

بهبود روش شناسایی وب سایت فیشینگ با استفاده از داده‌کاوی روی

صفحات وب

مهديه بهارلو* علیرضا یاری**

*دانشجوی کارشناسی ارشد دانشگاه آزاد اسلامی، واحد علوم و تحقیقات تهران
**استادیار پژوهشگاه ارتباطات و فناوری اطلاعات

تاریخ دریافت: ۱۳۹۹/۰۲/۰۱ تاریخ پذیرش: ۱۳۹۹/۰۷/۳۰

نوع مقاله: پژوهشی

چکیده

فیشینگ یک نوع حمله اینترنتی در سطح وب است که هدف آن سرقت مشخصات فردی کاربران برای سرقت آنلاین است. فیشینگ دارای اثر منفی در از بین بردن اعتماد بین کاربران در کسب‌وکارهای الکترونیکی است؛ بنابراین در این تحقیق سعی بر بررسی روش‌های تشخیص وب‌سایت‌های فیشینگ با استفاده از داده‌کاوی شده‌است. شناسایی ویژگی‌های اصلی از صفحات وب فیشینگ یکی از پیش‌شرط‌های مهم در طراحی یک سیستم تشخیص فیشینگ دقیق است. در پژوهش حاضر، برای افزایش کارایی سامانه تشخیص فیشینگ، یک روش ترکیبی برای کاهش ویژگی‌های وب سایت‌های فیشینگ پیشنهاد شده است. این روش ترکیبی از روش‌های انتخاب و کاهش ویژگی است که در دو مرحله انجام می‌شود. برای پیاده‌سازی و ارزیابی این روش پیشنهادی، بعد از کاهش ویژگی‌ها دسته بندی داده‌ها از طریق روش‌های درخت تصمیم‌گیری J48، جنگل تصادفی و بی‌زین ساده مورد بررسی قرار گرفت. نتایج نشان می‌دهند دقت مدل ایجاد شده برای تعیین وب سایت‌های فیشینگ با استفاده از کاهش ویژگی دو مرحله‌ای مبتنی بر پوششی و الگوریتم تحلیل مؤلفه اصلی در روش جنگل تصادفی که به میزان ۹۶٫۵۸٪ است، نسبت به سایر روش‌ها نتیجه مناسب‌تری را دارد.

واژگان کلیدی: حملات اینترنتی، فیشینگ، داده‌کاوی، انتخاب ویژگی، استخراج ویژگی

۱- مقدمه

تهدیدات و حملات مختلف در آن به وجود آید که ممکن است باعث خسارت مالی، سرقت هویت، از دست دادن اطلاعات خصوصی، آسیب شهرت نام تجاری و از دست دادن اعتماد مشتریان در تجارت الکترونیک شود [۱]. عواملی که

امروزه اینترنت به یک جزء ضروری از زیرساخت‌های اجتماعی و اقتصادی روزمره مردم تبدیل شده است. حجم بالای اطلاعات محرمانه و امنیتی اینترنت باعث می‌شود انواع

نویسنده مسئول: علیرضا یاری a_yari@itrc.ac.ir

تلفن، حمله بایوزی، فاپیدن تب، فارمینگ، فیشینگ موتور جستجو و فیشینگ مبتنی بر بدافزار دسته بندی کرد [۶،۷]. برای تشخیص صفحات فیشینگ از روش‌های متعددی استفاده می‌شود که هر یک به دسته‌ای از ویژگی‌های صفحات می‌پردازند. در این مقاله روشی جدیدی پیشنهاد شده است که برای دسته‌بندی داده‌های وب‌سایت‌های فیشینگ، از کاهش دمرحله‌ای ویژگی‌ها استفاده می‌نماید. بدین صورت که در مرحله اول بر اساس یک روش انتخاب ویژگی، زیرمجموعه‌ای مفید از ویژگی‌ها انتخاب شده و در مرحله دوم با کاهش بیشتر ترکیبی از ویژگی‌های مفید با ابعاد کمتر حاصل می‌شود.

در مقاله جاری، در ابتدا ابزارها و روش‌های مقابله با فیشینگ در قسمت پیشینه تحقیقات معرفی خواهند شد. سپس در بخش ۳ روش انجام کار، ویژگی‌های که توسط فیشرها برای ایجاد وب سایت‌های جعلی استفاده می‌شوند و همچنین روش‌های کاهش ویژگی و روش پیشنهادی کاهش ویژگی در این مقاله معرفی می‌گردد. در ادامه در بخش ۴ داده‌های حاصل از انجام آزمایشات مورد تجزیه و تحلیل قرار می‌گیرند و نتیجه با کارهای مشابه گذشته مورد مقایسه قرار می‌گیرد. در نهایت در بخش نتیجه گیری، نتیجه حاصل از انجام تحقیق بررسی خواهد شد.

۲- پیشینه پژوهش

تکنیک‌های فیشینگ روزبه‌روز در حال افزایش است و درعین حال پیچیده‌تر می‌شود. در نتیجه نیاز فوری به یافتن راه‌حل‌های مناسب برای مبارزه با حملات فیشینگ وجود دارد. تاکنون، راه‌حل‌های مختلفی در پاسخ به حملات فیشینگ پیشنهاد شده است. این راه‌حل‌ها با توجه به شکل (۱) به سه روش مقابله با فیشینگ شامل ابزارهای ضد فیشینگ، راه‌حل‌های فنی و غیر فنی تقسیم می‌شوند [۸]. از جمله ابزار ضد فیشینگ می‌توان افزونه مرورگر WOT [۸]، سایت ضد فیشینگ فیش تانک [۹]، سازمان جیو تراست، موزیلا تاندر برد [۱۰] را نام برد؛ اما در ادامه به توضیح راه‌حل‌های فنی و غیر فنی خواهیم پرداخت.

تهدید و حمله را در یک شبکه اینترنتی به وجود می‌آورند عبارتند از: دسترسی بدون محدودیت به اینترنت، گمنامی افراد، سرعت بالای انتشار، عدم ارتباط چهره به چهره، دسترسی آزاد به خدمات و محتویات ارزشمند و همچنین عدم وجود قوانین و توافقات مناسب [۲]؛ بنابراین، مناسب بودن اینترنت به‌عنوان یک کانال برای انجام معاملات تجاری مطرح می‌شود.

در اوایل ۱۹۹۰، با محبوبیت رو به رشد اینترنت، ما شاهد تولد یک نوع جدید از جرائم اینترنتی بودیم؛ که فیشینگ نام دارد [۳].

برخلاف سایر روش‌های هک و ورود به سیستم، در روش فیشینگ معمولاً هیچ نفوذی انجام نشده و از رخنه‌ها و آسیب‌پذیری‌ها استفاده نمی‌شود. بلکه خود کاربر است که با استفاده از روش‌های گوناگون فریب‌خورده و اطلاعاتی نظیر نام کاربری، رمز عبور، اطلاعات حساب بانکی را در اختیار حمله‌کننده که به اصطلاح فیشر نامیده می‌شود، قرار می‌دهد [۴]. طبق بررسی انجام شده توسط "گروه کاری ضد فیشینگ" تعداد حملات فیشینگ سراسر جهان در سه ماهه چهارم سال ۲۰۱۹ کاهش یافته و به میانگین نزدیکتر شده است. البته در همین سال در کشور برزیل تا ۲۳۲ درصد افزایش یافته است. حملات فیشینگ که کاربران وب، ایمیل و سرویس‌های نرم افزاری را هدف قرار می‌دهد، همچنان بزرگترین گروه حملات فیشینگ است. تقریباً سه چهارم از همه سایت‌های فیشینگ اکنون از حفاظت SSL استفاده می‌کنند، بالاترین میزان ثبت شده از اوایل سال ۲۰۱۵، و این نشانگر این است که کاربران نمی‌توانند به تنهایی به SSL اعتماد کنند و برای درک درست نیاز به ویژگی‌های بیشتری دارند [۵].

برای وب‌سایت فیشینگ تعاریف زیادی ارائه شده است که می‌توان تمام تعاریف را به صورت جامع و کامل در یک جمله بیان نمود: "وب‌سایت فیشینگ عمل ایجاد یک کپی از یک وب‌سایت قانونی و استفاده از مهارت‌های اجتماعی برای فریب قربانی برای ارسال اطلاعات شخصی او است" [۳].

انواع حملات فیشینگ را می‌توان فیشینگ سرنیزه یا فیشینگ هدفمند، کلون فیشینگ، صید نهنگ، فیشینگ

جعلی؛ لیست سفید در واقع لیست صفحاتی هست که قانونی هستند و اما در مقابل در لیست سیاه سعی می‌شود تمام صفحات جعلی شناسایی و پوشش داده شود. در این روش شناسایی تمام صفحات و به‌روز رسانی اطلاعات کار بسیار دشوار و زمان بر هست، چراکه بطور دائم صفحات زیادی ایجاد شده و یا از بین می‌روند [۱۲].

برخی از ارائه‌دهندگان لیست سیاه مانند کاوش ایمن گوگل [۱]، فیش نت، لیست سیاه مبتنی بر DNS [۶] و نرم‌افزارهای ضد فیشینگ نت‌کرفت، وب‌سن و کلودمارک [۱۳] با استفاده از این روش مانع از حملات فیشینگ می‌شوند.

۲-۲-۲- روش‌های اکتشافی

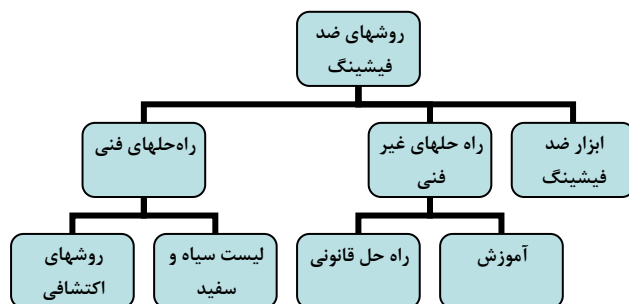
در این روش از ویژگی‌های وب‌سایت‌های فیشینگ برای تشخیص وب‌سایت‌های فیشینگ جدید استفاده می‌شود. در واقع قدرت اصلی تشخیص‌دهنده‌های فیشینگ بر اساس روش اکتشافی این است که آن‌ها قادر به تشخیص وب‌سایت‌های فیشینگ جدید هستند [۱۰].

۲-۲-۲-۱- روش‌های مبتنی بر الگوریتم‌های فازی

یک روش که توسط ابارس و همکاران در سال ۲۰۱۰ به‌کار گرفته شده است، بر مبنای الگوریتم‌های دسته‌بندی مبتنی بر قانون فازی برای تشخیص وب‌سایت‌های فیشینگ بانکداری الکترونیکی است [۱۳]. سعیدی در سال ۱۳۹۴ از روش دسته‌بندی مجموعه‌های فازی و روش ترکیبی تصمیم‌گیری AHP_TOPSIS برای تشخیص سریع‌تر و کارآمدتر وب‌سایت‌های فیشینگ در بانکداری الکترونیکی استفاده کرده است [۱۴]. عبدالحمید و همکاران در سال ۲۰۱۴ متد خاص دسته‌بندی انجمنی به نام دسته‌بند چند برچسب بر مبنای دسته‌بندی انجمنی را ارائه دادند [۱]. هادی و همکاران در سال ۲۰۱۶ از روش جدید دسته‌بندی انجمنی به نام الگوریتم دسته‌بندی انجمنی سریع استفاده کرده‌اند [۱۵].

۲-۲-۲-۲- روش‌های مبتنی بر یادگیری ماشین

باست و همکاران در سال ۲۰۱۲ از الگوریتم‌های یادگیری ماشین از جمله بیزین ساده، رگرسیون لجستیک و جنگل



شکل (۱): دسته‌بندی روش‌های شناسایی فیشینگ

۲-۱- راه‌های غیر فنی

۲-۱-۱- آموزش

آموزش مشتریان برای افزایش سطح آگاهی آن‌ها از جرائم آنلاین ضروری است تا با درک شاخص‌های امنیتی داخل وب‌سایت‌ها بتوانند به مقابله با فیشینگ بپردازند [۱۱]. اما این کار بسیار مشکل است چراکه کاربران باید زمان زیادی را صرف فراگیری متدهای فیشینگ کنند، علاوه بر آن فیشرها در ایجاد تکنیک‌های جدید هوشمندانه‌تر عمل می‌کنند [۳].

۲-۱-۲- راه‌های قانونی

این راه‌حل‌ها در کشورهای زیادی استفاده شده‌اند. ایالت متحده اولین کشوری بود که قانونی را در مورد این وب‌سایت‌ها وضع کرد و همین امر موجب گردید، فعالیت‌های زیادی توقیف گردد و عاملین آن به زندان انداخته شوند. با این حال، راه‌حل‌های قانونی نمی‌توانند مانع عمل وب‌سایت‌های جعلی شوند، چراکه ردیابی آن‌ها به دلیل مخفی شدن سریع آن‌ها در دنیای سایبری مشکل است [۱].

۲-۲- راه‌های فنی

علاوه بر روش‌های غیر فنی نظیر آموزش و قانون، برای مبارزه با فیشینگ، عموماً، دو متد فنی رایج در مبارزه با حملات فیشینگ، لیست سیاه و سفید و روش‌های اکتشافی است [۶].

۲-۲-۱- لیست سیاه و سفید

در روش لیست سیاه و سفید تمام URL‌های درخواستی با URL‌های موجود در لیست سیاه و سفید فیشینگ مقایسه می‌شوند تا مشخص شود وب‌سایت موردنظر قانونی است یا

مجموعه داده سایت داده کاوی UCI مورد بررسی قرار گرفته است. این مجموعه داده حاوی ۳۰ ویژگی مبتنی بر ۴ معیار زیر است [۲۶]:

- نوار آدرس
- غیرطبیعی بودن
- جاوا اسکریپت و HTML
- دامنه

در این قسمت با این ویژگی‌ها که در آزمایش‌ها مورد استفاده قرار گرفته‌اند آشنا می‌شویم. ویژگی‌ها به صورت صفر و یک و منفی یک کدگذاری شده‌اند.

۳-۱-۱- ویژگی‌های مبتنی بر نوار آدرس

- آدرس مبتنی بر IP
- URL های طولانی
- استفاده از خدمات کوتاه کننده URL
- استفاده از نماد @ در URL
- تغییر مسیر با "///"
- وجود نماد "-" در دامنه
- تعداد نقطه‌ها در دامنه
- HTTPS
- مدت ثبت دامنه
- فاوآیکون
- استفاده از درگاه غیراستاندارد
- وجود توکن HTTPS در دامنه

۳-۱-۲- ویژگی‌های مبتنی بر غیرطبیعی بودن

- درخواست URL
- تگ انکر
- پیوندها در تگ‌های <Script>, <Meta> و <Link>
- فرم هندلر در سرور
- ارسال اطلاعات به ایمیل
- URL غیرطبیعی

۳-۱-۳- ویژگی‌های مبتنی بر جاوا اسکریپت و

HTML

- ارسال وبسایت
- سفارشی‌سازی نوار وضعیت
- غیرفعال کردن کلیک راست
- استفاده از پنجره پاپ آپ
- تغییر مسیر با تگ آیفریم

۳-۱-۴- ویژگی‌های مبتنی بر دامنه

- طول عمر دامنه
- رکورد DNS
- ترافیک وبسایت
- رتبه‌بندی صفحه
- شاخص گذاری گوگل
- تعداد پیوندهای اشاره کننده به صفحه
- ویژگی‌های مبتنی بر گزارش‌های آماری

۳-۲- روش‌های کاهش ویژگی

روش‌های کاهش ابعاد داده به دودسته روش‌های مبتنی بر استخراج ویژگی و همچنین روش‌های مبتنی بر انتخاب ویژگی تقسیم می‌شوند [۲۷]. روش‌های انتخاب ویژگی سعی می‌کنند با انتخاب زیرمجموعه‌ای از ویژگی‌های اولیه، ابعاد داده‌ها را کاهش دهند. برخلاف روش‌های مبتنی بر استخراج ویژگی، این روش‌ها معنای اصلی ویژگی‌ها را بعد از کاهش حفظ می‌کنند. روش‌های مبتنی بر انتخاب ویژگی خود به سه روش فیلتر، پوششی و جاسازی شده تقسیم می‌شوند [۲۸]. هرکدام از این سه روش انتخاب ویژگی حاوی الگوریتم‌هایی برای اجرا هستند. در این قسمت به توضیح الگوریتم‌هایی که قرار است در این پژوهش مورد استفاده قرار گیرند می‌پردازیم.

۳-۲-۱- روش پوششی

در این روش انتخاب زیرمجموعه ویژگی با استفاده از الگوریتم یادگیری انجام می‌شود. الگوریتم دسته‌ای از ویژگی‌ها را برای یادگیری انتخاب می‌نماید و نهایتاً آن دسته

داشته و بر پایه احتمال وقوع یا عدم وقوع یک پدیده شکل می‌گیرد [۲۹].

• جنگل تصادفی

جنگل تصادفی مجموعه‌ای از درخت‌های تصمیم است که داده‌های آموزشی جهت ساخت هر درخت از روش انتخاب تصادفی با جایگذاری استفاده می‌کنند. هر درخت یک دسته‌بندی را می‌دهد که گفته می‌شود آن درخت به آن دسته رأی داده است. در انتها، دسته‌بندی که بیشترین رأی را داشته باشد انتخاب می‌شود. درخت‌ها هرس نمی‌شوند و در هر گره تعدادی ویژگی به‌طور تصادفی از مجموعه‌ی کل ویژگی‌ها برای انشعاب بررسی می‌شوند [۳۰].

• J48

این روش از معیار شاخص جینی جهت انتخاب ویژگی استفاده می‌کند [۲۷]. از میان ویژگی‌ها، هرکدام که مقدار شاخص جینی آن کوچک‌تر است، برای گروه جاری درخت تصمیم در نظر گرفته می‌شود.

۴-۳- روش پیشنهادی

همانطوریکه در شکل ۲ آمده است، روش پیشنهادی شامل سه مولفه هست: مولفه پیش‌پردازش، مولفه مدل‌سازی و مولفه آزمون و ارزیابی.

در این پژوهش در مولفه پیش‌پردازش داده‌ها برای کاهش ویژگی‌ها از روش ترکیبی جدیدی برای کاهش دومرحله‌ای ویژگی‌ها استفاده شده است. در این روش سعی شده است از پیچیدگی مسئله تا حد امکان کاسته شود و ویژگی‌هایی انتخاب گردد که در دسته‌بندی آن‌ها دارای بالاترین کارایی باشند و خطای دسته‌بندی نمونه‌ها را حتی‌الامکان کاهش دهد. بدین‌صورت که در مرحله اول بر اساس یک روش انتخاب ویژگی، زیرمجموعه‌ای مفید از ویژگی‌ها انتخاب می‌شود. برای این منظور با استفاده از روش‌های انتخاب ویژگی مبتنی بر پوششی، CFS و IG زیرمجموعه‌ای مهم از

⁠Gini coefficient

از ویژگی‌ها که دقت بالاتری دارند، انتخاب می‌شود. الگوریتمی که کار ارزیابی زیرمجموعه ویژگی‌ها و انتخاب بهترین زیرمجموعه را انجام می‌دهد، خود به‌عنوان بخشی از تابع ارزیابی، کار جستجو برای انتخاب بهترین مدل را انجام می‌دهد [۷].

• روش CFS

CFS مقدار همبستگی بین ویژگی‌ها و کلاس‌هایشان و همچنین همبستگی بین خود ویژگی‌ها را اندازه‌گیری می‌کند. ایده کلی این است که زیرمجموعه ویژگی‌های خوب، همبستگی زیادی با کلاس‌ها دارند، اما نباید با یکدیگر همبستگی داشته باشند [۱۷]. در الگوریتم CFS، هیوریستیکی برای ارزیابی ارزش یا شایستگی یک زیرمجموعه ویژگی وجود دارد [۱۸].

• روش IG

تفاوت بین آنترپی $H(S)$ از مجموعه داده S و آنترپی مشروط $H(S|F)$ از مجموعه داده که پس از جداسازی توسط ویژگی F ساخته شده، به دست می‌آید. آنترپی روش اندازه‌گیری ناخالصی در یک مجموعه داده است و اگر یک مجموعه داده تعداد مساوی از نمونه‌ها برای هر کلاس داشته باشد، مقدار آن حداکثر در نظر گرفته می‌شود [۱۷].

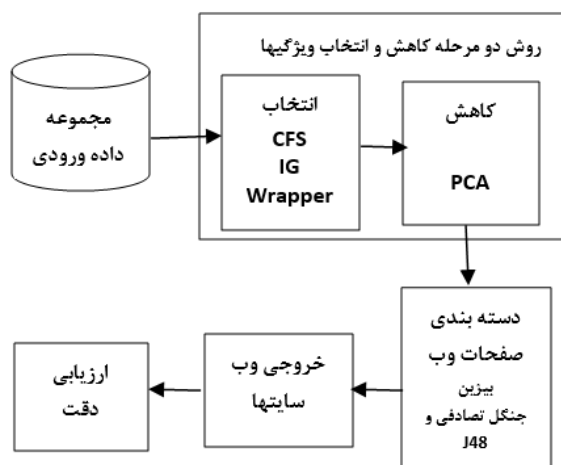
• روش PCA

یکی از عمومی‌ترین و شناخته‌شده‌ترین روش‌های آماری برای کاهش ویژگی‌ها است. هدف این روش به دست آوردن یک تبدیل تصویری است که از طریق آن بتوان با ترکیب خطی ویژگی‌های اصلی، تعداد کمتری ویژگی‌های جدید تولید نمود [۲۸].

۳-۳- الگوریتم‌های مدل‌سازی

• بیزین ساده

در روش بیزین ساده، دسته‌بندی بر پایه احتمالات و با فرض استقلال متغیرهای تصادفی ساخته می‌شود. این روش از ساده‌ترین الگوریتم‌های دسته‌بندی است که دقت قابل قبولی



شکل (۲): الگوریتم شناسایی فیشینگ

۴- آزمایش و ارائه نتیجه

در این قسمت از مقاله ابتدا به معرفی مجموعه داده و ابزار مورد استفاده در پژوهش می پردازیم؛ سپس به تحلیل نتایج حاصل از آزمایش ها پرداخته و در آخر نتایج حاصل را با چند روش دیگر مقایسه می کنیم.

۴-۱- دادگان آزمون

داده های مورد استفاده در این پژوهش تعدادی وبسایت جعلی و واقعی است که از سایت داده کاوی UCI [۲۶] استخراج شده است. نوع داده های این پژوهش از نوع دسته ای است. هر وبسایت در مجموعه داده دارای یک ویژگی کلاس یا هدف است که دسته آن را نشان می دهد. یک بودن ویژگی نشان دهنده وبسایت جعلی، صفر بودن آن نشان دهنده مشکوک و منفی یک نشان دهنده وبسایت واقعی است. در این مجموعه داده ۳۰ ویژگی در دسترس است که در بخش سوم تشریح گردید.

این مجموعه داده شامل ۱۱۰۵۵ وبسایت فیشینگ و قانونی متعلق به سال ۲۰۱۵ است که شامل ۴۸۹۸ وبسایت قانونی و ۶۱۵۷ وبسایت فیشینگ است.

۴-۲- ابزارهای آزمون

پایه سازی این پژوهش به کمک نرم افزار وکا نسخه ۳،۸،۱؛ یک نرم افزار منبع باز از دانشگاه Waikato [۱۲] است که

ویژگی ها انتخاب می شوند. سپس در مرحله بعد با اعمال روش استخراج ویژگی PCA بر روی ویژگی های باقیمانده، ترکیبی از این ویژگی ها با ابعاد کمتر به دست می آید. بدین صورت با کاهش دومرحله ای ویژگی ها، ترکیبی از ویژگی های مفید با ابعاد کمتر حاصل می شود. تکنیک های کاهش ویژگی دومرحله ای پیشنهادی، به صورت زیر نام گذاری شدند:

- ۱) کاهش ویژگی با استفاده از روش انتخاب ویژگی مبتنی بر CFS و اعمال روش استخراج ویژگی PCA بر روی آن (CFS+PCA)
- ۲) کاهش ویژگی با استفاده از روش انتخاب ویژگی مبتنی بر IG و اعمال روش استخراج ویژگی PCA بر روی آن (IG+PCA)
- ۳) کاهش ویژگی با استفاده از روش انتخاب ویژگی مبتنی بر پوششی و اعمال روش استخراج ویژگی PCA بر روی آن (Wrapper+PCA)

در مولفه دوم برای مدل سازی از الگوریتم های مختلف داده کاوی نظیر J48، جنگل تصادفی و بیژین ساده برای داده های آموزشی و آزمون استفاده شده است. از این مدل ها همچنین برای تشخیص دسته ی وبسایت استفاده شده است. با استفاده از قوانین حاصل از این مدل ها می توان وبسایت های ورودی را دسته بندی کرده و در دسته فیشینگ یا غیر فیشینگ جای داد.

در بخش آخر هم ارزیابی و اندازه گیری دقت دسته بندی با توجه به داده های آزمون صورت می گیرد.

جدول (۳): تعداد ویژگی‌ها پس از اعمال روش‌های کاهش ویژگی

تعداد ویژگی			نام روش
۳۰			داده‌های نرمال بدون کاهش ویژگی
۵			CFS + PCA
۸			IG + PCA
بیزین ساده	جنگل تصادفی	J48	Wrapper + PCA
۵	۸	۸	

پس از کاهش ویژگی‌ها، روش‌های دسته‌بندی درخت تصمیم J48، جنگل تصادفی و بیزین ساده به ترتیب با استفاده از تابع‌های J48، RandomForest و NaiveBayes در نرم‌افزار وکا اجرا شدند و معیارهای ارزیابی با استفاده از روش اعتبار سنجی تقاطعی با ۱۰ تکرار محاسبه و مورد ارزیابی قرار گرفتند تا بتوان به بهترین روش کاهش ویژگی و بهترین روش دسته‌بندی دست یافت.

جدول (۴): نتایج بررسی روش‌های دسته‌بندی با استفاده از دقت

روش‌های مدل‌سازی			روش‌های کاهش ویژگی
بیزین ساده	جنگل تصادفی	J48	داده‌های نرمال بدون کاهش ویژگی
۹۲,۹۸۰۶	۹۷,۲۵۹۲	۹۵,۹۷۴۷	CFS + PCA
۹۱,۷۰۵۱	۹۴,۶۳۵۹	۹۴,۰۱۱۸	IG + PCA
۹۱,۳۷۹۵	۹۶,۵۷۱۷	۹۵,۱۷۸۷	Wrapper + PCA
۹۱,۱۳۵۲	۹۶,۵۸۰۷	۹۵,۲۷۸۲	

در جداول (۴) و (۵) سه الگوریتم داده‌کاوی را با توجه به روش‌های کاهش ویژگی به ترتیب از نظر دقت و شاخص F1 مقایسه می‌کنیم تا به بهترین الگوریتم دسته‌بندی و بهترین روش کاهش ویژگی برسیم.

جدول (۴) نشان می‌دهد از نظر دقت که مهم‌ترین معیار ارزیابی است، الگوریتم‌های دسته‌بندی جنگل تصادفی و J48 با تعداد ۸ ویژگی و بیزین ساده با ۵ ویژگی به ترتیب دارای بهترین دقت هستند. از نظر بهترین روش کاهش

در محیط سیستم‌عامل ویندوز ۷ و پردازشگر Intel Core 2 Duo و RAM 2GB انجام گرفته است.

در این پژوهش از روش اعتبار سنجی تقاطعی با ۱۰ تکرار استفاده شده است که موجب می‌شود نتیجه به دست آمده دقیق‌تر باشد. اعتبار سنجی تقاطعی با ۱۰ تکرار موجب ارزیابی منطقی از مدل‌ها و کاهش سرریز می‌شود [۱۵]، این روش به این شکل است که داده‌ها به k قسمت تقسیم شده و $k-1$ قسمت آن به عنوان آموزش و ۱ قسمت به عنوان تست استفاده می‌شود؛ این عمل k مرتبه تکرار می‌شود. ارزیابی دقت نهایی برابر با میانگین k دقت محاسبه می‌شود [۳۱].

۳-۴- نتیجه آزمون

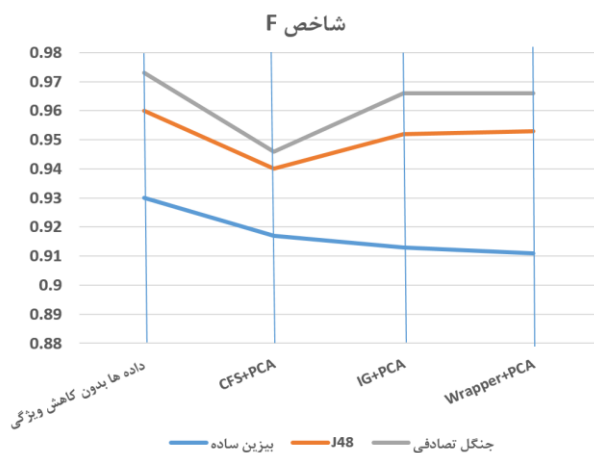
مدل آموزش یافته و آزمون شده با یک میزان دقت مشخص، می‌تواند جهت تشخیص کلاس یا اشیایی که برچسب کلاس آن‌ها ناشناخته است، مورد استفاده قرار گیرد [۳۲]. یکی از اهداف انتخاب ویژگی این است که یک زیرمجموعه از ویژگی‌ها برای افزایش دقت تشخیص، انتخاب شوند. به عبارت دیگر کاهش اندازه ساختار بدون کاهش قابل ملاحظه در دقت تشخیص دسته‌بندی که با استفاده از ویژگی‌های مدل به دست می‌آید، صورت گیرد [۱۵].

پس از اعمال روش‌های کاهش ویژگی بر روی داده‌ها، تعداد ویژگی‌ها به صورت نشان داده شده در جدول (۳) کاهش یافت. همان‌گونه که در جدول (۳) مشخص شده است، تعداد ویژگی‌ها با استفاده از تمام روش‌های کاهش ویژگی دو مرحله‌ای به طور قابل توجهی کاهش یافت. روش CFS+PCA و بیزین ساده در Wrapper+PCA به کمترین مقدار با تعداد ۵ ویژگی دست یافته‌اند؛ و مابقی روش‌ها به تعداد ۸ ویژگی رسیده‌اند.

همچنین براساس نتایج از ۴ معیار ذکر شده، معیار مبتنی بر جاوا اسکریپت و HTML با ۵ ویژگی متعلق به آن، کمترین تأثیر را در تشخیص وبسایت‌های فیشینگ دارد.

همان گونه که در جدول (۵) می بینیم با استفاده از کاهش ویژگی دو مرحله ای Wrapper+PCA و IG+PCA الگوریتم جنگل تصادفی بهترین عملکرد نسبت به سایر الگوریتم ها داشته و به مقدار یکسان ۰,۹۶۶ دست یافته است. پس از الگوریتم جنگل تصادفی الگوریتم J48 و بیزین ساده به ترتیب به بهترین نتیجه دست یافته اند.

مقدار F1 در بهترین نتیجه در الگوریتم جنگل تصادفی و J48 به میزان ۰,۰۰۷ و در بیزین ساده به میزان ۰,۰۱۳ نسبت به داده ها بدون انتخاب ویژگی کاهش یافته اند؛ که نشان می دهد الگوریتم جنگل تصادفی و J48 کمترین کاهش F1 را نسبت به مقدار اولیه داشته اند و از این نظر در یک سطح هستند. در نمودار (۲) شاخص F برای روش های دسته بندی در مقایسه با یکدیگر آمده است.



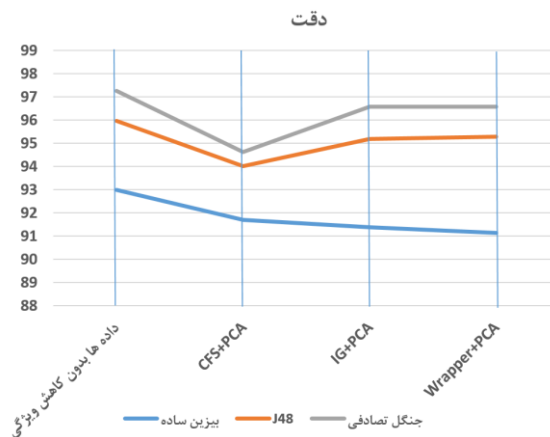
نمودار (۲): نتایج دسته بندی با استفاده از شاخص F

کاهش نامحسوس در مقدار دقت و F1 پس از کاهش ویژگی ها، حاصل افزایش هر دو مقدار نرخ FP و نرخ FN^۷ است. ولی همواره نرخ FN نسبت به نرخ FP پایین تر بوده و این امر بدین معنی است که احتمال اینکه یک وبسایت فیشینگ به عنوان وبسایت قانونی در نظر گرفته شود کمتر است. در نتیجه احتمال گیر افتادن در دام فیشینگ کاهش می یابد.

^۷False Positive

^۷False Negative

ویژگی ها در دو الگوریتم جنگل تصادفی و J48 به ترتیب روش های Wrapper+PCA، IG+PCA و CFS+PCA دارای بهترین نتیجه هستند؛ اما در الگوریتم بیزین ساده این ترتیب بالعکس است. در نمودار (۱) دقت روش های دسته بندی در مقایسه با یکدیگر آمده است.



نمودار (۱): نتایج دسته بندی با استفاده از شاخص دقت

مقدار دقت در الگوریتم جنگل تصادفی، J48 و بیزین ساده در بهترین نتیجه به ترتیب به میزان ۰,۶۷۸۵، ۰,۶۹۶۵ و ۱,۲۷۵۵ نسبت به داده ها بدون انتخاب ویژگی کاهش یافته اند؛ که نشان می دهد الگوریتم جنگل تصادفی کمترین کاهش را نسبت به دو الگوریتم دیگر داشته است. همچنین بالاترین دقت با مقدار ۹۶,۵۸۰۷ مربوط به الگوریتم جنگل تصادفی با روش کاهش ابعاد Wrapper+PCA است؛ در حالی که تعداد ویژگی ها کاهش محسوسی نداشته اند.

جدول (۵): نتایج بررسی روش های دسته بندی با استفاده از F1

روش های مدل سازی			روش های کاهش ویژگی
بیزین ساده	جنگل تصادفی	J48	
۰,۹۳	۰,۹۷۳	۰,۹۶	داده های نرمال بدون کاهش ویژگی
۰,۹۱۷	۰,۹۴۶	۰,۹۴	CFS + PCA
۰,۹۱۳	۰,۹۶۶	۰,۹۵۲	IG + PCA
۰,۹۱۱	۰,۹۶۶	۰,۹۵۳	Wrapper + PCA

۳-۴- مقایسه با روش‌های مشابه

در ادامه برای ارزیابی هرچه بیشتر رویکرد فوق، روش کاهش ویژگی دومرحله‌ای پیشنهادی با چند روش دیگر مقایسه می‌شود. از آنجاکه در اغلب روش‌های مشابه از معیار دقت برای ارزیابی استفاده کرده‌اند ما برای مقایسه روش خود با روش‌های دیگر از این معیار استفاده کردیم. جدول (۶) دقت به‌دست‌آمده از تحقیقات گذشته را با دقت به‌دست‌آمده از روش پیشنهادی نشان می‌دهد.

جدول (۶): مقایسه روش پیشنهادی با روش‌های مشابه

ردیف	روش	تعداد ویژگی‌ها	تعداد ویژگی‌های باقی‌مانده	دقت	روش پیشنهادی
	روش پیشنهادی	۱۱۰۵۵	۸	۳۰	۹۶,۵۸
۲۰۱۱	[۱۹]	۳۰۰۰	۷	۷	۹۳
۲۰۱۲	[۲۳]	۲۵۰۰	۱۲	۲۳	۹۷,۶
۲۰۱۴	[۱]	۱۳۵۳	۹	۱۶	۹۴,۴
۲۰۱۵	[۲۵]	۱۰۰۰	۱۵	۱۵	۹۶,۵۶
۲۰۱۵	[۱۸]	۲۴۵۶	۱۵	۳۰	۹۷,۴۷
۲۰۱۶	[۳۳]	۱۱۰۵۵	۱۲	۳۰	۹۲,۵
۲۰۱۹	[۲۰]	۱۱۰۵۵	۶	۳۰	۹۴,۶

نتایج مطالعات ذکرشده محققان در مقایسه با روش جدید کاهش ویژگی دومرحله‌ای پیشنهادی، حاکی از برتری روش پیشنهادی است چراکه تعداد ویژگی‌های ما در بهترین عملکرد (۸ ویژگی) از تمامی آن‌ها کمتر بوده درحالی‌که به‌دقت بالاتری نسبت به آن‌ها دست‌یافته‌ایم. مجموعه داده استفاده‌شده در این پژوهش با مجموعه داده استفاده‌شده در [۳۳] یکسان بوده که در آن از یک روش طبقه‌بندی مشارکتی^۱ با الگوریتم FACA استفاده‌شده است.

در خصوص تحقیق‌های جدول ۶ که دقت بالاتری نسبت به روش پیشنهادی دارند، می‌بایست به تعداد ویژگی‌ها و همچنین تعداد نمونه به‌کاربرده شده توجه نمود.

تحقیق [۲۰] نیز با کاهش دو مرحله‌ای به تعداد ویژگی ۲۰ درصدی دست یافته‌است.

۵- نتیجه‌گیری

با توجه به مشکلات و پیچیدگی‌های زیاد و سایر چالش‌ها در تشخیص صفحات فیشینگ، بررسی و ارائه روشی هوشمندانه برای تشخیص فیشینگ در صفحات وب ضروری است. در این پژوهش برای ساده‌سازی، کاهش دومرحله‌ای ویژگی‌ها پیشنهاد شده که در آن ابتدا با استفاده از روشهای پوششی، CFS و IG و ویژگیها انتخاب شده و سپس با اعمال روش استخراج ویژگی PCA از میان آنها بهترین ویژگی‌ها انتخاب می‌شوند که نهایتاً لیست خیلی کوتاهی از ویژگیهای صفحات وب حاصل می‌شود.

نتایج حاصل از اجرای روش پیشنهادی بر روی مجموعه داده UCI، در بهترین حالت به‌دقت ۹۶,۵۸٪ با روش دومرحله‌ای Wrapper+PCA و الگوریتم جنگل تصادفی دست‌یافت. در این روش تعداد ویژگی‌ها از ۳۰ ویژگی به ۸ ویژگی کاهش یافت.

مزیت مدل پیشنهادی نسبت به سایر سامانه‌های مشابه، دستیابی به کمترین تعداد ویژگی‌ها پس از کاهش ویژگی است که این امر موجب ساده‌سازی و کاهش پیچیدگی مدل می‌شود، مزیت دیگر استفاده از تعداد نمونه بیشتر است که موجب می‌شود نتایج به واقعیت نزدیک‌تر باشند.

به‌عنوان کارهای آتی به محققان پیشنهاد می‌شود زمان محاسباتی الگوریتم‌های یادگیری را موردبررسی و تحلیل قرار دهند. همچنین پیشنهاد می‌شود از مجموعه داده حاوی وبسایت‌های فارسی برای تشخیص وبسایت‌های فیشینگ استفاده نمایند.

^۱Associative classification

مراجع

(ICACCS -2016), Jan. 22 – 23, 2016, Coimbatore, INDIA, Available: IEEE Xplore, <http://www.ieee.org>.

[۱۱] محمدی، شهریار، غروی، عرفانه، "کاربرد تکنیک‌های داده‌کاوی جهت تشخیص آدرس‌های فیشینگ"، کنگره ملی مهندسی برق، کامپیوتر و فناوری اطلاعات، مشهد: موسسه آموزش عالی خیام، ۱۳۹۲.

[۱۲] Sanglerdsinlapachai, N., Rungsawang, A., "Using Domain Top-page Similarity Feature in Machine Learning-based Web Phishing Detection", *Third International Conference on Knowledge Discovery and Data Mining*, IEEE, pp. 17-190, 2010.

[۱۳] Aburrous, M., Hossain, M. A., Keshav, D., Thabtah, F., "Predicting Phishing Websites using Classification Mining Techniques with Experimental Case Studies", *IEEE Seventh International Conference on Information Technology*, pp. 176-181, 2010.

[۱۴] سعیدی، پریسا، "بررسی سیستم‌های هوشمند تشخیص وب‌سایت فیشینگ در بانکداری الکترونیکی به روش منطق فازی"، نخستین کنفرانس بین‌المللی فناوری اطلاعات، تهران: مرکز همایش‌های توسعه ایران، ۱۳۹۴.

[۱۵] حاتمی خواه، نفیسه، "بررسی روش‌های مبتنی بر انتخاب ویژگی"، تهران، دانشگاه صنعتی مالک اشتر، ۱۳۹۲.

[۱۶] Basnet, R. B., Sung, A.H., Liu, Q., "Feature Selection for Improved Phishing Detection", *international conference on Industrial Engineering and Other Applications of Applied Intelligent Systems*, pp 252-261, 2012, Available: <https://link.springer.com>.

[۱۷] Khonji, M., Jones, A., Iraqi, Y., "A Study of Feature Subset Evaluators and Feature Subset Searching Methods for Phishing Classification", *Proceedings of the 8th Annual Collaboration, Electronic messaging, Anti-Abuse and Spam Conference*, pp.135-144, ACM, 2011.

[۱۸] Singh, P., Jain, N., Maini, A., "Investigating the Effect Of Feature Selection and

[۱] Abdelhamid, N., Ayes, A., Thabtah, F., "Phishing detection based Associative Classification data mining", *Expert Systems with Applications* 41 5948–5959, 2014.

[۲] معاونی، مسعود، "تشخیص حملات در بانکداری الکترونیکی با استفاده از سیستم ترکیبی فازی-راف (Fuzzy_rough)" گروه کامپیوتر دانشگاه امام رضا (ع)، ۱۳۹۴.

[۳] Mohammad, R. M., Thabtah, F., McCluskey, L., "Tutorial and critical analysis of phishing websites methods", *Computer Science Review* 17 (2015) 1-24.

[۴] Chaudhry, J. A., Rittenhouse, R. G., "Phishing: Classification and Countermeasures", *7th International Conference on Multimedia, Computer Graphics and Broadcasting*, pp. 28-31, IEEE, 2015.

[۵] Anti Phishing Working Group, Phishing activity trends report, http://www.antiphishing.org/resources/apwg-reports/apwg_trends_report_q4_2019.pdf.

[۶] Buber, E., Demir, Ö., Sahingoz, O.K., "Feature Selections for the Machine Learning based Detection of Phishing Websites", *International Artificial Intelligence and Data Processing Symposium (IDAP) IEEE*, 2017.

[۷] Kohavi, R., John, G. H., "Wrappers for feature subset selection", *Artificial Intelligence*, Vol. 97, pp. 273-324, 1997.

[۸] Abur-rous, M. R. M., "Phishing Website Detection Using Intelligent Data Mining Techniques", Ph.D, dissertation, Dept. Computing, Bradford Univ, Bradford, 2010.

[۹] PhishTank.<http://www.phishtank.com>, 2017.

[۱۰] Aravindhana, R., Shanmugalakshmi, Dr.R., Ramya, K., Dr.Selvan C, "Certain Investigation on Web Application Security: Phishing Detection and Phishing Target Discovery", *2016 3rd International Conference on Advanced Computing and Communication Systems*

[۲۶] Mohammad, R. M., Thabtah, F., McCluskey, L., Phishing Website Dataset, <https://archive.ics.uci.edu/ml/datasets/Phishing+websites>, 2015.

[۲۷] اسماعیلی، مهدی، مفاهیم و تکنیک‌های داده‌کاوی، کاشان: سوره، ۱۳۹۲.

[۲۸] ورسلیز، کارلو، هوش تجاری داده‌کاوی و بهینه‌سازی برای تصمیم‌گیری، ترجمه‌ی احمدی، عباس، محبی، آزاده، ویرایش دوم، تهران، نشر دانشگاه صنعتی امیرکبیر (پلی تکنیک تهران)، زمستان ۱۳۹۲.

[۲۹] H.John, George, and pat Langley, "Estimating Continuous Distribution in Bayesian Classifiers", In Proceeding of the Eleventh Conference on Uncertainty in Artificial Intelligence. Morgan Kaufman, 1995.

[۳۰] Breiman, Leo. "Random Forests", Machine Learning, Kluwer Academic Publishers. Manufactured in The Netherlands, Statistics Department University of California Berkeley, 45, 5-32, 2001.

[۳۱] Kohavi, Ron, "A Study of Cross-Validation and Bootstrap for Accuracy Estimation and Model Selection", *Proceedings of the 14th international joint conference on Artificial intelligence (IJCAI)*, pp. 1137-1143, ACM, 1995.

[۳۲] Lakhita, Yadav, S., Bohra, B., Pooja, "A Review on Recent Phishing Attacks in Internet", *IEEE International Conference on Green Computing and Internet of Things (ICGCIoT)*, pp. 1312-1315, 2015.

[۳۳] Hadi, W., Aburub, F., Alhawari, S., "A new fast associative classification algorithm for detecting phishing websites", *Applied Soft Computing* 48 (2016) 729-734.

"Dimensionality Reduction On Phishing Website Classification Problem", *1st International Conference on Next Generation Computing Technologies (NGCT) Dehradun, India, IEEE*, pp. 388-393, 2015.

[۱۹] rahmi A. H., isredza, Abawajy, J., "Phishing Email Feature Selection Approach", *10th International Joint Conference of IEEE TrustCom.*, pp. 916-921, 2011.

[۲۰] K. L. Chiew, C. L. Tan, K. Wong, K. S. Yong, and W. K. Tiong, "A new hybrid ensemble feature selection framework for machine learningbased phishing detection system," *Information Sciences*, vol. 484, pp. 153-166, 2019.

[۲۱] M. Almseidin, A. A. Zuraiq, M. Alkasassbeh, and N. Alnidami, "Phishing detection based on machine learning and feature selection methods," *International Journal of Interactive Mobile Technologies (IJIM)*, vol. 13, no. 12, pp. 171-183, 2019.

[۲۲] Meenu, Sunila godara, "Phishing Detection using Machine Learning Techniques", *International Journal of Engineering and Advanced Technology (IJEAT)*, Volume-9 Issue-2, December, 2019.

[۲۳] Pandey, M., Ravi, V., "Detecting phishing e-mails using Text and Data mining", *IEEE International Conference on Computational Intelligence and Computing Research (ICCIC)*, 2012.

[۲۴] Pandey, M., Ravi, V., "Text and Data Mining to Detect Phishing Websites and Spam Emails", *Proceedings of the 4th International Conference on Swarm, Evolutionary, and Memetic Computing*, Vol. 8298, pp.559-573, 2013.

[۲۵] لنگری، نفیسه، عبدالرزاق نژاد، مجید، "شناسایی وبگاه فیشینگ در بانکداری اینترنتی با استفاده از الگوریتم بهینه‌سازی صفحات شیب‌دار"، *مجله پدافند الکترونیکی و سایبری*. شماره ۱، صفحه ۴۰-۲۹، ۱۳۹۴.

An efficient method for detecting phishing websites using data mining on web pages

Abstract

Phishing is regarded as a kind of internet attack on the web which aimed to steal the users' personal information for online stealing. Phishing plays a negative role in reducing the trust among the users in the business network based on the E-commerce framework. therefore, in this research, we tried to detect phishing websites using data mining. The detection of the outstanding features of phishing is regarded as one of the important prerequisites in designing an accurate detection system. Therefore, in order to detect phishing features, a list of 30 features suggested by phishing websites was first prepared. A new idea based on two steps: feature selection and feature extraction, has been proposed. To evaluate the proposed method, the performance of decision tree J48, random forest, naïve Bayes methods were evaluated on the reduced features. The results indicated that accuracy of the model created to determine the phishing websites by using the two-stage feature reduction-based Wrapper and Principal Component Analysis (PCA) algorithm in the random forest method of 96.58%, which is a desirable outcome compared to other methods.

Keywords: Internet attack, Phishing, Data Mining, Feature Selection, Feature Extraction.