

انتخاب ویژگی با الگوریتم بهینه‌سازی حاصلخیزی زمین‌های کشاورزی برای تشخیص صفحات وب هرز

محمد سخی دل هوسین

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

فرهاد سلیمانان قره چپق*

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

چکیده

FFANB کاهش ویژگی‌ها به منظور افزایش صحت با استفاده از الگوریتم حاصلخیزی زمین‌های کشاورزی می‌باشد که از مجموعه داده WEBSpam-UK2007 که از معتبرترین مجموعه داده در زمینه شناسایی صفحات وب هرز می‌باشد استفاده شده است. این مجموعه داده شامل سه دسته ویژگی با عناوین ویژگی‌های مبتنی بر محتوا (۹۶ ویژگی)، ویژگی‌های مبتنی بر پیوند (۴۱ ویژگی) و ویژگی‌های مبتنی بر پیوند تبدیل یافته (۱۳۸ ویژگی) می‌باشد که تعداد کل ویژگی‌ها برابر با ۲۷۵ ویژگی است. نتایج ارزیابی‌های صورت گرفته بر روی مدل FFANB نشان دهنده درصد دقت ۰/۹۲۴۱ و صحت ۰/۹۵۸۴ می‌باشند که حاکی از برتری مدل FFANB در مقایسه با بسیاری از روش‌های پیشین می‌باشد. واژه‌های کلیدی: صفحات وب هرز، طبقه‌بندی، الگوریتم حاصلخیزی زمین‌های کشاورزی، الگوریتم بیز ساده، انتخاب ویژگی

۱. مقدمه

هر وبگاه شامل مجموعه‌ای از صفحات وب متصل به هم است که شامل اطلاعات خاصی در مورد یک مبحث

در فضای اینترنت، امکان به‌کارگیری انواع سرویس‌ها و خدمات متعدد برای کاربران مهیا شده است. همزمان با رشد و گسترش استفاده از اینترنت، تعداد هرزنویسان وب افزایش یافته است. صفحات وب هرز به اشکال مختلفی چون تبلیغات تجاری و ویروس‌هایی نهان شده در صفحات وب جایگذاری می‌شود. صفحات وب هرز علاوه بر تهدید امنیت کاربران در وب، موجب هدر رفتن منابع سیستم و ایجاد ترافیک مخرب نیز می‌گردند؛ لذا ارایه راهکارهایی جهت مقابله با وب هرز ضروری به نظر می‌رسد. یکی از روش‌های شناسایی و مقابله با صفحات وب هرز، طبقه‌بندی صفحات با استفاده از الگوریتم‌های یادگیری ماشینی است. در این مقاله، مدلی جدید بر مبنای الگوریتم حاصلخیزی زمین‌های کشاورزی و بیز ساده با عنوان 'FFANB' برای تشخیص صفحات وب هرز پیشنهاد شده است. در مدل FFANB از الگوریتم حاصلخیزی زمین‌های کشاورزی برای انتخاب ویژگی و بیز ساده برای طبقه‌بندی نمونه‌ها استفاده شده است. هدف مدل

* نویسنده مسئول

1- Farmland Fertility Algorithm Naive Bayes (FFANB)

می‌باشند و معمولاً روی یک کارساز قرار دارند. وبگاه‌ها معمولاً انبوهی از اطلاعات را در قالب متن، تصویر، صدا، و فیلم در اختیار کاربران قرار می‌دهند. مهم‌ترین چالش اصلی برای موتورهای جستجوگر در جستجوی مطالب و ارائه خدمات کارآمد به کاربران، صفحات هرز می‌باشند که نقشی تخریب‌کننده در نتایج موتورهای جستجوگر از نظر دقت، سرعت و فضای ذخیره‌سازی دارند. هرز در صفحات وب چالش بزرگی است که هم بر روی کاربران و هم بر روی تامین‌کنندگان سرویس تاثیرگذار است. کاربران زمانی که در حال انجام فعالیت‌های خود مانند جستجوی مبتنی بر وب هستند با وب هرز مواجه می‌شوند و تشخیص اطلاعات واقعی برای آن‌ها مشکل می‌گردد. در واقع موتورهای جستجوگر هرچه قدر هم که قوی و کارآمد باشند به دلیل بالا بودن حجم و تنوع موضوعی اطلاعات و منابع موجود در شبکه اینترنت قادر به سازماندهی تمامی منابع این شبکه نیستند و لذا تشخیص وب هرز در میان انبوهی از اطلاعات سخت است. موتورهای جستجوگر باید صفحات وب را پیدا کنند و محتوای آن‌ها را بر مبنای کلمات وب هرز تجزیه کنند تا صفحات وب هرز را تشخیص دهند [۱،۲].

هدف از راه‌اندازی موتورهای جستجوگر، کمک به کاربران در پیدا کردن اطلاعات موجود در صفحات وب می‌باشد. در سیستم موتورهای جستجوگر با وارد کردن کلمات کلیدی توسط کاربران، لیستی از صفحات وبگاه‌ها ظاهر می‌شوند که به موضوع مورد علاقه کاربر مرتبط می‌باشند. موتورهای جستجوگر در واقع این اطلاعات را در میان صفحات وبگاه‌ها جستجو کرده و در مدت زمان کوتاه و با بیشترین سرعت در اختیار کاربران قرار می‌دهند. اما نکته مهم این است که همزمان با اطلاعات استخراج شده یکسری از صفحات وب هرز بازایی می‌شوند که باید شناسایی شوند و صفحات آن‌ها باز نشود. برای این کار باید کلمات هرز تشخیص داده شوند و اگر یک صفحه حاوی آن کلمات باشد از بازایی آن اجتناب شود [۳،۴]. مسئله انتخاب ویژگی در بسیاری از کاربردها (مانند

طبقه‌بندی) اهمیت ویژه‌ای دارد، زیرا در این کاربردها تعداد زیادی ویژگی وجود دارد، که بسیاری از آن‌ها یا بی‌استفاده هستند و یا این‌که بار اطلاعاتی چندانی ندارند. حذف نکردن این ویژگی‌ها مشکلی از لحاظ اطلاعاتی ایجاد نمی‌کنند ولی بار محاسباتی را افزایش می‌دهند. با افزایش تعداد ویژگی‌ها، فضای ویژگی‌ها نیز افزایش می‌یابد، تجزیه و تحلیل داده‌ها و طبقه‌بندی نیز به‌طور قابل توجهی سخت‌تر می‌شود. به‌علاوه، داده‌ها به‌طور فزاینده‌ای در فضایی که اشغال کرده‌اند پراکنده می‌شوند که منجر به مشکلات بزرگی هم برای الگوریتم‌های بانظارت و هم بدون نظارت می‌شود. این پدیده به عنوان مشکل ابعاد شناخته شده است و براساس این واقعیت است که اغلب کارکردن با داده‌های با ابعاد بالا مشکل است. تعداد زیاد ویژگی‌ها می‌تواند اختلال در داده‌ها را افزایش دهد و در نتیجه خطای الگوریتم یادگیری نیز افزایش می‌یابد، بخصوص اگر تعداد نمونه‌ها در مقایسه با تعداد ویژگی‌ها کم باشد [۵،۶].

در مقالات پیشین در زمینه تشخیص ایمیل هرزنامه از الگوریتم‌های مختلف فرا ابتکاری و طبقه‌بندی استفاده کردیم [۷،۸]. برای طبقه‌بندی نمونه‌ها از الگوریتم‌هایی مانند k نزدیکترین همسایه و آدابوست استفاده کردیم. دقت تشخیص مدل‌های پیشین که مبتنی بر الگوریتم‌های فرا ابتکاری بودند بسیار بالا و دقیق بود. لذا در این مقاله هم از الگوریتم فرا ابتکاری به دلیل قدرت بالا در یافتن راه حل‌های بهینه جهت تشخیص صفحات وب هرز استفاده می‌کنیم و مدلی جدید بر روی مجموعه داده وب هرز ارزیابی می‌شود. طبق نتایج مشخص شد که الگوریتم بهینه‌سازی حاصلخیزی زمین‌های کشاورزی که از دسته الگوریتم‌های فرا ابتکاری است در تشخیص صفحات وب هرز کارایی خوبی دارد. در تشخیص ایمیل هرزنامه از مجموعه داده معتبر Spambase با ۵۷ ویژگی استفاده کردیم. اما در این مقاله از مجموعه داده‌های وب هرز که بر مبنای ویژگی‌های وب طراحی شده‌اند استفاده می‌کنیم. در مجموعه داده‌های وب هرز، تعداد ویژگی‌ها، نوع داده‌ها و مقدار داده‌ها متفاوت از ایمیل هرزنامه است.

در این مقاله، مدلی جدید بر مبنای الگوریتم بهینه‌سازی حاصلخیزی زمین‌های کشاورزی [۹] و بیز ساده [۱۰] با عنوان FFANB برای تشخیص صفحات وب هرز بر روی مجموعه داده WEBSpam-UK2007 [۱۱] پیشنهاد می‌شود. در مدل FFANB از الگوریتم بهینه‌سازی حاصلخیزی زمین‌های کشاورزی برای انتخاب ویژگی‌ها و از الگوریتم بیز ساده برای طبقه‌بندی نمونه‌ها استفاده می‌کنیم. الگوریتم بیز ساده یکی از روش‌های یادگیری ماشین برای طبقه‌بندی به شمار می‌آید. در این روش رده‌های مختلف، هر کدام به شکل یک فرضیه دارای احتمال در نظر گرفته می‌شوند. هر داده آموزشی جدید، احتمال درست بودن فرضیه‌های پیشین را افزایش و یا کاهش می‌دهد و در نهایت، فرضیاتی که دارای بالاترین احتمال باشند، به عنوان یک رده در نظر گرفته شده و برچسبی بر آن‌ها تخصیص داده می‌شود [۱۲].

۲. تحقیقات مرتبط

پژوهشگران برای تشخیص صفحات وب هرز از روش Smart-BT استفاده کرده‌اند [۱۳]. آن‌ها در این روش برای انتخاب ویژگی و کشف وابستگی بین ویژگی‌ها از آزمون کای‌دو و برای طبقه‌بندی از الگوریتم بیز ساده استفاده کرده‌اند. آزمون کای‌دو، میزان ارتباط یا وابستگی بین ویژگی‌ها را اندازه‌گیری می‌کند. آزمایش مربعی کای‌دو بر پایه یک آزمایش آماری (χ^2) می‌باشد. در این روش ارزش هر ویژگی با استفاده از محاسبه مقدار آماری کای‌دو و با توجه به رده مد نظر ارزیابی می‌گردد. این مقدار با استفاده از آزمون کای‌دو جزء آزمون‌های غیر پارامتری است محاسبه می‌شود. نتایج نشان داده که مدل Smart-BT در مقایسه با ماشین بردار پشتیبان، درخت تصمیم‌گیری و شبکه عصبی مصنوعی دقت بیشتری دارد. همچنین انتخاب ویژگی با استفاده از کای دو از رتبه بالاتری در مقایسه با الگوریتم‌های ژنتیک، و بهینه‌سازی اجتماع ذرات بهره‌مند است.

برای تشخیص صفحات وب هرز از روش‌های یادگیری ماشین استفاده شده است [۱۴]. ارزیابی بر روی مجموعه

داده‌های WEBSpam-UK2006 و WEBSpam-UK2007 انجام شده است. بیشترین درصد صحت متعلق به الگوریتم آدابوست است که درصد صحت آن برای مجموعه داده‌های WEBSpam-UK2006 و WEBSpam-UK2007 به ترتیب برابر با ۰/۹۳۷۰ و ۰/۸۵۲۰ است. همچنین ماشین بردار پشتیبان، شبکه عصبی مصنوعی چندلایه و بیز ساده از دقت تشخیص مناسب برخوردار هستند.

مدلی بر مبنای ترکیب الگوریتم بهینه‌سازی اجتماع ذرات و انتخاب ویژگی بر مبنای همبستگی برای تشخیص صفحات وب هرز پیشنهاد شده است [۱۵]. در مدل ترکیبی از انتخاب ویژگی بر مبنای همبستگی برای بهبود الگوریتم بهینه‌سازی اجتماع ذرات استفاده شده است. ارزیابی بر روی مجموعه داده WEBSpam-UK2006 انجام شده است. نتایج نشان داده است که بهترین مقدار F-Measure برابر با ۸۸ درصد است.

برای تشخیص صفحات وب هرز از روش‌های بیز ساده، درخت C5.0، ماشین بردار پشتیبان، بگینگ، آدابوست و فیلتر استفاده شده است [۱۶]. ارزیابی بر روی مجموعه داده WEBSpam-UK2007 انجام شده است. نتایج نشان داده است که درصد صحت مدل فیلتر برابر با ۰/۸۱۹۰ درصد است که در مقایسه با روش‌های بیز ساده، درخت C5.0، ماشین بردار پشتیبان، بگینگ و آدابوست از کارایی بهتری برخوردار است. شبکه باور عمیق که برگرفته از شبکه‌های عصبی مصنوعی است برای تشخیص صفحات وب هرز پیشنهاد شده است [۱۷]. ایده اصلی شبکه باور عمیق، تکرار یادگیری شبکه به منظور آموزش بهینه و افزایش دقت تشخیص است. ارزیابی بر روی مجموعه داده WEBSpam-UK2007 انجام شده است. نتایج نشان داده است که درصد صحت برابر با ۰/۹۱۲۶ درصد است.

مدلی بر مبنای شبکه عصبی مصنوعی چندلایه و ژنتیک برای طبقه‌بندی صفحات وب هرز بر روی مجموعه داده WEBSpam-UK2007 پیشنهاد شده است. در مدل ترکیبی از شبکه عصبی مصنوعی چندلایه برای آموزش نمونه‌ها و

از ژنتیک برنامه‌نویسی شده برای طبقه‌بندی استفاده شده است. نتایج نشان داده که درصد صحت برابر با ۰/۹۲ درصد است که در مقایسه با مدل‌های درصد بالاتری دارد [۱۸]. مدلی برمبنای ماشین بردار پشتیبان برای تشخیص صفحات وب هرز پیشنهاد شده است. ابتدا ماشین بردار پشتیبان با N نقاط داده آموزش داده شده که هرکدام قبلاً به عنوان وب هرز یا وب غیرهرز طبقه‌بندی شده‌اند. در این مقاله از ترکیب ویژگی‌ها استفاده شده است. درصد نتایج با ۲۹۶ ویژگی برابر با ۰/۸ است [۱۹].

۳. مدل FFANB

هدف وب هرزنویسان این است که رتبه صفحات وب را به روش وب هرز به‌طور غیرقانونی افزایش دهند. لذا هدف این است که این صفحات هرز شناسایی شوند و سپس حذف شوند. اگر این محتوا شناسایی نشود و نتایج جستجو فیلتر نشود، صفحات وب هرز می‌توانند دقت موتورهای جستجو را به شدت کاهش دهند. به‌طور کلی، وب هرزها قصد دارند به صورت غیرقانونی مرتباً رتبه صفحات هرز خود را از طریق موتورهای جستجوگر افزایش دهند که ممکن است منجر به جعل، تخریب اطلاعات و تحریف نتایج جستجو شود و از این طریق کل فرآیند جستجوی اطلاعات تحت تأثیر قرار می‌گیرد. با توجه به موفقیت ابزارهای ضد وب هرز ایمیل توسط الگوریتم‌های یادگیری ماشین، ما در این مقاله برای شناسایی صفحات وب هرز از روش‌های یادگیری ماشین استفاده می‌کنیم. به‌طور معمول، دقت تشخیص بالا و ضریب مثبت کاذب پایین، اصلی‌ترین ویژگی‌های لازم برای شناسایی وب هرز براساس روش‌های یادگیری ماشین می‌باشد. این امر به ویژه برای شناسایی صفحات وب هرز و اطمینان از عدم جعل وبگاه‌ها خیلی مهم و ضروری است.

مدل FFANB، ترکیبی از الگوریتم حاصلخیزی زمین‌های کشاورزی و بیز ساده است. در مدل FFANB از الگوریتم حاصلخیزی زمین‌های کشاورزی برای انتخاب ویژگی و

از بیز ساده برای طبقه‌بندی نمونه‌ها استفاده می‌شود. در این مقاله از مجموعه داده وب هرز WEBSAM-UK2007 استفاده شده است. این مجموعه داده توسط پژوهشگران برچسب‌گذاری شده است و هر صفحه با عنوان «صفحه هرز» یا «صفحه غیرهرز» مشخص شده است. در مدل FFANB در ابتدا مجموعه داده‌ها در قالب فایل متنی خوانده می‌شود. سپس با استفاده از روش‌های پیش پردازش داده، مقدار کمینه و بیشینه ویژگی‌ها به بازه ۰ و ۱ تبدیل می‌شود. برای توسعه مدل، در مرحله اول مقادیر هریک از ویژگی‌ها جهت طبقه‌بندی و افزایش دقت، نرمال‌سازی می‌شوند. جهت نرمال‌سازی به عنوان یکی از گام‌های پیش‌پردازش داده از روش کمینه-بیشینه طبق معادله (۱) استفاده می‌شود [۲۰].

$$x_n = \frac{(x_r - x_{min})}{(x_{max} - x_{min})} \quad (1)$$

در معادله (۱) x_n ، x_r ، x_{min} و x_{max} به ترتیب نشان‌دهنده مقادیر واقعی، استاندارد شده، حداکثر و حداقل داده‌های تحت بررسی هستند. قالب بردارها برمبنای مقدار ویژگی‌ها ایجاد می‌شود و توسط الگوریتم حاصلخیزی زمین‌های کشاورزی بهترین ویژگی‌ها در میان بردارها انتخاب می‌شوند.

۳-۱ انتخاب ویژگی

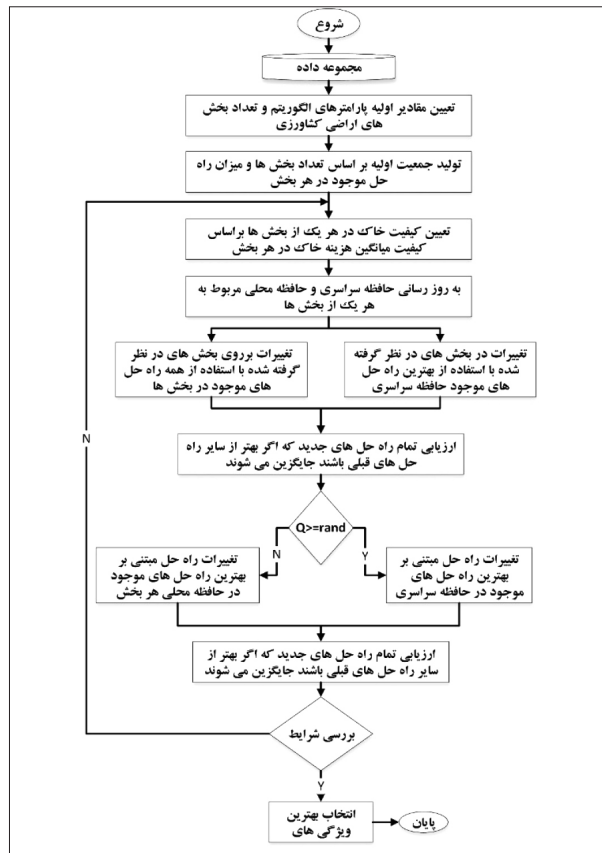
در الگوریتم حاصلخیزی زمین‌های کشاورزی در ابتدا فضا را به چهار بخش تقسیم می‌کنیم. مجموعه داده‌ها را در میان چهار بخش تقسیم می‌کنیم. یکی از بخش‌ها را انتخاب می‌کنیم و ویژگی‌ها مربوط به آن بخش را برمبنای میانگین فاصله ویژگی‌های یک رکورد ارزیابی می‌کنیم. ویژگی‌های انتخاب شده از یک رکورد برمبنای معیار فاصله انجام می‌شود. اگر میانگین فاصله در یک رکورد کمتر باشد به منزله این است که برزندگی ویژگی‌های انتخاب شده دقت بیشتری خواهند داشت. هر کدام از بخش‌ها شامل یک حافظه محلی و یک حافظه سراسری هستند که در واقع حافظه سراسری همان بهترین میانگین به دست آمده

در مدل FFANB ابتدا در مرحله اول مجموعه داده با تمام‌های ویژگی‌ها استخراج شده و جهت انتخاب بهترین ویژگی‌ها از الگوریتم حاصلخیزی زمین‌های کشاورزی استفاده می‌شود. هر موقعیت معادل یک آرایه به تعداد تمام ویژگی‌ها است. ابعاد فضای جستجو در موقعیت هر بخش از اراضی برابر با تعداد ویژگی‌ها است. هر خانه از آرایه می‌تواند عدد صفر و یا یک داشته باشد که به ترتیب به معنی آن است که یک ویژگی انتخاب نشده و یا انتخاب شده است. به طور مثال اگر ۱۵ ویژگی برای بررسی به الگوریتم حاصلخیزی زمین‌های کشاورزی داده می‌شود، فضای جستجو برای موقعیت هر بخش از اراضی ۱۵ بعدی می‌شود. هر بعد متناظر با یک ویژگی است. موقعیت هر بخش از اراضی با توجه به تابع سیگموئید به محدوده ۰ و ۱ نگاشت داده می‌شود. پس مقدار هر بعد از موقعیت با تابع سیگموئید به مقدار ۰ و ۱ تبدیل می‌شود که مقدار ۰ به معنی حذف ویژگی متناظر و مقدار ۱ به معنی حفظ ویژگی متناظر با آن بعد است.

$$S(x) = \frac{1}{1 + e^{-x}} \quad (2)$$

$$x_i = \begin{cases} 1, & \text{if } S(x) \geq rand \\ 0, & \text{if } S(x) < rand \end{cases} \quad (3)$$

برای تبدیل اعداد پیوسته به دودویی از تابع سیگموئید طبق معادله (۲)، استفاده می‌شود. خروجی تابع سیگموئید در یک محدوده عددی خاصی (عموماً بین صفر و یک) قرار می‌گیرد. در این تابع جواب ۰ یا ۱ نخواهد بود بلکه مجموعه اعدادی بین صفر و یک است. لذا در الگوریتم حاصلخیزی زمین‌های کشاورزی برای تبدیل حالت پیوسته به گسسته موقعیت هر بخش از اراضی x_i طبق معادله (۳)، تعریف می‌شود. اگر موقعیت هر بخش از اراضی به صورت معادله (۴) تعریف شود آنگاه موقعیت جدید هر بخش از اراضی بعد از اعمال سیگموئید به صورت معادله (۵) خواهد بود. سپس ویژگی‌های متناظر با ستون‌های غیرصفر، به طبقه‌بند اعمال می‌شوند.

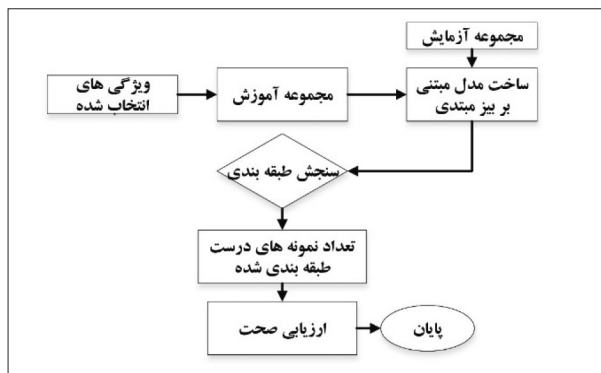


شکل ۱: روندنمای مدل FFANB به منظور انتخاب ویژگی

می‌باشد و حافظه محلی دربرگیرنده ارزیابی رکوردها به منظور نزدیک شده به حافظه سراسری می‌باشد.

در شکل (۱) روندنمای مدل FFANB بر مبنای انتخاب ویژگی توسط الگوریتم حاصلخیزی زمین‌های کشاورزی نشان داده شده است.

برای ارزیابی مدل FFANB، ابتدا لازم است مجموعه داده‌ها به دو بخش آموزش و آزمایش تقسیم شوند. داده‌های بخش آموزش مدل را تولید می‌کنند و داده‌های بخش آزمایش با کمک تعدادی رکورد، مدل تولید شده را آزمایش و برچسب مربوط به رکوردهای مذکور را تعیین می‌نمایند. نرخ صحت یک روش طبقه‌بندی بر روی مجموعه داده‌های آموزشی، درصد مشاهداتی از مجموعه آموزش است که به درستی توسط روش مورد استفاده طبقه‌بندی شده است. برای محاسبه این معیار داده‌های آزمون استفاده می‌شوند. همچنین می‌توان نرخ خطا را بر اساس معیار صحت محاسبه کرد.



شکل ۲: روندنمای مدل FFANB به منظور طبقه بندی

یادگیری طبقه بندی کمک می کند تا فرآیند یادگیری را به سمت رده اقلیت سوق دهد.

روش سطح داده: این روش با باز نمونه گیری از فضای داده باعث تغییر توزیع داده ها می شود به طوری که تغییری در الگوریتم یادگیری ایجاد نمی شود و تلاش می کند در مرحله پیش پردازش تأثیرات ناشی از ناتوانی را برطرف کند.

روش یادگیری حساس به هزینه: این روش مابین روش الگوریتمی و داده ای قرار دارد. به طوری که هم در سطح داده و هم در سطح الگوریتم تغییر ایجاد خواهد کرد. مهم ترین نقطه ضعف این رویکرد تعریف هزینه طبقه بندی نادرست می باشد که عموماً در مجموعه داده وجود ندارند. در قسمت طبقه بندی از الگوریتم بیز ساده برای طبقه بندی نمونه های صفحات وب هرز استفاده می شود. در شکل (۲) روندنمای مدل FFANB بر مبنای انتخاب ویژگی توسط الگوریتم بیز ساده نشان داده شده است.

یک مجموعه داده از نمونه های آموزشی با برچسب رده (نوع رده مشخص است) و یک مجموعه داده نمونه های آزمایش E با n مقدار ویژگی $(a_1, a_2, \dots, a_n | c)$ در نظر گرفته می شود و طبقه بندی بیز برای دسته بندی E به صورت معادله (۷) تعریف می شود.

$$c(E) = a \quad c \in C, \max P(c). \quad (7)$$

$$P(a_1, a_2, \dots, a_n | c)$$

فرض پایه الگوریتم بیز ساده این است که در هر رده مقادیر ویژگی ها از یکدیگر مستقل باشند. بنابراین با استفاده از قانون احتمالی استقلال رابطه (۸) تعریف می شود.

$$Pos1 = [3.10, 3.56, -3.12, 2.15, 3.10, -2.26, 0.52, -1.84, 2.80, -3.15] \quad (4)$$

$$Pos1 = [1, 1, 0, 1, 1, 0, 1, 0, 1, 0] \quad (5)$$

بعد از انتخاب ویژگی ها، تابع برازندگی برای ارزیابی کارایی ویژگی های انتخاب شده استفاده می شود. تابع برازندگی در مدل FFANB طبق معادله (۶) تعریف شده است [۲۱].

$$f(x_i) = \sigma \times E_{xi} + (1 - \sigma) \times \frac{|x_i|}{d} \quad (6)$$

در معادله (۶) مقدار خطای طبقه بندی، $|x_i|$ تعداد ویژگی های انتخاب شده، d تعداد کل ویژگی ها. پارامتر σ یک مقدار تصادفی در بازه ۰ و ۱ است که به منظور توازن بین دقت طبقه بندی و تعداد ویژگی های انتخاب شده استفاده می شود.

۳-۲ طبقه بندی

وجود رده های نامتوازن در طبقه بندی به یکی از چالش های بزرگ در این زمینه تبدیل شده است. مجموعه داده نامتوازن بر اساس تعریف عبارت است از یک مجموعه داده ای که تعداد نمونه های متعلق به یک رده در آن با تعداد نمونه های رده دیگر به طور مساوی توزیع نشده باشد. رده با تعداد داده های بیشتر را رده اکثریت و رده با داده های کمتر را رده اقلیت می نامند. در الگوریتم های طبقه بندی استاندارد، توزیع رده ها متوازن در نظر گرفته می شود و این دسته از الگوریتم ها در مواجهه با مجموعه داده های نامتوازن عملکرد مناسبی را از خود ارایه نمی دهند؛ چرا که الگوریتم های معمول طبقه بندی به سمت نمونه های آموزشی رده بزرگتر متمایل می شوند که این موضوع باعث افزایش خطا در شناسایی نمونه های اقلیت می شود. به منظور ارزیابی به مسائل مربوط به مجموعه داده های نامتوازن تکنیک های متعددی معرفی شده اند که در سه دسته زیر طبقه بندی می شوند:

روش سطح الگوریتم: این روش به الگوریتم های

$$p(a_1, a_2, \dots, a_n | c) = p(a_1 | c) p(a_2 | c) \dots p(a_n | c) \quad (8)$$

$$p(a_1, a_2, \dots, a_n | c) = \prod_{i=1}^n p(a_i | c)$$

با جایگذاری معادله (۷) در (۸)، نتیجه به صورت معادله (۹) نشان داده می‌شود که طبقه‌بند بی‌ساده نامیده می‌شود. در این رابطه، E طبقه‌بند بی‌ساده روی نمونه آزمایش E است. با استفاده از معادله (۹) تمامی احتمالات می‌توانند مستقیماً از روی داده‌های آموزشی تعیین شوند.

$$c_N(E) = a \quad c \in C, \max P(c) \prod_{i=1}^n p(a_i | c) \quad (9)$$

۳-۳ معیارهای ارزیابی

در این مقاله، جهت ارزیابی کارایی طبقه‌بندی از چهار معیار دقت، بازخوانی، F-Measure، و صحت استفاده شده است. دقت نشان دهنده تعداد نمونه‌های درست به تعداد کل نمونه‌ها است. معیار بازخوانی نشان می‌دهد که اگر نمونه‌ای، نوع A تشخیص داده شد با چه احتمالی نوع A می‌باشد. هرچه میزان بازخوانی بالاتر باشد بیانگر این است که قابلیت شناسایی درست رده‌ها بیشتر است. در مقابل هرچه میزان دقت بیشتر باشد، بیانگر این است که الگوریتم اشتباه‌های کمتری در شناسایی رده‌ها دارد. درصد صحت یک روش طبقه‌بندی بر روی مجموعه داده‌های آموزشی، درصد مشاهداتی از مجموعه آموزش است که به درستی توسط روش مورد استفاده طبقه‌بندی شده است [۲۲].

$$\text{دقت} = \frac{TP}{TP + FP} \quad (10)$$

$$\text{بازخوانی} = \frac{TP}{TP + FN} \quad (11)$$

$$F - \text{Measure} = \frac{2 * \text{Precision} * \text{Recall}}{\text{Precision} + \text{Recall}} \quad (12)$$

$$\text{صحت} = \frac{TP + TN}{TP + TN + FP + FN} \quad (13)$$

۴. ارزیابی و نتایج

در این بخش مدل FFANB که ترکیبی از الگوریتم

2- Precision
3- Recall
4- Accuracy

حاصلخیزی زمین‌های کشاورزی و بی‌ساده است به منظور طبقه‌بندی بر روی مجموعه داده‌های WEBSpam-UK2007 ارزیابی می‌گردد. این مجموعه داده شامل ۳۷۷۶ نمونه وب غیرهرز و ۲۲۲ نمونه وب هرز در بخش آموزش و شامل ۱۹۳۳ نمونه وب غیرهرز و ۱۲۲ نمونه وب هرز در بخش آزمایش می‌باشد. به منظور انجام آزمایش‌ها، از نمونه‌های با برچسب نامشخص صرف‌نظر شده است. همچنین به صورت تصادفی ۷۵ درصد از داده‌ها به عنوان داده‌های آموزش برای ساخت مدل و ۲۵ درصد از آن به عنوان داده‌های آزمون در نظر گرفته شده است.

در فرآیند اجرای مدل FFANB تعداد تکرار یا نسل می‌تواند ثابت یا براساس معیار توقف همچون خطا باشد. در این مقاله تعداد تکرار مدل FFANB در راستای طبقه‌بندی داده‌ها متغیر و در بازه ۱۰۰ تا ۵۰۰ مرتبه در نظر گرفته شده است. در این تحقیق جهت مقایسه عملکرد مدل FFANB با سایر مقالات، تعداد جمعیت برابر و یکسان در نظر گرفته شده است.

الگوریتم‌های فراابتکاری از ماهیت و رفتاری تصادفی در فرآیند جستجو و کشف بهینه سراسری برخوردار می‌باشند، لذا براساس رفتار تصادفی در فرآیند ارائه تأثیر ماهیت تصادفی نشان داده می‌شود. در راستای ارزیابی ثبات و پایداری و عملکرد الگوریتم‌های فراابتکاری، این الگوریتم‌ها چندین مرتبه به طور مستقل اجرا می‌شوند و سپس بر اساس تحلیل آماری عملکرد آن‌ها مورد ارزیابی و سنجش قرار می‌گیرد. در این مقاله، در راستای ارزیابی و سنجش عملکرد و کارایی مدل FFANB در راستای طبقه‌بندی داده‌ها برای هر مجموعه داده به طور مجزا ۱۵ مرتبه فرآیند طبقه‌بندی مورد آزمایش و شبیه‌سازی قرار گرفته است.

در جدول (۱) نتایج مدل FFANB بر مبنای تکرارهای مختلف نشان داده شده است. درصد صحت در تکرار ۱۰۰ و ۵۰۰ به ترتیب برابر با ۰/۸۱۲۳ و ۰/۹۵۸۴ است. طبق جدول (۱) می‌توان اظهار کرد که درصد صحت با افزایش

جدول ۱: نتایج مدل FFANB بر مبنای تکرارهای مختلف

تعداد تکرار	دقت	بازخوانی	F-Measure	صحت
۱۰۰	۰/۸۰۱۶	۰/۸۲۶۵	۰/۸۱۳۹	۰/۸۱۲۳
۲۰۰	۰/۸۴۵۹	۰/۸۵۲۰	۰/۸۴۸۹	۰/۸۵۴۹
۳۰۰	۰/۸۸۱۳	۰/۸۹۴۶	۰/۸۸۷۹	۰/۸۹۷۲
۴۰۰	۰/۹۱۶۸	۰/۹۲۰۸	۰/۹۱۸۸	۰/۹۲۳۵
۵۰۰	۰/۹۲۴۱	۰/۹۳۱۱	۰/۹۲۷۶	۰/۹۵۸۴

تکرارها، بیشتر شده و همچنین درصد معیارهای دقت و بازخوانی و F-Measure هم افزایش یافته است.

در جدول (۲) نتایج مدل FFANB بر مبنای تعداد ویژگی‌ها بر روی دسته مبتنی بر محتوا نشان داده شده است. تعداد ویژگی‌های دسته مبتنی بر محتوا برابر با ۹۶ ویژگی می‌باشد. اگر تعداد ویژگی‌ها کمتر باشد درصد صحت بیشتر خواهد بود. درصد صحت با ۱۰ ویژگی برابر با ۰/۹۸۱۳، با ۴۰ ویژگی برابر با ۰/۹۶۳۹ و با ۹۶ ویژگی برابر با ۰/۹۳۵۶ است. درصد دقت با ۱۰ ویژگی برابر با ۰/۹۷۹۱، با ۴۰ ویژگی برابر با ۰/۹۶۰۵ و با ۹۶ ویژگی برابر با ۰/۹۳۰۸ است. درصد بازخوانی با ۱۰ ویژگی برابر با ۰/۹۸۰۵، با ۴۰ ویژگی برابر با ۰/۹۶۴۷ و با ۹۶ ویژگی برابر با ۰/۹۳۷۸ است.

در جدول (۳) نتایج مدل FFANB بر مبنای تعداد ویژگی‌ها بر روی دسته مبتنی بر پیوند نشان داده شده است. تعداد ویژگی‌های دسته مبتنی بر پیوند برابر با ۴۱ ویژگی می‌باشد. درصد صحت با ۸ ویژگی برابر با ۰/۹۷۶۳، با ۳۰ ویژگی برابر با ۰/۹۶۴۹ و با ۴۱ ویژگی برابر با ۰/۹۵ است. درصد دقت با ۸ ویژگی برابر با ۰/۹۶۲۵، با ۳۰ ویژگی برابر با ۰/۹۴۷۲ و با ۴۱ ویژگی برابر با ۰/۹۳۹۱ است. درصد بازخوانی با ۸ ویژگی برابر با ۰/۹۶۳۲، با ۳۰ ویژگی برابر با ۰/۹۴۸۱ و با ۴۱ ویژگی برابر با ۰/۹۴۱۳ است.

در جدول (۴) نتایج مدل FFANB بر مبنای تعداد ویژگی‌ها بر روی دسته مبتنی بر پیوند تبدیل یافته نشان داده شده است. تعداد ویژگی‌های دسته مبتنی بر پیوند تبدیل یافته برابر با ۱۳۸ ویژگی می‌باشد. درصد صحت با ۱۰ ویژگی

جدول ۲: نتایج مدل FFANB بر مبنای تعداد ویژگی‌ها بر روی دسته مبتنی بر محتوا

تعداد ویژگی	دقت	بازخوانی	F-Measure	صحت
۱۰	۰/۹۷۹۱	۰/۹۸۰۵	۰/۹۷۹۸	۰/۹۸۱۳
۱۵	۰/۹۷۸۶	۰/۹۸۲۱	۰/۹۸۰۳	۰/۹۸۰۰
۲۰	۰/۹۷۵۶	۰/۹۷۸۹	۰/۹۷۷۲	۰/۹۷۸۱
۲۵	۰/۹۷۴۵	۰/۹۸۰۲	۰/۹۷۷۳	۰/۹۷۵۵
۳۰	۰/۹۶۱۳	۰/۹۶۶۵	۰/۹۶۳۹	۰/۹۶۷۳
۴۰	۰/۹۶۰۵	۰/۹۶۴۷	۰/۹۶۲۶	۰/۹۶۳۹
۴۵	۰/۹۶۰۹	۰/۹۶۸۲	۰/۹۶۴۵	۰/۹۶۱۶
۵۵	۰/۹۵۱۶	۰/۹۵۷۳	۰/۹۵۴۴	۰/۹۵۴۹
۶۰	۰/۹۵۰۶	۰/۹۵۴۶	۰/۹۵۲۶	۰/۹۵۱۸
۷۰	۰/۹۴۴۴	۰/۹۵۰۶	۰/۹۴۷۵	۰/۹۴۷۵
۷۵	۰/۹۴۷۱	۰/۹۴۹۳	۰/۹۴۸۲	۰/۹۴۳۶
۸۰	۰/۹۴۲۳	۰/۹۴۶۹	۰/۹۴۴۶	۰/۹۴۰۵
۸۵	۰/۹۳۵۴	۰/۹۳۴۱	۰/۹۳۴۷	۰/۹۳۸۲
۹۶	۰/۹۳۰۸	۰/۹۳۷۸	۰/۹۳۴۳	۰/۹۳۵۶

برابر با ۰/۷۹۰۰، با ۳۰ ویژگی برابر با ۰/۷۷۳۴، با ۶۰ ویژگی برابر با ۰/۷۵۴۹، با ۹۰ ویژگی برابر با ۰/۷۴۲۳ و با ۱۳۸ ویژگی برابر با ۰/۶۹۵۲ است. درصد دقت با ۱۰ ویژگی برابر با ۰/۷۹۱۱، با ۳۰ ویژگی برابر با ۰/۷۷۴۵، با ۶۰ ویژگی برابر با ۰/۷۵۱۲، با ۹۰ ویژگی برابر با ۰/۷۴۶۴ و با ۱۳۸ ویژگی برابر با ۰/۶۹۵۷ است.

۴-۲ مقایسه و ارزیابی

در جدول (۵) مقایسه مدل FFANB با مدل‌های دیگر نشان داده شده است. درصد صحت الگوریتم ژنتیک، شبکه عصبی مصنوعی و ماشین بردار پشتیبان به ترتیب برابر با ۰/۹۲۴۱، ۰/۹۳۶۶ و ۰/۹۴ است. درصد دقت الگوریتم بهینه‌سازی کلونی مورچه برابر با ۰/۸۴۴۰ است. درصد دقت و F-Measure در ماشین بردار پشتیبان برابر با ۰/۹۱ و ۰/۹۵ است. درصد دقت، بازخوانی و F-Measure در شبکه باور عمیق به ترتیب برابر با ۰/۹۶۷۴، ۰/۹۶۶۶ و ۰/۹۶۶۶ است. درصد دقت، بازخوانی و F-Measure در شبکه باور عمیق با آدابوسست به ترتیب برابر با ۰/۹۲۷۶، ۰/۹۳۲۷ و ۰/۹۳۰۱ است. درصد دقت، بازخوانی و

جدول ۳: نتایج مدل FFANB بر مبنای تعداد ویژگی‌ها بر روی دسته مبتنی بر پیوند

صحت	F-Measure	بازخوانی	دقت	تعداد ویژگی	مبتنی بر پیوند
۰/۹۷۲۹	۰/۹۶۱۸	۰/۹۶۲۶	۰/۹۶۱۱	۱۰	
۰/۹۷۰۳	۰/۹۵۸۳	۰/۹۵۸۰	۰/۹۵۸۷	۱۵	
۰/۹۶۸۲	۰/۹۵۷۶	۰/۹۵۹۶	۰/۹۵۵۶	۲۲	
۰/۹۶۴۹	۰/۹۴۷۶	۰/۹۴۸۱	۰/۹۴۷۲	۳۰	
۰/۹۶۲۱	۰/۹۴۴۶	۰/۹۴۵۵	۰/۹۴۳۸	۳۵	
۰/۹۵۵۵	۰/۹۴۲۸	۰/۹۴۴۲	۰/۹۴۱۴	۳۹	
۰/۹۵۰۰	۰/۹۴۰۲	۰/۹۴۱۳	۰/۹۳۹۱	۴۱	

جدول ۵: مقایسه مدل FFANB با مدل‌های دیگر

صحت	F-Measure	بازخوانی	دقت	مدل‌ها	منابع
۰/۹۳۴۱	-	-	-	الگوریتم ژنتیک	[۲۳]
۰/۹۳۶۶	-	-	-	شبکه عصبی مصنوعی	
۰/۹۴۰۰	-	-	-	ماشین بردار پشتیبان	
-	-	-	۰/۸۴۴۰	الگوریتم بهینه‌سازی کلونی مورچه	[۲۴]
-	۰/۹۶۶۶	۰/۹۶۷۴	۰/۹۶۵۹	شبکه باور عمیق	[۱۷]
-	۰/۹۳۰۱	۰/۹۳۲۷	۰/۹۲۷۶	شبکه باور عمیق + آدابوست	
-	۰/۹۲۳۱	۰/۹۲۴۶	۰/۹۱۴۵	ماشین بردار پشتیبان + آدابوست	
-	۰/۹۵۰۰	-	۰/۹۱۰۰	ماشین بردار پشتیبان	[۱۹]
۰/۹۵۸۴	۰/۹۲۷۶	۰/۹۳۱۱	۰/۹۲۴۱	مدل FFANB	-

جدول ۴: نتایج مدل FFANB بر مبنای تعداد ویژگی‌ها بر روی دسته مبتنی بر پیوند تبدیل یافته

صحت	F-Measure	بازخوانی	دقت	تعداد ویژگی	مبتنی بر پیوند تبدیل یافته
۰/۷۹۰۰	۰/۷۹۲۱	۰/۷۹۳۲	۰/۷۹۱۱	۱۰	
۰/۷۸۵۲	۰/۷۸۴۲	۰/۷۸۵۸	۰/۷۸۲۶	۲۰	
۰/۷۷۳۴	۰/۷۷۶۷	۰/۷۷۸۹	۰/۷۷۴۵	۳۰	
۰/۷۶۰۰	۰/۷۶۵۳	۰/۷۶۹۱	۰/۷۶۱۵	۴۰	
۰/۷۵۸۲	۰/۷۵۸۰	۰/۷۵۹۸	۰/۷۵۶۳	۵۰	
۰/۷۵۴۹	۰/۷۵۵۶	۰/۷۶۰۰	۰/۷۵۱۲	۶۰	
۰/۷۵۰۷	۰/۷۵۶۹	۰/۷۵۹۰	۰/۷۵۴۹	۷۰	
۰/۷۴۶۹	۰/۷۴۲۶	۰/۷۴۳۵	۰/۷۴۱۸	۸۰	
۰/۷۴۲۳	۰/۷۴۸۰	۰/۷۴۹۷	۰/۷۴۶۴	۹۰	
۰/۷۲۱۶	۰/۷۲۷۲	۰/۷۲۹۳	۰/۷۲۵۱	۱۰۰	
۰/۷۱۸۹	۰/۷۲۱۸	۰/۷۲۲۵	۰/۷۲۰۹	۱۱۰	
۰/۷۰۳۵	۰/۷۰۲۹	۰/۷۰۴۷	۰/۷۰۱۱	۱۲۰	
۰/۷۰۱۳	۰/۷۰۷۴	۰/۷۰۹۲	۰/۷۰۵۶	۱۳۰	
۰/۶۹۵۲	۰/۶۹۷۵	۰/۶۹۹۴	۰/۶۹۵۷	۱۳۸	

F-Measure در ماشین بردار پشتیبان با آدابوست به ترتیب برابر با ۰/۹۱۴۵، ۰/۹۲۴۶ و ۰/۹۲۳۱ است.

درصد دقت در ماشین بردار پشتیبان [۱۹] برابر با ۰/۹۱ و معیار F-Measure برابر با ۰/۹۵ است. در میان مدل‌های مقایسه‌ای، شبکه باور عمیق از درصد دقت بیشتری برخوردار است. یادگیری عمیق برگرفته از شبکه عصبی مصنوعی است که شامل مجموعه‌ای از الگوریتم‌ها است که در تلاشند تا مفاهیم انتزاعی و پیچیده را در سطوح و لایه‌های مختلف مدل کنند. ایده اصلی شبکه‌های

مبتنی بر یادگیری عمیق، نزدیک شدن به بهترین راه حل‌ها بر مبنای آموزش و آزمایش می‌باشد. برتری درصد صحت مدل FFANB در مقایسه با مدل‌های دیگر به دلیل این است که الگوریتم حاصلخیزی زمین‌های کشاورزی در چندین منطقه اقدام به کشف جواب‌های بهینه می‌نماید و سپس در بین آن‌ها، مقدار سراسری برای جواب مسئله انتخاب می‌شود.

- Semi-supervised Technique, Progress in Intelligent Computing Techniques: Theory, Practice, and Applications, pp. 281-286, 2018.
5. A. Sharaff, Spam Detection in SMS Based on Feature Selection Techniques, Emerging Technologies in Data Mining and Information Security, pp. 555-563, 2018.
 6. M.K. Sohrabi, F. Karimi, A Feature Selection Approach to Detect Spam in the Facebook Social Network, Arabian Journal for Science and Engineering, Vol. 43, Issue 2, pp. 949-958, 2018.
 7. N. Karimpour, F.S. Gharehchopogh, A New Approach to Detect Spam Emails Using the Hybrid Model of Ant Colony and Firefly Algorithms, Computing Science Journal (CSJ), Vol. 14, pp. 1-17, 2019
 8. F.S. Gharehchopogh, M. Vafadar, M. Motamanfar, Improving the Invasive Weed Optimization algorithm with the nearest neighbor in the classification of spam emails, Computing Science Journal (CSJ), 10: 5464, 2018
 9. H. Shayanfar, F.S. Gharehchopogh, Farmland fertility: A new metaheuristic algorithm for solving continuous optimization problems, Vol. 71, pp. 728-746, 2018.
 10. A. McCallum, K. Nigam, A Comparison of Event Models for Native Bayes Text Classification, In AAAI-98 workshop on learning for text categorization, Vol. 752, pp. 41-48, 1998.
 11. Webpage Spam data set, <http://chato.cl/webspam/datasets/uk2007/>, [last available 2019.09.09]
 12. J. Wu, A generalized tree augmented naive Bayes link prediction model, Journal of Computational Science, Vol. 27, pp. 206-217, 2018.
 13. F. Asdaghi, A. Soleimani, An effective feature selection method for web spam detection, Knowledge-Based Systems, Vol. 166, 15 pp. 198-206, 2019.
 14. K.L. Goh, A.K. Singh, Comprehensive Literature Review on Machine Learning Structures for Web Spam Classification, Procedia Computer Science, Vol. 70, pp. 434-441, 2015.
 15. S. Singh, A.K. Singh, Web-Spam Features Selection Using CFS-PSO, Procedia Computer Science, Vol. 125, pp. 568-575, 2018.
 16. J.F. Glez, D. Ruano-Ordas, J.R. MEndez, F. Fdez-Riverola, R. PavOn, A dynamic model for integrating simple web spam classification techniques, Expert Systems with Applications, Vol. 42, Issue 21, pp. 7969-7978, 2015.
 17. Y. Li, X. Nie, R. Huang, Web spam classification method based on deep belief networks, Expert Systems with Applications, Vol. 96, pp. 261-270, 2018
 18. A.H. Keyhanipour, B. Moshiri, Designing a web spam classifier based on feature fusion in the Layered Multi-population Genetic Programming framework, Proceedings of the 16th International Conference on Information Fusion, IEEE, pp. 53-60, 2013
 19. R.C. Patil and D.R. Patil, Web spam detection using SVM classifier, 2015 IEEE 9th International Conference on Intelligent Systems and Control (ISCO), Coimbatore, IEEE, pp. 1-4, 2015.
 20. S. Jain, S. Shukla, R. Wadhvani, Dynamic selection of normalization techniques using data complexity measures, Expert Systems with Applications, Vol. 106, pp. 252-262, 2018
 21. M.A. Elaziz, A.A. Ewees, R.A. Ibrahim, S. Lu, Opposition-based moth-flame optimization improved by differential evolution for feature selection, Mathematics and Computers in Simulation, In press, corrected proof, Available online 2 August 2019
 22. P. Galdi, R. Tagliaferri, Data Mining: Accuracy and Error Measures for Classification and Prediction, Encyclopedia of Bioinformatics and Computational Biology, Vol. 1, pp. 431-436, 2019
 23. S.H.R. Mohammadi, M.A.Z. Chahooki, Web Spam Detection Using Multiple Kernels in Twin Support Vector Machine, Information Retrieval, pp. 1-10, 2016
 24. B. Manaskasemsak, J. Jiarpakdee, A. Rungsawang, Adaptive Learning Ant Colony Optimization for Web Spam Detection, International Conference on Computational Science and Its Applications, ICCSA 2014: Computational Science and Its Applications – ICCSA 2014, Vol. 8584, pp. 642-653, 2014

صفحات وب در فضای اینترنت به فضایی محبوب برای تبادل اطلاعات، ثبت و مشاهده اطلاعات تبدیل شده است. این میزان محبوبیت و جامعیت، هدف مناسبی جهت فعالیت‌های مخرب و وب هرزنویسان شده است. امنیت و نظارت بر صفحات وب یکی از مهم‌ترین ملزومات فضای اینترنت به دلیل جلوگیری از نفوذ بدافزارها می‌باشد. به‌طور کلی بدافزارها پس از ورود به سیستم می‌توانند اقداماتی نظیر سرقت اطلاعات، تخریب اطلاعات و افشای اطلاعات را انجام دهند. بنابراین ایجاد روشی که بتواند به صورت کارا به شناسایی و جلوگیری از وب هرز بپردازد، همواره مورد نیاز خواهد بود. پس از بیان کارهای انجام شده در مورد تشخیص صفحات وب هرز، یک مدل مبتنی بر یادگیری ماشین برای تشخیص وب هرز پیشنهاد شد. مدل FFANB ترکیبی از الگوریتم حاصلخیزی زمین‌های کشاورزی و بیض ساده بود. به منظور افزایش صحت الگوریتم طبقه‌بندی از انتخاب ویژگی استفاده شد. انتخاب ویژگی در مجموعه داده‌های اصلی در ابتدا با تعداد ویژگی کم شروع شد و به مرور تعداد ویژگی‌ها افزایش یافت. اگر تعداد ویژگی‌ها کمتر باشد آنگاه درصد صحت، دقت، بازخوانی و F-Measure بیشتر خواهد شد. نتایج ارزیابی نشان داد که مدل FFANB در مقایسه با الگوریتم‌های دیگر از درصد صحت بالاتری بهره‌مند است.

منابع

1. M. Yu, J. Zhang, J. Wang, J. Gao, T. Xu, R. Yu, The Research of Spam Web Page Detection Method Based on Web Page Differentiation and Concrete Cluster Centers, International Conference on Wireless Algorithms, Systems, and Applications, WASA 2018: Wireless Algorithms, Systems, and Applications, pp. 820-826, 2018.
2. X.Y. Lu, M.S. Chen, J.L. Wu, P.C. Chang, M.H. Chen, A novel ensemble decision tree based on under-sampling and clonal selection for web spam detection, Pattern Analysis and Applications, Vol. 21, Issue 3, pp. 741-754, 2018.
3. K. Hans, L. Ahuja, S.K. Muttoo, Performance Evaluation of Neural Network Training Algorithms in Redirection Spam Detection, Nature Inspired Computing pp 177-183, 2018.
4. R. Narayan, J.K. Rout, S.K. Jena, Review Spam Detection Using