
Fake Websites Detection Improvement Using Multi-Layer Artificial Neural Network Classifier with Ant Lion Optimizer Algorithm

Farhang Padidaran Moghadam*, Mahshid Sadeghi Bajgiran**

*Assistant Professor, Computer Department, Esfarayen Higher Education Technical Complex, Esfarayen, Iran

**Master's degree, Eshraq Institute of Higher Education, Bojnourd, Iran

Abstract

In phishing attacks, a fake site is forged from the main site, which looks very similar to the original one. To direct users to these sites, Phishers or online thieves usually put fake links in emails and send them to their victims, and try to deceive users with social engineering methods and persuade them to click on fake links. Phishing attacks have significant financial losses, and most attacks focus on banks and financial gateways. Machine learning methods are an effective way to detect phishing attacks, but this is subject to selecting the optimal feature. Feature selection allows only important features to be considered as learning input and reduces the detection error of phishing attacks. In the proposed method, a multilayer artificial neural network classifier is used to reduce the detection error of phishing attacks, the feature selection phase is performed by the ant lion optimization (ALO) algorithm. Evaluations and experiments on the Rami dataset, which is related to phishing, show that the proposed method has an accuracy of about 98.53% and has less error than the multilayer artificial neural network. The proposed method is more accurate in detecting phishing attacks than BPNN, SVM, NB, C4.5, RF, and kNN learning methods with feature selection mechanism by PSO algorithm.

Keywords: Phishing Attacks, Feature Selection, Ant Lion Optimization Algorithm, Fake Pages, Fake Links

بهبود تشخیص وبگاه های جعل شده با استفاده از طبقه بندی کننده شبکه عصبی مصنوعی چند

لایه با الگوریتم بهینه سازی شیر مورچه

فرهنگ پدیداران مقدم*، مهشید صادقی باجگیران**

*استادیار گروه کامپیوتر، مجتمع آموزش عالی فنی مهندسی اسفراین، اسفراین، ایران

**دانش آموخته کارشناسی ارشد، موسسه آموزش عالی اشراق، بجنورد، ایران

تاریخ پذیرش: ۱۴۰۱/۰۹/۱۹

تاریخ دریافت: ۱۴۰۱/۰۵/۱۷

نوع مقاله: پژوهشی

چکیده

در حملات فیشینگ یک وبگاه جعلی از روی وبگاه اصلی جعل می گردد که ظاهر بسیار شبیه به وبگاه اصلی دارد. فیشر یا سارق آنلاین برای هدایت کاربران به این وبگاهها، معمولا لینکهای جعلی را در ایمیل قرار داده و برای قربانیان خود ارسال نموده و با روشهای مهندسی اجتماعی سعی در فریب کاربران و مجاب نمودن آنها برای کلیک روی لینکهای جعلی دارد. حملات فیشینگ زیان مالی قابل توجهی دارند و بیشتر روی بانکها و درگاههای مالی متمرکز هستند. روشهای یادگیری ماشین یک روش موثر برای تشخیص حملات فیشینگ است اما این مشروط به انتخاب بهینه ویژگی است. انتخاب ویژگی باعث می شود فقط ویژگیهای مهم به عنوان ورودی یادگیری در نظر گرفته شوند و خطای تشخیص حملات فیشینگ کاهش داده شود. در روش پیشنهادی برای کاهش دادن خطای تشخیص حملات فیشینگ یک طبقه بندی کننده شبکه عصبی مصنوعی چند لایه استفاده شده که فاز انتخاب ویژگی آن با الگوریتم بهینه سازی شیر مورچه انجام می شود. ارزیابی و آزمایشها روی مجموعه داده Rami که مرتبط با فیشینگ است نشان می دهد روش پیشنهادی دارای دقتی در حدود ۹۸٫۵۳٪ است و نسبت به شبکه عصبی مصنوعی چند لایه خطای کمتری دارد. روش پیشنهادی در تشخیص حملات فیشینگ از روشهای یادگیری SVM، BPNN، NB، C4.5، RF و kNN با سازوکار انتخاب ویژگی توسط الگوریتم PSO دقت بیشتری دارد.

واژگان کلیدی: حملات فیشینگ، انتخاب ویژگی، الگوریتم بهینه سازی شیر مورچه، صفحات جعلی، لینکهای جعلی

۱. مقدمه

وبگاه‌های جعلی^۱ به عنوان یکی از تهدیدهای مهم در فناوری اطلاعات و تجارت الکترونیک به شمار می‌روند زیرا این صفحات بسیار شبیه صفحات قانونی بوده و اطلاعات کاربران را مورد سرقت قرار می‌دهند. حملات فیشینگ^۲ با استفاده از صفحات جعلی فضای وب را برای کاربران اینترنت ناامن نموده است و در این نوع حملات که می‌تواند مبتنی بر مهندسی اجتماعی یا مبتنی بر فریب توسط بدافزار می‌باشد یک کاربر به صورت خودکار یا توسط لینک‌های جعلی به سمت صفحات وب جعلی هدایت شده و اطلاعات خود را در این صفحات وارد می‌نماید. فیشر با دریافت اطلاعات کاربران می‌تواند از آنها برای سرقت اطلاعات استفاده نماید. صفحات جعلی در اینترنت و حملات فیشینگ دارای مجموعه‌ای از ویژگی‌ها است که می‌تواند برای تشخیص صفحات جعلی استفاده شود. در بیشتر صفحات جعلی عمر دامنه اندک است زیرا این صفحات با سرعت ایجاد و سریع نیز شناسایی و حذف می‌شوند پس می‌توان از این ویژگی برای تشخیص صفحات جعلی و حملات فیشینگ استفاده نمود. اطلاعات مرتبط با دامنه فقط برای تشخیص حملات فیشینگ مهم نبوده بلکه اطلاعات مرتبط با لینک و آدرس نیز مهم می‌باشند و مشاهده می‌شود در بیشتر صفحات جعلی طول آدرس بیش از اندازه طولانی است [۱]. وجود کاراکترهای خاص در صفحات جعلی یا آدرس وبگاه می‌تواند نشانه فیشینگ باشد و به عنوان نمونه استفاده از کاراکتر @ در آدرس یک وبگاه یک نشانه مهم در تشخیص صفحات جعلی است زیرا هرگز با این کاراکتر اطلاعات مهم کاربران را برای خود ایمیل می‌نماید. اطلاعات و ویژگی‌های مرتبط با کد منبع وبگاه هم یک عامل مهم در تشخیص صفحات جعلی است به گونه‌ای که وجود کدهای خاص جاوا اسکریپت مانند عدم کلیک راست می‌تواند نشانه مهم فیشینگ باشد. اطلاعات و اعتبار سنجی صفحات وب نیز توسط موتورهای جستجو مانند آکسا و گوگل هم برای تشخیص صفحات جعلی مهم است زیرا در صفحات قانونی و معتبر به علت سابقه و تعداد لینک‌های ورودی دارای اعتبار بیشتری بوده و توسط موتورهای جستجوگر مانند گوگل شاخص گذاری می‌شوند و دارای رتبه مناسبی

می‌باشند.

مجموعه ویژگی‌های بکار رفته برای تشخیص صفحات جعلی متنوع و زیاد است و نیاز است که برای تشخیص صفحات جعلی و حملات فیشینگ توسط روش‌های کشف دانش مانند یادگیری ماشین از مرحله انتخاب ویژگی بخوبی استفاده شود تا فقط یادگیری بر روی ویژگی‌های مهم انجام شود تا الگوی صفحات فیشینگ شناسایی شود. مسئله انتخاب ویژگی در تشخیص حملات فیشینگ یک مسئله بهینه‌سازی است که نیاز است بردار ویژگی با دقت بالایی انتخاب شود و می‌توان برای حل آن از روش‌های مبتنی و طبیعت و الگوریتم-های فراابتکاری استفاده نمود [۲]. برای تشخیص حملات فیشینگ تاکنون روش‌های مختلفی توسعه داده شده است که بیشتر آنها بر اساس سه استراتژی لیست سیاه [۳]، روش‌های اکتشافی^۳ [۴] و روش‌های کشف دانش [۵] متمرکز است. روش‌های کشف دانش به دو دسته روش‌های یادگیری ماشین [۶] و یادگیری عمیق [۷] طبقه-بندی می‌شود. در روش‌های یادگیری ماشین زمان پردازش به مراتب کمتر از روش‌های یادگیری عمیق است اما فاقد مرحله انتخاب ویژگی می‌باشند از این جهت نیاز به سازوکار انتخاب ویژگی دارند تا با دقت بالا حملات فیشینگ را تشخیص دهند.

با توجه به اینکه حملات فیشینگ سالانه رو به افزایش است و تعداد زیادی از حملات در حوزه مالی انجام می‌شود لذا برآورد می‌شود که زیان آنها نیز قابل توجه است و از طرفی نیز تعداد حملات فیشینگ و چالش آن حتی از ویروس‌ها بیشتر است و از این جهت شناسایی آنها اهمیت بالایی دارد.

اهداف این مقاله، شناسایی صفحات جعلی و شناسایی صفحات فیشینگ در اینترنت با خطای حداقلی، مبارزه با سرقت‌های آنلاین و فیشینگ در فضای مجازی و بهبود روش‌های داده‌کاوی برای تشخیص فیشینگ و لینک‌های جعلی با استفاده از الگوریتم‌های فراابتکاری هوشمند نظیر بهینه‌سازی شیرمورچه است.

روش‌های یادگیری ماشین برای تشخیص حملات فیشینگ نیاز دارند تا ویژگی‌های مهم صفحات وب را دریافت نموده و از این ویژگی‌ها برای آموزش و یادگیری استفاده نمایند. به عبارت بهتر در برخی مطالعات فقط بر ویژگی‌های زبانی و محتوی صفحات وب برای

³ Heuristic

⁴ Deep learning

¹ Fake websites

² Phishing

حملات فیشینگ در صدر حملات به شبکه‌های کامپیوتری است و حتی سهم بیشتری از بدافزار و ویروس را به خود اختصاص می‌دهد.

در پژوهش [۱۲]، برای انتخاب ویژگی در حملات فیشینگ یک روش مبتنی بر یادگیری و تئوری ریاضیات Rough ارائه نمودند. نتایج ارزیابی آنها نشان می‌دهد که شاخص F-measure در روش پیشنهادی آنها برای تشخیص فیشینگ در حدود ۹۵٪ است. همچنین آنها نشان دادند که ۹ ویژگی توسط روش پیشنهادی آنها یا FRS بر روی تمامی سه مجموعه داده بکار رفته در تشخیص فیشینگ مهم می‌باشند. چالش‌های مهمی که در استفاده از روش‌های ریاضی در دو بعد می‌تواند مطرح باشد این است که اولاً این روش‌ها از هوشمندی روش‌های فراابتکاری برخوردار نبوده و از طرفی دیگر انعطاف‌پذیری بالایی ندارند.

در پژوهش [۱۳]، برای تشخیص لینک‌های جعلی مبتنی بر فیشینگ از یک روش جدید مبتنی بر یادگیری ماشین استفاده نمودند. در این پژوهش یک مجموعه داده جدید ساخته می‌شود و نتایج آزمایش بر روی آن آزمایش می‌شود. با توجه به نتایج تجربی و مقایسه‌ای می‌توان دریافت الگوریتم Random Forest با ویژگی‌های مبتنی بر NLP دارای دقتی بالا برای شناسایی URL های فیشینگ است.

در پژوهش [۱۴]، برای تشخیص صفحات جعلی و حملات فیشینگ یک روش اکتشافی غیرخطی مبتنی بر رگرسیون را انتخاب نمودند. نتایج پیاده‌سازی آنها نشان می‌دهد رگرسیون غیرخطی براساس الگوریتم جستجوی هارمونی و ماشین بردار پشتیبان با دو فاز انتخاب ویژگی درخت تصمیم‌گیری و بسته‌بندی ویژگی‌ها به ترتیب دارای دقت ۹۴،۱۳٪ و ۹۲،۸۰٪ در تشخیص حملات فیشینگ هستند در نتیجه، مطالعه نشان می‌دهد که الگوریتم جستجوی هارمونی مبتنی بر رگرسیون غیر خطی باعث عملکرد بهتر در مقایسه با ماشین بردار پشتیبان می‌شود.

تشخیص فیشینگ تاکید شده است اما برای تشخیص دقیق‌تر صفحات فیشینگ نیاز است که ویژگی‌های مهم این صفحات مورد بررسی قرار گرفته شود و ویژگی‌های مهم آن انتخاب شود زیرا انتخاب ویژگی باعث می‌شود که یادگیری فقط بر روی ویژگی‌های انجام شود که اهمیت بیشتری دارند و خروجی مدل را دقیق‌تر می‌نمایند و حال آنکه در پژوهش‌های مانند پژوهش [۸]، مسئله انتخاب ویژگی در ترکیب با روش‌های داده‌کاوی در نظر گرفته نشده است. در روش پیشنهادی برای تشخیص صفحات جعلی یک روش ترکیبی استفاده می‌شود و همزمان ویژگی‌های مرتبط با دامنه صفحات وب، ویژگی‌های مرتبط با آدرس وبگاه و ویژگی‌های مرتبط با کد صفحات وب در نظر گرفته شده و به کمک انتخاب ویژگی، بهترین ویژگی‌ها برای آموزش یک شبکه عصبی مصنوعی ارزیابی می‌گردد. این مقاله دارای چند بخش است در بخش دوم مقاله پیشینه تحقیق در مورد حملات فیشینگ ارزیابی می‌شود. در بخش سوم، روش پیشنهادی برای تشخیص حملات فیشینگ ارزیابی می‌گردد. در بخش چهارم نیز آزمایشها و تجزیه تحلیل و در نهایت در بخش پنجم نتیجه‌گیری تحقیق و پیشنهادات آتی بحث شده است.

۲. پیشینه پژوهش

برآوردها نشان می‌دهد تعداد حملات فیشینگ در سال ۲۰۲۰ قابل توجه است و فقط در ماه ژوئیه حدود ۱۰۰ هزار وبگاه فیشینگ در دنیا شناسایی شده است و این تعداد دارای یک روند صعودی است. برآوردها نشان می‌دهد که حملات فیشینگ در سال ۲۰۲۰ تا حدود ۹۰ هزار حمله گزارش شده است اما این حملات به حدود ۲۵۰ هزار مورد در سال ۲۰۲۱ رسیده است [۹ و ۱۰]. بررسی‌ها نشان می‌دهد حدود ۲۴،۹٪ از صفحات فیشینگ و جعلی برای سرقت‌های مالی ایجاد می‌شود و این موضوع باعث می‌شود زیان این حملات در این حوزه نیز قابل توجه شود. برآوردها نشان می‌دهد زیان حملات فیشینگ سالانه به میلیون‌ها دلار بالغ می‌شود و یکی از دلایل آن وجود حملات بر علیه زیرساخت‌های مالی مانند درگاه‌های پرداخت اینترنتی است. اهمیت شناسایی حملات فیشینگ فقط به تعداد حملات یا زیان مالی آنها محدود نمی‌شود. برآوردها نشان می‌دهد

بهبود تشخیص وبگاه های جعل شده با استفاده از طبقه بندی کننده شبکه عصبی مصنوعی چند لایه با الگوریتم بهینه سازی شیرمورچه

در پژوهش [۱۸]، جهت تشخیص صفحات جعلی یک روش یادگیری ماشین مبتنی بر سه تکنیک ماشین بردار پشتیبان، آنتروپی بیشینه و یادگیری عمیق ارائه نمودند. در روش پیشنهادی آنها اطلاعات مبتنی بر محتوی صفحات و کلمات کلیدی آنها برای آموزش و یادگیری مورد استفاده قرار گرفته شد. در این روش محتوی مورد نظر از صفحات وب استخراج شده و نرخ تکرار کلمات نیز به عنوان ویژگی در یادگیری این روش ها استفاده شده است. نتایج آزمایشها آن ها نشان می دهد دقت روش پیشنهادی توسط این تکنیک بالا بوده و چالش عمده این روش ها، زمان اجرای آن در نظر گرفته می شود.

۳. مراحل انجام طرح پیشنهادی

در روش پیشنهادی یک بردار ویژگی یک شیر مورچه یا مورچه است و یک عضو جمعیت الگوریتم بهینه سازی شیرمورچه فرض می شود. در مرتبه اول یک جمعیت تصادفی از بردارهای ویژگی در قالب جمعیت الگوریتم شیر مورچه ایجاد می شود که دارای مقادیر تصادفی صفر و یک می باشند. نقش بردارهای ویژگی آموزش دادن شبکه عصبی با انتخاب ویژگی های ورودی است و نقش الگوریتم بهینه سازی شیرمورچه به روزرسانی ویژگی های انتخاب شده و یافتن ویژگی مهم برای یادگیری در تشخیص فیشینگ است. در هر تکرار، الگوریتم بهینه سازی شیر مورچه بر روی بردارهای ویژگی یا جمعیت بردارهای ویژگی اعمال شده و آنها را به روزرسانی می نماید.

یک بردار ویژگی برای ارزیابی به دو عامل مهم ذیل نیاز دارد تا شایستگی آن مشخص شود:

- کاهش یافتن خطای تشخیص صفحات جعلی از اصلی
- کاهش دادن تعداد ویژگی انتخاب شده

هر بردار ویژگی که بتواند این مقادیر را کمینه تر نماید به عنوان بردار ویژگی بهینه در نظر گرفته می شود. در تکرار آخر تلاش می شود بهینه ترین بردار ویژگی برای آموزش شبکه عصبی مصنوعی انتخاب و بر اساس آن حملات فیشینگ تشخیص داده شود. در هر مرحله ارزیابی می توان کیفیت خروجی شبکه عصبی را با خطای طبقه بندی صفحات جعلی از اصلی تشخیص داد.

در شکل (۱)، سازوکار ترکیب الگوریتم شیر مورچه و شبکه عصبی

در پژوهش [۱۵]، برای تشخیص لینک های جعلی و آدرس های فیشینگ یک چارچوب مبتنی بر یادگیری ماشین با توجه به ویژگی بکار رفته در این صفحات ارائه نمودند. نتایج ارزیابی آنها نشان می دهد که تکنیک جنگل تصادفی نسبت به سایر روش ها میزان خطای کمتری برای تشخیص وبگاه های جعلی دارد.

در پژوهش [۱۶]، جهت تشخیص صفحات جعلی از یک رویکرد یادگیری مبتنی بر شبکه عصبی مصنوعی استفاده نمودند. در تکنیک پیشنهادی آنها خروجی مدل شبکه عصبی توسط تکنیک مونت کارلو بهبود داده می شود تا دقت آن افزایش یابد نتایج ارزیابی آنها بر روی صفحات اینترنتی فیشینگ و قانونی نشان داده که روش پیشنهادی نسبت به تکنیک های رگرسیون، ماشین بردار پشتیبان خطی، شبکه بیزین، نزدیکترین همسایه و ماشین بردار پشتیبان شعاعی دقت بیشتری در شناسایی حملات فیشینگ دارد.

در پژوهش [۱۷]، از یک سازوکار یادگیری مبتنی بر تجزیه و تحلیل لینک و پیوندهای موجود در سورس کد صفحات برای تشخیص جعلی بودن آنها و به عبارت بهتر اعتبار آنها استفاده نمودند. در این تکنیک، یک صفحه اینترنتی در وب در نظر گرفته می شود و اطلاعات درون سورس کد آن نظیر پیوند و لینک ها گردآوری شده تا بر اساس آنها اعتبار وبسایت مشخص شود. در این روش مجموعه ای از ویژگی های مرتبط با لینک ها نظیر لینک های داخلی و لینک های خارجی بکار رفته در سورس کد، بردار ویژگی را ایجاد نموده و این بردار ویژگی برای یادگیری توسط تکنیک یادگیرنده مورد استفاده قرار گرفته می شود تا صفحات به دو دسته غیرقانونی و قانونی طبقه بندی شوند. در این روش ۱۲ ویژگی مرتبط با لینک های موجود در سورس کد صفحات، بردار ویژگی را ایجاد نموده و سپس این بردار ویژگی برای یادگیری تکنیک های مختلف مانند رگرسیون، درخت تصمیم گیری و ماشین بردار پشتیبان بکار برده می شود. نتایج پیاده سازی آنها نشان می دهد تکنیک رگرسیون برای این روش بیشترین دقت ممکن و تکنیک ماشین بردار پشتیبان کمترین دقت ممکن را در بین این روش ها دارد.

شبکه عصبی مصنوعی تنظیم می‌شود و یک بردار ویژگی به عنوان یک عضو الگوریتم ALO کدگذاری می‌شود که دارای الگوی صفر و یک است. هر مولفه صفر نشان دهنده عدم انتخاب ویژگی و هر مولفه یک نشان دهنده انتخاب ویژگی است.

ایجاد یک جمعیت اولیه از بردارهای ویژگی به صورت تصادفی به عنوان جمعیت الگوریتم بهینه‌سازی شیر مورچه:

انتخاب یک بردار ویژگی از بردارهای ویژگی یا عضو الگوریتم شیر مورچه برای به روزرسانی با استفاده از مراحل الگوریتم انتخاب ویژگی

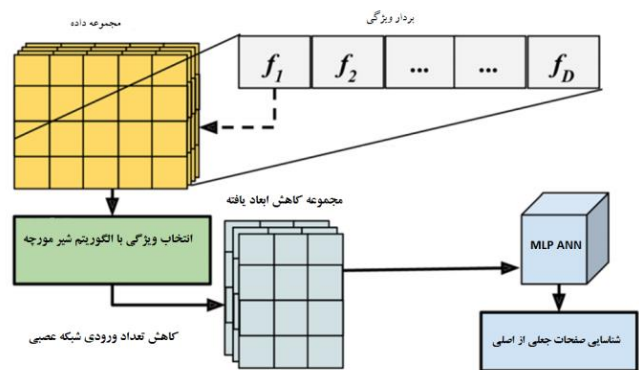
آموزش شبکه عصبی با الگوی صفر و یک بردار ویژگی
 ارزیابی بردار ویژگی بر اساس خطای شبکه عصبی و تعداد بردار ویژگی و تعیین شایستگی یک بردار ویژگی. هر چقدر مقدار تابع انتخاب ویژگی توسط یک بردار ویژگی کمینه‌تر شود آنگاه بردار ویژگی بهینه‌تر است.

به روزرسانی بردارهای ویژگی در تشخیص فیشینگ با الگوریتم ALO و با استفاده از حرکت بردارهای ویژگی از نوع مورچه به سمت بردارهای ویژگی از نوع شیرمورچه

اگر شمارنده تکرار بیشینه نشده مراحل قبلی تکرار شود و اگر شمارنده تکرار بیشینه شده است بردار ویژگی بهینه در تکرار آخر برای آموزش شبکه عصبی مصنوعی مورد استفاده قرار گرفته شود.

شبکه عصبی مصنوعی با بردار ویژگی بهینه آموزش داده می‌شود و ارزیابی می‌گردد.

مصنوعی چند لایه برای تشخیص حملات فیشینگ نمایش داده شده است. با توجه به شکل مورد نظر در ابتدا یک مجموعه داده با همه ویژگی‌ها در نظر گرفته می‌شود و یک بردار ویژگی با الگوی صفر و یک با D ویژگی به عنوان یک عضو الگوریتم شیر مورچه تعیین و کدگذاری می‌شود. الگوریتم مورد نظر در هر تکرار تعدادی بردار ویژگی دارد که روی مجموعه داده نگاشت داده می‌شوند. با توجه به شکل الگوریتم شیر مورچه یا ALO در هر تکرار سعی در انتخاب بردار ویژگی بهینه برای کاهش دادن ابعاد مجموعه داده و انتخاب ویژگی‌های بهینه دارد.



شکل ۱: سازوکار ترکیب الگوریتم شیر مورچه و شبکه عصبی مصنوعی

ویژگی‌های بهینه در تکرار آخر الگوریتم ALO محاسبه شده و از این ویژگی‌ها برای کاهش دادن ابعاد مجموعه داده اصلی استفاده می‌شود و شبکه عصبی مصنوعی از این مجموعه داده کاهش ابعاد به عنوان ورودی استفاده می‌کند. در روش پیشنهادی شبکه عصبی مصنوعی با دریافت این مجموعه داده به عنوان ورودی سعی در طبقه‌بندی صفحات وب به دو دسته جعلی و اصلی دارد.

چارچوب روش پیشنهادی برای تشخیص حملات فیشینگ با استفاده از انتخاب ویژگی و یادگیری و طبقه‌بندی با شبکه عصبی مصنوعی چند لایه در شکل (۲)، به تصویر کشیده شده است. در روش پیشنهادی در ابتدا ۷۰٪ از داده‌ها و صفحات وب از نوع داده آموزشی و ۳۰٪ دیگر از نوع آزمون است. در روش پیشنهادی داده‌های آموزشی در فاز انتخاب ویژگی و یادگیری شبکه عصبی استفاده می‌شود. در مرحله ارزیابی از داده‌های آزمون نیز برای تحلیل مدل پیشنهادی در تشخیص حملات فیشینگ استفاده می‌گردد. در چارچوب پیشنهادی برای تشخیص حملات فیشینگ از یک روال چند مرحله‌ای استفاده شده است که مراحل آن در ذیل آورده شده است:

پارامترهای شبکه عصبی مانند تعداد لایه‌ها و نورون‌های

بهبود تشخیص وبگاه های جعل شده با استفاده از طبقه بندی کننده شبکه عصبی مصنوعی چند لایه با الگوریتم بهینه سازی شیرمورچه

رویکرد هوش گروهی در نظر گرفت که فضای مسئله توسط مورچه ها و شیر مورچه ها مورد پیمایش قرار میگیرد [۱۹]. در روش پیشنهادی هر بردار ویژگی به عنوان یک عضو الگوریتم بهینه سازی شیرمورچه کدگذاری و تعریف می شود و این بردارهای ویژگی مرتباً به روزرسانی می شوند تا بهینه شوند. در این الگوریتم فراابتکاری، مورچه ها و شیرمورچه ها که بردارهای ویژگی می باشند در ابتدا به شکل تصادفی در فضای جستجوی مسئله پراکنده و توزیع می شوند. در این الگوریتم یک بردار ویژگی از نوع مورچه می تواند با حرکت تصادفی به سمت یک بردار ویژگی از نوع شیرمورچه حرکت نماید. برای به روزرسانی یک مورچه به سمت بردار ویژگی از نوع شیرمورچه ها از رابطه (۲)، استفاده می شود [۱۹]:

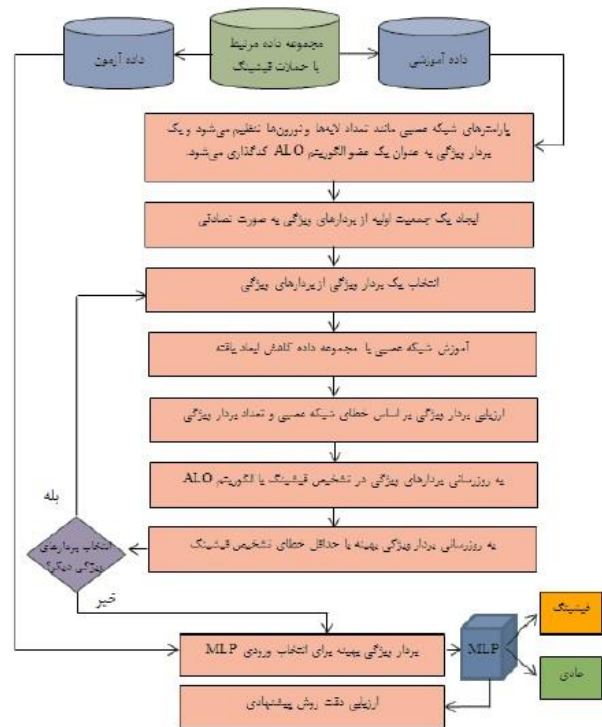
$$X(t) = [0, \text{cumsum}(2r(t_1) - 1), \text{cumsum}(2r(t_2) - 1), \dots, \text{cumsum}(2r(t_n) - 1)] \quad (2)$$

در رابطه فوق، cumsum تابع تجمعی، $X(t)$ بردار جابجایی یک بردار ویژگی از نوع مورچه در حرکت به سمت یک حفره یا یک ناحیه از فضای مسئله است و n حداکثر تعداد تکرار الگوریتم پیشنهادی است. در این رابطه، $r(t)$ یک تابع تصادفی برای تولید عدد صفر و یک است که مطابق رابطه (۳)، تعریف می شود [۱۹]:

$$r(t) = \begin{cases} 1 & \text{if rand} > 0.5 \\ 0 & \text{if rand} \leq 0.5 \end{cases} \quad (3)$$

برای فرموله کردن و مدل سازی الگوریتم بهینه سازی شیرمورچه در مرحله اول یک جمعیت از بردارهای ویژگی از نوع مورچه و بردار ویژگی از نوع شیرمورچه ها را به صورت ماتریس های جداگانه می توان تعریف نمود. در رابطه (۴) و (۵) به ترتیب دو ماتریس M_{Ant} و M_{OA} برای ذخیره بردارهای ویژگی از نوع مورچه ها و ماتریس شایستگی آنها با استفاده از تابع هدف انتخاب ویژگی نمایش داده شده است.

$$M_{Ant} = \begin{bmatrix} A_{1,1} & A_{1,2} & \dots & A_{1,d} \\ A_{2,1} & A_{2,2} & \dots & A_{2,d} \\ \vdots & \vdots & \vdots & \vdots \\ A_{n,1} & A_{n,2} & \dots & A_{n,d} \end{bmatrix} \quad (4)$$



شکل ۲: چارچوب روش پیشنهادی برای تشخیص حملات فیشینگ

هر بردار ویژگی در ارتباط با حملات فیشینگ نیاز به ارزیابی دارد و برای ارزیابی هر بردار ویژگی دو عامل خطای طبقه بندی صفحات وب نرمال و فیشینگ و تعداد بردار ویژگی مهم است. یک تابع هدف می تواند به صورت ترکیبی از این دو عامل مانند رابطه (۱)، باشد:

$$f = \alpha \cdot \frac{1}{n} \sum_{i=1}^n |\bar{Y}_i - Y_i| + \beta \cdot \frac{F}{A} \quad (1)$$

در این رابطه، n تعداد نمونه های بکار رفته برای ارزیابی است. متغیر F نشان دهنده تعداد ویژگی انتخاب شده به کل ویژگی ها یا A است. \bar{Y}_i شماره کلاس تخمین زده شده برای صفحات فیشینگ است و از طرفی Y_i برابر شماره واقعی یک کلاس از نظر نوع فیشینگ یا نرمال است. در این رابطه، α و β به ترتیب دو ضریب وزنی می باشند که مقدار آنها بین صفر و یک است و مجموع آنها نیز برابر یک است که اهمیت عامل خطا و تعداد ویژگی انتخاب شده را تعیین می کنند. هر بردار ویژگی که بتواند این تابع را کمینه نماید به عنوان بردار ویژگی بهینه برای آموزش شبکه عصبی مصنوعی استفاده می شود.

الگوریتم بهینه سازی شیرمورچه را می توان یک الگوریتم فراابتکاری با

$$c_i^t = Antlion_j^t + c^t \quad (۹)$$

$$d_i^t = Antlion_j^t + d^t \quad (۱۰)$$

در این رابطه، $Antlion_j^t$ موقعیت یک بردار ویژگی از نوع شیر مورچه مانند J در تکرار t -ام مسئله است. حرکت بردار ویژگی از نوع مورچه به سمت بردار ویژگی از نوع شیرمورچه در واقع یک جستجو برای یافتن بهینه‌های سراسری در فضای مسئله است. در الگوریتم بهینه‌سازی شیرمورچه نیاز است که با نزدیک شدن یک بردار ویژگی از نوع مورچه به سمت حفره مرتباً فاصله یک بردار ویژگی از نوع مورچه از شیرمورچه‌ای که در انتهای حفره است مطابق رابطه (۱۱) و (۱۲) کاهش یابد [۱۹]:

$$c^t = \frac{c^t}{I} \quad (۱۱)$$

$$d^t = \frac{d^t}{I} \quad (۱۲)$$

در روابط فوق، I نرخ کاهش شعاع، c^t و d^t به ترتیب کمینه و بیشینه شعاع عملیاتی یک بردار ویژگی از نوع شیرمورچه می‌باشد. در الگوریتم بهینه‌سازی شیرمورچه نرخ کاهش شعاع را می‌توان به صورت رابطه (۱۳)، محاسبه نمود [۱۹]:

$$I = 10^w \frac{t}{T} \quad (۱۳)$$

در این رابطه، t شماره تکرار فعلی الگوریتم بهینه‌سازی شیر مورچه و T حداکثر تعداد تکرار الگوریتم شیر مورچه و w یک ثابت است.

۴. شبیه‌سازی روش پیشنهادی

برای تشخیص حملات فیشینگ نیاز به مجموعه داده استاندارد است که یک نمونه آن مجموعه داده رامی و همکاران است. این مجموعه داده را می‌توان از پایگاه داده آنلاین UCI و Kaggle بارگذاری و دانلود نمود. این مجموعه داده دارای حدود ۱۱ هزار رکورد است و می‌توان یک بخش از آن را برای یادگیری استفاده نمود. این مجموعه داده دارای ۳۰ ویژگی است و هر کدام از این ویژگی‌ها می‌تواند معرف وبگاه‌های فیشینگ باشد و ویژگی شماره ۳۱ نیز نشان دهنده آن است که وبگاه جعلی یا اصلی است [۲۰]. برای پیاده‌سازی از متلب نسخه ۲۰۱۸ استفاده شده است. تعداد بردار ویژگی برابر ۲۰ تنظیم شده است که برابر تعداد کل مورچه‌ها و شیرمورچه‌ها است. تعداد اعضای جمعیت اولیه بردارهای ویژگی از نوع مورچه‌ها که در اینجا

$$M_{OA} = \begin{bmatrix} f([A_{1,1}, A_{1,2}, \dots, A_{1,d}]) \\ f([A_{2,1}, A_{2,2}, \dots, A_{2,d}]) \\ \vdots \\ f([A_{n,1}, A_{n,2}, \dots, A_{n,d}]) \end{bmatrix} \quad (۵)$$

در روابط فوق، d و n به ترتیب تعداد ابعاد مسئله برای انتخاب ویژگی و اندازه جمعیت بردارهای ویژگی از نوع مورچه‌ها است. علاوه بر ماتریس فوق می‌توان بردارهای ویژگی از نوع شیرمورچه‌ها و میزان شایستگی آنها را نیز به ترتیب مطابق رابطه (۶) و (۷) کدگذاری نمود [۱۹]:

$$M_{Antlion} = \begin{bmatrix} AL_{1,1} & AL_{1,2} & \dots & AL_{1,d} \\ AL_{2,1} & AL_{2,1} & \dots & AL_{2,d} \\ \vdots & \vdots & \vdots & \vdots \\ AL_{n,1} & AL_{n,2} & \dots & AL_{n,d} \end{bmatrix} \quad (۶)$$

$$M_{OAL} = \begin{bmatrix} f([AL_{1,1}, AL_{1,2}, \dots, AL_{1,d}]) \\ f([AL_{2,1}, AL_{2,2}, \dots, AL_{2,d}]) \\ \vdots \\ f([AL_{n,1}, AL_{n,2}, \dots, AL_{n,d}]) \end{bmatrix} \quad (۷)$$

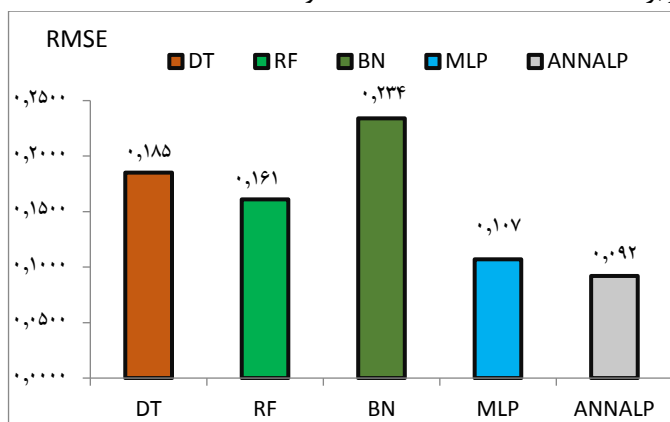
با کدگذاری اعضای جمعیت اولیه بردارهای ویژگی از نوع مورچه‌ها و بردارهای ویژگی از نوع شیرمورچه‌ها در قالب دو جمعیت ماتریس مانند، می‌توان مراحل الگوریتم بهینه‌سازی شیرمورچه را بر روی آنها به اجرا گذاشت. بردار ویژگی از نوع مورچه‌ها با حرکت‌های تصادفی و مبتنی بر گام می‌تواند به سمت یک بردار ویژگی از نوع شیرمورچه انتخاب شده حرکت نمایند و برای این حرکت می‌تواند از رابطه (۸)، استفاده نماید [۱۹]:

$$X_i^t = \frac{(X_i^t - a_i) \times (d_i - c_i^t)}{(d_i^t - a_i)} \quad (۸)$$

در رابطه فوق، X_i^t موقعیت بردار ویژگی از نوع مورچه i -ام در تکرار t -ام، a_i مقدار کمینه بعد i -ام، b_i مقدار بیشینه بعد i -ام، d_i^t و c_i^t به ترتیب بیشینه و کمینه بعد i -ام یک بردار ویژگی از نوع مورچه در تکرار t -ام است. با حرکت بردار ویژگی از نوع مورچه‌ها به سمت بردار ویژگی از نوع شیر مورچه آنها در گودال بردار ویژگی از نوع شیر مورچه گرفتار شده و توسط شیر مورچه شکار می‌شود. برای شبیه‌سازی این حالت موقعیت یک بردار ویژگی از نوع شیرمورچه را در نظر گرفته و به فاصله و شعاع c^t و d^t از آن می‌توان یک حفره تعریف نمود که ضابطه آن در رابطه (۹) و (۱۰) بیان شده است [۱۹]:

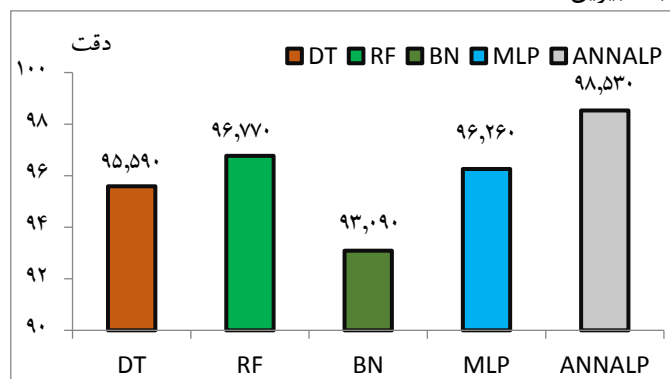
بهبود تشخیص وبگاه های جعل شده با استفاده از طبقه بندی کننده شبکه عصبی مصنوعی چند لایه با الگوریتم بهینه سازی شیرمورچه

مصنوعی چند لایه در تشخیص حملات فیشینگ به ترتیب برابر ۰,۰۹۲، ۰,۱۸۵، ۰,۱۶۱، ۰,۲۳۴ و ۰,۱۰۷ است.



شکل ۳: مقایسه خطای تشخیص روش پیشنهادی با چند طبقه بندی کننده

ارزیابی‌ها نشان می‌دهد دقت روش پیشنهادی، درخت تصمیم‌گیری، جنگل تصادفی، شبکه بیزین و شبکه عصبی مصنوعی چند لایه در تشخیص حملات فیشینگ به ترتیب برابر ۹۸,۵۳٪، ۹۵,۵۹٪، ۹۶,۷۷٪، ۹۳,۰۹٪، ۹۶,۲۶٪ است. در نمودار شکل (۴) به ترتیب مقایسه روش پیشنهادی با چند طبقه‌بندی کننده در شاخص دقت نمایش داده شده است. ارزیابی‌ها و مقایسه‌ها نشان می‌دهد دقت روش پیشنهادی از روش‌های درخت تصمیم‌گیری، جنگل تصادفی، شبکه بیزین و شبکه عصبی مصنوعی چند لایه در تشخیص حملات فیشینگ بیشتر است و بدترین عملکرد در بین این روشها مرتبط با شبکه بیزین است.



شکل ۴: مقایسه دقت تشخیص روش پیشنهادی با چند طبقه بندی کننده

دقت شبکه عصبی چند لایه در تشخیص حملات فیشینگ در حدود ۹۶,۲۶٪ است و دقت جنگل تصادفی نیز ۹۶,۷۷٪ است و بعد از

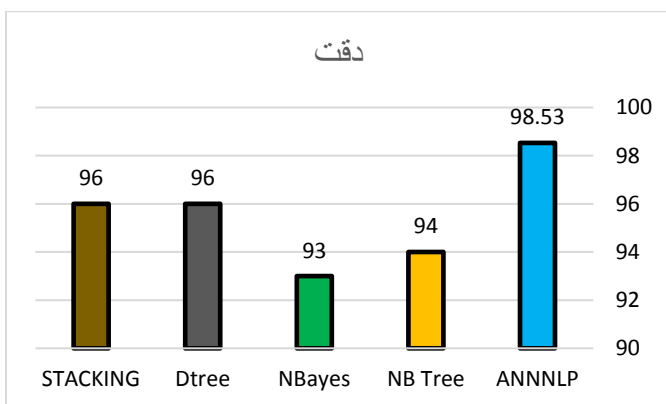
برابر ۱۰ و نصف کل جمعیت اولیه مورچه و شیرمورچه‌ها است. حداکثر تکرار الگوریتم پیشنهادی برای انتخاب ویژگی که برابر ۳۰ تنظیم شده است. تعداد لایه‌های شبکه عصبی که برابر ۲، ۳ و ۴ تنظیم شده است. تعداد نورون‌های لایه پنهان که برابر ۳۰ است. برای ارزیابی‌ها در این پژوهش از شاخص خطا و شاخص دقت استفاده شده است که به ترتیب در رابطه (۱۴) و (۱۵) فرموله شده است:

$$rmse = \sqrt{\frac{1}{m} \sum_{i=1}^m (\bar{Y}_i - Y_i)^2} \quad (14)$$

$$Accuracy = \frac{TP + TN}{TP + TN + FP + FN} \quad (15)$$

در این رابطه، \bar{Y}_i شماره کلاس پیش بینی شده برای نمونه i -ام و Y_i مقدار واقعی کلاس یک نمونه از نظر فیشینگ یا عادی است. در اینجا i یک شماره نمونه است و تعداد نمونه‌ها نیز برابر m است.

برای محاسبه شاخص دقت در ابتدا باید مفاهیم شاخص‌های مثبت واقعی، منفی واقعی، مثبت کاذب و منفی کاذب که به ترتیب با TP ، TN ، FP و FN نشان داده می‌شوند. تجزیه و تحلیل آماری نمودارها نشان می‌دهد با افزایش تعداد لایه‌های پنهان از ۲ به ۴، تابع هدف انتخاب ویژگی از ۰,۰۹۲ به ۰,۰۸۱ کاهش خواهد یافت و این کاهش فقط در حدود ۱۱,۹۵٪ است. با افزایش تعداد لایه‌های پنهان خطای $RMSE$ در تشخیص حملات فیشینگ از ۰,۲۱۴ به ۰,۱۹۷ کاهش خواهد یافت و این کاهش فقط در حدود ۷,۹۴٪ است. شاخص دقت بر حسب افزایش تعداد لایه‌های پنهان افزایش داشته است و این افزایش از ۹۸,۳۲٪ به حدود ۹۸,۵۳٪ است و افزایش فقط در حدود ۰,۲۱٪ بوده است. روش پیشنهادی برای تشخیص حملات فیشینگ از شبکه عصبی چند لایه به عنوان ابزار طبقه‌بندی کننده استفاده می‌کند. در این بخش روش پیشنهادی در ابزار وکا با چند طبقه‌بند مطرح مانند درخت تصمیم‌گیری، جنگل تصادفی، شبکه بیزین و شبکه عصبی مصنوعی چند لایه در شاخص دقت و $RMSE$ با هم مقایسه شده که مقایسه این روشها در شکل (۳)، ذکر شده است. ارزیابی‌ها نشان می‌دهد خطای روش پیشنهادی، درخت تصمیم‌گیری، جنگل تصادفی، شبکه بیزین و شبکه عصبی

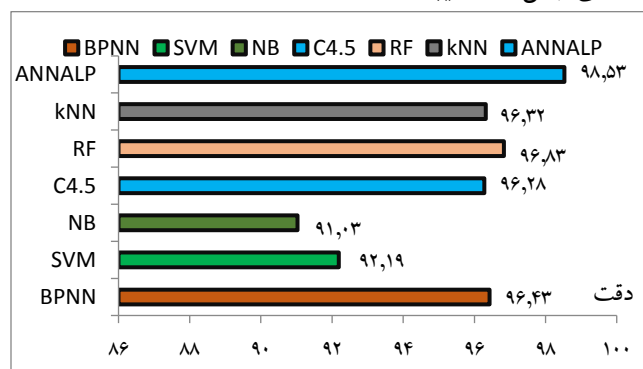


شکل ۶: مقایسه دقت تشخیص روش پیشنهادی با مقاله [۲۲]

۵. نتیجه گیری و پیشنهادات

در روش پیشنهادی بردار ویژگی بهینه به عنوان ورودی اصلی شبکه عصبی مصنوعی به عنوان طبقه‌بندی کننده صفحات جعلی از اصلی استفاده می‌شود. آزمایشها انجام شده روی بیش از ۱۱ هزار نمونه جعلی و اصلی در متلب انجام شده است. آزمایشها نشان داد روش پیشنهادی در تشخیص حملات فیشینگ نسبت به شبکه عصبی چند لایه خطایی کمتر و در حدود ۰.۱۴٪ دارد. روش پیشنهادی در تشخیص حملات فیشینگ نسبت به جنگل تصادفی دقتی در ۱.۷۶٪ بیشتر دارد و نسبت به شبکه عصبی چند لایه دقت آن در حدود ۲.۲۷٪ بیشتر شده است. روش پیشنهادی نسبت به ترکیب انتخاب ویژگی با الگوریتم ذرات و روشهای الگوریتم بهینه‌سازی ذرات با روشهای BPNN, SVM, NB, C4.5, RF, kNN و ANNNLP دارای دقت بیشتری است. همچنین با مقایسه روش پیشنهادی با روشهای ارایه شده در پژوهش [۲۲] مشاهده می‌گردد که روش پیشنهادی دارای دقت بالاتری می‌باشد. از پیشنهادات آتی ما استفاده از شبکه‌های عصبی یادگیری عمیق مانند LSTM در فاز یادگیری و طبقه‌بندی در تشخیص حملات فیشینگ است.

روش پیشنهادی، بهترین طبقه‌بند در تشخیص حملات فیشینگ جنگل تصادفی است و در مرتبه بعدی شبکه عصبی چند لایه است. روش پیشنهادی نسبت به جنگل تصادفی و شبکه عصبی چند لایه به ترتیب شاخص دقت را در حدود ۱.۷۶٪ و ۲.۲۷٪ بهبود داده است. روش پیشنهادی برای تشخیص حملات فیشینگ یک روش مبتنی بر انتخاب ویژگی است و از این جهت آن را با روش پژوهش [۲۱]، که سال ۲۰۲۰ انجام شده مورد مقایسه قرار داده‌ایم. در این پژوهش از ترکیب الگوریتم ذرات یا PSO با چند طبقه‌بندی کننده استفاده شده و صفحات جعلی از اصلی را مورد طبقه‌بندی قرار داده است. در نمودار شکل (۵)، یک مقایسه بین روش پیشنهادی و ترکیب انتخاب ویژگی الگوریتم بهینه‌سازی ذرات با روشهای SVM, BPNN, NB, C4.5, RF, kNN و انجام شده است. ارزیابی‌ها نشان دهنده آن است که روش پیشنهادی دارای دقتی برابر ۹۸.۵۳٪ است و این در حالی است که دقت روش RF, C4.5, NB, SVM, BPNN و kNN به ترتیب برابر ۹۶.۴۳٪، ۹۲.۱۹٪، ۹۱.۰۳٪، ۹۶.۲۸٪، ۹۶.۸۳٪ و ۹۶.۳۲٪ است. روش انتخاب ویژگی ذرات در ترکیب با جنگل تصادفی دارای بیشترین دقت است و این دقت در حدود ۹۶.۸۳٪ است و اما نسبت به روش پیشنهادی دقت آن کمتر است. در پژوهش [۲۲]، از پایگاه داده مشابهی برای ارزیابی چند روش تشخیص فیشینگ به نامهای Stacking, Dtree, NBayes, NB-Tree استفاده شده و روشها مقایسه شده اند، طبق شکل (۶) مشاهده می‌گردد روش پیشنهادی دارای دقت بالاتری در تشخیص وبگاه‌های جعل شده می‌باشد



شکل ۵: مقایسه خطای تشخیص روش پیشنهادی با چند روش

- techniques". *Telecommunication Systems*, 76(1), pp. 139-154, 2021.
- [10] https://docs.apwg.org/reports/apwg_trends_report_q1_2021.pdf
- [11] O. K. Sahingoz, E. Buber, O. Demir & B. Diri, "Machine learning based phishing detection from URLs". *Expert Systems with Applications*, 117, pp. 345-357, 2020.
- [12] M. H. Alkawaz, J. S. Stephanie, F. M. Omar, and Md. Johar, "Identification and analysis of phishing website based on machine learning methods." In *2022 IEEE 12th Symposium on Computer Applications & Industrial Electronics (ISCAIE)*, pp. 246-251, 2022.
- [13] M. M. Uddin, K. A. Islam, M. Mamun, V. K. Tiwari & J. Park, "A Comparative Analysis of Machine Learning-Based Website Phishing Detection Using URL Information". In *2022 5th International Conference on Pattern Recognition and Artificial Intelligence (PRAI)*, pp. 220-224, 2022.
- [14] A. Chawla, "Phishing website analysis and detection using Machine Learning." *International Journal of Intelligent Systems and Applications in Engineering* 10, no. 1, pp.10-16, 2022.
- [15] R. S. Rao & A. R. Pais, "Detection of phishing websites using an efficient feature-based machine learning framework". *Neural Computing and Applications*, pp. 1-23, 2019.
- [16] M. Gori, J. Visumathi, M. Mahdal and M. Elangovan, "An effective and secure mechanism for phishing attacks using a machine learning approach." *Processes* 10, no. 7, pg. 1356, 2022.
- [17] L. Tang, & Q. H. Mahmoud "A survey of machine learning-based solutions for phishing website detection". *Machine Learning and Knowledge Extraction*, 3(3), 672-694.
- [18] B. Altay, T. Dokeroglu, & A. Cosar, "Context-sensitive and keyword density-based supervised machine learning techniques for malicious webpage detection". *Soft Computing*, 23(12), pp. 4177-4191, 2019.
- [19] S. Mirjalili, "The ant lion optimizer". *Advances in engineering software*, 83, pp. 80-98, 2015.
- [20] R. Mohammad, F. Thabtah & T. L. McCluskey, "Phishing websites dataset", 2015.
- [21] S. Priya, S. Selvakumar, and R. Leela Velusamy. "PaSOFuAC: Particle Swarm
- [1] M. Mohammed, and Kh. Swayeb, "Phishing Website Detection Using a Hybrid Approach Based on Support Vector Machine and Ant Colony Optimization." In *2023 IEEE 3rd International Maghreb Meeting of the Conference on Sciences and Techniques of Automatic Control and Computer Engineering (MI-STA)*, pp. 402-406, 2023.
- [2] M. A. Tawhid, & A. M. Ibrahim, "Hybrid Binary Particle Swarm Optimization and Flower Pollination Algorithm Based on Rough Set Approach for Feature Selection Problem". In *Nature-Inspired Computation in Data Mining and Machine Learning*, pp. 249-273, Springer, Cham, 2020.
- [3] Lakshmana Rao K., Srinivasa Rao R., Abraham A., and Gabralla L.A., "Multilayer stacked ensemble learning model to detect phishing websites." *IEEE Access* 10 (2022): 79543-79552.
- [4] R. S. Rao, A. R. Pais & P. Anand, "A heuristic technique to detect phishing websites using TWSVM classifier". *Neural Computing and Applications*, 33(11), pp. 5733-5752, 2021.
- [5] M. Raj, M. and A. Jothi, "Website Phishing Detection Using Machine Learning Classification Algorithms." In *International Conference on Applied Informatics*, pp. 219-233. Cham: Springer International Publishing, 2022.
- [6] A. Alhogail & A. Alsabih, "Applying Machine Learning and Natural Language Processing to Detect Phishing Email". *Computers & Security*, 102414, 2021.
- [7] L. Lakshmi, M. P. Reddy, C. Santhaiiah & U. J. Reddy, "Smart Phishing Detection in Web Pages using Supervised Deep Learning Classification and Optimization Technique ADAM". *Wireless Personal Communications*, 118(4), pp. 3549-3564, 2021.
- [8] M. Korkmaz, E. Koçyiğit, Ö. Şahingöz, and B. Diri, "A Hybrid Phishing Detection System by Using Deep Learning-Based URL and Content Analysis." *Elektronika ir Elektrotehnika* 28, no. 5, 2022.
- [9] A. Basit, M. Zafar, X. Liu, A. R. Javed, Z. Jalil & K. Kifayat, "A comprehensive survey of AI-enabled phishing attacks detection

فرهنگ پدیداران مقدم و مهشید صادقی باجگیران، دوفصلنامه فناوری اطلاعات و ارتباطات ایران، سال پانزدهم، شماره های ۵۵ و ۵۶، بهار و تابستان ۱۴۰۲،
صفحه ۲۹۹ الی ۳۱۰

Optimization Based Fuzzy Associative Classifier
for Detecting Phishing Websites." *Wireless
Personal Communications* 125, no. 1, pp. 755-784,
2022.