

تاریخ دریافت مقاله: ۹۶/۱۰/۰۳

تاریخ پذیرش مقاله: ۹۷/۰۳/۱۹

بهبود الگوریتم بهینه‌سازی علف‌های هرز مهاجم با K نزدیک‌ترین همسایه در طبقه‌بندی ایمیل هرزنامه

فرهاد سلیمانیان قره چیق*

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

مهدی وفادار

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

محسن موتمن فر

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

چکیده

هرزنامه پیشنهاد شده است. در مدل پیشنهادی با استفاده از الگوریتم بهینه‌سازی علف‌های هرز مهاجم، ویژگی‌های مهم برای تشخیص ایمیل هرزنامه را انتخاب می‌کنیم و سپس توسط K نزدیک‌ترین همسایه، نمونه‌ها را طبقه‌بندی می‌کنیم. انتخاب ویژگی باعث می‌شود که درصد دقت تشخیص بیشتر باشد. ارزیابی مدل پیشنهادی بر روی مجموعه داده Spambase انجام شده است. نتایج نشان داده که مقدار صحت در مدل پیشنهادی برابر ۹۱،۱۱٪ است که در مقایسه با مدل‌های دیگر دقت تشخیص بهتری دارد. **واژه‌های کلیدی:** تشخیص ایمیل هرزنامه، الگوریتم بهینه‌سازی علف‌های هرز مهاجم، K نزدیک‌ترین همسایه

امروزه ایمیل که یک وسیله ارتباطی سریع و کم‌هزینه می‌باشد به یکی از مهم‌ترین ابزارهای مدرن برای ارتباط در محیط‌های برخط تبدیل شده است. اما مشکلی اصلی این است که ارسال هرزنامه باعث آزار کاربران، اتلاف زمان، هزینه، منابع شبکه و پهنای باند می‌شود. بنابراین ایمیل هرزنامه به یک مشکل جدی برای کاربران و سازمان‌هایی که مدت زمان زیادی از وقت خود را با ایمیل و کامپیوتر کار می‌کنند تبدیل شده است. در ایمیل هرزنامه بیشتر اوقات، کلاهبرداران با فرستادن هرزنامه قصد نفوذ به کامپیوتر یا کلاهبرداری را دارند که این گونه هرزنامه‌ها اغلب حاوی نرم‌افزارها و پیوندهای مخرب هستند که در درون خود انواع بدافزارها و ویروس‌ها را دارند. در این مقاله یک مدل بهبود یافته مبتنی بر الگوریتم بهینه‌سازی علف‌های هرز مهاجم و K نزدیک‌ترین همسایه برای تشخیص ایمیل

۱. مقدمه

پرکاربردترین و مهم‌ترین ابزار ارتباطی در جهان امروز که برای همه کاربران ساده و در دسترس می‌باشد ایمیل

* نویسنده مسئول

است. پیشرفت قابل توجه در عرصه فناوری‌های ارتباطات جهانی و همچنین ایمیل این امکان را به کاربران می‌دهد که در تمام نقاط جهان و با دسترس بودن اینترنت به راحتی می‌توانید ایمیل‌های خود را به سرتاسر جهان ارسال نمایید و بدون آن‌که زمان زیادی را تلف کنید فقط با چند ثانیه ایمیل خود را به گیرنده ارسال نمایید [۱]. از طرفی دیگر، دسترسی آسان، سریع و بی‌شمار به ایمیل توسط کاربران موجب شده است تا حجم وسیعی از محتوا و اطلاعات با ساختار و بدون ساختار به ایمیل کاربران ارسال گردد. از آنجایی که امکان تحقیق و نیز فرصت کافی برای بررسی اعتبار آن‌ها در زمان کوتاه وجود ندارد، لذا به یک تهدید امنیتی تبدیل می‌شوند و در بیشتر موارد به شکل ایمیل هرزنامه برای کاربران ارسال می‌شوند [۲]. هرزنامه برای اهداف مختلف بدون اجازه و تمایل دریافت کننده، ارسال می‌شود به طوری که کاربر پس از دریافت آن ایمیل (که می‌تواند با اهداف خرابکارانه یا با هدف تبلیغات ارسال شود) نمی‌تواند جلوی دریافت مجدد آن را بگیرد [۳]. در نتیجه، تهدیدهای امنیتی در محیط اینترنت استفاده از ایمیل را آسیب‌پذیر کرده‌اند که یکی از مهم‌ترین تهدیدها هرزنامه است [۴]. امروزه هرزنامه مانند سیلی عظیم، و با کپی‌های فراوان به ایمیل‌های کاربران ارسال می‌شود. در بیشتر مواقع هرزنامه‌ها فقط باعث مزاحمت کاربرانی که آن را دریافت می‌کند می‌شوند اما گاهی قضیه فراتر از این است و نفوذگرها با فرستادن هرزنامه قصد نفوذ به کامپیوتر کاربران یا کلاهبرداری را دارند که این گونه هرزنامه‌ها اغلب حاوی نرم‌افزارها و پیوندهای مخربی هستند که در درون خود انواع بدافزارها، ویروس‌ها و باج افزارها را دارند [۵]. تاکنون روش‌های مختلفی برای مبارزه با هرزنامه ایجاد شده که یکی از راهکارها استفاده از نرم‌افزارهای فیلترینگ است که در این برنامه‌ها با استفاده از کلید واژه‌های خاصی، موضوع پیام و ایمیل‌ها بررسی می‌شود و در صورت شناسایی آن‌ها ایمیل دریافت شده حذف می‌شود [۶]. کاربرانی که در ایمیل خود هرزنامه

دریافت می‌کنند، می‌توانند نشانی ایمیل مورد نظر را در قسمت ایمیل‌های مسدود شده قرار دهند تا از دریافت مجدد ایمیل از آن نشانی جلوگیری شود.

استفاده از فناوری اطلاعات و ارتباطات در حوزه‌های مختلف در حال رشد و گسترش بوده و اطلاعات کاربران برخط با ایمیل هرزنامه در خطر حمله نفوذگرها قرار دارد. به همین دلیل شناسایی و تشخیص ایمیل هرزنامه یک راهکار مفید برای جلوگیری از نفوذ افراد ناشناس به حساب کاربری کاربران می‌باشد [۷]. در این مقاله با بررسی ویژگی‌ها و مقادیر مربوط به مجموعه داده Spambase [۸] با استفاده از ترکیب الگوریتم بهینه‌سازی علف‌های هرز مهاجم [۹] و K نزدیک‌ترین همسایه [۱۰] روشی موثر برای دقت تشخیص ایمیل هرزنامه پیشنهاد می‌کنیم. در مدل پیشنهادی از الگوریتم بهینه‌سازی علف‌های هرز مهاجم برای انتخاب ویژگی و از K نزدیک‌ترین همسایه برای طبقه‌بندی نمونه‌ها استفاده می‌شود. الگوریتم بهینه‌سازی علف‌های هرز مهاجم یکی از الگوریتم‌های قابلیت تطابق‌پذیری و تولید مثل، روابط بین مقادیر را پیدا می‌کند. این الگوریتم با الهام از تکثیر و رشد علف‌های هرز ابداع شده است. طبق تعریف، علف‌های هرز به گیاهانی ناخواسته اطلاق می‌شود که دارای رفتار تهاجمی برای رشد بوده و به عنوان تهدیدی برای دیگر گیاهان زراعی مفید می‌باشند و جلوی رشد آن‌ها را می‌گیرند. این الگوریتم در عین سادگی، در یافتن نقاط بهینه بسیار مؤثر و سریع می‌باشد. در واقع این الگوریتم براساس ویژگی‌های طبیعی علف‌های هرز مانند تولید بذر، رشد و تنازع برای بقا عمل می‌کند.

هدف در انتخاب ویژگی در واقع پیدا کردن کوچک‌ترین زیرمجموعه از ویژگی‌های ورودی با بیشترین خاصیت پیشگویانه یا بیشترین اطلاعات جداکننده در کل مجموعه داده‌ها است. با انتخاب ویژگی‌های مناسب از میان

1- Invasive Weed Optimization (IWO)
2- K-Nearest Neighbor (KNN)

ویژگی‌های فراوان استخراج شده، می‌توان علاوه بر دست یافتن به نرخ دقت دسته‌بندی بالا، هزینه‌های محاسباتی را کاهش داده و از استخراج ویژگی‌های غیر ضروری خودداری کرد.

ساختار کلی مقاله را به شرح زیر سازماندهی کردیم: در بخش دوم، کارهای قبلی را توضیح می‌دهیم. در بخش سوم، مدل پیشنهادی و توضیحات لازم در مورد الگوریتم بهینه‌سازی علف‌های هرز مهاجم و K نزدیک‌ترین همسایه را توضیح خواهیم داد. در بخش چهارم، ارزیابی و نتایج مدل پیشنهادی و مقایسه‌ها را توضیح می‌دهیم و نهایتاً در بخش پنجم به نتیجه‌گیری و کارهای آینده خواهیم پرداخت.

۲. تحقیقات انجام شده

در دوده اخیر، با رشد چشمگیر استفاده از ایمیل و از طرفی حجم ایمیل‌های ناخواسته و مزاحم تحت عنوان هرزنامه باعث انگیزه محققان برای طراحی و پیاده‌سازی سیستم‌هایی جهت فیلترسازی ایمیل‌های هرزنامه با الهام از تکنیک‌های طبقه‌بندی شد.

مدل شبکه عصبی مصنوعی بر پایه داده‌گردانی با روش گروهی^۲ [۱۱] برای طبقه‌بندی ایمیل هرزنامه پیشنهاد شده است. این شبکه عصبی مصنوعی یکی از پرکاربردترین شبکه‌های عصبی مصنوعی است که از توانایی بالایی در مدل‌سازی داده‌های پیچیده برخوردار است. این روش یک روش مدل‌سازی آماری رده‌یک نیست؛ بلکه فرآیندی منظم برای غلبه بر ضعف‌های آماری و شبکه‌های عصبی مصنوعی است. این نوع شبکه عصبی، حاوی مجموعه‌ای از نرون‌ها است که از پیوند جفت‌های مختلف از طریق یک چندجمله‌ای درجه دوم به وجود می‌آیند. شبکه با ترکیب چند جمله‌ای‌های درجه دوم حاصل از تمامی نرون‌ها برای یک مجموعه از ورودی‌ها با کمترین خطای خروجی تعریف می‌شود. ارزیابی شبکه عصبی مصنوعی بر پایه داده‌گردانی با روش گروهی بر روی مجموعه داده جهانی و معتبر Spambase انجام شده است. مجموعه داده

Spambase شامل ۵۷ ویژگی و ۶۶۰۱ نمونه است. نمونه‌ها در دو رده هرزنامه و غیر هرزنامه طبقه‌بندی شده‌اند. نتایج با انتخاب ویژگی‌های مختلف نشان داده که مقدار معیار صحت در شبکه عصبی مصنوعی بر پایه داده‌گردانی با روش گروهی در مقایسه با مدل‌های پرسپترون چندلایه و بیز ساده بیشتر است. بیشترین مقدار صحت در مدل شبکه عصبی مصنوعی بر پایه داده‌گردانی با روش گروهی برابر ۹۲٫۴ و در مدل‌های پرسپترون چندلایه و بیز ساده^۴ به ترتیب برابر ۹۱٫۷ و ۷۵٫۴ است.

مدل ترکیبی بهینه‌سازی اجتماع ذرات-الگوریتم انتخاب منفی^۵ به منظور طبقه‌بندی ایمیل هرزنامه پیشنهاد شده است [۱۲]. ارزیابی بر روی مجموعه داده Spambase انجام شده است. در مدل ترکیبی بهینه‌سازی اجتماع ذرات-الگوریتم انتخاب منفی از مدل بهینه‌سازی اجتماع ذرات برای جستجوی مقدار ویژگی‌ها در فضای مسئله و از مدل الگوریتم انتخاب منفی برای انتخاب ویژگی‌ها استفاده شده است. روش پیشنهادی دارای دو فاز آموزش و آزمایش است. در فاز آموزش ابتدا با مدل الگوریتم انتخاب منفی از طریق فیلترینگ نمونه‌ها تعدادی ویژگی انتخاب می‌شوند. در فاز آزمایش، براساس مدل بهینه‌سازی اجتماع ذرات-الگوریتم انتخاب منفی تصمیم مناسبی در راستای هرزنامه یا غیر هرزنامه بودن یک نمونه انجام می‌شود. نتایج نشان داده که مدل بهینه‌سازی اجتماع ذرات-الگوریتم انتخاب منفی در مقایسه با مدل‌های بیز ساده، انتخاب ویژگی متمایز-ماشین بردار پشتیبان^۶ و الگوریتم انتخاب منفی دقت تشخیص بالاتر و در مقایسه با مدل ماشین بردار پشتیبان دقت تشخیص کمتری دارد. مقدار معیار صحت در مدل بهینه‌سازی اجتماع ذرات-الگوریتم انتخاب منفی برابر ۸۳٫۲۰ و در مدل‌های بیز ساده، انتخاب ویژگی متمایز-ماشین بردار پشتیبان و الگوریتم انتخاب منفی به ترتیب برابر ۷۸٫۸، ۷۱ و ۶۸٫۸۶ است.

مدل QUANT [۱۳] که ترکیبی از شبکه عصبی مصنوعی

4- Naive Bayes

5- Particle Swarm Optimization-Negative Selection Algorithm

6- Distinguishing Feature Selection and Support Vector Machine

3- Group Method of Data Handling

و درخت تصمیم‌گیری است برای تشخیص ایمیل 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 ایمیل 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 SpamAssassin و Corpus 2006 انجام شده است. نتایج بر روی مجموعه داده SpamAssassin نشان داده که مقدار صحت در مدل QUANT برابر ۸۹،۱۵ و در مدل‌های بیز ساده، بهینه‌سازی حداقل متوالی^۷ و C4.5 به ترتیب برابر ۸۱،۰۸، ۸۸،۶۲ و ۷۳،۰۸ است. و همچنین بر روی مجموعه داده Corpus 2006 مقدار صحت در مدل QUANT برابر ۹۰،۸۷ و در مدل‌های بیز ساده، بهینه‌سازی حداقل متوالی و C4.5 به ترتیب برابر ۸۸،۱۵ و ۸۹،۷۹ است.

مدل ترکیبی شبکه عصبی مصنوعی-بهینه‌سازی اجتماع ذرات به منظور تشخیص ایمیل هرزنانه پیشنهاد شده است [۱۴]. از الگوریتم بهینه‌سازی اجتماع ذرات برای انتخاب ویژگی‌ها و از مدل پرسپترون چندلایه برای آموزش و آزمایش داده‌ها و طبقه‌بندی استفاده شده است. در مدل شبکه عصبی مصنوعی-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. تعداد لایه‌های مخفی در مدل پرسپترون چندلایه بین ۱۵-۳ لحاظ شده و تکرار الگوریتم بهینه‌سازی اجتماع ذرات برای انتخاب ویژگی برابر ۲۰۰ است. ارزیابی بر روی مجموعه داده LingSpam با ۴۸۱ ایمیل هرزنانه و ۲۱۷۱ ایمیل غیر هرزنانه و مجموعه داده SpamAssassin با ۶۰۰۰ نمونه ایمیل انجام شده است. ارزیابی بر روی مجموعه داده SpamAssassin و LingSpam نشان داده که مقدار معیار صحت در مدل شبکه عصبی مصنوعی-بهینه‌سازی اجتماع ذرات به ترتیب برابر

7- Sequential Minimal Optimization

۹۹،۹۸ و ۹۹،۷۹ است. مقایسه‌ها نشان داده که مدل شبکه عصبی مصنوعی-بهینه‌سازی اجتماع ذرات در مقایسه با مدل‌های ماشین بردار پشتیبان با تابع هسته، ماشین بردار پشتیبان با تابع شعاعی پایه^۸ و شبکه عصبی مصنوعی تابع شعاعی پایه دقت تشخیص بهتری دارد.

مدل‌های درخت تصمیم‌گیری، ماشین بردار پشتیبان و شبکه عصبی مصنوعی تابع شعاعی پایه و ترکیب آن‌ها بر روی دو مجموعه داده با ۱۴ ویژگی آزمایش و اجرا شده‌اند [۱۵]. مجموعه داده اولی شامل ۵۰۴ ایمیل (۳۳۶ ایمیل و ۱۶۸ ایمیل هرزنانه) و مجموعه داده دومی شامل ۶۵۷ ایمیل (۲۸۷ ایمیل و ۲۷۰ ایمیل هرزنانه) است. در مدل درخت تصمیم‌گیری از آنتروپی، مدل ماشین بردار پشتیبان از تابع کرنل و شبکه عصبی مصنوعی تابع شعاعی پایه از خطای میانگین استفاده شده است. نتایج بر روی مجموعه داده اولی نشان داده که مقدار صحت در مدل ترکیبی برابر ۹۱،۰۷ و در مدل‌های درخت تصمیم‌گیری، ماشین بردار پشتیبان و شبکه عصبی مصنوعی تابع شعاعی پایه به ترتیب برابر ۸۹،۸۸، ۸۸،۶۹ و ۸۹،۸۸ است. و بر روی مجموعه داده دومی مقدار صحت در مدل ترکیبی برابر ۹۱،۷۸ و در مدل‌های درخت تصمیم‌گیری، ماشین بردار پشتیبان و شبکه عصبی مصنوعی تابع شعاعی پایه به ترتیب برابر ۹۰،۸۷، ۹۰،۸۷ و ۸۹،۰۴ است.

ادریس و همکارانش [۱۶] مدل ترکیبی تکامل تفاضلی-الگوریتم انتخاب منفی^۹ را برای تشخیص ایمیل هرزنانه پیشنهاد کرده‌اند. در این مدل از الگوریتم تکاملی تفاضلی برای جستجوی ویژگی‌ها در فضای مسئله استفاده شده است. الگوریتم تکاملی تفاضلی یکی از جدیدترین روش‌های جستجو است. الگوریتم تکاملی تفاضلی به عنوان روشی قدرتمند و سریع برای مسائل بهینه‌سازی در فضاهای پیوسته معرفی شده است. ارزیابی مدل تکامل تفاضلی-الگوریتم انتخاب منفی بر روی مجموعه داده Corpus با ۱۸۱۳ ایمیل هرزنانه و ۲۷۸۸ ایمیل غیرهرزنانه طبق

8- Radial Basis Function

9- Differential Evolution-Negative Selection Algorithm

معادله (۱) انجام شده است. معادله (۱) به منظور ضریب همبستگی^{۱۰} نمونه‌ها استفاده شده است.

$$CC = \frac{[(TP)(TN) - (FP)(FN)]}{\sqrt{(TP + FN)(TP + FP)(TN + FP)(TN + FN)}} \quad (1)$$

پارامترهای درست مثبت (TP)، درست منفی (TN)، کاذب مثبت (FP)، کاذب منفی (FN) از پارامترهای اصلی برای معیار ضریب همبستگی هستند. خروجی دقت معیار ضریب همبستگی در معادله (۱) بر مبنای پارامترهای ذکر شده می‌باشد. پارامتر درست مثبت، بیانگر تعداد رکوردهایی است که دسته واقعی آن‌ها مثبت بوده و الگوریتم دسته‌بندی نیز دسته آن‌ها را به درستی مثبت تشخیص داده است. پارامتر درست منفی، بیانگر تعداد رکوردهایی است که دسته واقعی آن‌ها را به درستی منفی تشخیص داده است. پارامتر کاذب مثبت، بیانگر تعداد رکوردهایی است که دسته واقعی آن‌ها منفی بوده و الگوریتم دسته‌بندی، دسته آن‌ها را به اشتباه مثبت تشخیص داده است. پارامتر کاذب منفی، بیانگر تعداد رکوردهایی است که دسته واقعی آن‌ها مثبت بوده و الگوریتم دسته‌بندی دسته آن‌ها را به اشتباه منفی تشخیص داده است.

نتایج مدل ترکیبی تکامل تفاضلی-الگوریتم انتخاب منفی نشان داده که مقدار ضریب همبستگی در مدل الگوریتم انتخاب منفی و تکامل تفاضلی-الگوریتم انتخاب منفی به ترتیب برابر ۶۸٫۸۶ و ۸۰٫۶۶ است. و مقدار ضریب همبستگی در مدل الگوریتم انتخاب منفی و تکامل تفاضلی-الگوریتم انتخاب منفی به ترتیب برابر ۳۶٫۰۶ و ۶۰٫۰۸ است. مدل بهبود داده شده بهینه‌سازی اجتماع ذرات-الگوریتم انتخاب منفی [۱۷] برای تشخیص طبقه‌بندی ایمیل 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. در مدل بهینه‌سازی اجتماع ذرات-الگوریتم انتخاب منفی برای تغییر موقعیت ذرات و سرعت ذرات از یک مدل

ریاضی جدید طبق معادله (۲) استفاده شده است.

(۲)

$$dv_i = \left\{ -\alpha \sum_{j \neq i} \left[\left[\frac{r}{x_{ij}} \right]^p - \left[\frac{r}{x_{ji}} \right]^q \right] x_{ij} - \beta \sum_{j \neq i} \left[\left[\frac{r}{x_{ij}} \right]^p - \left[\frac{r}{x_{ji}} \right]^q \right] V_{ij} + f_i(t, x_i, v_i) \right\} dt$$

در معادله (۲)، x موقعیت ذرات، v سرعت ذرات، r پارامتر یادگیری، t تعداد تکرار، مقدار p و q در بازه $[0, 1]$ می‌باشد. هدف معادله (۲)، همگرایی به جواب بهینه است. همگرایی یعنی این‌که الگوریتم مقادیر مناسب و بهینه را پیدا کند.

نتایج نشان داده که مدل بهینه‌سازی اجتماع ذرات-الگوریتم انتخاب منفی در مقایسه با مدل‌های بهینه‌سازی اجتماع ذرات، الگوریتم انتخاب منفی، بیز ساده، ماشین بردار پشتیبان، انتخاب ویژگی متمایز-ماشین بردار پشتیبان^{۱۱}، شبکه عصبی مصنوعی و فازی دقت تشخیص بهتری دارد.

در جدول (۱)، مزایا و معایب هر مدل نشان داده شده است. مدل‌های مختلفی برای تشخیص ایمیل هرزنامه پیشنهاد شده است. اما هر مدل می‌تواند مزایا و معایبی داشته باشد.

۳. مدل پیشنهادی

در این بخش، مدل پیشنهادی که ترکیبی از الگوریتم بهینه‌سازی علف‌های هرز مهاجم و K نزدیک‌ترین همسایه است را توضیح می‌دهیم. از الگوریتم بهینه‌سازی علف‌های هرز مهاجم برای انتخاب ویژگی و از K نزدیک‌ترین همسایه برای طبقه‌بندی داده‌ها استفاده می‌کنیم. در شکل (۱)، روندنمای کلی مدل پیشنهادی نشان داده شده است.

در ابتدا، خواندن مجموعه داده‌ها که شامل ۶۰۱ نمونه می‌باشد انجام می‌گیرد. نمونه‌ها سطر به سطر وارد یک ماتریس 57x4601 می‌شوند. در ابتدا ۵۷ ویژگی از مجموعه داده Spambase به عنوان ورودی برای الگوریتم بهینه‌سازی علف‌های هرز مهاجم در نظر گرفته می‌شوند.

11- Distinguishing Feature Selection and Support Vector Machine

10- Correlation Coefficient

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

| مراجع | مدل‌ها | مزایا | معایب |
|-------|--|--|--|
| [۱۱] | مدل شبکه عصبی مصنوعی بر پایه داده‌گردانی با روش گروهی | *مدل شبکه عصبی مصنوعی بر پایه داده‌گردانی با روش گروهی، شبکه‌ای خودسازمانده و یک سوپره است که از چندین لایه و هر لایه از چندین نرون تشکیل یافته است و لذا آموزش تا کمترین خطا انجام می‌گیرد. | *سرعت همگرایی آن خیلی کند است و به نوبه‌های موجود در مجموعه داده‌های ورودی و خروجی که جهت آموزش شبکه عصبی به کار می‌روند بسیار حساس است. |
| [۱۲] | بهینه‌سازی اجتماع ذرات- الگوریتم انتخاب منفی | *انتخاب ویژگی بر مبنای نزدیک‌ترین فاصله *ماتریسی جدید ویژگی‌ها که تولید شده با ماتریس اولیه مقایسه می‌شود و اگر ماتریس جدید دارای دقت بیشتری باشد جایگزین ماتریس اولیه می‌شود. | *گیر افتادن در راه حل غیربهینه |
| [۱۳] | ترکیبی از شبکه عصبی مصنوعی و درخت تصمیم‌گیری | *استفاده از قوانین درخت تصمیم‌گیری *کاهش حجم ویژگی‌ها با استفاده از درخت تصمیم‌گیری *آموزش و آزمایش نمونه‌ها *کاهش خطا در هر مرحله از آموزش | *افزایش زمان محاسباتی |
| [۱۴] | شبکه عصبی مصنوعی - بهینه‌سازی اجتماع ذرات | *آموزش و آزمایش نمونه‌ها بر مبنای مقایسه *جستجو در فضا برای انتخاب بهترین ویژگی‌ها | *گیر افتادن در راه حل غیربهینه |
| [۱۵] | درخت تصمیم‌گیری ماشین بردار پشتیبان شبکه عصبی مصنوعی تابع شعاعی پایه | *استفاده از قوانین تصمیم‌گیری *پیدا کردن مرز بهینه بین رده‌های هرزنامه و غیرهرزنامه | *افزایش زمان محاسباتی |
| [۱۶] | تکامل تفاضلی- الگوریتم انتخاب منفی | *کشف ویژگی‌های جدید بعد از هر دور مقایسه *احتمال تشخیص درست ویژگی‌ها به دلیل روش ادغام | *افزایش زمان محاسباتی |
| [۱۷] | بهینه‌سازی اجتماع ذرات- الگوریتم انتخاب منفی | * انتخاب ویژگی بر مبنای نزدیک‌ترین فاصله | *گیر افتادن در راه حل غیربهینه |

تولید مثل به عنوان راهکاری برای نزدیک شدن به ویژگی‌ها بهینه استفاده می‌شود.

لذا، فضای مسئله را براساس ویژگی‌های مجموعه داده ایجاد می‌کنیم. سپس عملیات پیش پردازش بر روی نمونه‌های ماتریس انجام می‌گیرد. هدف پیش پردازش این است که نمونه‌هایی با مقدار تهی و مقدار نامعتبر شناسایی شوند و یک مقدار مناسب با توجه به بازه ویژگی‌ها برای آن‌ها تعیین شود. به عبارتی نرمال‌سازی بر روی مقدار نمونه‌ها انجام گیرد.

۳-۱- انتخاب ویژگی

مجموعه داده Spambase دارای ویژگی‌های گمراه‌کننده، نامرتب، زائد و نوبه‌دار می‌باشد که در صحت عملکرد و زمان محاسباتی تأثیر منفی دارند. برای حل

جمعیت اولیه بر مبنای نمونه‌های مجموعه داده Spambase ایجاد می‌شوند. مراحل الگوریتم بهینه‌سازی علف‌های هرز برای انتخاب ویژگی را می‌توان به طور خلاصه به صورت زیر بیان کرد:

الف) تولید جمعیت اولیه (بر مبنای مجموعه داده) و ارزیابی تابع هدف آن‌ها. یک جمعیت اولیه در فضای حل مسئله بر مبنای نمونه‌های اولیه پراکنده و سپس ارزیابی می‌شوند. از محاسبه فاصله بین ویژگی‌ها برای ارزیابی استفاده می‌گردد. ویژگی‌هایی که فاصله کمتری دارند انتخاب می‌شوند.

ب) تولید مثل بر مبنای شایستگی و بروزسانی انحراف معیار. در مراحل بعدی ویژگی‌هایی انتخاب می‌شوند که به ویژگی‌های انتخاب شده در مراحل اولیه نزدیک‌تر باشند و همچنین مقدار برازندگی آن‌ها کمتر باشد. با این روش از

بیشترین تعداد تغییر می‌کند. تعداد دانه‌هایی که هر علف هرز می‌تواند تولید کند طبق معادله (۳) محاسبه می‌شود. در معادله (۱)، پارامتر f_{min}^k حداقل تابع هزینه کلونی در تکرار k است و پارامتر f_{max}^k حداکثر تابع هزینه کلونی در تکرار k است.

(۳)

$$seed_i^k = \left[S_{max} - (S_{max} - S_{min}) \frac{f_i^k - f_{min}^k}{f_{max}^k - f_{min}^k} \right]$$

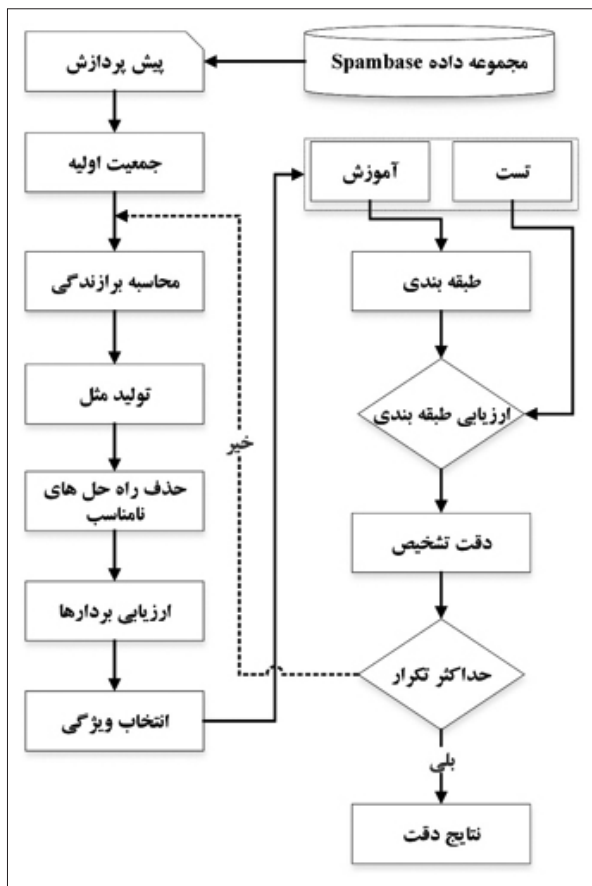
طبق معادله (۳)، برای هر کدام از ویژگی‌های مجموعه داده Spambase بر مبنای فضای مسئله، فاصله بهینه پیدا می‌شود. در ابتدا یک ویژگی انتخاب می‌شود و سپس مقدار آن ویژگی بر مبنای معادله (۳)، با مقدار ویژگی‌های دیگر مقایسه می‌گردد. در انتها در بین مقادیر به دست آمده از محاسبه فاصله بین ویژگی‌ها، یک ویژگی به عنوان نزدیک‌ترین ویژگی از لحاظ مقدار، انتخاب می‌شود. در شکل (۲)، بردار تولید دانه‌ها و انتخاب دانه‌های بهینه از بردار دانه‌ها نشان داده شده است.

از الگوریتم بهینه‌سازی علف‌های هرز مهاجم برای انتخاب بهترین ویژگی‌ها به منظور بالابردن دقت تشخیص و کاهش زمان محاسباتی استفاده شده است. به منظور ارزیابی دانه‌ها در الگوریتم بهینه‌سازی علف‌های هرز مهاجم از معیار برازندگی استفاده می‌کنیم. معیار صحت اصلی‌ترین معیار ارزیابی مربوط به طبقه‌بندی و تشخیص نمونه‌ها در فاز آموزش است. سازماندهی نمونه‌های هر رده به درستی و یا نادرستی دسته‌بند بستگی دارد که به صورت درست مثبت، درست منفی، کاذب منفی و کاذب مثبت تعریف می‌شود.

$$accuracy = \frac{TP + TN}{TP + TN + FN + FP} \quad (۵)$$

۳-۲ طبقه‌بندی: K نزدیک‌ترین همسایه

بعد از انتخاب ویژگی‌ها، نمونه‌ها را به دو مجموعه آموزش و آزمایش تقسیم می‌کنیم. برازندگی دانه‌ها با استفاده از داده‌های آزمایش ارزیابی می‌شود. در طی



شکل ۱: روندنمای مدل پیشنهادی

این مشکل می‌توان از انتخاب زیرمجموعه ویژگی بهینه استفاده نمود. انتخاب ویژگی مدل را ساده‌تر می‌کند، این کار باعث می‌شود که هزینه محاسبات کاهش یابد. با حذف ویژگی‌های غیرمفید، مدل شفاف‌تر و جامع‌تر می‌گردد. همچنین باعث تسریع در فرآیند یادگیری، کاهش فضای ذخیره‌سازی و بهبود کارایی مانند صحت می‌شود. در نتیجه وجود الگوریتم انتخاب ویژگی برای کاهش ابعاد داده‌ها در داده‌های با ابعاد بالا ضروری است.

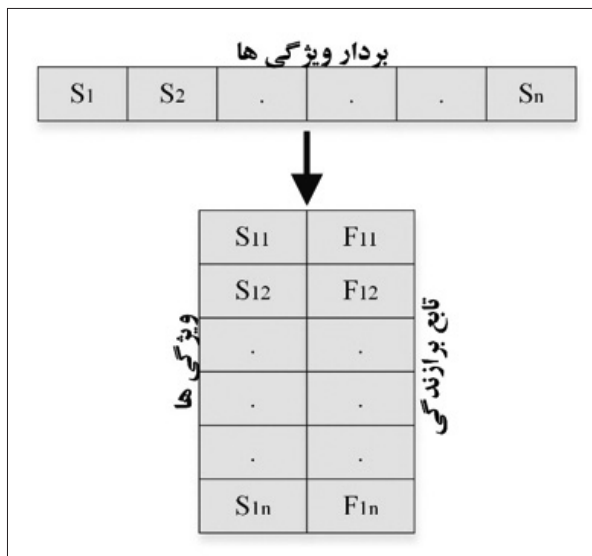
الگوریتم بهینه‌سازی علف‌های هرز مهاجم یکی از الگوریتم‌های بهینه‌سازی جدید و توانمند است که از قابلیت تطابق‌پذیری و تصادفی بودن کلونی علف‌های هرز، الهام گرفته شده است. هر عضو از جمعیت بر طبق توانایی‌اش می‌تواند بین دو مقدار تعیین شده حداکثر (S_{max}) و حداقل (S_{min}) تولید دانه کند. تعداد دانه‌هایی که هر گیاه می‌تواند تولید کند به طور خطی از کمترین تعداد دانه ممکن تا

اقلیدسی بین نقاط محاسبه شد، با مرتب‌سازی عناصر برحسب فاصله اقلیدسی، از میان k همسایه، برچسبی که از میان k همسایه دارای اکثریت است به نمونه ناشناخته داده می‌شود.

نمونه‌های آزمایشی، نمونه‌هایی هستند که از قبل به مدل داده نمی‌شوند، بلکه هدف نمونه‌های آزمایشی، ارزیابی نمونه‌های آموزشی است. برای مثال طبق جدول (۲)، فرض کنید که ۱۰ نمونه آموزشی و یک نمونه آزمایشی داریم. هدف این است که برای نمونه آزمایش یک رده بر مبنای فاصله اقلیدسی و با استفاده از الگوریتم K نزدیک‌ترین همسایه مشخص شود. در مرحله اول باید فاصله بین نمونه‌ها محاسبه شود و سپس عمل رتبه‌بندی بر مبنای مقدار فاصله انجام گیرد. سپس نمونه آزمایش بر مبنای تعداد k به رده Y یا N اختصاص داده می‌شود. اگر $k=3$ باشد در این حالت $[Y=1]$ ، $[Y=2]$ و $[Y=3]$ است. لذا رده نمونه آزمایشی برابر Y است. زیرا تعداد Y بیشتر از N است.

۴. ارزیابی و نتایج

در این بخش ارزیابی و نتایج مدل پیشنهادی بر روی مجموعه داده Spambase با ۴۶۰۱ نمونه و ۵۷ ویژگی که شامل ۳۹٫۴٪ ایمیل هرزنامه و ۶۰٫۶٪ ایمیل غیرهرزنامه است در محیط برنامه‌نویسی ویتوال سی شارپ ۲۰۱۷ انجام شده است. مجموعه داده Spambase شامل ۵۷ ویژگی عددی است. بعضی از این ویژگی‌ها مقدار نزدیک به همدیگر دارند و اگر ویژگی‌های مشابه به الگوریتم طبقه‌بندی داده شوند، دقت تشخیص بیشتر خواهد بود. زیرا دسته‌ها، شامل نمونه‌های دقیق‌تری خواهند بود. منظور این است که فاکتورهای درست مثبت، درست منفی، کاذب مثبت و کاذب منفی در دسته‌های متعلق به خودشان قرار می‌گیرند. ویژگی‌ها توسط الگوریتم بهینه‌سازی علف‌های هرز مهاجم انتخاب می‌شوند و ویژگی‌هایی که در تکرارهای برنامه بیشترین دقت صحت را داشته باشند



شکل ۲: انتخاب دانه بهینه از بردار دانه‌ها

فرآیند تکامل الگوریتم، دانه‌ها به سمت نقاط بهینه همگر می‌شوند. با استفاده از این روش تعداد ویژگی کاهش می‌یابد و بهینه‌ترین ویژگی‌ها را انتخاب می‌کنیم. در آخرین تکرار بهترین دانه از نظر معیار نزدیکی به عنوان دانه نهایی که نشان دهنده بهترین دانه است انتخاب بهینه شود. مدل K نزدیک‌ترین همسایه در اغلب موارد برای طبقه‌بندی به کار می‌رود، هرچند که می‌توان از آن برای تخمین و پیش‌بینی نیز استفاده نمود. در مدل K نزدیک‌ترین همسایه یک نمونه طبقه‌بندی نشده ممکن است به‌سادگی با مقایسه آن با شبیه‌ترین نمونه‌ها در مجموعه آموزشی یافت شود. بنابراین لازم است معیاری را برای تعیین فاصله بین نمونه‌ها مشخص نماییم. اگر یک بردار ویژگی به صورت $\langle a_1(x), a_2(x), \dots, a_n(x) \rangle$ داشته باشیم از فاصله اقلیدسی طبق معادله (۴) برای به‌دست آوردن فاصله بین دو ویژگی x_i و x_j استفاده می‌کنیم.

$$d(x_i, x_j) = \sqrt{\sum_{r=1}^n (a_r(x_i) - a_r(x_j))^2} \quad (4)$$

در روش K نزدیک‌ترین همسایه، یک دسته باید شامل K نمونه از مجموعه نمونه‌های آموزشی باشد، به طوری که مشابه‌ترین نمونه‌ها در یک دسته باشند و براساس برتری دسته یا برچسب مربوط به آن‌ها در مورد نمونه آزمایشی جدید تصمیم‌گیری می‌شود. در نتیجه، پس از آن‌که فاصله

جدول ۲: نمایش نمونه‌های آموزش و آزمایش

| رتبه‌بندی | فاصله | رده | b | a | ردیف | داده‌ها |
|-----------|-------------------------------------|-----|---|----|------|-----------------|
| ۱ | $\text{SQRT}(12-11)^2+(4-5)^2=1.41$ | Y | ۵ | ۱۱ | ۱ | داده‌های آموزشی |
| ۲ | $\text{SQRT}(12-12)^2+(4-6)^2=1.41$ | Y | ۶ | ۱۲ | ۲ | |
| ۷ | $\text{SQRT}(12-13)^2+(4-7)^2=3.16$ | N | ۷ | ۱۳ | ۳ | |
| ۵ | $\text{SQRT}(12-15)^2+(4-4)^2=3$ | N | ۴ | ۱۵ | ۴ | |
| ۴ | $\text{SQRT}(12-14)^2+(4-5)^2=2.23$ | N | ۵ | ۱۴ | ۵ | |
| ۸ | $\text{SQRT}(12-12)^2+(4-8)^2=4$ | Y | ۸ | ۱۲ | ۶ | |
| ۳ | $\text{SQRT}(12-13)^2+(4-5)^2=1.41$ | Y | ۵ | ۱۳ | ۷ | |
| ۱۰ | $\text{SQRT}(12-16)^2+(4-6)^2=4.47$ | N | ۶ | ۱۶ | ۸ | |
| ۶ | $\text{SQRT}(12-12)^2+(4-7)^2=3$ | N | ۷ | ۱۲ | ۹ | |
| ۹ | $\text{SQRT}(12-13)^2+(4-8)^2=4.12$ | Y | ۸ | ۱۳ | ۱۰ | |
| - | - | Y | ۴ | ۱۲ | ۱ | داده آزمایش |

جدول ۳: ارزیابی مدل پیشنهادی با انتخاب ویژگی‌های مختلف

| تعداد ویژگی | ویژگی‌های انتخاب شده | مقدار صحت |
|-------------|---|-----------|
| ۲۸ | ۱۳، ۱۷، ۲۶، ۱۱، ۳، ۲۷، ۲۴، ۵، ۵۶، ۱۹، ۱۶، ۲، ۳۰، ۱۸، ۶، ۵۵، ۲۰، ۸، ۵۷، ۳۷، ۲۵، ۲۳، ۲۱، ۵۳، ۵۲، ۷، ۳۵، ۹ | ۸۹، ۰۳ |
| ۳۵ | ۲۲، ۵۰، ۴۸، ۱۰، ۲۸، ۴، ۱۵، ۱۲، ۳۵، ۹، ۲، ۳۰، ۱۸، ۲۱، ۵۳، ۵۲، ۷، ۲۶، ۱۱، ۶، ۵۵، ۲۰، ۸، ۵۷، ۳۷، ۱۷، ۲۵، ۲۳، ۳، ۲۷، ۲۴، ۵، ۵۶، ۱۹، ۱۶ | ۹۰، ۲۳ |
| ۴۲ | ۵۶، ۱۹، ۱۶، ۲۵، ۲۳، ۲۱، ۵۳، ۵۲، ۲۲، ۷، ۳۳، ۱۴، ۴۵، ۳۲، ۴۸، ۳۱، ۱۳، ۱۰، ۴۶، ۲۹، ۱، ۲۸، ۱۵، ۱۲، ۳۵، ۹، ۲، ۳۰، ۱۸، ۶، ۵۵، ۲۰، ۸، ۵۷، ۳۷، ۲۴، ۵ | ۹۰، ۰۶ |
| ۴۷ | ۳۳، ۴۲، ۴۰، ۳۲، ۳۱، ۱۳، ۱۰، ۴۶، ۲۹، ۱، ۲۸، ۱۵، ۱۲، ۳۵، ۵۲، ۷، ۳۹، ۳۶، ۲۲، ۱۴، ۴، ۵۴، ۹، ۲، ۳۰، ۱۸، ۶، ۵۵، ۲۰، ۸، ۵۷، ۳۷، ۱۷، ۲۶، ۱۱، ۳، ۲۷، ۲۴، ۵، ۵۶، ۱۹، ۲۶، ۲۵، ۲۳، ۲۱، ۵۳، ۴۵ | ۹۲، ۸۷ |
| ۵۲ | ۵، ۳۹، ۳۶، ۲۲، ۱۴، ۴، ۲۵، ۲۳، ۲۱، ۵۳، ۵۲، ۷، ۴۹، ۴۸، ۴۴، ۴۳، ۴۱، ۵۴، ۴۵، ۳۳، ۴۲، ۴۰، ۳۲، ۳۱، ۱۳، ۱۰، ۴۶، ۲۹، ۱، ۵۶، ۲۸، ۱۹، ۱۶، ۱۵، ۱۲، ۳۵، ۹، ۲، ۳۰، ۱۸، ۶، ۵۵، ۲۰، ۸، ۵۷، ۳۷، ۱۷، ۲۶، ۱۱، ۳، ۲۷، ۲۴ | ۹۱، ۳۷ |

الگوریتم الگوریتم بهینه‌سازی علف‌های هرز مهاجم به انتخاب می‌شود. تعداد ۲۰۰ تکرار اجرا می‌گردد. همان‌طور که مشاهده می‌کنید تعداد تکرار الگوریتم بهینه‌سازی علف‌های هرز مهاجم در دقت تشخیص خیلی موثر است و با افزایش تعداد تکرار مقدار معیار صحت بیشتر شده است. دلیل این‌که با تعداد تکرار، مقدار صحت افزایش می‌یابد این است که الگوریتم بهینه‌سازی علف‌های هرز مهاجم در یافتن ویژگی‌ها در فضای جستجو دقت بیشتری حاصل می‌کند و از گیرافتادن در بهینگی محلی فرار می‌کند.

در شکل (۳)، نمودار مقایسه تعداد تکرار در الگوریتم بهینه‌سازی علف‌های هرز مهاجم و مقدار K در مدل K نزدیک‌ترین همسایه نشان داده شده است. در شکل (۳)،

حداکثر تعداد تکرار در الگوریتم بهینه‌سازی علف‌های هرز مهاجم برابر ۲۰۰ است و مقدار K در K نزدیک‌ترین همسایه با سه مقدار مختلف تست شده است. در ابتدا مدل پیشنهادی را با انتخاب ویژگی‌های مختلف ارزیابی کردیم و با هر تعداد ویژگی نتایج متفاوتی حاصل شده است. در جدول (۳)، ارزیابی مدل پیشنهادی با انتخاب ویژگی‌ها مختلف بر مبنای معیار صحت نشان داده شده است.

در جدول (۴)، مدل پیشنهادی با تعداد تکرار الگوریتم بهینه‌سازی علف‌های هرز مهاجم و مقدار K در مدل K نزدیک‌ترین همسایه برای ۵۷ ویژگی ارزیابی شده است.

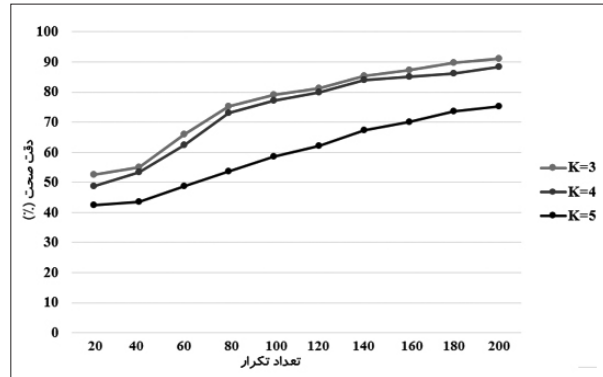
جدول ۴: ارزیابی مدل پیشنهادی با فاکتورهای تعداد تکرار و مقدار K

| تعداد تکرار | مقدار K | | |
|-------------|---------|-------|-------|
| | K=۵ | K=۴ | K=۳ |
| ۲۰ | ۴۲,۳۱ | ۴۸,۶۷ | ۵۲,۶۴ |
| ۴۰ | ۴۳,۴۰ | ۵۳,۴۱ | ۵۵,۰۲ |
| ۶۰ | ۴۸,۶۲ | ۶۲,۴۰ | ۶۵,۹۴ |
| ۸۰ | ۵۳,۶۱ | ۷۳,۰۱ | ۷۵,۱۳ |
| ۱۰۰ | ۵۸,۵۰ | ۷۷,۲۵ | ۷۹,۰۶ |
| ۱۲۰ | ۶۲,۰۴ | ۷۹,۸۴ | ۸۱,۳۵ |
| ۱۴۰ | ۶۷,۲۸ | ۸۳,۹۷ | ۸۵,۳۰ |
| ۱۶۰ | ۷۰,۰۰ | ۸۵,۰۹ | ۸۷,۱۸ |
| ۱۸۰ | ۷۳,۴۹ | ۸۶,۰۴ | ۸۹,۶۲ |
| ۲۰۰ | ۷۵,۳۲ | ۸۸,۳۲ | ۹۱,۱۱ |

جدول ۵: مقایسه مدل پیشنهادی با مدل‌های دیگر بر مبنای صحت

| مدل‌ها | مقدار صحت |
|--|-----------|
| بهینه‌سازی اجتماع ذرات-الگوریتم انتخاب منفی [۱۷] | ۹۱,۲۲ |
| بهینه‌سازی اجتماع ذرات [۱۷] | ۸۱,۳۲ |
| الگوریتم انتخاب منفی [۱۷] | ۶۸,۸۶ |
| بیز ساده [۱۸] | ۷۹,۳ |
| ماشین بردار پشتیبان [۱۹] | ۹۰ |
| انتخاب ویژگی متمایز-ماشین بردار پشتیبان [۲۰] | ۷۱ |
| شبکه عصبی مصنوعی [۲۱] | ۸۶ |
| مجموعه فازی [۲۲] | ۸۶,۹ |
| مدل پیشنهادی | ۹۱,۱۱ |

هرزنامه، مدل‌هایی طراحی و پیاده‌سازی کرد که توانایی تشخیص مفید و نامناسب بودن اطلاعات را داشته باشند. معمولاً شاهد بوده‌ایم که الگوریتم‌های یادگیری ماشین در تشخیص و طبقه‌بندی کارایی و دقت بالایی داشته‌اند. در این مقاله برای طبقه‌بندی ایمیل هرزنامه از مدل ترکیبی الگوریتم بهینه‌سازی علف‌های هرز مهاجم و K نزدیک‌ترین همسایه استفاده شد. ما نشان دادیم که در مدل پیشنهادی با تعداد ویژگی‌های مختلف، مقدار صحت نتایج متفاوتی دارد. برای مدل پیشنهادی مقدار صحت برابر ۹۱,۱۱٪ به دست آمد و در مقایسه با مدل‌هایی مانند ماشین بردار پشتیبان



شکل ۳: نمودار مقایسه تعداد تکرار در الگوریتم بهینه‌سازی علف‌های هرز مهاجم و مقدار K در مدل K نزدیک‌ترین همسایه

نمایان است که مقدار K=۳ در مقایسه با K=۴ و K=۵ بیشتر است. دلیل این‌که درصد معیار صحت در K=۳ بیشتر است برای این‌که فاصله بین ویژگی‌ها برای طبقه‌بندی بهتر شناسایی می‌شود.

در جدول (۵)، مقایسه مدل پیشنهادی با مدل‌های دیگر نشان داده شده است. همان‌طور که مشاهده می‌کنید مدل پیشنهادی در مقایسه با مدل‌های بهینه‌سازی اجتماع ذرات، الگوریتم انتخاب منفی، بیز ساده، ماشین بردار پشتیبان، انتخاب ویژگی متمایز-ماشین بردار پشتیبان، شبکه عصبی مصنوعی و فازی مقدار صحت بیشتری دارد. در میان مدل‌های جدول (۵)، مقدار صحت در مدل ماشین بردار پشتیبان به مدل پیشنهادی نزدیک‌تر و صحت در مدل بهینه‌سازی اجتماع ذرات-الگوریتم انتخاب منفی با اندکی جزئی بیشتر از مدل پیشنهادی است.

در جدول (۴) نشان دادیم که مقدار صحت با مقادیر مختلف K نتایج متفاوتی دارد. طبق نتایج این مقاله، به این نتیجه دست یافتیم که تعداد ویژگی‌ها، مهم‌ترین فاکتور در افزایش دقت صحت می‌باشند. انتخاب ویژگی‌های موثر به الگوریتم K نزدیک‌ترین همسایه کمک می‌کنند تا بر مبنای فاصله بتواند دسته‌بندی دقیقی انجام دهد.

۵. نتیجه‌گیری و کارهای آینده

در دنیای امروزی با پیشرفت علوم کامپیوتری باید برای موارد ناشناخته در محیط‌های برخط و هجوم ایمیل

its application to spam e-mail detection, *Expert Systems with Applications*, Vol. 37, 7976-7985, 2010

A.R. Behjat, A. Mustapha, H. Nezamabadi-pour, Md. Nasir Sulaiman, and N. Mustapha, A PSO-Based Feature Subset Selection for Application of Spam /Non-spam Detection, Springer-Verlag Berlin Heidelberg, M-CAIT 2013, CCIS 378, pp. 183-193, 2013.

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.

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, 97-110, 2014

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

Y. Zhang, H.Y. Li, M. Niranjan, P. Rockett, Applying cost-sensitive multi objective genetic programming to feature extraction for spam e-mail filtering, *Lecture Notes in Computer Science, Genetic Programming, Berlin/Heidelberg, Springer*, Vol. 4971, pp. 325-336, 2008

T. Fagbola, S. Olabiyisi, A. Adigun, Hybrid GA-SVM for efficient feature selection in e-mail classification, *Comput. Eng. Intell. Syst.*, 3 (3), 2012.

A.K. Uysal, S. Gunal, A novel probabilistic feature selection method for text classification, *Knowl.-Based Syst.*, Vol. 36, pp. 226-235, 2012.

L. Ozgur, T. Gungor, F. Gurgen, Spam mail detection using artificial neural network and Bayesian filter, in: Z. Yang, H. Yin, R. Everson (Eds.), *Intelligent Data Engineering, and Automated Learning- IDEAL 2004*, Springer, Berlin/Heidelberg, 2004, pp. 505-510, 2004.

R. Ariaeinejad, A. Sadeghian, Spam detection system: A new approach based on interval type-2 fuzzy sets, in: 24th Canadian Conference on Electrical and Computer Engineering (CCECE, 2011), 2011.

و بیز ساده دقت بیشترى داشت. مهم‌ترین مشکل در تشخیص ایمیل هرزنامه، انتخاب ویژگی‌های مناسب از میان انبوهی ویژگی‌ها است که در آینده امیدواریم این مشکل را با الگوریتم‌های متنوع فرابتکاری بهینه‌تر کنیم و به دقت ۱۰۰٪ دست یابیم.

مراجع

Q. Fu, B. Feng, D. Guo, Q. Li, Combating the evolving spammers in online social networks, *Computers & Security*, Vol. 72, pp. 60-73, 2018.

Y. Cohen, D. Hendler, A. Rubin, Detection of malicious web-mail attachments based on propagation patterns, *Knowledge-Based Systems*, in press, corrected proof, Available online 11 November 2017.

S. Rathore, V. Loia, J.H. Park, SpamSpotter: An efficient spammer detection framework based on intelligent decision support system on Facebook, *Applied Soft Computing*, In press, corrected proof, Available online 22 September 2017.

L. Nguyen; A.Q. Tran; L.T. 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.

S. Salehi, A. Selamat; M. Bostanian, Enhanced genetic algorithm for spam detection in email, *IEEE 2nd International Conference on Software Engineering and Service Science*, pp. 594-597, 2011

M. Takesue, Cascaded Simple Filters for Accurate and Lightweight Email-Spam Detection, *Fourth International Conference on Emerging Security Information, Systems and Technologies*, pp. 160-165, 2010.

W. Ma, D. Tran; D. Sharma, A Novel Spam Email Detection System Based on Negative Selection, *Fourth International Conference on Computer Sciences and Convergence Information Technology*, IEEE, pp. 987-992, 2009

<https://archive.ics.uci.edu/ml/datasets/Spambase>

A.R. Mehrabian, C. Lucas, A novel numerical optimization algorithm inspired from weed colonization, *Ecol. Inform.* 1(4): 355-366, 2006

Martin, *Instance-Based Learning: Nearest Neighbor with Generalization*, Doctoral dissertation, University of Waikato, 1995

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.

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

Mu-Chun Su, Hsu-Hsun Lo, Fu-Hau Hsu, A neural tree and