

تاریخ دریافت: ۱۳۹۷/۰۸/۱۳

تاریخ پذیرش: ۱۳۹۸/۰۱/۰۸

رویکرد ترکیبی نوین برای تشخیص هرزنامه با استفاده از الگوریتم‌های کلونی مورچه و کرم شب‌تاب

ناصر کریم‌پور

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

فرهاد سلیمانان قره چیق*

استادیار، گروه مهندسی کامپیوتر، واحد ارومیه، دانشگاه آزاد اسلامی، ارومیه، ایران
پست الکترونیکی: bonab.farhad@gmail.com

چکیده:

شب‌تاب برای طبقه‌بندی ایمیل هرزنامه پیشنهاد می‌شود. در مدل پیشنهادی از بهینه‌سازی کلونی مورچه به منظور انتخاب ویژگی و از الگوریتم کرم شب‌تاب برای آموزش و آزمایش نمونه‌ها استفاده شده است. نتایج نشان می‌دهد که درصد صحت مدل پیشنهادی بر روی مجموعه داده Spambase با ۲۰۰ بار تکرار و انتخاب همه ویژگی‌ها برابر ۹۲،۰۵ است و همچنین درصد صحت مدل پیشنهادی در مقایسه با بهینه‌سازی اجتماع ذرات-الگوریتم انتخاب منفی، الگوریتم انتخاب منفی، ترکیب تکاملی تفاضلی-الگوریتم انتخاب منفی، درخت تصادفی، شبکه بیزین، ماشین بردار پشتیبان و الگوریتم بیزین ساده بیشتر است.

واژه‌های کلیدی: تشخیص ایمیل هرزنامه، طبقه‌بندی، بهینه‌سازی کلونی مورچه، الگوریتم کرم شب‌تاب

۱. مقدمه

اگر دنیای دیروز را دنیای فناوری بدانیم، دنیای امروز را می‌توان دنیای ابزارها فرض کرد. گسترش فناوری و کاهش محدودیت دسترسی به اطلاعات، عملاً باعث گردیده

با گسترش اینترنت و شبکه‌های کامپیوتری، دسترسی به اطلاعات سریع‌تر انجام می‌شود. اما از طرفی دیگر، دنیای اینترنت با وجود نفوذگرها و هرزنامه نویسان محیط امنی برای کاربران نیست. هر لحظه ممکن است خرابکاری‌هایی از طرف مهاجمان اینترنتی رخ بدهد. میزان آسیب‌ها و نفوذهای گزارش شده به سیستم‌های کامپیوتری در سرتاسر جهان روز به روز در حال افزایش است. با گذشت زمان ابزارها و روش‌های نفوذ به شبکه‌های کامپیوتری ساده و ساده‌تر می‌شوند و نفوذگرها با حداقل ابزارها، مجال نفوذ را می‌یابند. چند سالی است که هرزنامه در صندوق ایمیل‌های کاربران خیلی زیاد شده است و حتی در برخی مواقع کاربران را مجبور به ترک استفاده از ایمیل می‌کند. هرزنامه یکی از جنبه‌های منفی و آزار دهنده استفاده از خدمات ایمیل و در بسیاری از مواقع یکی از راه‌های آلوده شدن سیستم‌های کامپیوتری به انواع ویروس است. در این مقاله، مدل ترکیبی بر مبنای بهینه‌سازی کلونی مورچه و الگوریتم کرم

* نویسنده مسئول

است که سطح فناوری در فرهنگ‌ها و جوامع مختلف، مشابه و یا حداقل قابل مقایسه باشد. آنچه که امروزه می‌تواند رشد و توسعه یا توقف و پیشرفت یک فرهنگ و جامعه را مشخص نماید، ابزارسازی و ابزارشناسی در آن جامعه است. ایمیل در میان ابزارهای مورد استفاده در جامعه امروزی جایگاه ویژه‌ای دارد که ابزاری بسیار ارزان قیمت با کمترین هزینه اطلاع رسانی، تبلیغات و روابط عمومی می‌باشد.

در دهه‌های اخیر به دلایل استفاده ارزان و راحت و عدم اتلاف زمان ارسال پیام از طریق پست الکترونیکی، بنابر نیازهای جوامع مختلف از جمله اهداف علمی، تبلیغاتی به شدت افزایش پیدا کرده است که این امر باعث پر شدن صندوق‌های پست الکترونیکی کاربران اینترنت از ایمیل‌های هرز گردیده است [۱]. هرزنامه به معنای پیام یا نامه الکترونیکی است که بدون درخواست گیرنده و به صورت ناخواسته در حجم انبوه برای کاربران بی‌شماری فرستاده می‌شود. امروزه هرزنامه تبدیل به یک معضل برای کاربران و شرکت‌های حوزه فناوری اطلاعات و سازمان‌های اجتماعی، اقتصادی و دولتی شده است، زیرا در این زمینه مقدار بزرگی از پهنای باند به هدر می‌رود و در سیستم فرستادن ایمیل سرریز رخ می‌دهد [۲]. همچنین از نظر فنی ارسال هرزنامه تقریباً بدون هزینه است و این مسئله باعث شده است که شرکت‌های بازاریابی به سمت آن حرکت کنند.

در برخی از مواقع هرزنامه‌ها حاوی اطلاعاتی مخرب برای کاربران می‌باشند که با اهداف خاصی طراحی می‌شوند و می‌توانند شامل ویروس یا نرم‌افزارهای مخرب باشد که منجر به خرابی در کامپیوترها و شبکه‌ها و سرقت اطلاعات شخصی، سرقت داده، دزدی هویت کاربران و دستاوردهای فکری دیگران می‌گردد [۳]. در بیشتر مواقع، فضایی از صندوق ایمیل کاربران را اشغال می‌کنند و از کاربر می‌خواهند که ایمیل را به بقیه افراد ارسال نماید و یا وبگاه شخصی را ملاقات کنند که در جامعه اینترنتی چنین

ایمیل‌هایی، نمونه بارز هرزنامه می‌باشند. هرزنامه وقت و انرژی زیادی از کاربران را هدر داده و می‌تواند سرچشمه کلاهبرداری و سرقت باشد و همچنین تاثیرات نامطلوب فرهنگی نیز دارند. در صورتی که هر کاربر اینترنت دارای یک نشانی پست الکترونیکی عمومی باشد، ممکن است صدها پیام ناخواسته را با ظاهری کاملاً متفاوت و نامعقول و غیرمنطقی دریافت نماید [۴].

بیشتر تولیدکننده‌های هرزنامه سال‌های زیادی برای پنهان کردن منشأ تولید هرزنامه فعالیت می‌کنند و معمولاً هرزنامه ایمیل‌ها را از طریق کارسازهای پراکسی ناامن ارسال می‌کنند. هرزنامه‌نویسان مکرراً از نام‌ها، نشانی‌ها، شماره تلفن‌ها و دیگر اطلاعات جعلی برای پرداخت حساب‌ها استفاده می‌کنند که این کار به آن‌ها اجازه می‌دهد سریعاً از حسابی به حساب دیگر در صورت برملا شدن حرکت کنند [۵]. هدف از ارسال هرزنامه می‌تواند جنبه‌های مختلفی داشته باشد که عبارتند از: اطلاع رسانی، تبلیغاتی، ایجاد ترافیک در شبکه و دستیابی به رمزهای کاربران. شاید این سؤال به وجود آید که آیا همه پیام‌هایی که در ایمیل هرزنامه قرار می‌گیرند مخرب یا ویروسی هستند؟ می‌توان گفت خیر، همه پیام‌ها در زمره پیام‌های مخرب و یا ویروسی نیستند، بلکه بعضی از آن‌ها می‌توانند مخرب و یا ویروسی باشد. بنابراین، راه‌هایی برای مقابله و جلوگیری از قرار نگرفتن پیام‌ها در هرزنامه وجود دارد. اما این راهکارها هم نمی‌تواند به طور کامل کارگشا باشند و جلوی همه آن‌ها را بگیرد، بلکه می‌توان تعداد پیام‌های هرزنامه را کاهش داد. ما در این مقاله می‌خواهیم ایمیل هرزنامه را با استفاده از مدل ترکیبی الگوریتم بهینه‌سازی کلونی مورچه [۶] و الگوریتم کرم شب‌تاب [۷] که هر دو از الگوریتم‌های فراابتکاری می‌باشند، تشخیص دهیم. در این مقاله مجموعه داده‌های Spambase [۸] با ۶۰۱ نمونه و ۵۷ ویژگی را مورد ارزیابی قرار می‌دهیم. مجموعه Spambase شامل دو رده هرزنامه با ۱۸۱۳ نمونه (۳۹٫۴٪) و غیرهرزنامه با ۲۷۸۸ نمونه (۶۰٫۶٪) است.

الگوریتم بهینه‌سازی کلونی مورچه [۶]، یکی از جدیدترین و کارآمدترین الگوریتم‌های تکاملی می‌باشد که از زندگی مورچه‌ها الهام گرفته است. الگوریتم بهینه‌سازی کلونی مورچه برای اولین بار توسط دوریگو و همکاران به عنوان یک راه حل چند عامله برای حل مسائل مشکل بهینه‌سازی مثل فروشنده دوره‌گرد ارائه شد. مورچه‌ها هنگام راه رفتن از خود ردی از ماده شیمیایی فرمون به جای می‌گذارند، این ماده قدرت تبخیر بالایی دارد ولی در کوتاه مدت به عنوان رد مورچه بر سطح زمین باقی می‌ماند. یک رفتار پایه‌ای ساده در مورچه‌های وجود دارد. آن‌ها هنگام انتخاب بین دو مسیر به صورت احتمالاتی مسیری را انتخاب می‌کنند که فرمون بیشتری داشته باشد یا به عبارت دیگر مورچه‌های بیشتری قبلاً از آن عبور کرده باشند. الگوریتم کرم شب‌تاب [۷]، یکی از الگوریتم‌های الهام گرفته شده از طبیعت است. کرم‌های شب‌تاب، درک پیرامون خود را از طریق حسگرها انجام می‌دهند و مبتنی بر هوش جمعی می‌باشند.

الگوریتم بهینه‌سازی کلونی مورچه و کرم شب‌تاب به دلیل قابلیت‌هایی مانند تعداد پارامترهای کم، معادلات آسان، فرار از بهینه‌های محلی و قابلیت جستجوی ناحیه اطراف جواب بهینه یافته شده از کیفیت بالایی در انتخاب ویژگی و طبقه‌بندی بهره‌مند هستند. این الگوریتم‌ها می‌توانند احتمال به دام افتادن در کمینه محلی را کاهش دهند و همچنین دارای نرخ همگرایی (کشف راه حل بهینه) بالایی هستند. در تمام الگوریتم‌های فرا ابتکاری از جمله الگوریتم بهینه‌سازی کلونی مورچه و کرم شب‌تاب دو ویژگی ذاتی باید در نظر گرفته شود: توانایی الگوریتم برای جستجوی تمام بخش‌های فضای جستجو (اکتشاف) و توانایی در بهره‌برداری از بهترین راه حل. در الگوریتم بهینه‌سازی کلونی مورچه، توانایی اکتشاف با استفاده از به‌روزرسانی فرمون و توانایی بهره‌برداری با استفاده از احتمال حرکت انجام می‌گیرد. در الگوریتم کرم شب‌تاب، با انتخاب مقدار مناسب برای پارامترهای اولیه، اکتشاف

می‌تواند تضمین شود و حرکت کرم‌ها بر مبنای فاصله می‌تواند توانایی بهره‌برداری را تضمین کند.

ساختار کلی مقاله به شرح زیر سازماندهی شده است: در بخش دوم کارهای قبلی را توضیح می‌دهیم. در بخش سوم، مدل پیشنهادی و مراحل آن را به تفکیک توضیح می‌دهیم. در بخش چهارم، ارزیابی و مقایسه مدل پیشنهادی را توضیح می‌دهیم و نهایتاً در بخش پنجم نتیجه‌گیری و کارهای آینده را توضیح خواهیم داد.

۲. کارهای قبلی

تاکنون تکنیک‌های متعددی برای تشخیص هرزنامه ارائه شده است که شامل تکنیک‌های رگرسیون، داده کاوی، شبکه‌های احتمالاتی، شبکه‌های عصبی مصنوعی، منطق فازی و الگوریتم‌های هوش جمعی می‌باشد که در این بخش برخی از آن‌ها را بررسی و تحلیل می‌کنیم.

مدل RAN-LSH [۹] که ترکیبی از شبکه عصبی مصنوعی تابع شعاعی پایه و ماشین بردار پشتیبان است برای تشخیص ایمیل هرزنامه پیشنهاد شده است. ارزیابی بر روی دو مجموعه داده که در سال ۲۰۱۳ گردآوری شده‌اند و هر کدام شامل ۸۹۴۸ (۱۳۱۱) ایمیل هرزنامه و ۷۶۳۷ (غیرهرزنامه) و ۲۰۹۰۵ (۸۸۶۳) و ۱۲۰۴ ایمیل غیرهرزنامه) نمونه هرزنامه هستند انجام شده است. مدل RAN-LSH از تابع درهم‌ساز به منظور جستجو و ذخیره نمونه‌ها استفاده می‌کند. در مدل RAN-LSH [۹]، با استفاده از جدول درهم‌ساز یک گروه از نمونه‌های مشابه به یک سطر خاص در جدول درهم‌ساز اختصاص داده می‌شوند که می‌توانند با سرعت با استفاده از مقدار نمایه بازیابی شوند. همچنین افزودن داده‌های جدید در جدول درهم‌ساز به زمان کمی نیاز دارد. لذا هدف از تابع درهم‌ساز ایجاد یک جدول برای نمونه‌های هرزنامه و غیرهرزنامه است. جدول درهم‌ساز شامل وزن نمونه‌ها و تعداد نمونه‌ها است و ارزیابی آن بر مبنای معیارهای دقت و بازخوانی انجام

1- Resource Allocating Network with Locality Sensitive Hashing (RAN-LSH)

شده است. نتایج نشان داده که مدل RAN-LSH در پیش‌بینی دقت مناسبی دارد. به دلیل این‌که مجموعه داده‌ها شامل ویژگی‌های زیادی هستند از انتخاب ویژگی استفاده شده است. در مدل RAN-LSH برای انتخاب ویژگی از ماشین بردار پشتیبان و برای آموزش و آزمایش از شبکه عصبی مصنوعی تابع شعاعی پایه استفاده شده است. همچنین برای انتخاب ویژگی‌های براننده از لایه وسطی مدل RAN-LSH استفاده می‌شود.

مدل شبکه عصبی مصنوعی نوع روش گروهی مدیریت داده‌ها [۱۰] برای طبقه‌بندی هرزنانه پیشنهاد شده است. شبکه عصبی نوع روش گروهی داده‌گردانی^۲ یکی از پرکاربردترین شبکه‌های عصبی مصنوعی است که از توانایی بالایی در مدل‌سازی داده‌های پیچیده برخوردار است. ارزیابی مدل روش گروهی داده‌گردانی بر روی مجموعه داده جهانی و معتبر Spambase انجام شده است. مجموعه داده Spambase شامل ۵۷ ویژگی و ۴۶۰۱ نمونه است. نمونه‌ها در دو رده هرزنانه و غیرهرزنانه طبقه‌بندی شده‌اند. نتایج با انتخاب ویژگی‌های مختلف نشان داده که درصد صحت در مدل روش گروهی داده‌گردانی در مقایسه با مدل‌های پرسپترون چندلایه و بیزین ساده بیشتر است. بیشترین درصد صحت در مدل روش گروهی داده‌گردانی برابر ۹۲٫۴ و در مدل‌های پرسپترون چندلایه و بیزین ساده به ترتیب برابر ۹۱٫۷ و ۷۵٫۴ است.

مدل ترکیبی بهینه‌سازی اجتماع ذرات-الگوریتم انتخاب منفی به منظور طبقه‌بندی هرزنانه پیشنهاد شده است [۱۱]. ارزیابی بر روی مجموعه داده Spambase انجام شده است. در مدل ترکیبی بهینه‌سازی اجتماع ذرات-الگوریتم انتخاب منفی از بهینه‌سازی اجتماع ذرات برای جستجوی مقدار ویژگی‌ها در فضای مسئله و از مدل الگوریتم انتخاب منفی برای انتخاب ویژگی استفاده شده است. روش پیشنهادی دارای دو مرحله آموزش و آزمایش است. نتایج نشان داده که مدل ترکیبی بهینه‌سازی اجتماع ذرات-الگوریتم انتخاب منفی در مقایسه با مدل‌های بیزین ساده، ماشین

بردار پشتیبان-انتخاب ویژگی متمایز و الگوریتم انتخاب منفی دقت تشخیص بالاتر و در مقایسه با مدل ماشین بردار پشتیبان دقت تشخیص کمتری دارد. درصد صحت در مدل بهینه‌سازی اجتماع ذرات-الگوریتم انتخاب منفی برابر ۸۳٫۲۰ و در مدل‌های بیزین ساده، ماشین بردار پشتیبان-انتخاب ویژگی متمایز و الگوریتم انتخاب منفی به ترتیب برابر ۷۸٫۸، ۷۱ و ۶۸٫۸۶ است.

مدل MOEA^۳ [۱۲] با الهام از الگوریتم‌های تکاملی برای تشخیص هرزنانه بر روی دو مجموعه داده که شامل ۲۷۲ و ۴۲۶ ایمیل هستند اجرا و آزمایش شده است. در این مدل از قوانین برای تشخیص هرزنانه استفاده شده است. نتایج نشان داده که دقت تشخیص MOEA در مقایسه با NSGA^۴ و SOOA^۵ بیشتر است و همچنین خطای پیش‌بینی در مدل MOEA کمتر است. مدل QUANT^۵ [۱۳] که ترکیبی از شبکه عصبی مصنوعی و درخت تصمیم‌گیری است برای تشخیص هرزنانه پیشنهاد شده است. در این مدل داده‌ها با استفاده از شبکه عصبی مصنوعی آموزش داده و آزمایش می‌شوند و با استفاده از درخت C4.5 طبقه‌بندی می‌شوند. از الگوریتم C4.5 جهت تحلیل ویژگی‌های اصلی موثر بر هرزنانه و تشخیص استفاده شده است. در درخت C4.5 هر مسیر از ریشه به سمت یک گره، نمایانگر یک قانون طبقه‌بندی می‌باشد. ارزیابی بر روی دو مجموعه داده Spam-Assassin و Corpus 2006 انجام شده است. نتایج بر روی مجموعه داده Spam-Assassin نشان داده که درصد صحت در مدل QUANT برابر ۸۹٫۱۵ و در مدل‌های بیزین ساده، SMO و درخت C4.5 به ترتیب برابر ۸۱٫۰۸، ۸۸٫۶۲ و ۷۳٫۰۸ است. و همچنین بر روی مجموعه داده Corpus درصد صحت مدل QUANT برابر ۹۰٫۸۷ و در مدل‌های بیزین ساده، SMO و C4.5 به ترتیب برابر ۸۸٫۱۵، ۸۹٫۷۹ و ۸۸٫۱۵ است.

مدل‌های شبکه عصبی مصنوعی و درخت تصمیم‌گیری بر روی ۱۲۰۰ نوع ایمیل متفاوت آزمایش و اجرا شده‌اند

3- Multi-Objective Evolutionary Algorithms (MOEA)

4- Single Objective Optimization Algorithm (SOOA)

5- Quadratic-Neuron-based Neural Tree (QUANT)

2- Group Method of Data Handling (GMDH)

[۱۴]. در مدل درخت تصمیم‌گیری از قوانین کشف و در مدل شبکه عصبی مصنوعی از آموزش و آزمایش داده‌ها برای تشخیص استفاده شده است. در لایه‌های مختلف شبکه عصبی مصنوعی از معیار وزندهی برای انتخاب ویژگی استفاده شده است. در مدل درخت تصمیم‌گیری انتخاب یک ویژگی براننده به عنوان ریشه درخت به دقت تشخیص کمک می‌کند. نتایج نشان داده که مدل شبکه عصبی مصنوعی در مقایسه با درخت تصمیم‌گیری دقت بهتری دارد. الگوریتم سیستم ایمنی مصنوعی [۱۵] که یکی از الگوریتم‌های تکاملی است برای تشخیص ایمیل هرزنامه بر روی ۴۵۴۱۹ ایمیل از مجموعه TREC07 اجرا و آزمایش شده است. از مدل الگوریتم سیستم ایمنی مصنوعی برای انتخاب ویژگی و حذف ویژگی‌های نامرتبط استفاده شده است. ارزیابی با ۴ اجرا نشان داده که دقت تشخیص در الگوریتم سیستم ایمنی مصنوعی متفاوت است و دقت به حدود ۹۹ درصد نزدیک شده است.

مدل ترکیبی «شبکه عصبی مصنوعی چندلایه-بهینه‌سازی اجتماع ذرات» به منظور تشخیص هرزنامه پیشنهاد شده است [۱۶]. از الگوریتم بهینه‌سازی اجتماع ذرات برای انتخاب ویژگی و از مدل شبکه عصبی مصنوعی چندلایه برای آموزش و آزمایش داده‌ها و طبقه‌بندی استفاده شده است. در مدل شبکه عصبی مصنوعی چندلایه-بهینه‌سازی اجتماع ذرات از شبکه عصبی پرسپترون با تابع فعال‌سازی سیگموئید برای لایه پنهان و ۸۰ درصد داده‌ها برای آموزش و ۲۰ درصد برای آزمایش استفاده شده است. ارزیابی بر روی مجموعه داده Ling-Spam با ۴۸۱ هرزنامه و ۲۱۷۱ غیرهرزنامه و مجموعه داده Spam-Assassin با ۶۰۰۰ نمونه ایمیل انجام شده است. ارزیابی بر روی مجموعه داده Spam-Assassin و Ling-Spam نشان داده که درصد صحت در مدل شبکه عصبی مصنوعی چندلایه-بهینه‌سازی اجتماع ذرات به ترتیب برابر ۹۹٫۹۸ و ۹۹٫۷۹ است. مقایسه‌ها نشان داده که مدل شبکه عصبی مصنوعی

چندلایه-بهینه‌سازی اجتماع ذرات در مقایسه با مدل‌های ماشین بردار پشتیبان با تابع کرنل، ماشین بردار پشتیبان با تابع RBF و شبکه عصبی مصنوعی چندلایه دقت تشخیص بهتری دارد.

مدل‌های درخت تصمیم‌گیری، ماشین بردار پشتیبان و شبکه عصبی مصنوعی چندلایه و ترکیب آن‌ها بر روی دو مجموعه داده با ۱۴ ویژگی آزمایش و اجرا شده‌اند [۱۷]. مجموعه داده اولی شامل ۵۰۴ ایمیل (۳۳۶ ایمیل و ۱۶۸ هرزنامه) و مجموعه داده دومی شامل ۶۵۷ ایمیل (۳۸۷ ایمیل و ۲۷۰ هرزنامه) است. در مدل درخت تصمیم‌گیری از آنتروپی، مدل ماشین بردار پشتیبان از تابع کرنل و شبکه عصبی مصنوعی چندلایه از خطای میانگین استفاده شده است. نتایج بر روی مجموعه داده اولی نشان داده که درصد صحت در مدل ترکیبی برابر ۹۱٫۰۷ و در مدل‌های درخت تصمیم‌گیری، ماشین بردار پشتیبان و شبکه عصبی مصنوعی چندلایه به ترتیب برابر ۸۹٫۸۸، ۸۸٫۶۹ و ۸۹٫۸۸ است. و بر روی مجموعه داده دومی درصد صحت در مدل ترکیبی برابر ۹۱٫۷۸ و در مدل‌های درخت تصمیم‌گیری، ماشین بردار پشتیبان و شبکه عصبی مصنوعی چندلایه به ترتیب برابر ۹۰٫۸۷، ۹۰٫۸۷ و ۸۹٫۰۴ است.

ادریس و همکاران [۱۸] مدل ترکیبی تکاملی تفاضلی-الگوریتم انتخاب منفی را برای تشخیص ایمیل هرزنامه پیشنهاد کرده‌اند. در این مدل از الگوریتم تکامل تفاضلی برای جستجوی ویژگی‌ها در فضای مسئله استفاده شده است. الگوریتم تکامل تفاضلی یکی از جدیدترین روش‌های جستجو است. الگوریتم تکامل تفاضلی به عنوان روشی قدرتمند و سریع برای مسائل بهینه‌سازی در فضاهای پیوسته معرفی شده است. ارزیابی مدل تکاملی تفاضلی-الگوریتم انتخاب منفی بر روی مجموعه داده Corpus با ۱۸۱۲ هرزنامه و ۲۷۸۸ ایمیل غیرهرزنامه انجام شده است. نتایج نشان داده که درصد صحت در الگوریتم انتخاب منفی و تکاملی تفاضلی-الگوریتم انتخاب منفی به ترتیب برابر ۶۸٫۸۶ و ۸۰٫۶۶ است. و مقدار ضریب همبستگی در

الگوریتم انتخاب منفی و تکاملی تفاضلی-الگوریتم انتخاب منفی به ترتیب برابر ۳۶,۰۰۶ و ۶۰,۰۰۸ است.

مدل بهبود داده شده بهینه‌سازی اجتماع ذرات-الگوریتم انتخاب منفی [۱۹] برای تشخیص طبقه‌بندی هرزنامه پیشنهاد شده است. ارزیابی مدل بهینه‌سازی اجتماع ذرات-الگوریتم انتخاب منفی بر روی مجموعه داده Corpus با ۱۸۱۳ ایمیل هرزنامه و ۲۷۸۸ ایمیل غیرهرزنامه انجام شده است. نتایج نشان داده است که مدل بهینه‌سازی اجتماع ذرات-الگوریتم انتخاب منفی در مقایسه با مدل‌های بهینه‌سازی اجتماع ذرات، الگوریتم انتخاب منفی، بیزین ساده، ماشین بردار پشتیبان، ماشین بردار پشتیبان-انتخاب ویژگی متمایز، شبکه عصبی مصنوعی و فازی دقت تشخیص بهتری دارد.

طبقه‌بندی بیزین [۲۰] برای تشخیص هرزنامه بر روی سه مجموعه داده که شامل ۱۰۰۰، ۱۵۰۰ و ۲۰۰۰ ایمیل هستند اجرا و آزمایش شده است. در طبقه‌بندی بیزین از روابط احتمالاتی استفاده می‌شود. نتایج نشان داده دقت تشخیص برای سه مدل بیشتر از ۹۳ درصد است. برای ارزیابی از معیارهای صحت، دقت و بازخوانی استفاده شده است. درصد صحت به ترتیب برای سه مجموعه داده برابر ۹۳,۹۸، ۹۴,۸۵ و ۹۶,۴۶ است.

سیستم تشخیص هرزنامه با استفاده از ماشین بردار پشتیبان پیاده‌سازی شده است [۲۱]. این سیستم برای تشخیص ایمیل‌هایی که روزانه از طرف شبکه‌های اجتماعی ارسال می‌شوند طراحی شده است. به طور میانگین هر روز ۲۸۰۰۰۰ ایمیل با استفاده از ماشین بردار پشتیبان بررسی می‌شود. نتایج نشان داده که مدل ماشین بردار پشتیبان در زمان محاسباتی کم توانسته است حدود ۶۰ درصد ایمیل‌ها را بر مبنای هرزنامه و غیرهرزنامه تشخیص دهد. مدل بهینه‌سازی اجتماع ذرات بر مبنای واحد پردازش گرافیکی برای تشخیص هرزنامه بر روی ۸۴۸۱۰ ایمیل (۴۸۳۶۰ هرزنامه و ۳۶۴۵۰ غیرهرزنامه) آزمایش شده است [۲۲]. به منظور کاهش زمان محاسباتی و بهره‌وری مدل

بهینه‌سازی اجتماع ذرات از مدل واحد پردازش گرافیکی استفاده شده است. برای محاسبه تشخیص کلماتی که به عنوان هرزنامه تلقی می‌شوند از تعداد فراوانی کلمات استفاده شده است. نتایج نشان داده که مدل بهینه‌سازی اجتماع ذرات با به کارگیری واحد پردازش گرافیکی توانسته است دقت تشخیص را در مقایسه با بهینه‌سازی اجتماع ذرات در حالت معمولی افزایش دهد.

جدول (۱)، مزایا و معایب مدل‌های پیشنهاد شده برای تشخیص هرزنامه را نشان می‌دهد.

در میان مدل‌های پیشنهاد شده در جدول (۱) دلیل برتری نسبی شبکه‌های عصبی مصنوعی در طبقه‌بندی به دلیل فرایند یادگیری است. زیرا بر اساس فرایند یادگیری، تمام ارتباطات پیچیده غیرخطی و اثرهای متقابل میان ویژگی‌های مستقل فراگرفته می‌شوند. لذا قدرت طبقه‌بندی این مدل به طور قابل ملاحظه‌ای بالا است. با وجود این، در صورت عدم انتخاب معماری مناسبی برای شبکه عصبی مصنوعی و بخصوص انتخاب تعداد زیادی از لایه‌های پنهان و گره‌ها، ممکن است شبکه، خطای تصادفی موجود در داده‌های تحت یادگیری را نیز فراگیرد (بیش آموزش) و در نتیجه، با وجود دقت بالای شبکه در طبقه‌بندی داده‌های تحت یادگیری، در خصوص داده‌های تحت آزمون به خوبی عمل نکند. مهم‌ترین ضعف شبکه‌های عصبی مصنوعی، عدم تفسیرپذیری وزن‌ها، پیچیدگی و زمان بر بودن مقدار آستانه وزن‌ها است.

۳. مدل پیشنهادی

مدل پیشنهادی که ترکیبی از بهینه‌سازی کلونی مورچه و الگوریتم کرم شب‌تاب است بر روی مجموعه داده Spambase اجرا می‌شود (مدل پیشنهادی برای مجموعه داده Spambase طراحی شده است). مهم‌ترین چالش در مجموعه داده‌های بزرگ، ابعاد زیاد یا به عبارتی تعداد زیاد ویژگی‌ها (صفات) است. علاوه بر این، حذف ویژگی‌های اضافی سبب افزایش کارایی الگوریتم‌ها می‌شود و حتی

جدول ۱: مزایا و معایب مدل‌های پیشنهاد شده برای تشخیص هرزنامه

معايب	مزایا	مدل‌ها	رفرنس
<p>کاهش درصد صحت شبکه عصبی مصنوعی تابع شعاعی پایه با افزایش تعداد گره‌های میانی</p>	<p>انتخاب ویژگی‌های بارز با استفاده از ماشین بردار پشتیبان برای طبقه‌بندی</p> <p>آموزش و آزمایش نمونه‌ها بر مبنای کشف بهترین مقدار برای وزن نرون‌های شبکه عصبی مصنوعی</p>	<p>ترکیب شبکه عصبی مصنوعی تابع شعاعی پایه و ماشین بردار پشتیبان</p>	[۹]
<p>افزایش میانگین مربعات خطا به دلیل بیش آموزش</p>	<p>کشف رابطه بین داده‌های ورودی و خطای خروجی</p> <p>احتمال رسیدن به پاسخ صحیح، حتی اگر بخشی از گره‌های شبکه در فرایند وزن‌دهی مشکل یا عملکرد غلط داشته باشند.</p>	<p>شبکه عصبی مصنوعی از نوع روش گروهی مدیریت داده‌ها</p>	[۱۰]
<p>افزایش زمان محاسباتی</p>	<p>بهبود الگوریتم انتخاب منفی با استفاده از بهینه‌سازی اجتماع ذرات</p> <p>افزایش تنوع جمعیتی و کشف ویژگی‌های مناسب برای طبقه‌بندی</p>	<p>بهینه‌سازی اجتماع ذرات - الگوریتم انتخاب منفی</p>	[۱۱]
<p>افزایش زمان محاسباتی</p>	<p>پیدا کردن جواب‌های بهینه متعدد (سراسری و محلی) مرتبط با تابع برازندگی</p> <p>فضای جستجو به چندین زیر فضا تقسیم می‌شود، در هر زیر فضا عملیات جستجو به منظور کاهش خطای طبقه‌بندی انجام می‌شود.</p>	<p>الگوریتم تکاملی چندگانه</p>	[۱۲]
<p>افزایش پیچیدگی</p>	<p>از میان متغیرهای مستقل، متغیری که دارای بیشترین آنتروپی است برای طبقه‌بندی انتخاب می‌شود.</p>	<p>ترکیب شبکه عصبی مصنوعی چندلایه و درخت تصمیم C4.5</p>	[۱۳]
<p>افزایش عمق درخت تصمیم‌گیری</p>	<p>انتخاب یک ویژگی برازنده در درخت تصمیم‌گیری به عنوان ریشه درخت به دقت طبقه‌بندی کمک می‌کند.</p> <p>حداقل تعداد تکرار از شبکه عصبی مصنوعی به منظور آموزش و تست نمونه‌ها</p>	<p>ترکیب شبکه عصبی مصنوعی چندلایه و درخت تصمیم ID3</p>	[۱۴]
<p>گیر افتادن در بهینه محلی</p>	<p>الگوریتم از مدل ریاضی آسانی برخوردار است.</p> <p>انتخاب ویژگی بر مبنای نزدیکی مقدار ویژگی‌ها</p>	<p>الگوریتم سیستم ایمنی مصنوعی</p>	[۱۵]
-	<p>انعطاف پذیری در برابر مشکل بهینه محلی</p> <p>همگرایی بالا</p>	<p>شبکه عصبی مصنوعی چندلایه - بهینه‌سازی اجتماع ذرات</p>	[۱۶]
<p>حساس به مقدار پارامترهای اولیه</p>	<p>انتخاب ویژگی‌هایی که در طبقه‌بندی موثر هستند و دسته‌ی بیشتری توسط آن‌ها تشکیل می‌گردد.</p>	<p>مدل‌های درخت تصمیم‌گیری، ماشین بردار پشتیبان و شبکه عصبی مصنوعی چندلایه</p>	[۱۷]
<p>افزایش زمان اجرای الگوریتم ترکیبی</p> <p>به دست آوردن دقیق نقاط بهینه محلی و سراسری با احتمال کم</p>	<p>افزایش کارایی الگوریتم انتخاب منفی با استفاده از تکاملی تفاضلی</p>	<p>ترکیب تکاملی تفاضلی - الگوریتم انتخاب منفی</p>	[۱۸]
<p>افزایش زمان محاسباتی</p>	<p>افزایش تنوع جمعیتی و کشف ویژگی‌های مناسب برای طبقه‌بندی</p>	<p>بهینه‌سازی اجتماع ذرات - الگوریتم انتخاب منفی</p>	[۱۹]
<p>کاهش درصد صحت به دلیل انتخاب نمونه‌های مشابه</p>	<p>استفاده از احتمال به منظور کاهش خطای طبقه‌بندی</p>	<p>طبقه‌بندی بیزین</p>	[۲۰]
<p>تابع کرنل به منظور طبقه‌بندی نمونه‌ها باید دقیق باشد.</p>	<p>طبقه‌بندی نمونه‌ها به دو نمونه هرزنامه و غیرهرزنامه با تعداد تکرار کم</p> <p>آموزش نسبتاً ساده است.</p> <p>برای داده‌های با ابعاد بالا تقریباً خوب جواب می‌دهد.</p>	<p>ماشین بردار پشتیبان</p>	[۲۱]
<p>مشکل بهینه محلی دارد</p>	<p>استفاده از تعداد فراوانی کلمات</p> <p>استفاده از وزن‌دهی به منظور یافتن ویژگی‌های دقیق</p>	<p>بهینه‌سازی اجتماع ذرات</p>	[۲۲]

از بهینه‌سازی کلونی مورچه به منظور انتخاب ویژگی و از الگوریتم کرم شب‌تاب برای طبقه‌بندی داده‌ها بر مبنای آموزش و آزمایش نمونه‌ها استفاده می‌شود. در شکل (۱)، شمای کلی مدل پیشنهادی نشان داده شده است.

به دلیل کاهش ویژگی‌های نامرتب می‌توان به نتایج بهتری دست یافت. در مدل پیشنهادی در ابتدا نمونه‌ها از مجموعه داده Spambase خوانده می‌شوند و سپس عملیات نرمال‌سازی انجام می‌گیرد. در مدل پیشنهادی

۳-۱ نرمال‌سازی داده‌ها

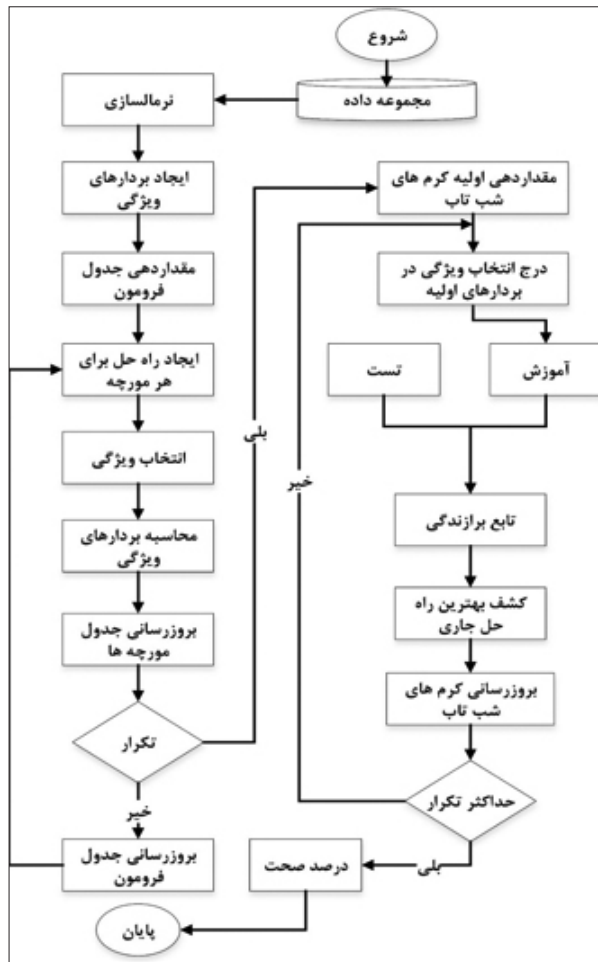
نرمال‌سازی داده‌ها گام مهمی در فرآیند شناسایی ویژگی‌های مهم می‌باشد. مجموعه داده‌ها ممکن است چندان کنترل شده نباشند و مقادیر پرت، تکراری و یا اشتباه می‌تواند منجر به خروجی نامعتبر شوند. وجود داده‌های پرت در اکثر موارد منجر به اختلال در دقت تشخیص خواهد شد. داده پرت معمولاً در فاصله دورتری از مقدار سایر داده‌ها قرار می‌گیرد. وجود داده‌های پرت اساساً به یکی از دلایل زیر است: (۱) غیر صحیح بودن اندازه‌گیری موارد ثبت شده یا وارد شده در رایانه، (۲) گردآوری اندازه‌گیری‌ها از جوامع مختلف. بنابراین انجام مجموعه فعالیت‌هایی قبل از استفاده از داده‌ها جهت تضمین کیفیت لازم است. اگر اطلاعات دارای افزونگی یا غیرمرتبط در فرآیند آموزش باشند، این مراحل بسیار مشکل خواهند شد. نرمال‌سازی داده‌ها طبق معادله (۱) انجام می‌شود [۲۳].

$$x' = (x_{\max} - x_{\min}) \times \frac{(x_i - x_{\min})}{(x_{\max} - x_{\min})} + x_{\min} \quad (1)$$

در معادله (۱)، پارامترهای x_{\max} و x_{\min} مربوط به بیشترین و کمترین مقدار ویژگی‌ها هستند.

۳-۲ انتخاب ویژگی

در بیشتر موارد برای یافتن دانشی که در میان داده‌ها نهفته است، انتخاب همه ویژگی‌ها مهم و حیاتی نیست. به همین دلیل در بسیاری از زمینه‌ها کاهش ابعاد داده یکی از مباحث قابل توجه است. لذا، انتخاب ویژگی برای کاهش فضای ویژگی و در افزایش کارایی طبقه‌بندی موثر است. هدف انتخاب ویژگی، انتخاب زیرمجموعه‌ای از ویژگی‌های مناسب از بین مجموعه ویژگی‌های اولیه، برای بهبود دقت طبقه‌بندی است. در مدل پیشنهادی با استفاده از مورچه‌ها، ویژگی‌ها انتخاب می‌شوند. ویژگی‌های انتخاب شده از هر تکرار با تکرار بعدی مقایسه می‌شوند تا بهترین مسیر فرمون ایجاد شود. هر مورچه دارای یک حافظه است که



شکل ۱: روندنمای مدل پیشنهادی

بهترین راه حل به دست آمده را در آن ذخیره می‌کند. هر مورچه از این حافظه به منظور ساخت فرمون استفاده می‌کند. فرآیند بهینه‌سازی با تعدادی راه حل تصادفی شروع می‌گردد. در هر تکرار از الگوریتم، هر مورچه با نمونه‌برداری از یک توزیع نرمال که ویژگی‌های آن فرمون اصلاح شده توسط مورچه‌ها می‌باشد، یک راه حل جدید ایجاد می‌کند. حافظه مورچه‌ها بر مبنای راه حل‌های ساخته شده در هر تکرار بروزرسانی می‌شود. در صورتی که راه حل جدید بهتر از راه حل موجود در حافظه مورچه باشد جایگزین آن می‌شود. این فاز همان بروزرسانی فرمون است که با هدف حرکت به سمت بهترین راه حل انجام می‌گیرد.

در بهینه‌سازی کلونی مورچه باید فضای جستجو برای انتخاب ویژگی به صورت یک گراف وزن دار بدون

جهت نمایش داده شود. این گراف به صورت $G=(F,E)$ که $F=\{F_1, F_2, \dots, F_n\}$ یک مجموعه از ویژگی‌هایی اصلی و نمایش گره‌های گراف و $E=\{(F_i, F_j): F_i, F_j \in F\}$ یال‌های گراف هستند. وزن هر یال $(F_i, F_j) \in E$ طبق معادله (۲) تعریف می‌شود. همبستگی بین دو بردار F_i و F_j طبق معادله (۲) تعریف می‌شود.

$$w_{ij} = \left| \frac{\sum_p (x_i - \bar{x}_i)(x_j - \bar{x}_j)}{\sqrt{\sum_p (x_i - \bar{x}_i)^2} \sqrt{\sum_p (x_j - \bar{x}_j)^2}} \right| \quad (2)$$

در معادله (۲)، x_i و x_j بردارهای ویژگی و متغیرهای \bar{x}_i و \bar{x}_j میانگین مقادیر بردارهای x_i و x_j و تعداد p نمونه‌ها است. احتمال این‌که مورچه k ام ویژگی بعدی را انتخاب کند طبق معادله (۳) تعریف می‌شود.

$$P_k(i, j) = \begin{cases} \frac{[\tau_i] \cdot [\eta_1(F_j)]^\alpha [\eta_2(F_i, F_j)]^\beta}{\sum_{u \in J_i^k} [\tau_u] [\eta_1(F_u)]^\alpha [\eta_2(F_i, F_u)]^\beta}, & \text{if } j \in J_i^k \\ 0 & \text{otherwise} \end{cases} \quad (3)$$

در معادله (۳)، پارامتر J_i^k مجموعه ویژگی‌های ملاقات نشده، پارامتر τ_u شدت فرمون و پارامترهای α و β به ترتیب تاثیر مقدار فرمون ریخته شده بر روی یال‌ها و کشف ویژگی‌ها هستند. ارتباط مقدار ویژگی‌ها با استفاده از η_1 و η_2 تعیین می‌شود.

۳-۳ طبقه‌بندی داده‌ها

مرحله طبقه‌بندی داده‌ها با استفاده از الگوریتم کرم شبتاب انجام می‌شود. در این مرحله، آموزش و آزمایش داده‌ها انجام می‌گیرد. در الگوریتم کرم شبتاب در ابتدا بردارهای ویژگی‌ها را بر مبنای الگوریتم بهینه‌سازی کلونی مورچه تولید می‌کنیم. سپس در مرحله آموزش بر مبنای نزدیکی و بازه ویژگی‌ها عمل نزدیکی و تشابه بین مقدار ویژگی‌ها انجام می‌شود. فاصله اقلیدسی بین دو کرم شبتاب به صورت معادله (۴) تعریف می‌شود.

$$r_{ij} = \|x_i - x_j\| = \sqrt{\sum_{k=1}^n (x_{i,k} - x_{j,k})^2} \quad (4)$$

در الگوریتم کرم شبتاب برای بروزرسانی بردارها از معادله (۵) استفاده می‌شود. برای این‌که طبقه‌بندی نمونه‌ها با دقت بالایی انجام شود از مقدار قبلی ویژگی‌ها با مقدار پیش‌بینی شده به منظور بروزرسانی بردارها استفاده می‌شود.

$$F = \frac{|x_i - \bar{x}_i|}{x_i} \times \left(rand - \frac{1}{2} \right) \quad (5)$$

از آنجایی که معیار شباهت در تابع هدف بر اساس فاصله تعریف می‌شود می‌توان از تعاریف مختلفی برای استفاده شود که در جدول (۲) چند نمونه از این توابع آورده شده است.

۳-۴ معیارهای ارزیابی

نتایج مدل پیشنهادی، باید در مرحله ارزیابی توسط معیارهای مهم ارزیابی مورد تحلیل قرار گیرد تا بتوان کارایی آن را تعیین نمود. این معیارها را می‌توان هم برای مجموعه داده‌های آموزشی در مرحله یادگیری و هم برای مجموعه رکوردهای آزمایشی در مرحله ارزیابی محاسبه نمود. در این مقاله مهم‌ترین معیارها را برای دقت پیش‌بینی انتخاب کرده‌ایم که درصد صحت مهم‌ترین معیار در دقت پیش‌بینی است [۲۴][۲۵].

$$\text{Precision} = \frac{TP}{TP + FP} \quad (6) \quad \text{دقت}$$

$$\text{Recall} = \frac{TP}{TP + FN} \quad (7) \quad \text{بازخوانی}$$

$$F\text{-Measure} = \frac{2 * \text{Precision} * \text{Recall}}{(\text{Precision} + \text{Recall})} \quad (8) \quad \text{اندازه‌گیری F}$$

$$\text{Accuracy} = \frac{(TP + TN)}{(TP + TN + FP + FN)} \quad (9) \quad \text{صحت}$$

$$\text{ErrorRate} = \frac{(FP + FN)}{(TP + TN + FP + FN)} = 1 - \text{Accuracy} \quad (10) \quad \text{نرخ خطا}$$

پارامتر TP (مثبت درست) نشان دهنده تعداد نمونه‌هایی که جزء دسته مثبت بوده و درست پیش‌بینی شده‌اند. پارامتر FP (مثبت نادرست) نشان دهنده تعداد نمونه‌هایی که نادرست به عنوان دسته مثبت پیش‌بینی شده‌اند. پارامتر

جدول ۲: محاسبه توابع فاصله برای پیدا کردن داده‌های مشابه

تابع فاصله	فرمول
فاصله اقلیدسی (Euclidean distance)	$d(x, y) = \sqrt{\sum_{i=1}^n (x_i - y_i)^2}$
فاصله همینگ (Hamming)	$d(x, y) = \sum_{i=1}^n x_i - y_i $
فاصله چیبشف (chebyshev)	$d(x, y) = \max_{i=1,2,\dots,n} x_i - y_i $
فاصله کنبرا (canberra)	$d(x, y) = \sum_{i=1}^n \frac{ x_i - y_i }{x_i + y_i}, x_i \text{ and } y_i \text{ are positive}$

جدول ۳: ارزیابی مدل‌ها بدون انتخاب ویژگی و با ۱۰۰ بار تکرار

معیارها	بهینه‌سازی کلونی مورچه	الگوریتم کرم شب‌تاب	مدل پیشنهادی
دقت	۶۹,۵۱	۶۷,۱۱	۷۳,۲۵
بازخوانی	۷۱,۳۲	۶۹,۰۷	۷۸,۱۵
اندازه‌گیری F	۷۰,۴۰	۶۸,۰۸	۷۵,۶۲
صحت	۷۷,۸۴	۷۶,۲۵	۸۹,۰۰
نرخ خطا	۰,۲۳۱۶	۰,۲۳۷۵	۰,۱۱۰۰

الگوریتم کرم شب‌تاب برای آموزش و آزمایش نمونه‌ها استفاده شده است.

۴-۲ ارزیابی مدل پیشنهادی بر مبنای انتخاب ویژگی و تعداد تکرار

در جدول (۴)، نتایج ارزیابی مدل پیشنهادی بر مبنای انتخاب ویژگی و با ۱۰۰ بار تکرار و با استفاده از معیارهای مختلف مورد ارزیابی قرار گرفته است.

در جدول (۵)، نتایج ارزیابی مدل پیشنهادی بر مبنای انتخاب ویژگی و با ۲۰۰ بار تکرار و با استفاده از معیارهای مختلف مورد ارزیابی قرار گرفته است.

در جدول (۶)، نتایج ارزیابی معیار صحت، بر مبنای انتخاب ویژگی‌های مختلف با ۱۰۰ بار تکرار نشان داده شده است. در جدول (۶) واضح است که با کاهش انتخاب ویژگی مقدار دقت افزایش یافته است. زیرا در مدل پیشنهادی برای طبقه‌بندی داده‌ها از ویژگی‌های موثر استفاده شده

FN (منفی نادرست) نشان دهنده تعداد نمونه‌هایی که نادرست به عنوان دسته منفی پیش‌بینی شده‌اند. پارامتر TN (منفی درست) نشان دهنده تعداد نمونه‌هایی که جز دسته منفی بوده و درست پیش‌بینی شده‌اند.

۴.۴ ارزیابی و نتایج

ارزیابی و نتایج بر روی مجموعه داده Spambase برگرفته از UCI با ۴۶۰۱ نمونه و ۵۷ ویژگی در محیط VC#.NET 2017 انجام شده است. به دلیل آن که الگوریتم‌های فرا ابتکاری از جستجوی تصادفی پیروی می‌کنند، از نظر سرعت همگرایی و رسیدن به جواب نمی‌توان تنها بر اجرای تکرار یکبار برنامه تکیه کرد. برای ارزیابی مدل پیشنهادی، مقدار پارامترهایی مانند جمعیت اولیه، تعداد تکرار، تعداد نسل و مقدار فرمون به ترتیب برابر ۱۰۰، ۱۰۰، ۱۰۰ و ۰٫۱ است و همچنین مرحله آموزش و آزمایش به ترتیب برابر با ۸۰ و ۲۰ درصد می‌باشند. همچنین طول بردار در الگوریتم کرم شب‌تاب به صورت متغیر تعیین شده است.

۴-۱ ارزیابی مدل‌ها بدون انتخاب ویژگی

در جدول (۳)، ارزیابی مدل‌ها بدون انتخاب ویژگی و با ۱۰۰ بار تکرار انجام شده‌اند. طبقه‌بندی در مدل پیشنهادی بر مبنای آموزش و آزمایش نمونه‌ها انجام شده است. از بهینه‌سازی کلونی مورچه برای انتخاب ویژگی و از

جدول ۴: ارزیابی مدل پیشنهادی بر مبنای انتخاب ویژگی و با ۱۰۰ بار تکرار

معیارها	انتخاب ویژگی							
	۲۰	۲۷	۳۲	۴۰	۴۹	۵۰	۵۳	۵۷
دقت	۹۰,۳۶	۹۲,۸۵	۹۱,۳۵	۹۴,۷۹	۹۵,۱۴	۹۲,۴۳	۹۰,۱۲	۸۹,۳۷
بازخوانی	۹۱,۰۶	۹۳,۶۴	۹۲,۰۸	۹۵,۴۸	۹۵,۶۶	۹۳,۱۰	۹۱,۵۲	۹۰,۳۲
اندازه‌گیری F	۸۹,۱۱	۹۱,۰۸	۸۸,۵۷	۹۲,۱۷	۹۲,۸۲	۸۹,۵۹	۸۹,۶۷	۸۸,۱۷
صحت	۹۲,۳۶	۹۱,۸۶	۹۱,۵۲	۹۰,۷۶	۸۹,۷۱	۸۸,۰۹	۹۰,۱۲	۸۹,۰۰
نرخ خطا	۷,۶۴	۸,۱۴	۸,۴۸	۹,۲۴	۱۰,۲۹	۱۱,۹۱	۹,۸۸	۱۱,۰۰

جدول ۵: ارزیابی مدل پیشنهادی بر مبنای انتخاب ویژگی و با ۲۰۰ بار تکرار

معیارها	انتخاب ویژگی							
	۲۰	۲۷	۳۲	۴۰	۴۹	۵۰	۵۳	۵۷
دقت	۹۵,۰۲	۹۴,۳۶	۹۳,۰۷	۹۶,۰۸	۹۵,۱۴	۹۲,۴۳	۹۰,۱۲	۸۹,۳۷
بازخوانی	۹۶,۱۶	۹۵,۹۷	۹۳,۹۱	۹۷,۹۶	۹۶,۹۴	۹۳,۴۵	۹۱,۰۱	۹۱,۹۰
اندازه‌گیری F	۹۰,۴۴	۹۰,۷۵	۹۰,۳۸	۹۳,۵۷	۹۲,۴۳	۹۱,۴۲	۸۹,۶۶	۹۰,۵۹
صحت	۹۵,۴۷	۹۴,۳۵	۹۴,۰۵	۹۳,۶۱	۹۳,۲۷	۹۳,۶۵	۹۲,۸۶	۹۲,۰۵
نرخ خطا	۴,۵۳	۵,۶۵	۵,۹۵	۶,۳۹	۶,۷۳	۶,۳۵	۷,۱۴	۷,۹۵

جدول ۶: ارزیابی معیار صحت با انتخاب ویژگی‌های مختلف با ۱۰۰ بار تکرار

تعداد ویژگی	ویژگی‌های انتخاب شده	درصد صحت
۲۰	۵۶, ۱۹, ۱۶, ۱۳, ۳۰, ۱۸, ۶, ۵۵, ۲۰, ۸, ۵, ۳۷, ۲۵, ۲۳, ۲۱, ۵۳, ۵۲, ۲۷, ۳۵, ۲	۹۲,۳۶
۲۷	۳۰, ۱۸, ۶, ۵۵, ۲۰, ۸, ۵۷, ۳۷, ۲۵, ۲۳, ۲۱, ۵۳, ۵۲, ۲۷, ۲۴, ۵, ۵۶, ۱۹, ۱۶, ۷, ۳۵, ۹, ۱۳, ۱۷, ۲۶, ۱۱, ۳	۹۱,۸۶
۳۲	۲, ۲۱, ۵۳, ۵۲, ۷, ۲۵, ۲۳, ۳, ۲۷, ۲۴, ۵, ۲۲, ۵۰, ۴۸, ۱۰, ۴, ۱۵, ۲۸, ۱۲, ۳۵, ۵۶, ۱۹, ۱۶, ۲۶, ۱۱, ۶, ۵۵, ۲۰, ۸, ۵۷, ۳۷, ۱۷	۹۱,۵۲
۴۰	۴۵, ۳۲, ۴۸, ۳۵, ۹, ۲, ۳۰, ۱۸, ۶, ۵۵, ۲۰, ۲۲, ۱۳, ۱۰, ۴۶, ۲۹, ۱, ۲۷, ۱۵, ۱۲, ۸, ۵۷, ۵, ۵۶, ۱۹, ۱۶, ۲۵, ۲۳, ۲۱, ۵۳, ۵۲, ۳۱, ۷, ۲۶, ۳, ۲۷, ۲۴, ۳۳, ۱۴	۹۰,۷۶
۴۹	۱۹, ۱۶, ۲۵, ۲۳, ۲۱, ۵۳, ۱۴, ۴, ۵۴, ۹, ۲, ۳۰, ۱۸, ۶, ۱۵, ۲۰, ۸, ۵۷, ۳۷, ۱۷, ۲۶, ۴۲, ۴۰, ۳۲, ۳۱, ۱۳, ۱۰, ۱۱, ۲۹, ۱, ۲۸, ۵۵, ۱۲, ۵, ۵۶, ۱۹, ۱۶, ۲۵, ۲۳, ۲۱, ۵۳, ۱۴, ۴, ۵۴, ۷, ۳۸, ۳۷, ۳۴, ۳۳, ۲۲, ۲۰, ۵۶, ۳۵	۸۹,۷۱
۵۰	۳۶, ۴۶, ۳, ۲۷, ۸, ۵۷, ۳۷, ۱۷, ۲۶, ۴۲, ۴۰, ۳۲, ۳۱, ۱۳, ۱۰, ۱۱, ۲۹, ۱, ۲۸, ۵۵, ۱۲, ۵, ۵۶, ۱۹, ۱۶, ۲۵, ۲۳, ۲۱, ۵۳, ۱۴, ۴, ۵۴, ۲۴, ۹, ۲, ۳۰, ۱۸, ۶, ۱۵, ۲۰, ۴۸, ۷, ۳۸, ۳۷, ۳۴, ۳۳, ۲۲, ۲۰, ۵۶, ۳۵	۸۸,۰۹
۵۳	۲۶, ۳۷, ۱۷, ۵۷, ۱۱, ۳۱, ۳۲, ۳, ۱۳, ۱۰, ۴۶, ۲۹, ۱, ۲۸, ۵۶, ۱۹, ۱۶, ۲۷, ۲۴, ۵, ۳۹, ۳۶, ۲۲, ۱۴, ۴, ۲۵, ۲۳, ۲۱, ۵۳, ۵۲, ۷, ۴۹, ۴۸, ۴۴, ۴۳, ۴۱, ۵۴, ۴۵, ۳۳, ۴۲, ۴۰, ۳۴, ۱۶, ۱۲, ۳۵, ۹, ۲, ۳۰, ۱۸, ۶, ۵۵, ۲۰, ۸	۹۰,۱۲
۵۷	انتخاب همه ویژگی‌ها	۸۹,۰۰

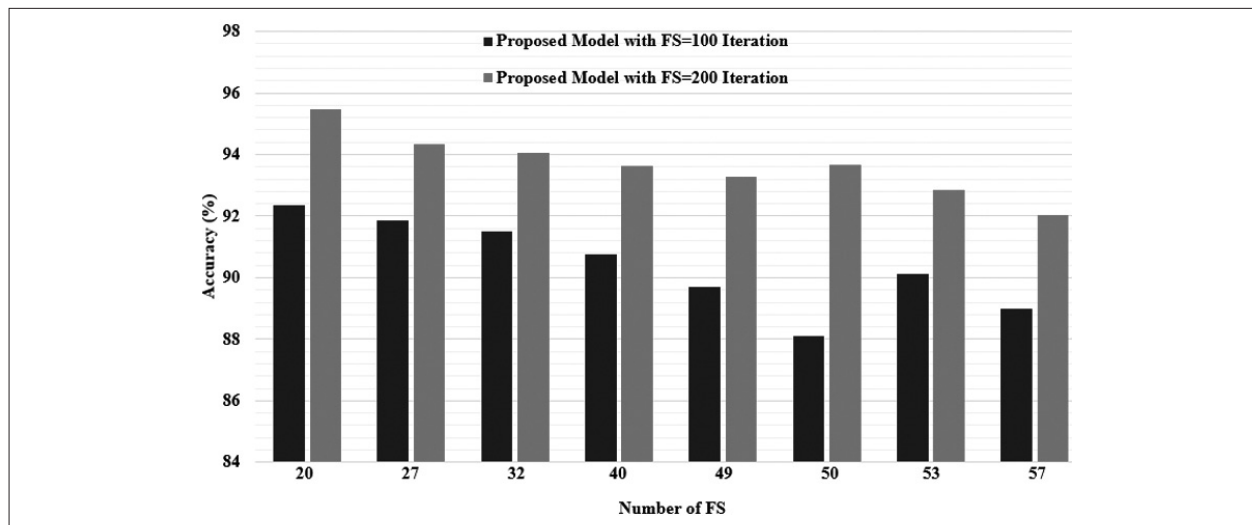
مختلف مقدار دقت متفاوتی دارند. لذا تعداد ویژگی در دقت مدل پیشنهادی تاثیرگذار بوده است.

در جدول (۷)، نتایج ارزیابی معیار صحت بر مبنای انتخاب ویژگی‌های مختلف با ۲۰۰ بار تکرار نشان داده شده است. در جدول (۷) واضح است که با کاهش انتخاب ویژگی

است. معیار صحت برای مقایسه انتخاب ویژگی‌ها در مدل پیشنهادی استفاده شده است. در بیشتر تحقیقات انجام شده از معیار صحت به عنوان معیار اصلی برای ارزیابی استفاده شده است. در جدول (۶) مقایسه دقت مدل پیشنهادی بر مبنای انتخاب ویژگی انجام شده و ویژگی‌های

جدول ۷: ارزیابی معیار صحت با انتخاب ویژگی‌های مختلف با ۲۰۰ بار تکرار

تعداد ویژگی	ویژگی‌های انتخاب شده	درصد صحت
۲۰	۵۶, ۱۹, ۱۶, ۱۳, ۳۰, ۱۸, ۶, ۵۵, ۲۰, ۸, ۵, ۳۷, ۲۵, ۲۳, ۲۱, ۵۳, ۵۲, ۲۷, ۳۵, ۲	۹۵,۴۷
۲۷	۳۰, ۱۸, ۶, ۵۵, ۲۰, ۸, ۵۷, ۳۷, ۲۵, ۲۳, ۲۱, ۵۳, ۵۲, ۲۷, ۲۴, ۵, ۵۶, ۱۹, ۱۶, ۷, ۳۵, ۹, ۱۳, ۱۷, ۲۶, ۱۱, ۳	۹۴,۳۵
۳۲	۲, ۲۱, ۵۳, ۵۲, ۷, ۲۵, ۲۳, ۳, ۲۷, ۲۴, ۵, ۲۲, ۵۰, ۴۸, ۱۰, ۴, ۱۵, ۲۸, ۱۲, ۳۵, ۵۶, ۱۹, ۱۶, ۲۶, ۱۱, ۶, ۵۵, ۲۰, ۸, ۵۷, ۳۷, ۱۷	۹۴,۰۵
۴۰	,۴۵, ۳۲, ۴۸, ۳۵, ۹, ۲, ۳۰, ۱۸, ۶, ۵۵, ۲۰, ۲۲, ۱۳, ۱۰, ۴۶, ۲۹, ۱, ۲۷, ۱۵, ۱۲, ۸, ۵۷, ۵, ۵۶, ۱۹, ۱۶, ۲۵, ۲۳, ۲۱, ۵۳, ۵۲, ۳۱, ۷, ۲۶, ۳, ۲۷, ۲۴, ۳۳, ۱۴	۹۳,۶۱
۴۹	,۱۹, ۱۶, ۲۵, ۲۳, ۲۱, ۵۳, ۱۴, ۴, ۵۴, ۹, ۲, ۳۰, ۱۸, ۶, ۱۵, ۲۰, ۸, ۵۷, ۳۷, ۱۷, ۲۶, ۴۲, ۴۰, ۳۲, ۳۱, ۱۳, ۱۰, ۱۱, ۲۹, ۱, ۲۸, ۵۵, ۱۲, ۵, ۵۶, ۱۹, ۱۶, ۲۵, ۲۳, ۲۱, ۵۳, ۱۴, ۴, ۵۴, ۷, ۳۸, ۳۷, ۳۴, ۳۳, ۲۲, ۲۰, ۵۶, ۳۵, ۳۶, ۴۶, ۳, ۲۷, ۲۴, ۵, ۵۶	۹۳,۲۷
۵۰	,۳۶, ۴۶, ۳, ۲۷, ۸, ۵۷, ۳۷, ۱۷, ۲۶, ۴۲, ۴۰, ۳۲, ۳۱, ۱۳, ۱۰, ۱۱, ۲۹, ۱, ۲۸, ۵۵, ۱۲, ۵, ۵۶, ۱۹, ۱۶, ۲۵, ۲۳, ۲۱, ۵۳, ۱۴, ۴, ۵۴, ۲۴, ۹, ۲, ۳۰, ۱۸, ۶, ۱۵, ۲۰, ۴۸, ۷, ۳۸, ۳۷, ۳۴, ۳۳, ۲۲, ۲۰, ۵۶, ۳۵	۹۳,۶۵
۵۳	,۲۶, ۳۷, ۱۷, ۵۷, ۱۱, ۳۱, ۳۲, ۳, ۱۳, ۱۰, ۴۶, ۲۹, ۱, ۲۸, ۵۶, ۱۹, ۱۶, ۲۷, ۲۴, ۵, ۳۹, ۳۶, ۲۲, ۱۴, ۴, ۲۵, ۲۳, ۲۱, ۵۳, ۵۲, ۷, ۴۹, ۴۸, ۴۴, ۴۳, ۴۱, ۵۴, ۴۵, ۳۳, ۴۲, ۴۰, ۳۴, ۱۶, ۱۲, ۳۵, ۹, ۲, ۳۰, ۱۸, ۶, ۵۵, ۲۰, ۸	۹۲,۸۶
۵۷	انتخاب همه ویژگی‌ها	۹۲,۰۵



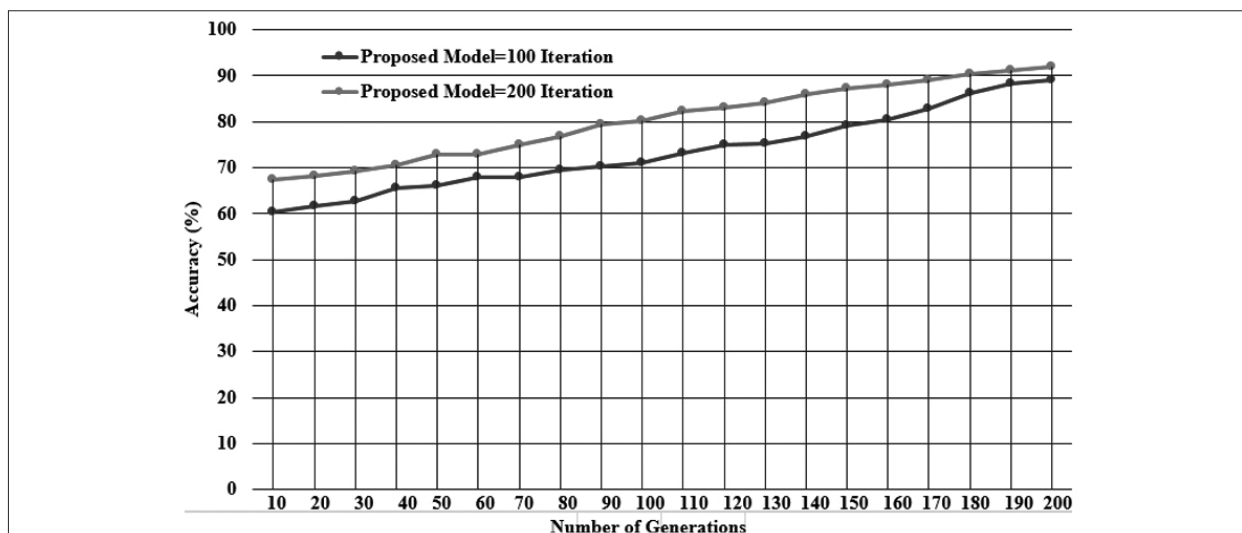
شکل ۲: نمودار مقایسه مدل پیشنهادی با انتخاب ویژگی بر مبنای تعداد تکرار

نسل نشان داده شده است. تعداد نسل بهینه می‌تواند دقت تشخیص را افزایش دهد. زیرا مدل پیشنهادی در هر چرخه از نسل می‌تواند ویژگی‌های متنوع و موثرتری را کشف کند و در نتیجه در طبقه‌بندی نمونه‌ها موثر است. الگوریتم‌های فرا ابتکاری بر مبنای تعداد نسل کاوش در محیط را انجام می‌دهند و تعداد نسل‌ها بهینه به مدل پیشنهادی کمک می‌کند تا محیط جستجو را به خوبی ارزیابی کند. در شکل (۳)، مقایسه بر مبنای ۱۰۰ و ۲۰۰ تکرار انجام شده است. بر مبنای اجراهای مختلف به این نتیجه رسیدیم که

مقدار دقت افزایش یافته است. زیرا در مدل پیشنهادی برای طبقه‌بندی داده‌ها از ویژگی‌های موثر استفاده شده است. در شکل (۲)، نمودار مقایسه مدل پیشنهادی با انتخاب ویژگی بر مبنای تعداد تکرار نشان داده شده است. همانطور که در شکل (۲)، مشاهده می‌کنید مدل پیشنهادی با انتخاب ویژگی و با ۲۰۰ بار تکرار، درصد صحت بیشتری دارد.

۳-۴ ارزیابی مدل پیشنهادی بر مبنای تعداد نسل

در شکل (۳) ارزیابی مدل پیشنهادی بر مبنای ۲۰۰



شکل ۳: ارزیابی مدل پیشنهادی بر مبنای تعداد نسل

جدول ۸: ارزیابی مدل پیشنهادی بر مبنای معیار فاصله همینگ و با ۱۰۰ بار تکرار

معیار فاصله	$d(x, y) = \sum_{i=1}^n x_i - y_i $							
	۲۰	۲۷	۳۲	۴۰	۴۹	۵۰	۵۳	۵۷
انتخاب ویژگی	۲۰	۲۷	۳۲	۴۰	۴۹	۵۰	۵۳	۵۷
دقت	۰.۹۶۲۳	۰.۹۵۱۲	۰.۹۴۶۳	۰.۹۳۴۹	۰.۹۲۵۸	۰.۹۱۶۸	۰.۹۱۳۲	۰.۹۰۱۰
بازخوانی	۰.۹۶۸۵	۰.۹۶۳۵	۰.۹۴۸۷	۰.۹۴۶۸	۰.۹۲۶۸	۰.۹۲۹۸	۰.۹۱۵۶	۰.۹۰۲۱
اندازه‌گیری F	۰.۹۶۵۴	۰.۹۵۷۳	۰.۹۴۷۵	۰.۹۴۰۸	۰.۹۲۶۳	۰.۹۲۳۳	۰.۹۱۴۴	۰.۹۰۱۵
صحت	۰.۹۶۴۲	۰.۹۶۸۴	۰.۹۵۰۸	۰.۹۴۸۲	۰.۹۳۲۱	۰.۹۲۸۷	۰.۹۱۶۵	۰.۹۰۴۶
نرخ خطا	۰.۰۳۵۸	۰.۰۳۱۶	۰.۰۴۹۲	۰.۰۵۱۸	۰.۰۶۷۹	۰.۰۷۱۳	۰.۰۹۳۵	۰.۰۹۵۴

از معیارهای مختلف مورد ارزیابی قرار گرفته است. در جدول (۹)، بیشترین درصد صحت برابر با ۰.۹۵۲۱ است.

در جدول (۱۰)، نتایج ارزیابی مدل پیشنهادی بر مبنای معیار فاصله کنبرا و با ۱۰۰ بار تکرار و با استفاده از معیارهای مختلف مورد ارزیابی قرار گرفته است. در جدول (۱۰)، بیشترین درصد صحت برابر با ۹۶.۹۴ است.

۴-۵ مقایسه و ارزیابی

مقایسه و ارزیابی به منظور کارایی و برتری مدل پیشنهادی در مقایسه با مدل‌های دیگر انجام شده است. مدل پیشنهادی در مقایسه با اغلب الگوریتم‌های فرا ابتکاری و داده کاوی از دقت بیشتری بهره‌مند است. یکی از دلایلی که باعث شده است که مدل پیشنهادی دقت بیشتری داشته

تکرارهای ۱۰۰ و ۲۰۰ برای مدل پیشنهادی حالت بهینه دارند. در بازه ۱۰۰ تا ۲۰۰ مدل پیشنهادی می‌تواند به درصد صحت برارنده دست یابد و مقادیر بیشتر از این بازه در درصد صحت تأثیر چندانی ندارند و فقط باعث می‌شوند که زمان محاسباتی افزایش یابد.

۴-۴ ارزیابی مدل پیشنهادی بر مبنای معیارهای فاصله

در جدول (۸)، نتایج ارزیابی مدل پیشنهادی بر مبنای معیار فاصله همینگ و با ۱۰۰ بار تکرار و با استفاده از معیارهای مختلف مورد ارزیابی قرار گرفته است. در جدول (۸)، بیشترین درصد صحت برابر با ۰.۹۶۴۲ است.

در جدول (۹)، نتایج ارزیابی مدل پیشنهادی بر مبنای معیار فاصله چبیشف و با ۱۰۰ بار تکرار و با استفاده

جدول ۹: ارزیابی مدل پیشنهادی بر مبنای معیار فاصله چبیشف و با ۱۰۰ بار تکرار

معیار فاصله	$d(x, y) = \max_{i=1,2,\dots,n} x_i - y_i $							
انتخاب ویژگی	۲۰	۲۷	۳۲	۴۰	۴۹	۵۰	۵۳	۵۷
دقت	۹۵,۱۶	۹۴,۹۱	۹۳,۲۱	۹۲,۸۴	۹۱,۳۶	۹۰,۵۴	۸۹,۶۴	۸۷,۹۵
بازخوانی	۹۶,۳۱	۹۵,۸۸	۹۳,۵۴	۹۳,۲۶	۹۱,۵۸	۹۱,۶۴	۸۹,۹۶	۸۹,۱۵
اندازه‌گیری F	۹۵,۷۳	۹۵,۳۹	۹۳,۳۷	۹۳,۰۵	۹۱,۴۷	۹۱,۰۹	۸۹,۸۰	۸۸,۵۵
صحت	۹۵,۲۱	۹۴,۶۸	۹۳,۰۸	۹۲,۶۳	۹۱,۳۵	۹۰,۴۸	۸۹,۱۳	۸۹,۵۱
نرخ خطا	۴,۷۹	۵,۳۲	۶,۹۲	۷,۳۷	۸,۶۵	۹,۵۲	۱۰,۸۷	۱۰,۴۹

جدول ۱۰: ارزیابی مدل پیشنهادی بر مبنای معیار فاصله کنبرا و با ۱۰۰ بار تکرار

معیار فاصله	$d(x, y) = \sum_{i=1}^n \frac{ x_i - y_i }{x_i + y_i}$							
انتخاب ویژگی	۲۰	۲۷	۳۲	۴۰	۴۹	۵۰	۵۳	۵۷
دقت	۹۷,۵۲	۹۶,۴۳	۹۴,۶۸	۹۳,۴۸	۹۲,۳۵	۹۱,۸۷	۹۱,۳۴	۹۰,۱۶
بازخوانی	۹۷,۸۲	۹۶,۵۶	۹۵,۱۳	۹۳,۶۷	۹۳,۶۵	۹۱,۹۰	۹۱,۸۴	۹۱,۵۴
اندازه‌گیری F	۹۷,۶۷	۹۶,۴۹	۹۴,۹۰	۹۳,۵۷	۹۳,۰۰	۹۱,۸۸	۹۱,۵۹	۹۰,۸۴
صحت	۹۶,۹۴	۹۵,۸۹	۹۴,۹۷	۹۳,۵۲	۹۲,۸۴	۹۱,۳۴	۹۰,۴۸	۹۰,۶۱
نرخ خطا	۳,۰۶	۴,۱۱	۵,۰۳	۶,۴۸	۷,۱۶	۸,۶۶	۹,۵۲	۹,۳۹

۵. نتیجه‌گیری و کارهای آینده

بر اساس بررسی انجام شده، هرزنامه بیش از چهل درصد از ترافیک ایمیل در اینترنت را شامل می‌شود. که این ترافیک باعث کاهش افت شدید سرعت اینترنت و همچنین کاهش کارایی کارسازها در پردازش داده‌ها می‌شوند. مهم‌ترین راهکار برای مقابله با هرزنامه، شناسایی و تشخیص نوع ایمیل‌های ارسالی از سوی هرزنامه‌نویسان است. در این مقاله، برای تشخیص هرزنامه از مدل ترکیبی بهینه‌سازی کلونی مورچه و الگوریتم کرم شبتاب بر روی مجموعه داده Spambase استفاده شد. نتایج نشان داد که مدل پیشنهادی در طبقه‌بندی دقت بالایی دارد. و همچنین کارایی مدل پیشنهادی بر مبنای انتخاب ویژگی و تکرار نشان داد که با انتخاب ویژگی‌های مهم می‌توان به دقت بیشتری دست یافت. درصد صحت مدل پیشنهادی با ۲۰۰ و ۱۰۰ بار تکرار برای همه ویژگی‌ها به ترتیب برابر ۸۹,۰۰ و ۹۲,۰۵ به دست آمد. زیرا با ۲۰۰ بار تکرار،

باشد استفاده از ترکیب دو مدل فرا ابتکاری و بهره‌مندی از عملگرهایی مانند بروزسانی، جستجو برای بهترین بردار ویژگی، انتخاب ویژگی بر مبنای تشابه و نزدیکی و تعداد نسل‌ها است. در جدول (۱۱)، مقایسه مدل پیشنهادی با مدل‌های دیگر نشان داده شده است.

در جدول (۱۱)، مدل‌های بهینه‌سازی اجتماع ذرات - الگوریتم انتخاب منفی، الگوریتم انتخاب منفی، بهینه‌سازی اجتماع ذرات، بوستینگ منطقی و بی‌زین ساده کمترین درصد صحت را در مقایسه با مدل پیشنهادی و مدل‌های دیگر دارند. و همچنین مدل پیشنهادی بر مبنای تعداد تکرار و تعداد ویژگی‌ها نتایج متفاوتی دارد. مدل پیشنهادی هم مانند مدل‌های دیگر بنابه دلایلی مانند عدم تنوع در جمعیت اولیه و گیرافتادن در بهینه محلی نمی‌توانند به طور صد در صد پاسخگوی سیستم تشخیص ایمیل هرزنامه باشد. اما درصد صحت مدل پیشنهادی با کاهش تعداد ویژگی‌ها، افزایش می‌یابد.

جدول ۱۱: مقایسه مدل پیشنهادی با مدل‌های دیگر

اندازه‌گیری F	بازخوانی	دقت	صحت	مدل‌ها	رفرنس‌ها
۹۰.۵	۸۸.۲	۹۳.۵	۹۱.۷	روش گروهی مدیریت داده‌ها	[۱۰]
۷۴.۹۵	۶۵.۹۹	۸۶.۷۱	۸۳.۲۰	بهینه‌سازی اجتماع ذرات-الگوریتم انتخاب منفی	[۱۱]
۳۶.۰۱	۲۲.۲۴	۹۴.۵۳	۶۸.۸۶	الگوریتم انتخاب منفی	[۱۸]
۶۹.۷۶	۵۶.۶۲	۹۰.۸۶	۸۰.۶۶	تکاملی تفاضلی-الگوریتم انتخاب منفی	
۴۳.۵۳	۳۲.۳۱	۸۱.۳۶	۹۱.۲۲	بهینه‌سازی اجتماع ذرات-الگوریتم انتخاب منفی	[۱۹]
۳۸.۹۴	۲۸.۶۴	۸۱.۷۷	۸۱.۳۲	بهینه‌سازی اجتماع ذرات	
۲۲.۰۹	۱۳.۶۳	۸۵.۰۲	۶۸.۸۶	الگوریتم انتخاب منفی	
۹۰.۸	۸۸.۶	۹۳.۱۱	۸۸.۵۶	شبکه بیزین	[۲۶]
۹۱.۶۶	۹۰.۴	۹۲.۸۶	۸۹.۷	بوستینگ منطقی	
۹۳.۰۰	۹۳.۲	۹۲.۷۵	۹۱.۵۴	درخت تصادفی	
۹۳.۷۱	۹۲.۹	۹۴.۴۷	۹۲.۳۲	JRIP	
۹۳.۷	۹۳.۵	۹۳.۹۰	۹۲.۳۴	درخت J۴۸	
۹۴.۴۵	۹۴.۵	۹۴.۳۳	۹۳.۲۸	شبکه عصبی مصنوعی چندلایه	
۹۴.۷۳	۹۳.۹	۹۵.۵۸	۹۳.۵۶	الگوریتم Kstar	
۹۵.۰۰	۹۴.۱	۹۵.۸۷	۹۳.۸۹	جنگل تصادفی	
۹۵.۳۲	۹۴.۵	۹۶.۱۲	۹۴.۲۸	کمینه تصادفی	
۵۹.۹۷	۹۲.۶۴	۴۳.۲۶	۷۹.۲۸	بیزین ساده	[۲۷]
۳۱.۷۷	۶۱.۰۸	۲۱.۴۷	۸۹.۸۰	شبکه بیزین	
۳۰.۴۸	۶۳.۳۵	۲۰.۰۷	۹۰.۴۱	ماشین بردار پشتیبان	
۲۳.۷۴	۵۰.۴۹	۱۵.۵۲	۹۳.۳۴	درخت توابع	
۲۷.۵۴	۵۲.۴۳	۱۸.۶۸	۹۲.۹۷	J۴۸	
۲۷.۲۵	۴۲.۲۳	۲۰.۱۲	۹۴.۸۲	جنگل تصادفی	
۲۸.۹۱	۶۱.۴۰	۱۸.۹۱	۹۰.۹۳	درخت تصادفی	
۳۱.۲۴	۵۳.۳۳	۲۲.۰۹	۹۲.۴۳	درخت نمونه CART	
۸۸.۱۷	۹۰.۳۲	۸۹.۳۷	۸۹.۰۰	انتخاب همه ویژگی‌ها و تکرار برابر با ۱۰۰	مدل پیشنهادی
۹۰.۵۹	۹۱.۹۰	۸۹.۳۷	۹۲.۰۵	انتخاب همه ویژگی‌ها و تکرار برابر با ۲۰۰	
۸۹.۱۱	۹۱.۰۶	۹۰.۳۶	۹۲.۳۶	انتخاب ۲۰ ویژگی و تکرار برابر با ۱۰۰	
۹۰.۴۴	۹۶.۱۶	۹۵.۰۲	۹۵.۴۷	انتخاب ۲۰ ویژگی و تکرار برابر با ۲۰۰	

منابع

1. F.S. Gharehchopogh, M. Vafadar, M. Motaman, Improve of Invasive Weed Optimization algorithm with K nearest neighbor for email spam classification, Computing Science Journal (CSJ), in press, March 2018.
2. H. Faris, A.M. Al-Zoubi, A.A. Heidari, I. Aljarah, H. Fujita, An intelligent system for spam detection and identification of the most relevant features based on evolutionary Random Weight Networks, Information Fusion, Vol. 48, pp. 67-83, 2019.
3. A.F. Colladon, P.A. Gloor, Measuring the impact of

جستجو در فضای محیط از تنوع و اکتشاف بیشتری بهره‌مند بوده است. برای جلوگیری و کاهش هرزنامه باید از الگوریتم‌های هوشمند استفاده شود. بدین‌گونه که قوانین هر هرزنامه شناسایی شود و از دریافت آن‌ها در دفعات بعدی جلوگیری شود. و همچنین بر مبنای آموزش باید مدل‌هایی ارائه دهیم که با استفاده از عنوان و محتوا هرزنامه را تشخیص دهد.

18. I. Idris, A. Selamat, S. Omatu, Hybrid email spam detection model with negative selection algorithm and differential evolution, *Engineering Applications of Artificial Intelligence*, Vol. 28, pp. 97-110, 2014
19. I. Idris, A. Selamat, Improved email spam detection model with negative selection algorithm and particle swarm optimization, *Applied Soft Computing*, Vol. 22, pp. 11-27, 2014.
20. S.B. Rathod, T.M. Pattewar, Content based spam detection in email using Bayesian classifier, *International Conference on Communications and Signal Processing (IC-CSP)*, IEEE, pp. 1257-1261, 2015.
21. Chi-Yao Tseng; Ming-Syan Chen, Incremental SVM Model for Spam Detection on Dynamic Email Social Networks, *2009 International Conference on Computational Science and Engineering*, Vol. 4, pp. 128-135, 2009.
22. M. Prilepok, T. Jezowicz, J. Platos, V. Snasel, Spam detection using compression and PSO, *Fourth International Conference on Computational Aspects of Social Networks (CASoN)*, pp. 263-270, 2012.
23. B.K. Singh, K. Verma, A.S. Thoke, Investigations on Impact of Feature Normalization Techniques on Classifier's Performance in Breast Tumor Classification, *International Journal of Computer Applications*, Vol. 116, No. 19, pp. 11-15, 2015.
24. R.S. Michalski, I. Bratko, M. Kubat, *Machine Learning, and Data Mining: Methods and Applications*, New York: Wiley, 1998.
25. D. Francois, Binary classification performances measure cheat sheet, 2009.
26. S. Sharma, A. Arora, Adaptive Approach for Spam Detection, *IJCSI International Journal of Computer Science Issues*, Vol. 10, Issue 4, No. 1, pp. 23-26, July 2013.
27. M. Rathi, V. Pareek, Spam Mail Detection through Data Mining-A Comparative Performance Analysis, *I.J. Modern Education and Computer Science*, Vol. 12, pp. 31-39, 2013.
- spammers on e-mail and Twitter networks, *International Journal of Information Management*, In press, corrected proof, Available online 24 September 2018
4. A. Heydari, M.A. Tavakoli, N. Salim, Z. Heydari, Detection of review spam: A survey, *Expert Systems with Applications*, Vol. 42, Issue 7, pp. 3634-3642, 2015
5. G. Dalkılıç, D. Sipahi, Spam filtering with sender authentication network, *Computer Communications*, Vol. 98, pp. 72-79, 2017.
6. M.Dorigo, L.M. Gambardella, The colony system: A cooperative learning approach to the traveling salesman problem, *IEEE Transactions on Evolutionary Computation*, Vol. 1, No.1, pp. 53-66, April 1997.
7. X.S. Yang, *Nature-Inspired Meta-heuristic Algorithms*, Luniver Press, 2008.
8. <https://archive.ics.uci.edu/ml/datasets/Spambase>
9. Siti-Hajar-Aminah Ali, S. Ozawa, J. Nakazato, T. Ban, J. Shimamura, An autonomous online malicious spam email detection system using extended RBF network, *International Joint Conference on Neural Networks (IJCNN)*, pp. 1-7, 2015.
10. El-Sayed M. El-Alfy, R.E. Abdel-Aal, Using GMDH-based networks for improved spam detection and email feature analysis, *Applied Soft Computing*, Vol. 11, Issue 1, pp. 477-488, 2011.
11. I. Idris, A. Selamat, N.T. Nguyen, S. Omatu, O. Krejcar, K. Kuca, M. Penhaker, A combined negative selection algorithm-particle swarm optimization for an email spam detection system, *Engineering Applications of Artificial Intelligence*, Vol. 39, pp. 33-44, 2015
12. L. Nguyen, A.Q. Tran, L. Thu Bui, DMEA-II and its application on spam email detection problems, *Seventh IEEE Symposium on Computational Intelligence for Security and Defense Applications (CISDA)*, IEEE, pp. 1-6, 2014
13. Mu-Chun Su, Hsu-Hsun Lo, Fu-Hau Hsu, A neural tree, and its application to spam e-mail detection, *Expert Systems with Applications*, Vol. 37, pp. 7976-7985, 2010.
14. S. Ergin, S. Isik, The investigation on the effect of feature vector dimension for spam email detection with a new framework, *9th Iberian Conference on Information Systems and Technologies (CISTI)*, pp. 1-4, 2014
15. W. Ma, D. Tran, D. Sharma, A Novel Spam Email Detection System Based on Negative Selection, *2009 Fourth International Conference on Computer Sciences and Convergence Information Technology*, pp. 987-992, 2009.
16. A.R. Behjat, A. Mustapha, H. Nezamabadi-pour, Md. Nasir Sulaiman, and N. Mustapha, A PSOBased Feature Subset Selection for Application of Spam /Non-spam Detection, *Springer-Verlag Berlin Heidelberg, M-CAIT 2013, CCIS 378*, pp. 183-193, 2013.
17. Kuo-Ching Ying, Shih-Wei Lin, Zne-Jung Lee, Yen-Tim Lin, An ensemble approach applied to classify spam e-mails, *Expert Systems with Applications*, Vol. 37, Issue 3, pp. 2197-2201, 2010.