که در اینجا ویژگی‌ها همان کلمات استخراج شده از متن پیام‌ها هستند)، نمایش داده خواهد شد.



شکل ۱. نمایش مراحل اصلی در سیستم‌های دسته‌بندی هرزنامه

1. Pre-processing
2. Redundancy
3. Stop terms
4. Word stemming

انتخاب ویژگی

پس از اینکه اطلاعات متن‌های اولیه ایمیل‌ها پالایش شده و اطلاعات غیرضروری تا حدودی حذف شدند، در مرحله بعد می‌توان به انتخاب ویژگی‌های برتر اقدام کرد. این کار معمولاً بهوسیله یک سری الگوریتم انتخاب ویژگی^۱ انجام می‌شود. هدف الگوریتم انتخاب ویژگی این است که ویژگی‌هایی را انتخاب کند که دارای اطلاعات بیشتری بوده و بتوانند نرخ خطای دسته‌بندی را تا حد ممکن کاهش دهنده و نرخ پیش‌بینی آن را بهبود بخشدند (Amini & Guiiping, 2021).

انتخاب ویژگی در کلیه مسائل دسته‌بندی^۲ از جمله تشخیص هرزنامه، یکی از موارد مهم به شمار می‌رود. اینکه چه مجموعه‌ای از ویژگی‌ها انتخاب شوند و تعداد این ویژگی‌ها چقدر باشد، یکی از مسائل مهم در مسائل دسته‌بندی به شمار می‌رود. در مسئله تشخیص هرزنامه، ویژگی‌ها همان کلمات یا لغات هستند که دسته‌بندی پیام‌های رسیده اغلب بر اساس آنها انجام می‌شود. بدلیل اینکه ممکن است تعداد ویژگی‌هایی که بتوان بر اساس آنها به دسته‌بندی پیام‌ها اقدام کرد، بسیار زیاد شود، بایستی با استفاده از یک الگوریتم کارآمد، ویژگی‌های برتر را برای عملیات دسته‌بندی انتخاب کرد. تاکنون برای انتخاب ویژگی، الگوریتم‌های مختلفی ارائه شده است و اکثر این الگوریتم‌ها سعی می‌کنند ویژگی‌هایی را انتخاب کنند که قادر باشند بهتر به دسته‌بندی کلاس‌ها اقدام کنند.

یکی از روش‌های ساده انتخاب ویژگی، تکرار سند (DF)^۳ است (Yang et al, 1997) است که عبارت است از تعداد سندهایی که آن ویژگی در آنها تکرار شده است. روش کار به این صورت است که ابتدا میزان DF برای تمام ویژگی‌ها محاسبه می‌شود و در پایان، ویژگی‌هایی که مقدار DF آنها از یک مقدار آستانه کمتر باشد، حذف می‌شوند. دلیل انجام کار نیز این است که ویژگی‌هایی که بهندرت در سندها اتفاق می‌افتد، برای دسته‌بند ما اطلاعات مفیدی ندارند. روش‌های دیگر انتخاب ویژگی، اطلاعات متقابل (MI)^۴ می‌افتد، برای دسته‌بند ما اطلاعات مفیدی ندارند. روش‌های دیگر TF-IDF است که از (Zorkadis et al, 2005) و نفع اطلاعاتی (IG)^۵ (Mitchell, 1996) هستند. روش دیگر TF-IDF است که از (Tکرار ویژگی^۶) ضرب در IDF (معکوس تکرار سند)^۷ به دست می‌آید (Wu, 2009).

در تمامی روش‌های گفته شده سعی شده تا بر اساس پارامتری خاصی به گزینش ویژگی‌های برتر اقدام شود. همین موضوع باعث می‌شود که هرزنامه‌نویس‌ها نیز به شناسایی این ویژگی‌ها پرداخته و بتوانند به راحتی با آن مقابله کنند. به همین دلیل، روش‌های دیگر انتخاب ویژگی وجود دارند که به جای انتخاب ویژگی‌هایی برتر، به صورت تصادفی ویژگی‌ها را انتخاب می‌کنند (DeBarr et al, 2012). این کار باعث افزایش مقاومت^۸ مدل نسبت به نویز می‌شود. برای انجام این کار، برای تبدیل بردار ورودی به بردار خروجی، از یک ماتریس با مقادیر تصادفی استفاده می‌شود. بیشتر مواقع برای تولید اعداد تصادفی از

1. Feature detection algorithms
2. Classification
3. Document frequency
4. Mutual information
5. Information gain
6. Term frequency
7. Inverse document frequency
8. Robust

توزیع گاوی استفاده می‌شود. سپس، جمع تمام مقادیر یک ستون بایستی برابر با یک (نرمالیزه) شود. بردارهای ورودی (همان بردار ویژگی‌ها) در این ماتریس ضرب شده و فضای بردار جدید را نتیجه می‌دهند. این کار در واقع، همان ترکیب خطی ویژگی‌ها است.

فرض کنید n نمونه آموزشی (پیام) وجود دارد که بایستی عملیات انتخاب ویژگی روی آنها اعمال شود. ابتدا پس از پیش‌پردازش، یک بردار از تمامی m ویژگی اولیه به شکل $F = \{f_1, f_2, f_3, \dots, f_m\}$ ساخته می‌شود که در آن \mathbf{z} نشان‌دهنده ویژگی (کلمه) زام است. سپس، هر پیام a به یک بردار منطقی به صورت $V_{ia} = \{v_{i1}, v_{i2}, \dots, v_{im}\}$ تبدیل می‌شود که در آن اگر v_{ij} برابر تعداد تکرار ویژگی \mathbf{z} در پیام a است، بنابراین برای مجموعه تمامی پیام‌های آموزشی، ماتریسی به صورت زیر خواهیم داشت:

$$\text{Input} = \begin{bmatrix} v_{11} & v_{12} & v_{13} & \cdots & v_{1m} \\ \vdots & \ddots & \vdots & & \vdots \\ v_{n1} & v_{n2} & v_{n3} & \cdots & v_{nm} \end{bmatrix}$$

این ماتریس حاوی ویژگی مدد نظر ورودی است. حال، می‌توان با ضرب کردن این ماتریس در یک ماتریس انتقال تعداد این ویژگی‌ها را کاهش داده و عملیات انتخاب ویژگی را انجام داد. ماتریس انتقال به صورت زیر خواهد بود:

$$M_{\text{Transform}} = \begin{bmatrix} t_{11} & t_{12} & t_{13} & \cdots & t_{1s} \\ \vdots & \ddots & \vdots & & \vdots \\ t_{m1} & t_{m2} & t_{m3} & \cdots & t_{ms} \end{bmatrix}$$

همان‌طور که در بالا مشاهده می‌کنید، تعداد سطرهای این ماتریس برابر تعداد ویژگی موجود در ماتریس ورودی (m) است و تعداد ستون‌ها (s) نیز برابر تعداد ویژگی‌های کاهش‌یافته پایانی است که مقدار آن توسط کاربر تعیین می‌شود. در واقع، با تنظیم مقدار s با هر عدد دلخواه می‌توان طول بردار نهایی ویژگی را تعیین کرد. این عملیات ضرب را می‌توان در زیر مشاهده کرد:

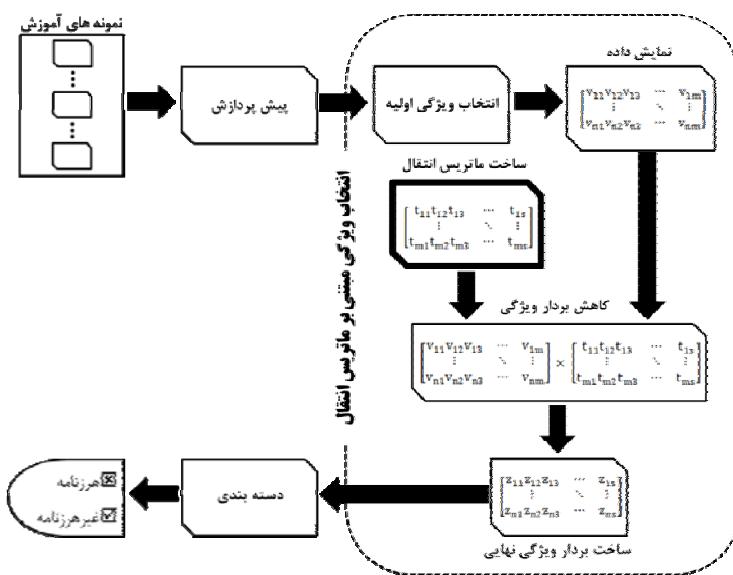
$$\begin{bmatrix} v_{11} & v_{12} & v_{13} & \cdots & v_{1m} \\ \vdots & \ddots & \vdots & & \vdots \\ v_{n1} & v_{n2} & v_{n3} & \cdots & v_{nm} \end{bmatrix} \times \begin{bmatrix} t_{11} & t_{12} & t_{13} & \cdots & t_{1s} \\ \vdots & \ddots & \vdots & & \vdots \\ t_{m1} & t_{m2} & t_{m3} & \cdots & t_{ms} \end{bmatrix} = \begin{bmatrix} z_{11} & z_{12} & z_{13} & \cdots & z_{1s} \\ \vdots & \ddots & \vdots & & \vdots \\ z_{n1} & z_{n2} & z_{n3} & \cdots & z_{ns} \end{bmatrix}$$

این موضوع باعث می‌شود که هر ویژگی نهایی از ترکیب خطی ویژگی‌های اولیه به دست آید. بنابراین، ماتریس انتقال در روند انتخاب ویژگی‌های قدرتمند و بهبود عملکرد دسته‌بندی نقش بسیار مهمی دارد و تمرکز اصلی پژوهش حاضر بر ساخت این ماتریس است. هدف روش پیشنهادی این است که این ماتریس انتقال را بر اساس الگوریتم ژنتیک و طی چند مرحله آموزش سیستم به نحوی بسازد که مقدار صحت سیستم به حداقل مقدار خود برسد. شکل ۱ فرایند پایه‌ای سیستم دسته‌بندی منفرد مبتنی بر انتخاب ویژگی با استفاده از ماتریس انتقال را نشان می‌دهد که در آن مرحله ساخت ماتریس انتقال با خط چین نشان داده شده است. روند کار به این صورت است که ابتدا عملیات پیش‌پردازش روی نمونه‌های ورودی انجام می‌گیرد. به دلیل اینکه ممکن است تعداد ویژگی‌ها بسیار زیاد باشد و به یک ماتریس انتقال بزرگ و محاسبات زیاد منجر شود، در ابتدا بایستی روی نمونه‌های آموزشی یک انتخاب

ویژگی اولیه انجام شده و تعداد کمتری از ویژگی‌ها انتخاب شوند. سپس، با ساخت یک ماتریس انتقال و ضرب کردن در ویژگی‌های اولیه بردار ویژگی نهایی ساخته می‌شود. در انتهای، عملیات دسته‌بندی بر اساس بردار ویژگی ساخته شده انجام می‌شود.

دسته‌بندی

پس از ساخت بردار ویژگی معادل هر پیام ایمیل، در مرحله بعد کافی است به یک الگوریتم دسته‌بند تحويل داده شود تا روی آنها عملیات دسته‌بندی انجام شود. در این مقاله، برای عملیات دسته‌بندی، از بیزین ساده استفاده شده است. اگرچه تاکنون، برای فیلترینگ هرزنامه‌ها، الگوریتم‌های یادگیری ماشینی فراوانی ارائه و حتی الگوریتم‌های مثل بوستینگ^۱ (Freund et al, 1999) و ماشین بردار پشتیبان^۲ (Drucker et al, 1999) نسبت به بیزین در زمینه دسته‌بندی متون دارای کارایی بالاتری هستند، اما از بیزین به طور خاص در کاربردهای تجاری و متن باز^۳ استفاده بسیاری می‌شود (Metsis et al, 2006). دلیل این موضوع شاید سادگی آن باشد که پیاده‌سازی آن را راحت می‌کند. پیچیدگی اجرایی خطی و دقت به نسبت بالای آن نیز باعث می‌شود که بتوان آن را با سایر الگوریتم‌های یادگیری ماشین مقایسه کرد (Androultsopoulos et al, 2004).



شکل ۲. انتخاب ویژگی مبنی بر ماتریس انتقال در یک سیستم دسته‌بندی منفرد

1. Boosting
2. SVM
3. Open-source

بر اساس تئوری بیزین، احتمال اینکه یک پیام با بردار ویژگی $\vec{x} = \langle x_1, \dots, x_m \rangle$ متعلق به کلاس c باشد، برابر است با (Metsis et al, 2006)

$$p(c|\vec{x}) = \frac{p(c).p(\vec{x}|c)}{p(\vec{x})} \quad \text{رابطه (1)}$$

که در این رابطه p برابر احتمال پیشین^۱ کلاس c ، $p(\vec{x}|c)$ برابر است با احتمال اینکه \vec{x} در کلاس c ظاهر شود و مقدار $p(\vec{x})$ نیز به صورت زیر محاسبه می‌شود:

$$p(c_s).p(\vec{x}|c_s) + p(c_h).p(\vec{x}|c_h) \quad \text{رابطه (2)}$$

که c_s کلاس هرزنامه و c_h کلاس غیرهرزنامه است.
از آنجا که عبارت مخرج برای هر دو کلاس برابر است، بنابراین پیام متعلق به کلاسی است که به‌ازای آن کلاس مقدار $p(\vec{x}|c).p(c)$ ماکریم شود. در واقع، $p(c)$ که برابر احتمال پیش‌بینی کلاس برای هر یک از کلاس‌های موجود است (در اینجا دو کلاس هرزنامه و غیرهرزنامه) و از تقسیم تعداد اعضای این کلاس در مجموعه آموزش به کل مجموعه آموزش به دست می‌آید.

روش بیز به کاررفته در پژوهش حاضر همان نیو بیز^۲ است که با فرض استقلال ویژگی‌ها (در اینجا کلمات) استفاده می‌شود. هرچند احتمال دارد در اینجا کلماتی که در هرزنامه‌ها ظاهر می‌شوند به طور کامل از هم مستقل نباشند، اما دومینگوس و پازانی^۳ (۱۹۹۶) نشان دادند تا زمانی که بتوان به‌طور صحیح احتمالات شرطی کلاس را مرتب کرد، وابستگی بین ویژگی‌ها بر عملکرد دسته‌بند تأثیری نخواهد داشت. اگرچه کلیت روش بیزین یک اصل ثابت است، اما تاکنون برای تعیین مقدار $p(\vec{x}|c)$ روش‌های مختلفی ارائه شده که در اکثر مقالات به این مسئله اشاره‌ای نشده است. اینکه مقدار $p(c)$ را به چه روشهای تعیین کنیم، ممکن است عملکرد الگوریتم دسته‌بند را نیز تحت تأثیر قرار دهد (Metsis et al, 2006). در این مقاله از نسخه‌ای از بیزین با عنوان روش «بیزین چندجمله‌ای با TF»^۴ استفاده می‌شود.

در این روش فرض می‌شود که یک پیام با بردار ویژگی به شکل $\langle x_1, \dots, x_m \rangle$ است که در آن هر x_i نشان‌دهنده میزان تکرار t_i در پیام مد نظر است. سپس هرزنامه بودن آن پیام بر اساس رابطه^۳ تعیین می‌شود:

$$\frac{p(c_s). \prod_{i=1}^m p(t_i|c)^{x_i}}{\sum_{c \in \{c_s, c_h\}} p(c). \prod_{i=1}^m p(t_i|c)^{x_i}} > T \quad \text{رابطه (3)}$$

-
1. Priori probabilities
 2. Naive Bayes
 3. Domingos & Pazzani
 4. Multinomial NB, TF attributes

در عبارت بالا مقدار $p(t|c)$ به صورت زیر محاسبه می‌شود:

$$p(t|c) = \frac{1 + N_{t,c}}{m + N_c} \quad \text{رابطه ۴}$$

که در آن $N_{t,c}$ برابر تعداد رخداد توکن t در پیام‌های مجموعه آموزشی کلاس c و $N_c = \sum_{i=1}^m N_{t_i,c}$ است.

الگوریتم ژنتیک

الگوریتم ژنتیک، یکی از متداول‌ترین جست‌وجویی‌های بر اساس اصول انتخاب طبیعی داروین است. در واقع، این الگوریتم نگاشتی از یک سیر طبیعی است که به صورت مسئله ریاضی فرموله شده و همین الهام گرفتن از طبیعت باعث شده است که در حل مسائل بهینه‌سازی و به خصوص مسئله جست‌وجوی کارایی بالایی داشته باشد. ایده اصلی این الگوریتم این است که خصوصیات موروثی توسط ژن^۱ از نسلی به نسل بعد منتقل می‌شوند. کارایی این الگوریتم در فضاهای جست‌وجوی وسیع مناسب بوده و در مقایسه با سایر الگوریتم‌ها احتمال کمتری دارد که در بهینه محلی گرفتار شود.

روش کار در این الگوریتم به این صورت است که طی چندین دوره^۲ پردازش کروموزوم‌ها و تولید جمعیت‌های پی‌درپی جواب بهینه را به دست می‌آورند. در ابتدا، مجموعه‌ای از راه حل‌های یک مسئله با عنوان یک جمعیت^۳ اولیه تشکیل می‌شود. در واقع، در این مرحله مسئله مدنظر به صورت ریاضی فرموله شده و ساختار هر کروموزوم تعیین می‌شود. در مرحله بعد، کیفیت و ارزش تمام اعضای جمعیت که کروموزوم^۴ نامیده می‌شوند، توسط یکتابع ارزش^۵ محاسبه شده و کروموزوم‌ها بر اساس ارزش به دست آمده مرتب می‌شوند. در مرحله بعد با استی تعدادی از کروموزوم‌ها برای ترکیب و تولید جمعیت بعدی انتخاب شوند. بعد از مشخص شدن کروموزوم‌هایی که باستی با هم ترکیب شوند، عملگر ادغام^۶ و جهش^۷ روی آنها انجام شده و کروموزوم‌های جدید به دست می‌آیند. عملگرهای جهش و ادغام، عملگرهای اصلی الگوریتم ژنتیک به شمار می‌روند. مهم‌ترین عملگر در الگوریتم ژنتیک، عملگر ادغام است که باعث جست‌وجوی کامل فضای جست‌وجو شده و در این فرایند دو کروموزوم برای بالا بردن ارزش خود به تبادل ژن‌های خود اقدام می‌کنند. عملگر جهش نیز می‌تواند از گیر افتادن کروموزوم‌ها در نقطه بهینه محلی جلوگیری به عمل آورد که برای پیاده‌سازی آن روش‌های متفاوتی وجود دارد. با استراتژی نخبه‌گزینی، نسل جدید جایگزین نسل قبل شده و جمعیت جدید را شکل می‌دهد. فرایند Goldberg et al, 1989; Davis et al, 1989) کند (

1. Genome
2. Iterative
3. Population
4. Chromosomes
5. Fitness
6. Crossover
7. Mutation

2012 & Stanimirović, 1991). این روند چندین دوره اجرا شده و برترین کروموزوم‌ها به عنوان جواب نهایی برگردانده می‌شوند.

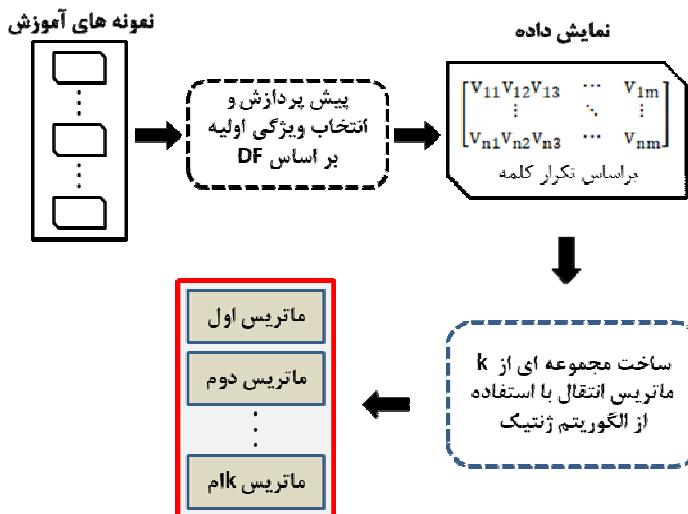
روش پیشنهادی برای انتخاب ویژگی و دسته‌بندی جمعی

روش پیشنهادی در این مقاله به عنوان یک راهکار مبتنی بر یادگیری ماشین با ناظر، دارای دو فاز آموزش (یادگیری) و آزمون (تست) است. در فاز آموزش، الگوریتم بر اساس نمونه‌های آموزشی یاد می‌گیرد که چگونه به دسته‌بندی نمونه‌های موجود اقدام کند و در فاز آزمون، الگوریتم برای دسته‌بندی نمونه‌های آزمایش به کار برده می‌شود. بدین ترتیب، میزان یادگیری الگوریتم امتحان و بررسی شده و کارایی آن مشخص می‌شود. در ادامه، روند هر یک از دو فاز نامبرده در روش پیشنهادی تشریح خواهد شد.

فاز آموزش در روش پیشنهادی

در این فاز، قوانینی که بایستی فرایند یادگیری بر اساس آنها انجام گیرد، ساخته خواهد شد. در شکل ۳ فاز آموزش به صورت کلی نمایش داده است. بر اساس این شکل، فاز آموزش شامل دو مرحله است که عبارت‌اند از:

- پیش‌پردازش و انتخاب ویژگی اولیه
 - ساخت مجموعه‌ای از k ماتریس انتقال با استفاده از الگوریتم ژنتیک
- بر اساس شکل ۳ هر یک از این مراحل دارای یک خروجی مربوط به خود هستند. خروجی مرحله انتخاب ویژگی مجموعه‌ای از بردار ویژگی‌های اولیه در قالب یک ماتریس ورودی است. خروجی مرحله ساخت ماتریس نیز مجموعه‌ای از ماتریس‌های انتقال است. در ادامه، هر یک از این مراحل به تفصیل بیان خواهد شد.



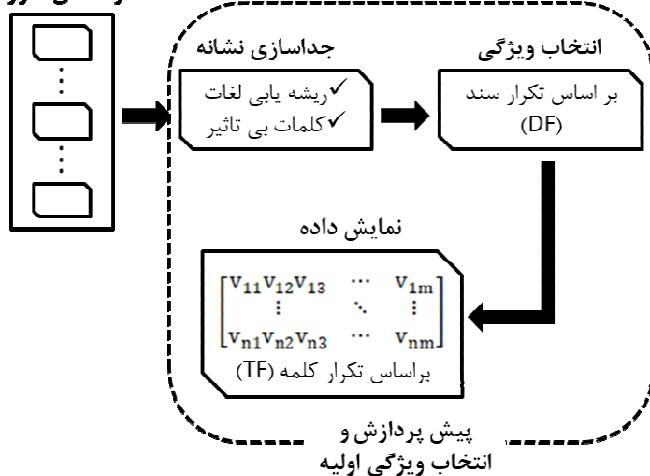
شکل ۳. فاز آموزش در روش پیشنهادی

پیش‌پردازش و انتخاب ویژگی اولیه

روند کلی این مرحله در شکل ۴ نشان داده شده است. همان‌طور که در این شکل مشاهده می‌شود، این روند شامل سه زیرفعالیت است که عبارت‌اند از:

۱. جداسازی نشانه‌ها و استخراج تمام ویژگی‌های موجود: در این مرحله پس از عملیات پیش‌پردازش (که خود شامل ریشه‌یابی لغات و کلمات بی‌تأثیر و غیره است) پایگاه داده اولیه‌ای از تمام ویژگی‌های موجود تشکیل می‌شود.
۲. انتخاب ویژگی اولیه: در این مرحله بر اساس الگوریتم تکرار سند (DF) از بین تمام ویژگی موجود در پایگاه داده اولیه مجموعه‌ای از ویژگی برتر انتخاب می‌شوند.
۳. نمایش داده: در این مرحله تمام نمونه‌های آموزشی (ایمیل‌ها) بر اساس بردار ویژگی خود بیان می‌شوند.

نمونه‌های آموزشی



شکل ۴. عملیات پیش‌پردازش و انتخاب ویژگی اولیه

طی گام‌های مطرح شده، نمونه‌های آموزشی از شکل تعدادی متن به یک ماتریس عددی تبدیل می‌شوند که تعداد سطرهای آن برابر تعداد نمونه‌ها و تعداد ستون‌های آن برابر تعداد ویژگی‌ها خواهد بود.

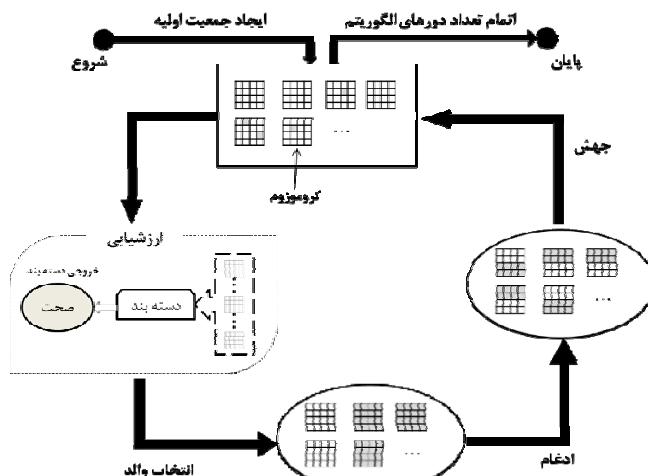
ساخت ماتریس انتقال با استفاده از الگوریتم ژنتیک

ورودی این مرحله، ماتریسی است که در مرحله انتخاب ویژگی اولیه به دست آمده بود. الگوریتم ژنتیک وظیفه دارد که بهینه‌ترین ماتریس انتقال را تولید کند. روند کار الگوریتم ژنتیک در شکل ۵ نشان داده است. چرخه این الگوریتم دارای مراحل مختلفی است و در هر مرحله فرایند خاصی انجام می‌شود و

بدین ترتیب امکان ادامه چرخه مقدور می‌شود. این مراحل عبارت‌اند از: ۱. ایجاد نسل اولیه؛ ۲. تابع ارزش؛ ۳. انتخاب؛ ۴. ادغام؛ ۵. جهش. در ادامه، هر یک از این مرحله‌ها به تفصیل بیان می‌شوند.

تولید جمعیت اولیه کروموزوم‌ها

نخستین مرحله، ایجاد نسل اولیه‌ای از کروموزوم‌ها است. همان‌طور که گفته شد، کروموزوم‌ها در الگوریتم ژنتیک همان راه حل‌های بالقوه مسئله هستند. در اینجا هدف به دست آوردن بهینه‌ترین ماتریس انتقال است، بنابراین کروموزوم‌ها همان ماتریس‌های انتقال هستند. تعداد سطرهای این ماتریس انتقال برابر تعداد ستون‌های ماتریس ویژگی ورودی و تعداد ستون‌های آن توسط کاربر تعیین می‌شود و تعیین‌کننده طول بردار ویژگی نهایی خواهد بود. مقدار المان‌های این ماتریس نیز بر اساس توزیع گوسی^۱ به صورت تصادفی مقداردهی می‌شوند. از آنجا که این اعداد تصادفی هستند، ممکن است شرط لازم برای این ماتریس که بایستی جمع المان‌های هر ستون برابر ۱ باشد رعایت نشود. به همین دلیل، بعد از اختصاص مقادیر تصادفی به المان‌های این ماتریس بایستی روی تمامی ستون‌های این ماتریس، عملیات نرمالیزه کردن را اعمال کرد. عملیات نرمالیزه کردن نیز بسیار ساده است و از تقسیم هر المان بر جمع مقادیر تمامی المان‌های ستونی که در آن قرار گرفته است، به دست می‌آید.



شکل ۵. فرایند ساخت ماتریس انتقال توسط الگوریتم ژنتیک

تابع ارزش

در مرحله بعد، بایستی کروموزوم‌ها بر اساس تابع ارزش مرتب شوند. تابع ارزش یکی از قسمت‌های مهم الگوریتم ژنتیک به شمار می‌رود و کارایی الگوریتم تا حدود زیادی در گروی انتخاب تابع ارزش مناسب

است، بنابراین باید این تابع طوری انتخاب شود که باعث افزایش کارایی الگوریتم شود. در روش پیشنهادی، الگوریتم ژنتیک در تعامل با دسته‌بند سعی در به دست آوردن ماتریسی دارد که بتواند بالاترین بازدهی را در دسته‌بندی پیام‌ها داشته باشد، بنابراین در روش پیشنهادی، الگوریتم ژنتیک محاسبه ارزش هر کروموزوم (ماتریس انتقال) را به الگوریتم دسته‌بند می‌سپارد و این الگوریتم دسته‌بند است که تعیین می‌کند کدام ماتریس برای دسته‌بندی پیام‌ها دارای ارزش بالاتری است. این ویژگی الگوریتم‌های رپر باعث می‌شود تا الگوریتم دسته‌بند نیز در فرایند انتخاب ویژگی تأثیر داشته باشد و بتواند برای خود ویژگی‌های کاراتر را انتخاب کند. تابع ارزش در نظر گرفته شده در الگوریتم ژنتیک در روش پیشنهادی پارامتر صحت است که به صورت زیر تعریف می‌شود (Androultsopoulos et al, 2000):

$$\text{رابطه (۵)} \quad \frac{\text{تعداد کل پیام‌هایی که به درستی دسته‌بندی شده‌اند}}{\text{تعداد کل پیام‌های مجموعه آموزش}} \times 100 = \text{صحت}$$

الگوریتم ژنتیک، هر کروموزوم را به الگوریتم دسته‌بند داده و بر حسب اینکه بر اساس این ماتریس انتقال میزان صحت دسته‌بند چقدر است، به مرتب‌سازی کروموزوم‌ها اقدام می‌کند. بر این اساس، بیشترین ارزش به کروموزومی (ماتریس انتقال) متعلق است که تابع دسته‌بند به‌ازای آن بیشترین مقدار صحت را داشته باشد.

انتخاب (والد)

در این قسمت باید بررسی شود که چه کروموزوم‌هایی بایستی در ساخت نسل بعدی دخیل باشند. در پژوهش حاضر، عملیات انتخاب بر اساس چرخ رولت^۱ و بر پایه تقاطع یکسان انجام می‌شود. همچنین نرخ انتخاب والد (بازترکیبی) برابر با $p_c = ۰/۸$ تنظیم شد.

ادغام (ترکیب)

در اینجا چون کروموزوم‌ها ماتریسی شکل و دو بعدی هستند، برای ترکیب دو کروموزوم نسبت به حالت تک بعدی، راه‌های بیشتری وجود دارد. ترکیب در نظر گرفته شده روش سنتی پخش است که به صورت یگانه پیاده‌سازی شده است. در حالت یگانه هر کروموزوم به دو قسمت شکسته شده و تکه اول کروموزوم اول با تکه دوم کروموزوم دوم و تکه دوم کروموزوم دوم با تکه اول کروموزوم دوم ترکیب می‌شود. برای شکستن هر ماتریس (کروموزوم‌ها) بایستی دو مورد مشخص شود:

- ماتریس به چه شکل شکسته شود. به صورت عرضی، طولی، مورب و ...
- تعیین نقطه‌ای برای شکستن ماتریس

در اینجا ماتریس‌ها به صورت عرضی شکسته می‌شوند و نقطه شکست وسط ماتریس در نظر گرفته شده است که باعث می‌شود اندازه دو قسمت شکسته شده با هم برابر شود. در شکل ۶ می‌توان این

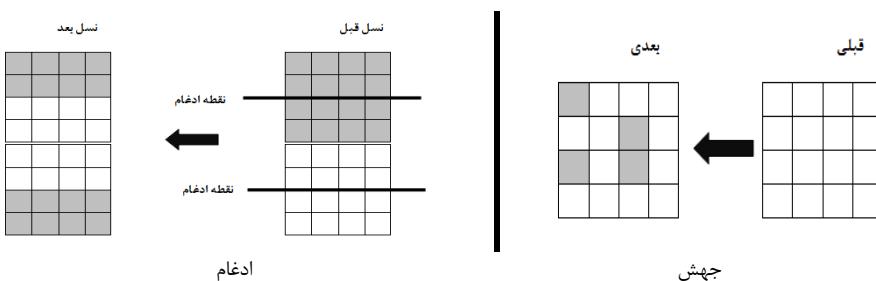
عملیات را به صورت گرافیکی مشاهده کرد. همان‌طور که در شکل مشاهده می‌شود، ماتریس‌ها به صورت عرضی و از وسط به دو نیم شکسته می‌شوند، سپس هر کروموزوم یک نیمة خود را با کروموزوم دیگر تعویض کرده و تشکیل دو کروموزوم جدید را می‌دهند. بعد از تشکیل کروموزوم‌های جدید بایستی شرط نرمال بودن روی آنها بررسی شده و در صورت لزوم به نرمال کردن آنها اقدام شود.

جهش

در انتخاب یک ماتریس به عنوان کروموزوم، المان‌های این ماتریس ژن‌های آن به حساب خواهند آمد. بر این اساس، روش پیاده‌سازی جهش به این شکل است که برای تمام المان‌های ماتریس (کروموزوم) یک عدد تصادفی بین صفر و ۱ تولید شده و در صورتی که عدد تولیدشده کوچک‌تر از مقدار ۰/۰۱ بود مقدار آن عضو تعییر می‌کند. معادله تعییر ژن‌ها به صورت زیر است:

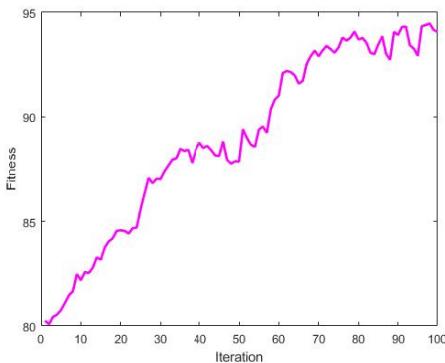
$$Gene [i,j] = (Gene [i,j] + 1.0) / 2.0 \quad (6)$$

برای تعییر ژن‌ها هر عبارت دیگری را نیز می‌توان انتخاب کرد که به نظر کاربر بستگی دارد. پس از انجام عملیات جهش روی اعضاء، بدلیل اینکه ممکن است جمع مقادیر یک ستون برابر ۱ نباشد، بایستی یک بار دیگر آن را نرمال کرد. در پژوهش حاضر، نرخ جهش p_m برابر با $2/0$ مقداردهی شد.



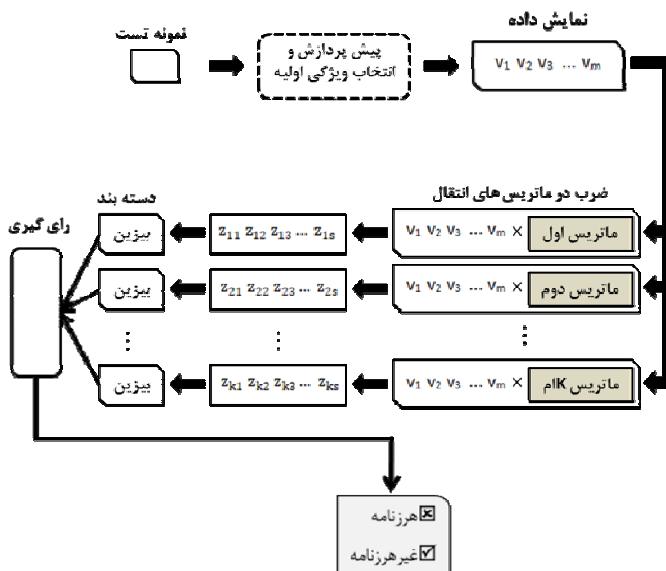
شکل ۶. شماتیک عملگرهای جهش و ادغام ماتریس‌های انتقال در الگوریتم ژنتیک

در پایان اجرای الگوریتم ژنتیک، همانند شکل ۵ مجموعه‌ای از بهترین کروموزوم‌های آخرین نسل (k) برگردانده می‌شوند. این مجموعه k ماتریس حاصل فرایند آموزش سیستم هستند. همان‌طور که در بالا گفته شد، پس از انجام عملگر جهش و ادغام، به نرمال کردن ستون‌های ماتریس انتقال اقدام می‌شود که موجب می‌شود اسکیماهای موجود در پاسخ‌های موفق دستخوش تعییر شده و ممکن است باعث گمراه شدن الگوریتم و عدم هم‌گرایی الگوریتم شود. بدین منظور، در شکل ۷ نمونه‌ای از تعییرات مقدار برازندگی (مقدار تابع ارزش) طی اجرای الگوریتم ژنتیک نشان داده شده است. همان‌طور که در این شکل مشاهده می‌شود، هم‌گام با پیشرفت الگوریتم، مقدار ارزش (Fitness) افزایش یافته و هم‌گرایی الگوریتم در این شکل مشهود است.



شکل ۷. نمونه‌ای از تغییرات مقدارتابع ارزش طی اجرای الگوریتم ژنتیک

فاز آزمون در روش پیشنهادی
کلیت فاز آزمون در روش پیشنهادی در شکل ۸ نمایش داده شده است. در این فاز، الگوریتم بر اساس آنچه در فاز آموزش فراگرفته است به دسته‌بندی پیام‌ها اقدام می‌کند، بنابراین در این مرحله هر سند ورودی (همانند فاز آموزش) به یک بردار ویژگی اولیه تبدیل می‌شود. سپس، به‌طور مجزا در k ماتریس انتقال (به‌دست‌آمده از فاز آموزش) ضرب شده و در نهایت، k بردار ویژگی نهایی را تولید می‌کند. همان‌طور که در شکل ۸ مشهود است، به‌ازای هر یک از این k ماتریس یک عملیات دسته‌بندی بر اساس دسته‌بند بیزین انجام شده و در نهایت، بین دسته‌بندهای مختلف رأی‌گیری شده و کلاسی که بیشترین رأی را داشته باشد به عنوان کلاس نمونه جدید معرفی می‌شود.



شکل ۸. فاز آزمون در روش پیشنهادی

ارزیابی روش پیشنهادی پایگاه داده استفاده شده

پایگاه داده‌های استفاده در این مقاله شامل PU1 و PU2 است¹ که مجموعه‌ای از ایمیل‌های شخصی هستند. به همین دلیل، برای اینکه محتوای آنها شخصی بماند، تمام کلمات با یک شناسه عددی (part1, ..., part10, unused) جایگزین شده‌اند. هر یک از این پایگاه داده‌ها شامل یازده زیرداده‌کتوري (part1, ..., part10, unused) هستند که طی آزمایش‌ها ده‌تای آنها استفاده می‌شوند. ارزیابی دسته‌بندی با یک رویه اعتبارسنجی ده‌گانه معروف اندازه‌گیری و گزارش شده است. در حقیقت در هر آزمایش، مجموعه داده به ده بخش مجزا تقسیم می‌شود و به‌ازای هر بخش الگوریتم یک بار اجرا می‌شود. هر بار یک بخش مختلف به عنوان مجموعه آزمایشی استفاده می‌شود و نه بخش دیگر با یکدیگر گروه‌بندی شده و به عنوان مجموعه آموزشی استفاده می‌شوند. پارامترهای دسته‌بندی (در مجموعه تست) در ده دور اجرا و سپس میانگین‌گیری شده و به عنوان صحت کل مجموعه داده برگردانده می‌شود. PU1 شامل ۱۰۹۹ پیام است که ۴۴ درصد آنها هرزنامه هستند و PU2 نیز شامل ۷۱۰ پیام است که ۲۰ درصد آنها هرزنامه هستند.

معیارهای کارایی

برای ارزیابی دسته‌بندهای هرزنامه معمولاً چندین پارامتر ارزیابی ارائه می‌شود (Mohammad et al. 2011). در پژوهش حاضر نیز برای ارزیابی روش پیشنهادی، از سه پارامتر مهم به صورت زیر استفاده می‌شود.

- معیار صحت

$$\text{Accuracy (\%)} = \frac{\text{TP} + \text{TN}}{\text{TP} + \text{TN} + \text{FP} + \text{FN}} \times 100 \quad (\text{رابطه ۷})$$

- معیار دقت

$$\text{Precision (\%)} = \frac{\text{TP}}{\text{TP} + \text{FP}} \times 100 \quad (\text{رابطه ۸})$$

- معیار بازخوانی

$$\text{Recall(\%)} = \frac{\text{TP}}{\text{TP} + \text{FN}} \times 100 \quad (\text{رابطه ۹})$$

که در آن داریم:

TN (منفی درست): برابر است با تعداد پیام‌هایی که دسته واقعی آنها منفی بوده و الگوریتم دسته‌بندی آنها را به درستی منفی تشخیص داده است.

TP (مثبت درست): برابر است با تعداد پیام‌هایی که دسته واقعی آنها مثبت بوده و الگوریتم دسته‌بندی آنها را به درستی مثبت تشخیص داده است.

FP (مثبت اشتباه): برابر است با تعداد پیام‌هایی که دسته واقعی آنها منفی بوده و الگوریتم دسته‌بندی آنها را به اشتباه مثبت تشخیص داده است.

FN (منفی اشتباه): برابر است با تعداد پیام‌هایی که دسته واقعی آنها مثبت بوده و الگوریتم دسته‌بندی آنها را به اشتباه منفی تشخیص داده است.

ارزیابی

در این بخش، روش پیشنهادی ارزیابی می‌شود. قبل از بیان نتایج به دست آمده، تنظیمات روش پیشنهادی تشریح می‌شوند:

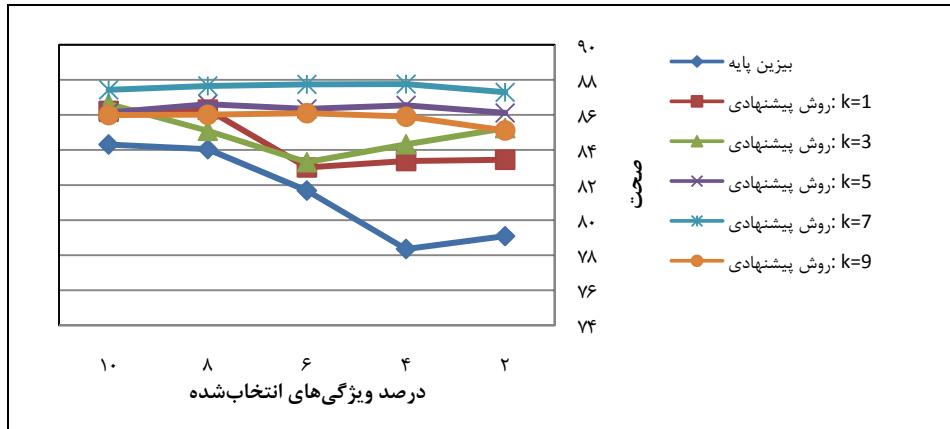
- در مرحله انتخاب ویژگی اولیه فقط ویژگی‌هایی انتخاب شدن که مقدار DF آنها از ۱۰ بیشتر باشد، یعنی در بیشتر از ده سند تکرار شده باشند.
- در تمامی روش‌ها شامل روش پیشنهادی، بیزین پایه و دو روش SVM و KNN، ویژگی‌های اولیه بر اساس DF مرتب شده و از بین آنها ویژگی‌های برتر انتخاب می‌شوند.
- در تمامی روش‌های پیشنهادی s (اندازه ستون ماتریس انتقال) برابر با ۱۰ مقداردهی شده است.
- الگوریتم ژنتیک استفاده شده دارای ۱۰۰ کروموزوم بوده و در ۱۰۰ دور اجرا شد.
- قبل از ارائه نتایج پژوهش، نتایج به دست آمده از پیش‌پردازش داده‌ها بیان می‌شود. تعداد ویژگی‌های (کلمات) هر یک از پایگاه داده‌ها پیش از عملیات پیش‌پردازش و همچنین تعداد ویژگی‌های کاهش‌یافته پس از مرحله پیش‌پردازش، در جدول ۱ نشان داده شده است. همان‌طور که مشاهده می‌شود، تعداد کلماتی که در پیام‌های غیرهرزنامه ظاهر می‌شوند بیشتر از پیام‌هایی است که در هرزنامه‌ها وجود دارند. همچنین پس از اعمال پیش‌پردازش، یک بردار ویژگی یکسان به‌ازای کل پایگاه داده ایجاد می‌شود که اندازه بردار ویژگی پایگاه داده PU1 بزرگ‌تر از PU2 است.

جدول ۱. تعداد ویژگی‌ها در مرحله پیش‌پردازش پایگاه داده‌ها

تعداد ویژگی‌ها		نوع کلاس	پایگاه داده
بعد از پیش‌پردازش	قبل از پیش‌پردازش		
۳۳۹۷	۱۲۲۱۳	هرزنامه	PU1
	۱۶۲۸۶	غیرهرزنامه	
	۲۳۳۷۰	مجموع	
۱۸۶۸	۵۰۷۸	هرزنامه	PU2
	۱۳۱۵۳	غیرهرزنامه	
	۱۵۳۳۸	مجموع	

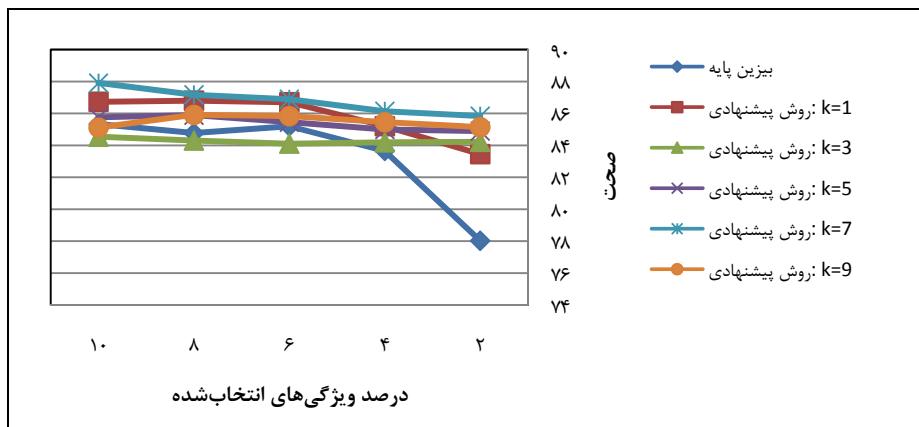
پس از مرحله پیش‌پردازش و ایجاد یک مجموعه اولیه از ویژگی‌ها، در مرحله بعد ویژگی‌های اولیه بر اساس DF مرتب شده و به منظور اعمال روش پیشنهادی و فرایند دسته‌بندی زیرمجموعه‌ای از آنها انتخاب می‌شود. در پژوهش حاضر، در کلیه روش‌ها از ۲ الی ۱۰ درصد از ویژگی‌های برتر انتخاب شده در مرحله پیش‌پردازش استفاده شده است (برای PU1 از بین کل ۳۳۹۷ ویژگی اولیه در مرحله پیش‌پردازش، بهازی ۲ درصد: ۶۹ ویژگی، ۴ درصد: ۱۳۶ ویژگی، ۶ درصد: ۲۰۴ ویژگی، ۸ درصد: ۲۷۲ ویژگی و ۱۰ درصد: ۳۴۰ ویژگی بر اساس معیار DF انتخاب شد، الگوریتم ژنتیک روی آنها اعمال شده و برای PU2 نیز به نحو مشابه اقدام شد).

روش‌های پیشنهادی بهازی ۵ مقدار مختلف k شامل ۱، ۳، ۵، ۷، ۹ ارزیابی شدند. دلایل انتخاب این اعداد این است که مشاهده شد، بهازی مقادیر k بالاتر از ۱۰ روش پیشنهادی دارای کارایی کاهشی است. نتایج مقایسه روش‌های پیشنهادی بر اساس سه پارامتر صحت، دقت و بازخوانی با روش بیزین پایه در دو پایگاه داده PU1 و PU2 در شکل‌های ۹ تا ۱۴ نشان داده شده‌اند.



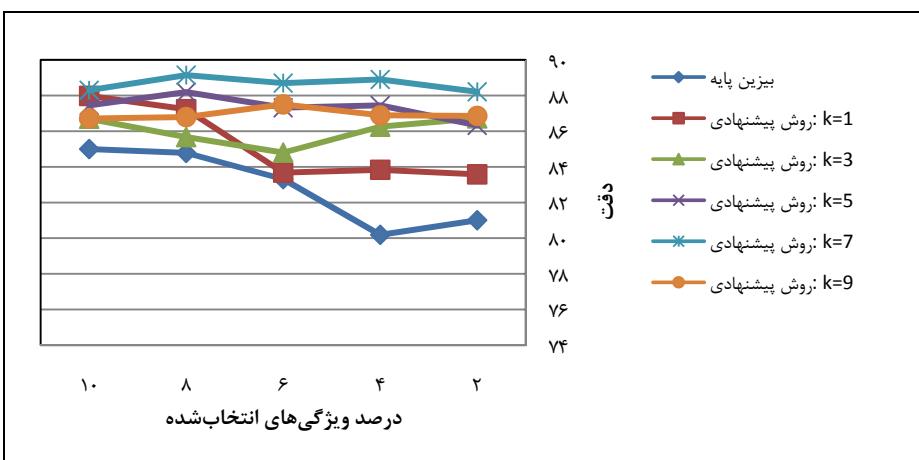
شکل ۹. نتایج پارامتر صحت (/) روی پایگاه داده PU1

نتایج بدستآمده از مقایسه میزان صحت روش‌ها در پایگاه داده‌های PU1 و PU2 در شکل‌های ۹ و ۱۰ نشان داده شده است. نتایج نشان می‌دهد که بر اساس پارامتر صحت، روش پیشنهادی ($k=7$) در هر دو پایگاه داده دارای عملکرد بهتری است. در شکل ۹ مشاهده می‌شود که تمامی روش‌های پیشنهادی از بیزین پایه بهتر عمل کرده‌اند و در شکل ۱۰ مشهود است که روش بیزین پایه فقط از روش پیشنهادی ($k=3$) بهتر بوده است و سایر روش‌ها همگی از بیزین پایه بهتر بوده‌اند. همچنین در روش‌های پیشنهادی با افزایش درصد ویژگی‌های انتخابی میزان صحت تغییر چندانی پیدا نمی‌کند، در حالی که در بیزین پایه با افزایش درصد ویژگی‌ها، عملکرد بهبود می‌یابد.

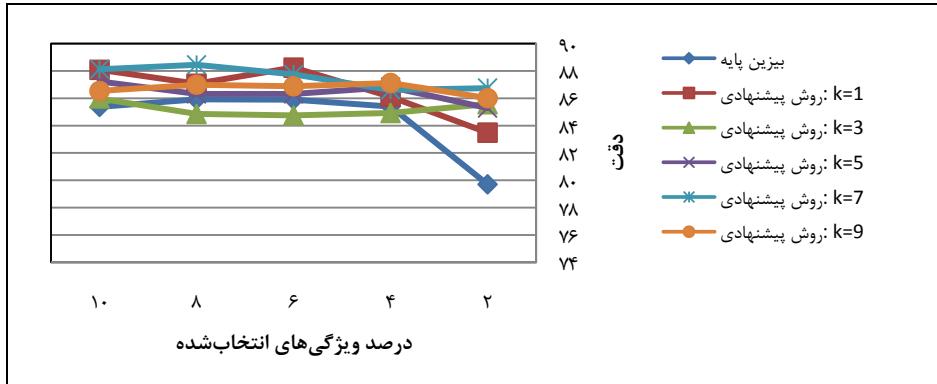


شکل ۱۰. نتایج پارامتر صحت (٪) روی پایگاه داده PU1

نتایج بهدست آمده از مقایسه میزان دقیق روش‌ها در پایگاه داده‌های PU1 و PU2 در شکل‌های ۱۱ و ۱۲ نشان داده شده است. در اینجا نیز روش پیشنهادی ($k = 7$) در بین تمامی روش‌های پیشنهادی بهترین کارایی را داشته است. در شکل ۱۱ مشاهده می‌شود که روش پیشنهادی ($k = 7$) در پایگاه داده PU1 دارای بهترین عملکرد است و بیزین پایه بدترین خروجی را داشته است. در شکل ۱۲ فاصله بین روش‌ها کمتر است، با وجود این، تمامی روش‌های پیشنهادی به غیر از ($k = 9$) از بیزین پایه بهتر عمل کرده‌اند.

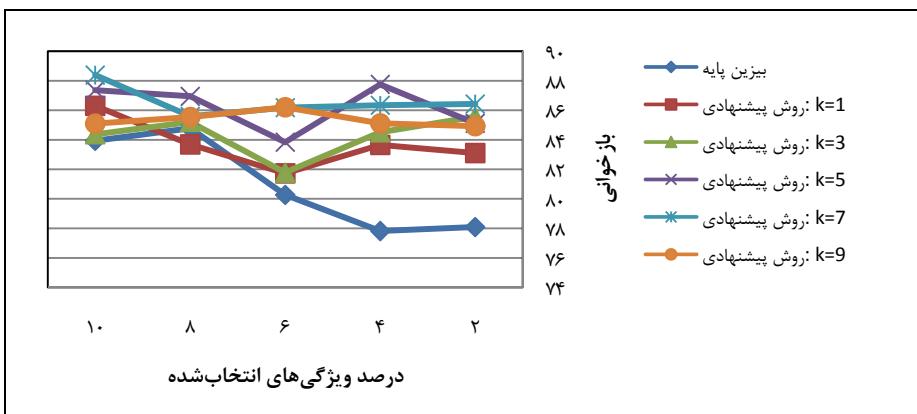


شکل ۱۱. نتایج پارامتر دقیق (٪) روی پایگاه داده PU1



شکل ۱۲. نتایج پارامتر دقت (%) روی پایگاه داده PU2

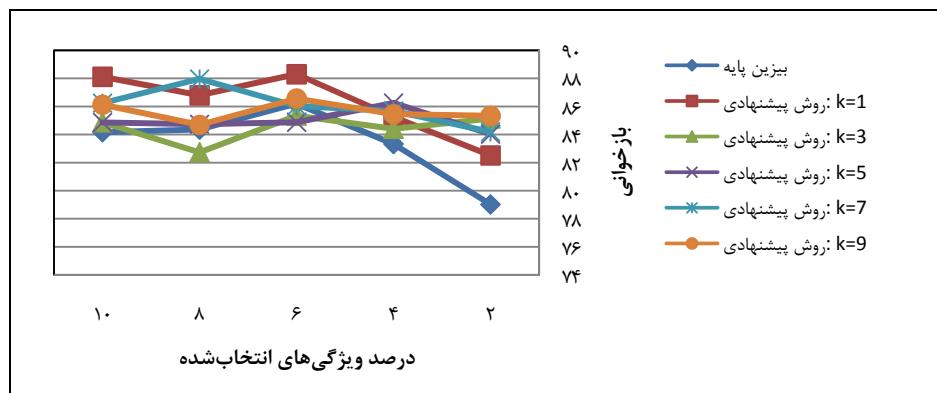
نتایج بدست آمده از مقایسه میزان بازخوانی روش‌ها در پایگاه داده‌های PU1 و PU2 در شکل‌های ۱۳ و ۱۴ نشان داده شده است. اگرچه در این دو شکل روش‌های پیشنهادی به طور کلی از بیزین پایه بهتر عمل کرده‌اند، اما برخلاف پارامترهای قبلی، هیچ یک از روش‌های پیشنهادی به صورت مطلق در مقایسه با سایر روش‌ها، دارای عملکرد بهتری نبوده‌اند.



شکل ۱۳. نتایج پارامتر بازخوانی (%) روی پایگاه داده PU1

پس از مقایسه روش‌های پیشنهادی با بیزین پایه و اثبات عملکرد بهتر آنها نسبت به روش بیزین پایه، در ادامه، عملکرد روش پیشنهادی با دو دسته‌بند پایه مهم و پرکاربرد SVM و KNN مقایسه می‌شود. برای الگوریتم SVM نسخه‌های مختلفی وجود دارد که در اینجا از SVM خطی استفاده شده است. در هر دو دسته‌بند پایه از انتخاب ویزگی پایه DF مشابه با روش‌های پیشنهادی استفاده شده است. نتایج بخش‌های قبل نشان داد که نتایج هر سه پارامتر صحت و دقت و بازخوانی روش‌های پیشنهادی

تقریباً دارای رفتار یکسانی است و روش پیشنهادی با مقدار $k=7$ در مقایسه با سایر روش‌ها دارای عملکرد بهتری است، بنابراین در این بخش به مقایسه پارامتر صحت روش پیشنهادی ($k=7$) با دو دسته‌بند مهم پرداخته می‌شود که نتایج در جدول ۲ نشان داده شده است.



شکل ۱۴. نتایج پارامتر بازخوانی (/) روی پایگاه داده PU۲

همان‌طور که در جدول ۲ مشاهده می‌شود، روش پیشنهادی در مقایسه با هر دو دسته‌بند دارای عملکرد بهتری است. اگرچه دسته‌بند SVM در مقایسه با KNN دارای نتایج بهتری است، اما هنوز هم در مقایسه با روش پیشنهادی دارای عملکرد پایین‌تری است.

جدول ۲. مقایسه پارامتر صحت روش پیشنهادی ($k=7$) با دو دسته‌بند پایه

درصد ویژگی‌های انتخابی					روش	پایگاه داده
۱۰	۸	۶	۴	۲		
۸۷/۴۳	۸۷/۶۵	۸۷/۷۴	۸۷/۷۶	۸۷/۳۹	روش پیشنهادی ($k=7$)	PU۱
۸۶/۸۹	۸۶/۶۶	۸۶/۲۳	۸۶/۴۴	۸۶/۰۱	SVM	
۸۴/۳۴	۸۴/۰۵	۸۴/۱۲	۸۴/۶۳	۸۴/۲۷	KNN	
۸۷/۹۱	۸۷/۱۷	۸۶/۸۹	۸۶/۱۴	۸۵/۸۳	روش پیشنهادی ($k=7$)	PU۲
۸۶/۳۴	۸۶/۱۱	۸۵/۹۴	۸۵/۸۳	۸۵/۶۹	SVM	
۸۶/۵۵	۸۴/۳۸	۸۴/۶۹	۸۴/۱۷	۸۳/۷۴	KNN	

بحث

در بخش قبل، پس از تجزیه و تحلیل روش پیشنهادی، عملکرد آن روی زیرمجموعه‌های مختلف ویژگی^۲ (۱۰ درصد) بررسی شد. مطالعه پیشینهٔ پژوهش نشان می‌دهد که برخی پژوهشگران روی پایگاه داده‌های پژوهش حاضر، روش‌هایی ارائه داده‌اند که نتایج خود را بر اساس درصدهای مختلف ویژگی‌های انتخاب شده گزارش کرده‌اند. رزم آرا و همکاران^۱ (۲۰۱۲) با استفاده از تمامی ویژگی‌ها (۱۰۰ درصد کل ویژگی‌ها) توانستند به نرخ صحت ۹۶/۲۶ و ۹۴/۲۳ در دو پایگاه داده یادشده دست یابند. مقایسهٔ پژوهش حاضر با این روش‌ها نشان می‌دهد که روش پیشنهادی با توجه به استفاده از درصد کمی از ویژگی‌ها (۱۰ درصد) توانسته است نتایج قابل قبولی کسب کند.

همچنین بر اساس نتایج پژوهش، در روش پیشنهادی، برخلاف روش بیزین ساده، با کاهش تعداد ویژگی‌ها کارایی افزایش می‌یابد که می‌تواند به عنوان یک مزیت برای الگوریتم دسته‌بند به شمار رود و می‌توان با استفاده از تعداد ویژگی‌های کمتر در مقایسه با بیزین پایه، به کارایی بالاتری دست یافت. دلیل کارآمدی روش پیشنهادی با تعداد ویژگی کم این است که انتخاب درصد کمتری از ویژگی‌ها باعث کاهش ابعاد ماتریس انتقال و در پی آن، کاهش فضای حل مسئله الگوریتم ژنتیک می‌شود. بدیهی است هرچقدر فضای جستجو کوچک‌تر باشد، احتمال پیدا کردن جواب بهینه افزایش می‌یابد و همین موضوع باعث می‌شود که با کاهش تعداد ویژگی‌های ورودی، پاسخ‌های به دست آمده دارای کارایی بالاتر بوده و در نهایت، باعث افزایش عملکرد دسته‌بندی شوند.

استفاده از یک ترکیب خطی از ویژگی‌ها در روش پیشنهادی، باعث می‌شود که هرزنامه‌نویس‌ها نتوانند به راحتی ویژگی‌ها را شناسایی کرده و تغییر دهنند. در واقع، در روش‌هایی که در آنها به دلیل حجم بسیار زیاد ویژگی‌ها از زیرمجموعه‌ای از ویژگی‌ها استفاده می‌شود، هرزنامه‌نویس‌ها ممکن است به شناسایی ویژگی‌های استفاده شده اقدام کرده و با تغییر یا حذف آنها سبب گمراهی و اشتباه سیستم دسته‌بند شوند. همچنین، استفاده از مجموعه‌ای از راه حل‌های برتر موجب استحکام و اعتماد پذیری بیشتر مدل دسته‌بندی می‌شود.

روش پیشنهادی به ازای مقادیر مختلف k بررسی، تجزیه و تحلیل شد. نتایج نشان داد که افزایش مقادیر k تا هر مقدار دلخواه به معنای افزایش کارایی دسته‌بندی نیست. به این معنا که با انتخاب تعداد بسیار زیاد ماتریس‌های انتقال برتر به عنوان راه حل‌های حاصل از آخرین گام الگوریتم ژنتیک، نمی‌توان به طور افزایشی کارایی را بهبود بخشید. بر اساس نتایج، افزایش تعداد ماتریس‌ها تا مقدار ۷ می‌تواند تأثیر مثبتی روی نتایج داشته باشد و پس از آن دیگر تغییرات مثبتی رخ نخواهد داد و حتی شاهد تأثیر منفی آن نیز خواهیم بود. دلیل این موضوع تا حدودی منطقی است. در الگوریتم ژنتیک راه حل‌ها بر اساس برآزندگی مرتب شده و به طور کلی چند پاسخ برتر الگوریتم ژنتیک راه حل‌های تکاملی حاوی راه حل بهینه بوده و راه حل‌های بعدی دارای برآزندگی کاهشی هستند. بنابراین، انتخاب تعداد مناسب راه حل‌های بهینه مرحله آخر، مسئله‌ای مهم است که در پژوهش حاضر مقدار آن برابر با ۷ تعیین شد.

نتیجه‌گیری

در این مقاله، برای شناسایی و تشخیص هرزنامه‌ها، روش نوینی ارائه شد. هدف روش پیشنهادی در این مقاله این بود که با ارائه یک روش انتخاب ویژگی رپر مبتنی بر دو روش مهم و پرکاربرد شامل الگوریتم ژنتیک و دسته‌بند بیزین، در گام نخست ویژگی‌هایی که به نحو بهتری قادر به تفکیک پیام‌های هرزنامه هستند را استخراج کرده و با توجه به عملکرد پذیرفته‌شده دسته‌بندهای جمعی نسبت به دسته‌بندهای منفرد، در گام بعد با اعمال یک روش دسته‌بندی جمعی بیزین به دسته‌بندی آنها اقدام کند.

کارایی روش پیشنهادی روی دو پایگاه داده PU1 و PU2 بررسی شد. نتایج نشان داد که روش‌های پیشنهادی با مقادیر مختلف k در مقایسه با بیزین پایه، عملکرد بهتری دارند. همچنین بر اساس نتایج، روش پیشنهادی با مقدار $k = 7$ در بین روش‌های پیشنهادی دارای بهترین کارایی بود. همچنین نتایج مقایسه روش پیشنهادی ($k = 7$) با دو دسته‌بند پرکاربرد KNN و SVM حاکی از برتری روش پیشنهادی بود. به عنوان کارهای آینده می‌توان سایر الگوریتم پردازش تکاملی را در ترکیب با سایر الگوریتم‌های دسته‌بندی همچون درخت تصمیم و روش‌های دسته‌بند جمعی به منظور انتخاب ویژگی به کار برد. همچنین می‌توان با تغییرات تابع ارزش الگوریتم ژنتیک و با انتخاب معیارهای بهینه‌تر از معیار صحت برای تابع ارزش، سعی در بهبود کارایی الگوریتم کرد.

فهرست منابع

- Amini, F., and Hu, G. (2021). A two-layer feature selection method using genetic algorithm and elastic net. *Expert Systems with Applications*, 166: 114072.
- Androutsopoulos, I., Koutsias, J., Chandrinou, K.V. and Spyropoulos, S.D. (2000). An experimental comparison of naive Bayesian and keyword-based anti-spam filtering with personal e-mail messages. In *Proceedings of the 23rd annual international ACM SIGIR conference on Research and development in information retrieval*, pp. 160-167.
- Androutsopoulos, I., Palioras, G. and Michelakis, E. (2004). *Learning to filter unsolicited commercial e-mail*. NCSR “Demokritos” Technical Report, No. 2004/2.
- Balamurugan, A.A., Rajaram, R., Pramala, S., Rajalakshmi, S., Jeyendran, C. and Surya Prakash, J.D. (2011). NB+: an improved naive Bayesian algorithm. *Knowledge-Based Systems*, 24(5), 563-569.
- Chinavle, D., Kolari, P., Oates, T. and Finin, T. (2009). Ensembles in adversarial classification for spam. In *Proceedings of the 18th ACM conference on Information and knowledge management*, pp. 2015-2018.
- Crawford, E., Koprinska, I. and Patrick, J. (2004). Phrases and Feature Selection in E-Mail Classification. Conference: *ADCS 2004, Proceedings of the Ninth Australasian Document Computing Symposium*, December 13.

- Dada, E. G., Bassi, G.S., Chiroma, H., Adetunmbi, A.O. and Ajibuwu, O.E. (2019). Machine learning for email spam filtering: review, approaches and open research problems. *Heliyon*, 5(6), e01802.
- Davis, L. (1991). *Handbook of genetic algorithms*. New York: Van Nostrand Reinhold.
- DeBarr, D. and Wechsler, H. (2012). Spam detection using random boost. *Pattern Recognition Letters*, 33(10), 1237-1244.
- Domingos, P., and Pazzani, M. (1996). Beyond independence: Conditions for the optimality of the simple bayesian classifier. In *Proc. 13th Intl. Conf. Machine Learning*. pp. 105-112.
- Drucker, H., Wu, D. and Vapnik, V.N. (1999). Support vector machines for spam categorization. *IEEE Transactions on Neural networks*, 10(5), 1048-1054.
- Faris, H., Al-Zoubi, A.M., Heidari, A.A., Aljarah, I., Mafarja, M., Hassonah, M.A. and Fujita, H. (2019). An intelligent system for spam detection and identification of the most relevant features based on evolutionary random weight networks. *Information Fusion* 48: 67-83.
- Freund, Y., Schapire, R. (1999). A short introduction to boosting. *Journal-Japanese Society for Artificial Intelligence*, 14(771-780), 1612.
- Goldberg, D. E. (1989). Genetic algorithms in search. *Optimization, and Machine Learning*. Reading, MA: Addison-Wesley.
- Hu, Y., Guo, C., Ngai, E. W. T., Liu, M. and Chen, Sh. (2010). A scalable intelligent non-content-based spam-filtering framework. *Expert systems with applications*, 37(12), 8557-8565.
- Huang, J., Cai, Y. and Xu, X. (2007). A hybrid genetic algorithm for feature selection wrapper based on mutual information. *Pattern recognition letters*, 28 (13), 1825-1844.
- Jadhav, S., He, H. and Jenkins, K. (2018). Information gain directed genetic algorithm wrapper feature selection for credit rating. *Applied Soft Computing*, 69: 541-553.
- Kohavi, R., and John, G.H. (1998). The wrapper approach. In: Liu H., Motoda H. (eds) *Feature Extraction, Construction and Selection*, pp. 33-50. The Springer International Series in Engineering and Computer Science, vol 453. Springer, Boston, MA.
- Kolcz, A., Chowdhury, A. and Alspector, J. (2004). The impact of feature selection on signature-driven spam detection. In *Proceedings of the 1st Conference on Email and Anti-Spam (CEAS-2004)*.
- Li, Ch.H. and Huang, J.X. (2012). Spam filtering using semantic similarity approach and adaptive BPNN. *Neurocomputing*, 92: 88-97.
- Metsis, V., Androutsopoulos, I. and Palioras, G. (2006). Spam filtering with naive bayes-which naive bayes?. *CEAS 2006 - The Third Conference on Email and Anti-Spam*, July 27-28, 2006, Mountain View, California, USA.
- Michie, D., Spiegelhalter, D.J. and Taylor, C.C. (1994). *Machine learning, neural and statistical classification*. *Technometrics*, 37(4). DOI:10.2307/1269742

- Mitchell, T. M. (1997). *Machine Learning*, McGraw-Hill Higher Education. New York.
- Mohammad, A.H., and Abu Zitar, R. (2011). Application of genetic optimized artificial immune system and neural networks in spam detection. *Applied Soft Computing*, 11(4), 3827-3845.
- Mohammadzadeh, H. and Soleimanian Gharehchopogh, F. (2021). A novel hybrid whale optimization algorithm with flower pollination algorithm for feature selection: Case study Email spam detection. *Computational Intelligence*, 37(1), 176– 209.
- Nadjate, S., Adi, K. and Allili, M.S. (2020). A semantic-based classification approach for an enhanced spam detection. *Computers & Security*, 94: 101716.
- Razmara, M., Asadi, B., Narouei, M. and Ahmadi, M. (2012). A novel approach toward spam detection based on iterative patterns. In 2012 2nd International eConference on Computer and Knowledge Engineering (ICCKE): pp 318-323
- Stanimirović, Z. (2012). A genetic algorithm approach for the capacitated single allocation p-hub median problem. *Computing and Informatics*, 29(1), 117-132.
- Su, M.C., Lo, H.H. and Hsu, F.H. (2010). A neural tree and its application to spam e-mail detection. *Expert Systems with Applications*, 37(12), 7976-7985.
- Wang, B., Jones, G. JF and Pan, W. (2006). Using online linear classifiers to filter spam emails. *Pattern analysis and applications*, 9(4), 339-351.
- Wu, C.H. (2009). Behavior-based spam detection using a hybrid method of rule-based techniques and neural networks. *Expert Systems with Applications*, 36(3), 4321-4330.
- Xu, H., and Yu, B. (2010). Automatic thesaurus construction for spam filtering using revised back propagation neural network. *Expert Systems with Applications*, 37(1), 18-23.
- Yang, Y., and Pedersen, J.O. (1997). A comparative study on feature selection in text categorization. *Proceedings of the Fourteenth International Conference on Machine*, 97(412-420), p. 35.
- Zorkadis, V., Karras, D.A. & Panayotou, M. (2005). Efficient information theoretic strategies for classifier combination, feature extraction and performance evaluation in improving false positives and false negatives for spam e-mail filtering. *Neural Networks*, 18(5), 799-807.

Ensemble Bayesian Classification Using Genetic Algorithm Wrapper Feature Selection in Spam Detection

Vahid Nosrati¹

Ph.D. Candidate, Computer Department, Faculty of Engineering, Arak University, Arak, Iran

Mohsen Rahmani

Associate Prof., Computer Department, Faculty of Engineering, Arak University, Arak, Iran

Abstract

The role of email in communication is seriously threatened by a phenomenon called spam. So far, many methods have been proposed to deal with this phenomenon, one of the most important of which is to classify emails based on their content into two categories: spam and non-spam. Content-based classification mechanisms use the words as features, where applying an efficient feature selection mechanism is critical due to the large number of features. Therefore, the main focus of this paper is to select useful features via proposing a wrapper feature selection approach based on a powerful genetic algorithm. We then apply a Bayesian classifier, which has demonstrated a high efficiency in text classification. The main steps of the proposed method is as follows: first, an initial feature vector is chosen, then it is optimized by multiplying the vector in a matrix called the transformation matrix made by the genetic algorithm, and finally, a set of k feature vectors is generated. An ensemble classification approach composed of k Bayesian classifiers is applied to the feature vectors, and the ultimate class label is determined by voting among ensemble members. The proposed method is implemented on two datasets PU1 and PU2. The results show that the classification accuracy of the proposed method with $k=7$ reaches 87.86 and 87.91 in PU1 and PU2, respectively. The results also indicate the efficiency of the proposed method compared to naïve Bayes and two well-known classifiers SVM and KNN.

Keywords: Email, Spam, Classification, Genetics Algorithm, Feature Selection, Transformation Matrix, Ensemble Learning

1. Corresponding Author: vh_nosratty@yahoo.com