

## مروری بر روش‌های شناسایی سایت‌های فیشینگ از سایت‌های قانونی

فرح مکی الهیجل<sup>۱</sup>، زهرا یعقوبی<sup>۲\*</sup>

<sup>۱</sup> دانشجوی کارشناسی ارشد، دانشکده فنی و مهندسی، دانشگاه بین‌المللی امام خمینی (ره)، قزوین، ایران

<sup>۲</sup> استادیار دانشکده فنی و مهندسی، دانشگاه بین‌المللی امام خمینی (ره)، قزوین، ایران

### چکیده

### مقاله پژوهشی

تاریخ دریافت:

۱۴۰۱/۱۱/۲۴

تاریخ پذیرش:

۱۴۰۲/۳/۳۰

کلیدواژه‌ها:

شناسایی سایت‌های فیشینگ، یادگیری ماشین، یادگیری عمیق

نویسنده مسئول:

z.yaghoubi@eng.ikiu.ac.ir

محبوبیت روزافزون اینترنت منجر به رشد چشمگیر تجارت الکترونیک شد. با این حال، چنین فعالیت‌هایی دارای چالش‌های امنیتی اساسی است که عمدتاً ناشی از کلاهبرداری‌های سایبری و سرقت هویت افراد است. از این رو، بررسی مشروعیت صفحات وب بازدید شده یک کار بسیار مهم برای ایمن کردن هویت مشتریان و جلوگیری از حملات فیشینگ است. استفاده از روش‌های یادگیری ماشین، و یادگیری عمیق به طور گسترده به عنوان یک راه حل امیدوارکننده شناخته شده است. تحقیقات، مملو از مطالعاتی است که از روش‌های یادگیری ماشین و یادگیری عمیق برای تشخیص فیشینگ وب سایت استفاده می‌کنند. با توجه به این که در تحقیقات قبلی، تشخیص سایت-های فیشینگ با استفاده از ویژگی‌های مختلف توسط شبکه‌های عمیق بررسی شده است، در تحقیق جاری، ما قصد داریم تشخیص سایت‌های فیشینگ را با استفاده از چندین نوع ویژگی به صورت ترکیبی با کمک شبکه‌های عمیق انجام دهیم. بنابراین در این مقاله ابتدا ما به مرور روش‌های شناسایی سایت‌های فیشینگ از سایت‌های قانونی می‌پردازیم و در انتها روش پیشنهادی خود را ارائه می‌دهیم.



: 10.22034/ABMIR.2022.2714



## ۱- مقدمه

طبقه‌بندی کننده‌های کارآمد، (۲) استفاده از ویژگی‌های متمایز، و (۳) مجموعه‌ای از نمونه‌های داده نماینده برای آموزش. سیستم‌های مبتنی بر یادگیری ماشین که برای تشخیص فیشینگ توسعه یافته‌اند، عمدتاً به دودسته اصلی طبقه‌بندی می‌شود: سیستم‌های مبتنی بر محتوا<sup>۶</sup> و سیستم‌های مبتنی بر منبع‌یاب یکنواخت<sup>۷</sup>. سیستم‌های مبتنی بر محتوا، فیشینگ توسط بررسی فعال یا غیرفعال بودن محتوای صفحات وب بازدید شده، شناسایی می‌شود و سیستم‌های مبتنی بر منبع‌یاب یکنواخت، فقط URL های صفحات وب بازدید شده مورد بررسی قرار می‌گیرند [۱].

روش‌های مختلفی برای شناسایی وب‌سایت فیشینگ پیشنهاد شده است که عبارت‌اند از: مبتنی بر فهرست<sup>۸</sup>، شباهت بصری<sup>۹</sup>، اکتشافی<sup>۱۰</sup>، یادگیری ماشین و یادگیری عمیق<sup>۹</sup> [۴]. جدول (۲) این روش‌ها را نشان می‌دهد.

جدول (۲): انواع روش‌های شناسایی سایت‌های فیشینگ

روش‌های شناسایی فیشینگ				
اکتشافی	شباهت بصری	مبتنی بر لیست	یادگیری بر ماشین	یادگیری عمیق

## ۲-۱ روش مبتنی بر فهرست

مرورگرهایی مانند مایکروسافت آج<sup>۱</sup>، فایرفاکس و گوگل کروم از روش‌های مبتنی بر فهرست برای شناسایی وب‌سایت‌های فیشینگ استفاده می‌کند. لیست سفید و لیست سیاه دو نوع از روش‌های مبتنی بر لیست هستند. لیست سفید حاوی فهرستی از URL های معتبر است که مرورگرها می‌توانند به آن‌ها دسترسی داشته باشند، به این معنی که اگر URL در لیست سفید باشد، مرورگر می‌تواند صفحه وب را دانلود کند. در همان زمان، پایگاه داده لیست سیاه شامل فیشینگ یا URL های جعلی است که مرورگرها از دانلود صفحات آن سایت‌ها متوقف شده‌اند. عیب عمده این است که یک تغییر جزئی در URL برای دور زدن روش‌های مبتنی بر فهرست کافی است و برای جلوگیری از URL های فیشینگ جدید، این فهرست‌ها باید مرتباً به‌روز شوند [۵].

فیشینگ<sup>۱</sup> به‌عنوان ساده‌ترین و گسترده‌ترین تهدید جرائم سایبری شناخته شده است. هکرها نیازی به کرک کردن کد رمز پیچیده و نقض فایروال سخت را ندارند. به‌جای آن، به‌سادگی ایمیل‌های احساسی، انتقادی یا معقول را به‌منظور درخواست از گیرندگان برای معرفی اعتبار شخصی خود با کلیک بر روی یک پیوند ارسال می‌کنند. سپس گیرندگان به صفحات وب جعلی هدایت می‌شوند که بسیار شبیه به وب‌سایت‌های معتبر مورد هدف است. در نتیجه، گیرندگان مانند ماهی در وب‌سایت‌های جعلی گرفتار می‌شوند [۱]. اخیراً، هکرها کار خود را بسیار حرفه‌ای انجام می‌دهند، روند اخیر گزارش شده [۲] از فعالیت فیشینگ نشان داد که ۷۸ درصد از تمام وب‌سایت‌های فیشینگ، از محافظت لایه سوکت‌های امن<sup>۲</sup> استفاده می‌کنند که منحصراً توسط وب‌سایت‌های معتبر استفاده می‌شود.

## ۲-۲ مروری بر روش‌های شناسایی

واندرا [۳] در گزارش چشم‌انداز تهدیدات موبایلی خود در سال ۲۰۲۰ اعلام کرد که هر ۲۰ ثانیه یک وب‌سایت جدید فیشینگ راه‌اندازی می‌شود. همه این حقایق از مطالعات تحقیقاتی عمیق در مورد شناسایی و پیشگیری از چنین حملات سایبری حمایت می‌کند. جدول (۱) انواع حملات فیشینگ را نشان می‌دهد.

جدول (۱): انواع حملات فیشینگ

انواع حملات فیشینگ						
ایمیل / اسپم	تحویل مبتنی بر وب	دست‌کاری لینک	فیشینگ تلفن	دزدی داده	فیشینگ مبتنی بر بدافزار	فیشینگ مبتنی بر سیستم نام دامنه

تشخیص حملات فیشینگ از دسته مسائل طبقه‌بندی محسوب می‌شود. بنابراین، روش‌های یادگیری ماشین<sup>۳</sup> به‌عنوان راه‌حل‌های امیدوارکننده در نظر گرفته می‌شوند. با این حال، سه جنبه اصلی باید در هنگام پذیرش چنین روش‌هایی در نظر گرفته شود: (۱) انتخاب

<sup>6</sup> Lists-Based

<sup>7</sup> Visual Similarity

<sup>8</sup> Heuristic

<sup>9</sup> Deep Learning

<sup>10</sup> Edge

<sup>1</sup> Phishing

<sup>2</sup> Secure Sockets Layer

<sup>3</sup> Machine Learning

<sup>4</sup> Content-based

<sup>5</sup> (Uniform Resource Locator) URL-based

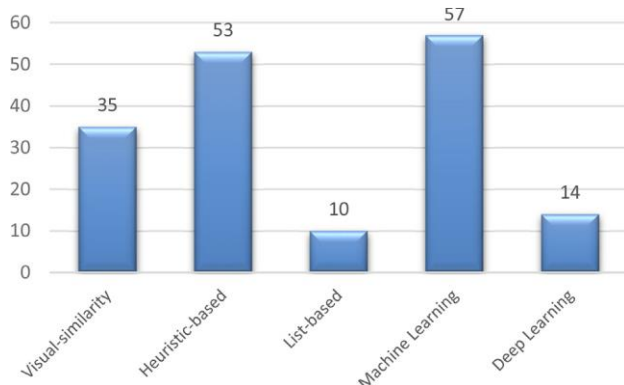
داده‌های حجیم دارای مجموعه داده‌هایی با سرعت، تنوع، حجم، ارزش و صحت بالا است. طبقه‌بندی‌کننده‌های مبتنی بر یادگیری ماشین بیش از ۹۹ درصد از دقت را به دست آوردند که ثابت کردند موثرترین روش‌ها بودند [۱۰].

## ۲-۵ روش یادگیری عمیق

با توجه به پیشرفت‌های اخیر در روش‌های یادگیری عمیق، شبکه‌های عصبی عمیق<sup>۶</sup>، باید عملکرد بهتری نسبت به روش‌های متداول یادگیری ماشین در تشخیص وب‌سایت‌های فیشینگ داشته باشند. برخی از الگوریتم‌های معروف یادگیری عمیق شامل، شبکه عصبی عمیق، شبکه عصبی مکرر<sup>۷</sup>، شبکه عصبی عمیق فید فوروارد<sup>۸</sup>، ماشین بولتزمن محدودشده<sup>۹</sup>، شبکه عصبی پیچشی<sup>۱۰</sup>، شبکه باور عمیق<sup>۹</sup> و رمزگذار خودکار عمیق<sup>۱۱</sup> است [۱۱].

## ۳-۳ - مروری بر کارهای انجام‌شده

شکل (۳) تعداد مقالات مربوط به انواع روش‌های شناسایی سایت‌های فیشینگ را نشان می‌دهد [۴].



شکل (۳): تعداد مقالات مربوط به انواع روش‌های شناسایی سایت‌های فیشینگ

شکل (۴) تعداد مقالات یافت شده برای هر الگوریتم مورد استفاده در شناسایی سایت‌های فیشینگ را نشان می‌دهد [۴].

## ۲-۲ روش شباهت بصری

این رویکرد مضمون را ارزیابی می‌کند و وب‌سایت‌های معتبر را بر اساس ویژگی‌های بصری مختلف احراز هویت می‌کند. زیرا به نظر می‌رسد که صفحه وب فیشینگ بسیار شبیه به صفحه قانونی<sup>۱</sup> خودش است، این ابزارها شباهت‌ها را مقایسه می‌کنند: این روش از برگه‌های سبک آبخاری<sup>۲</sup>، طرح‌بندی متن، کد منبع، آرم وب‌سایت، تصاویر صفحه وب، و سایر عناصر بصری استفاده می‌کند. از آن جایی که این روش‌ها صفحه وب مشکوک را با صفحات وب قبلاً بازدید شده یا ذخیره‌شده مقایسه می‌کنند آن‌ها نمی‌توانند حملات فیشینگ ساعت صفر<sup>۳</sup> را شناسایی کنند [۶].

## ۲-۳ روش اکتشافی

رویکرد اکتشافی از ویژگی‌های مشتق شده از وب‌سایت فیشینگ استفاده می‌کند. این استراتژی بر چندین ویژگی استوار است که می‌تواند یک وب‌سایت فیشینگ را از یک وب‌سایت واقعی متمایز کند. این روش‌ها داده‌ها را از منابع مختلف، مانند URL ها، محتوای متنی، سیستم نام دامنه، گواهی‌های دیجیتال و ترافیک وب‌سایت‌ها جمع‌آوری می‌کند. مجموعه ویژگی، نمونه‌های آموزشی و الگوریتم‌های طبقه‌بندی، همگی بر موفقیت این روش تأثیر می‌گذارد. یکی از مزایای این روش این است که می‌تواند حملات فیشینگ ساعت صفر را شناسایی کند [۷].

## ۲-۴ روش یادگیری ماشین

امروزه یادگیری ماشین روشی رایج برای شناسایی وب‌سایت‌های فیشینگ است [۸]. ویژگی‌های رایج مانند اطلاعات URL، ساختار وب‌سایت، و ویژگی‌های جاوا اسکریپت برای نشان دادن URL های فیشینگ و وب‌سایت‌های مرتبط جمع‌آوری می‌شوند. سپس، بر اساس آن ویژگی‌ها، مجموعه داده‌های فیشینگ به دست می‌آید. پس‌از آن، طبقه‌بندی‌کننده‌های یادگیری ماشین برای شناسایی وب‌سایت فیشینگ، بر اساس آن ویژگی‌ها آموزش داده می‌شوند [۹]. این روش‌ها با داده‌های حجیم بسیار خوب کار می‌کند.

<sup>6</sup> Feed-Forward Deep Neural Network

<sup>7</sup> Limited Boltzmann machine

<sup>8</sup> Convolutional Neural Network

<sup>9</sup> Deep belief network

<sup>10</sup> Deep auto-encoder

<sup>1</sup> Legitimate

<sup>2</sup> Cascading Style Sheets

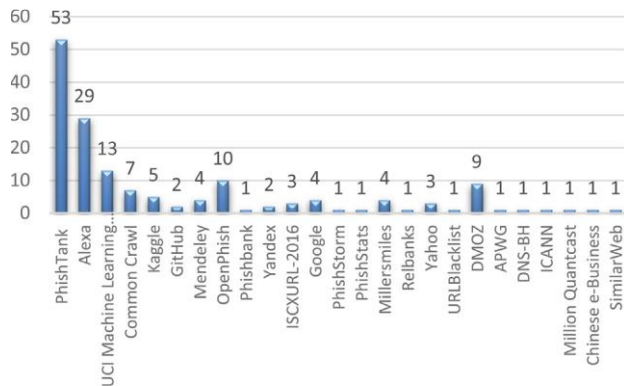
<sup>3</sup> Zero-hour phishing attacks

<sup>4</sup> Deep Neural Network

<sup>5</sup> Recurrent Neural Network

(یا ۱۶٪) در سال ۲۰۱۹، ۲ مقاله (یا ۲٪) در سال ۲۰۱۷ منتشر

شد [۴].



شکل (۶): تعداد مقالات یافت شده برای انواع مجموعه داده

مورد استفاده در شناسایی سایت‌های فیشینگ

#### ۴- روش پیشنهادی

با توجه به اینکه در کارهای قبلی، تشخیص فیشینگ با استفاده از ویژگی‌های مختلف توسط شبکه‌های عمیق بررسی شده است، در کار جاری قصد داریم تشخیص فیشینگ را با استفاده از چندین نوع ویژگی در کنار هم با کمک شبکه‌های عمیق انجام دهیم. ویژگی‌های ترکیبی به کاررفته عبارت‌اند از:

(۱) ویژگی‌های مبتنی بر URL مانند طول بخش‌های URL، آی پی، استفاده از کاراکترهای خاص در آدرس URL و غیره.  
 (۲) ویژگی‌های مبتنی بر محتوا مانند تعداد لینک‌های هر صفحه، نسبت لینک‌های داخلی به خارجی، نسبت لینک‌های خراب و غیره.

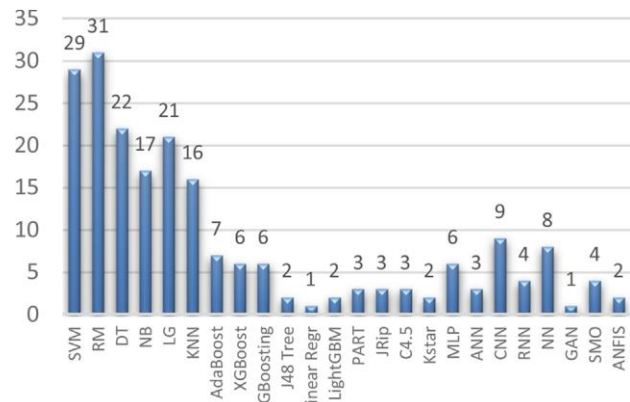
(۳) ویژگی‌های خارجی مانند دامنه ثبت شده در WHOIS، دوره اعتبار دامنه، عمر دامنه و غیره.

اطلاعات مربوط به ویژگی‌های فوق برای تعداد ۱۱۴۸۰ وبسایت، در دیتاست phishing\_data.csv موجود است. دیتاست ما از دیتا ست‌های آماده موجود در اینترنت است.

در ادامه، به توصیف دقیق ویژگی‌ها می‌پردازیم:

#### الف) ویژگی‌های مربوط به URL

۱- طول URL و بخش‌های آن: به منظور مخفی سازی دامنه و زیر دامنه‌های واقعی در سایت‌های فیشینگ از URL های طولانی استفاده می‌شود. به همین خاطر طول کل URL و طول نام

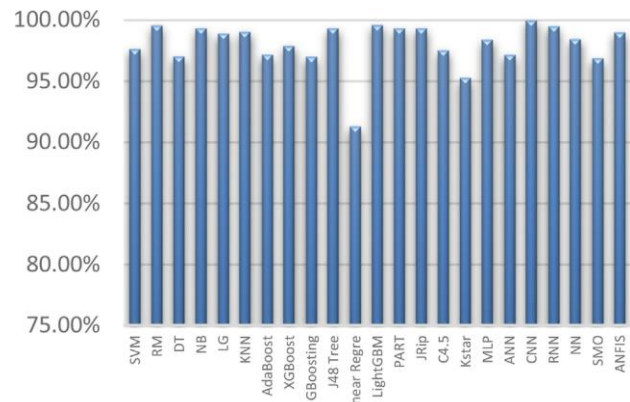


شکل (۴): تعداد مقالات یافت شده برای هر الگوریتم مورد استفاده در

شناسایی سایت‌های فیشینگ

شکل (۵) بیشترین دقت به دست آمده توسط هر الگوریتم را در

شناسایی سایت‌های فیشینگ نشان می‌دهد [۴].



شکل (۵): بیشترین دقت به دست آمده توسط هر الگوریتم در شناسایی

سایت‌های فیشینگ

شکل (۶) تعداد مقالات یافت شده برای انواع مجموعه داده

مورد استفاده در شناسایی سایت‌های فیشینگ را نشان می‌دهد [۴].

پس از جمع‌آوری ۵۳۷ مطالعه، فقط ۸۰ مطالعه، مرتبط با حملات فیشینگ بود. از این ۸۰ نشریه، ۳۰ مقاله که ۳۸ درصد از کل را

تشکیل می‌دهند، در مجله IEEE، ۲۰ (یا ۲۵٪) در اسپرینگر، ۱۷ (یا ۲۱٪) در الزویر، ۸ (یا ۱۰٪) و ۵ (یا ۶٪) در مجلات نمایه شده

Scopus یافت می‌شود. این نشان می‌دهد که IEEE در انتشار این موضوع پیشرو است. هم‌چنین از ۸۰ مقاله منتخب، ۳۴ مقاله در

سال ۲۰۲۰، که معادل ۴۳ درصد، و ۱۶ مقاله (یا ۲۰ درصد) در سال ۲۰۲۱، ۱۵ مقاله (یا ۱۹٪) در سال ۲۰۱۸، ۱۳ مقاله تحقیقاتی

مسیر یا در بخش زیردامنه قرار دارد، URL مشکوک به فیشینگ است و این ویژگی را با متغیر دودویی در نظر گرفتیم (دودویی).

۱۰- زیر دامنه‌های غیرعادی: ممکن است URL فیشینگ، به‌جای www از ترکیب w و اعداد و حالات مختلف آن استفاده کند. پس چنین URL هایی مشکوک به فیشینگ هستند و با متغیر دودویی مشخص می‌شوند.

۱۱- تعداد زیردامنه‌ها: URL های فیشینگ، معمولاً دارای تعداد زیردامنه بیشتری از سایت‌های قانونی هستند و تعداد آن‌ها را به‌عنوان یکی از ویژگی‌های فیشینگ در نظر گرفتیم (عدد صحیح).

۱۲- پیشوند-پسوند: پیشوندها و پسوندها که با علامت - از هم جدا می‌شوند، برای وانمود کردن ظاهری قانونی و بی‌خطر به کاربران در URL استفاده می‌شوند. به همین دلیل، وجود این علامت را (-) را با متغیر دودویی به‌عنوان ویژگی فیشینگ در نظر گرفتیم (دودویی).

۱۳- دامنه‌های تصادفی: URL های فیشینگ، از کلمات مرکب شامل کاراکترهای تصادفی تشکیل شده‌اند و این مسئله را با متغیر دودویی در نظر گرفته‌ایم (دودویی).

۱۴- سرویس‌های کوتاه‌سازی URL: نام‌های کوتاه شده‌ای که به دیگر URL های پیچیده و طولانی هدایت می‌شوند و در فیشینگ برای مخفی سازی نام هاست های واقعی قابل استفاده است. در صورت استفاده از این سرویس‌ها آن را با متغیر دودویی مشخص کرده‌ایم (دودویی).

۱۵- پسوند مسیر: اسکریپت‌های بدخواه را می‌توان به صفحات سالم اضافه کرد. برخی پسوند فایل‌های به‌کاررفته در URL ها به این منظور استفاده می‌شوند که در این کار، پسوند های txt و exe و js را در نظر گرفتیم (دودویی).

۱۶- هدایت مجدد (redirect): از هدایت مجدد یا تغییر مسیر URL، برای بازکردن صفحاتی غیر از صفحه انتخابی کاربر استفاده می‌شود که برای جلوگیری از دستیابی به لینک‌های خراب (پس از حذف شدن صفحه اصلی) از آن استفاده می‌شود. هدایت مجدد می‌تواند به صفحات داخلی سایت یا

هاست را به‌عنوان ویژگی در نظر گرفتیم (از نوع عدد صحیح).

۲- آدرس IP: برای مخفی سازی هویت وبسایت‌ها، گاهی در قسمت نام هاست از آدرس آی پی استفاده می‌شود. به همین دلیل، وجود یا عدم وجود IP در نام دامنه را به‌عنوان یک ویژگی در نظر گرفتیم (۱ است اگر وجود داشته باشد و ۰ است اگر وجود نداشته باشد).

۳- کاراکترهای خاص: برای آن‌که کاربر متوجه تفاوت بین دامنه واقعی و دامنه فیشینگ نشود از کاراکترهای خاص مانند نقطه و دیگر علائم نگارشی و غیره (شارپ، دلار و...) استفاده می‌شود، به همین دلیل تعداد این علائم در نام دامنه را به‌عنوان یک ویژگی در نظر گرفتیم (عدد صحیح).

۴- بخش‌های مشترک: اجزای مشترک مانند http و www و com. فقط یک‌بار در نام دامنه ظاهر می‌شوند، در غیر این صورت مشکوک بوده و تعداد آن در این ویژگی قرار می‌گیرد (عدد صحیح).

۵- توکن HTTPS: در مقایسه با وبسایت‌های قانونی، بیشتر وبسایت‌های فیشینگ، فاقد ابزار امنیتی هستند. به همین دلیل، وجود یا عدم وجود آن را در نظر گرفتیم (۰ و ۱).

۶- نسبت اعداد موجود در دامنه: وجود تعداد زیاد عدد در URL، مشکوک بوده و به‌عنوان نشانه فیشینگ در نظر گرفته شده است. نسبت اعداد در کل URL و در قسمت نام هاست را در نظر گرفتیم (اعشاری).

۷- Punycode: پانی کد، به کدهای یونیکد جایگزین شده با کدهای اسکی انگلیسی گفته می‌شود که دلیلی ندارد در سایت‌های قانونی از آن استفاده شود. به همین دلیل با یک متغیر دودویی، URL های دارای پانی کد را مشخص کرده‌ایم (دودویی).

۸- شماره پورت: بندرت از شماره پورت در URL های غیربدخواه استفاده می‌شود. به همین دلیل یک متغیر دودویی برای مشخص کردن این URL ها در نظر گرفتیم (دودویی).

۹- محل TLD: در URL های خوش‌فرم، TLD یا دامنه سطح بالا (مثل com)، قبل از مسیر ظاهر می‌شود. وقتی TLD در

- فیشینگ برعکس. از این رو، نسبت لینک‌های داخلی به خارجی را به‌عنوان یک ویژگی در نظر گرفتیم (عدد اعشاری).
- ۳- نسبت ابرلینک‌های تھی: برای تقلید از سایت‌های قانونی، لینک صفحات غیرفیشینگ به صورت تھی در سایت‌های فیشینگ قرار داده می‌شود (یعنی عملاً به‌جایی اشاره نمی‌کنند). نسبت لینک‌های تھی به سالم به‌عنوان ویژگی محتوایی دیگر در نظر گرفته شده است (اعشاری).
- ۴- تعداد فایل CSS خارجی: در وب‌سایت‌های قانونی، از یک سبک داخلی یا بیش از یک فایل CSS خالی استفاده می‌شود. اما در وب‌سایت‌های فیشینگ، فقط از یک فایل CSS خارجی استفاده شده که حاوی لینک‌هایی به فایل‌های CSS سایت هدف است. از این رو تعداد فایل‌های CSS خارجی را به‌عنوان ویژگی دیگر در نظر گرفتیم (عدد صحیح).
- ۵- تعداد تغییر مسیر (redirect) های داخلی/خارجی: لینک‌های صفحات فیشینگ، ممکن است به دیگر صفحات جعلی یا سالم، اشاره کنند. نسبت تغییر مسیرهای داخلی و خارجی برای تشخیص صفحات وب فیشینگ در نظر گرفته شده است (عدد صحیح).
- ۶- نسبت خطاهای داخلی/خارجی: معمولاً در صفحات فیشینگ، از لینک‌های جعلی استفاده می‌شود. بنابراین، همه لینک‌های صفحات وب، بازرسی شد و نسبت خطاهای اتصال داخلی و خارجی به‌عنوان ویژگی در نظر گرفته شده است (اعشاری).
- ۷- فرم‌های ورود (لاگین): فرم‌های لاگین، ابزار پرکاربرد دیگری برای ربودن اطلاعات کاربران وب‌سایت‌ها محسوب می‌شوند. وجود فرم‌های لاگین با لینک اکشن خارجی یا اکشن تھی به‌عنوان دیگر ویژگی فیشینگ در نظر گرفته شده است که با ۰ یا ۱ مشخص می‌شود (دودویی).
- ۸- فاو آیکن: برای تقلید از وب‌سایت‌های قانونی، در سایت‌های فیشینگ، از فاو آیکن‌های مشابه با سایت هدف استفاده می‌شود. به همین جهت، استفاده از فاو آیکن‌های خارجی به‌عنوان ویژگی در نظر گرفته شده است (دودویی).
- ۹- لینک در تگ: در وب‌سایت‌های قانونی، انتظار می‌رود تگ‌های <links>، از لینک‌هایی استفاده کنند که به صفحات

- کلاً مسیرهای خارجی استفاده شود که تعداد هر یک از این موارد را در یک ویژگی مشخص کرده‌ایم (عدد صحیح).
- ۱۷- ویژگی‌های NLP: برای تشخیص فیشینگ از ویژگی‌های NLP (پردازش زبان طبیعی) استفاده می‌شود که در این کار تعدادی از موارد زیر را به صورت ویژگی‌های جداگانه در نظر گرفتیم: تعداد کلمات، تکرار کاراکتر، کلمات کوتاه شده در URL و نام هاست و مسیر، طولانی‌ترین کلمه در URL و نام هاست و مسیر، میانگین طول کلمات در URL ها، نام هاست ها و مسیرها (عدد صحیح).
- ۱۸- Phish hints (سرنخ‌های فیشینگ): کلمات حساسی هستند که برای جلب اعتماد به صفحات بازدید شده به کار می‌روند. تعداد چنین کلماتی در URL، به‌عنوان یک ویژگی در نظر گرفته شد (عدد صحیح).
- ۱۹- دامنه‌های برند: URL های فیشینگ از نام دامنه‌های برند در بخش‌های مختلف استفاده می‌کنند. استفاده از نام برند در بخش دامنه، مجاز و استفاده از آن در زیردامنه یا مسیرها را غیرمجاز در نظر گرفتیم و این موارد را با متغیرهای دودویی نشان دادیم (دودویی).
- ۲۰- TLD های مشکوک: TLD ها بررسی شده‌اند و فهرستی از TLD های مشکوک در تحقیق جاری استفاده شده است که این ویژگی با متغیر دودویی نشان داده شده است (دودویی).
- ۲۱- گزارش آماری: دامنه‌های URL بررسی شدند و IP آن‌ها با دامنه‌های مهم فیشینگ، تطبیق داده شده است (دودویی).

#### ب) ویژگی‌های محتوایی

- ۱- تعداد ابرلینک‌ها: وب‌سایت‌های غیرفیشینگ، معمولاً تعداد صفحات بیشتری نسبت به سایت‌های فیشینگ دارند از این رو، تعداد لینک‌های موجود در محتوای صفحه را به‌عنوان یک ویژگی در نظر گرفتیم (عدد صحیح).
- ۲- نسبت لینک‌های داخلی/خارجی: وب‌سایت‌های غیرفیشینگ، بیشتر دارای لینک‌های داخلی هستند و وب‌سایت‌های



مشاهده سورس کد، از این قابلیت استفاده می‌کنند و آن را به‌عنوان یک ویژگی در نظر گرفتیم (دودویی).

۱۷- عنوان تھی: بیشتر وب‌سایت‌های قانونی، عنوان صفحه خود را در تگ <title> قرار می‌دهند. عدم وجود عنوان را در یک ویژگی در نظر گرفتیم (دودویی).

۱۸- وجود نام دامنه در عنوان: وب‌سایت‌های قانونی، از نام دامنه به‌عنوان بخشی از عنوان صفحه استفاده می‌کنند. وب‌سایت‌های فیشینگ، از دامنه‌های قانونی، در عنوان خود برای فریب دادن کاربران استفاده می‌کنند. از این رو، وجود دامنه URL به‌عنوان بخشی از عنوان صفحه را به‌عنوان ویژگی سایت‌های قانونی در نظر گرفتیم (دودویی).

۱۹- وجود دامنه در کپی‌رایت: وب‌سایت‌های قانونی، نام دامنه خود را در لوگوی کپی‌رایتشان نمایش می‌دهند که این امر در مورد وب‌سایت‌های فیشینگ درست نیست. این مسئله را به‌عنوان یک ویژگی در نظر گرفتیم (دودویی).

ج) ویژگی‌های به‌دست‌آمده از منابع خارجی:

۱- ثبت در whois: برخلاف سایت‌های قانونی، دامنه وب‌سایت‌های فیشینگ، معمولاً در WHOIS ثبت نشده است. این مسئله را به‌عنوان یک ویژگی در نظر گرفتیم (دودویی).

۲- طول عمر ثبت دامنه: سایت‌های فیشینگ معمولاً، عمر کوتاهی دارند (برخلاف سایت‌های قانونی). به همین دلیل تعداد سال تمدید دامنه را به‌عنوان یک ویژگی در نظر گرفتیم (عدد صحیح).

۳- عمر دامنه: عمر دامنه (فاصله زمانی تا انقضا)، نیز به‌عنوان ویژگی در نظر گرفته شد (زیرا سایت‌های فیشینگ، عمر کوتاهی دارند) (عدد صحیح).

۴- ترافیک وب: وب‌سایت‌های فیشینگ، معمولاً بازدیدکننده کمتری نسبت به سایت‌های قانونی دارند و به همین دلیل این مسئله را با استفاده از سایت الکسا، به‌عنوان یک ویژگی در نظر گرفتیم (عدد صحیح).

۵- رکورد DNS: سیستم نام دامنه، برای بازیابی آدرس IP و دستیابی به URL، ضروری است. از این رو، دامنه‌های URL

وب با دامنه مشابه با URL اشاره کنند. از این رو، نسبت لینک‌های داخلی در تگ‌های <link> به لینک‌های خارجی نیز در نظر گرفته شده است (اعشاری).

۱۰- ارسال به ایمیل: وب‌سایت‌های فیشینگ، ورودی‌های کاربران در وب فرم‌ها را به ایمیل‌های خاص ارسال می‌کنند. فرم‌های حاوی mailto: یا mail() به‌عنوان ویژگی در نظر گرفته شدند (دودویی).

۱۱- نسبت رسانه‌های داخلی/خارجی: وب‌سایت‌های قانونی، بیشتر از رسانه‌هایی که از عکس، صوت و فیلم استفاده می‌کنند در دامنه مشابه با URL ذخیره شده‌اند. به همین دلیل نسبت لینک‌های فایل رسانه‌ای داخلی و خارجی را نیز به‌عنوان ویژگی در نظر گرفتیم (عدد صحیح).

۱۲- SFH: معمولاً، پس از دریافت اطلاعات از کاربر، باید اقدامی (اکشن) اجرا شود. از همین رو فرم‌های با رشته تھی یا about:blank را به‌عنوان بخشی از سایت فیشینگ در نظر گرفتیم (دودویی).

۱۳- iFrame های غیرقابل مشاهده: از تگ‌های iFrame برای گنجانیدن دیگر صفحات وب در یک صفحه، استفاده می‌شود. وب‌سایت‌های فیشینگ معمولاً از این تگ‌ها با مرزهای (border) نامرئی استفاده می‌کنند تا کاربر فکر کند که در بخشی از صفحه واقعی است. این مسئله را نیز به‌عنوان ویژگی در نظر گرفتیم (دودویی).

۱۴- پنجره‌های باز شو (pop up): در سایت‌های قانونی، برای هشدار دادن به کاربران استفاده می‌شود و به‌ندرت برای جمع‌آوری و ارسال داده استفاده می‌شوند. وجود چنین پنجره‌هایی با فیلدهای متنی به‌عنوان نشانه فیشینگ استفاده شده است (دودویی).

۱۵- لنگر امن: از تگ <a>، برای لینک یک صفحه به صفحه دیگر استفاده می‌شود. تگ‌هایی با لینک‌های '#', '&', 'javascript', 'mailto' را به‌عنوان ناامن در نظر گرفته و تعداد آن‌ها را در یک ویژگی قراردادیم (عدد صحیح).

۱۶- کلیک راست: اسکریپت‌هایی برای غیرفعال سازی راست کلیک وجود دارد که سایت‌های فیشینگ برای جلوگیری از

مورد استفاده محققین در پنج سال گذشته در تشخیص وبسایت فیشینگ را گزارش می‌کند. پنج روش تشخیص فیشینگ به طور عمده توسط جامعه تحقیقاتی، مورد استفاده قرار گرفت که در میان آن‌ها، یادگیری ماشین بیشترین استفاده از رویکردهای مورد مطالعه بوده است. از ۸۰ مورد تحقیقاتی، ۵۷ مقاله، یا ۷۱٫۲۵٪ از مطالعات، از رویکردهای یادگیری ماشین در کار خود استفاده کردند. همچنین نظرسنجی‌ها نشان داد که عمدتاً از دو منبع برای تحلیل فیشینگ مجموعه داده‌ها، استفاده شد. ۵۳ یا ۶۶٫۲۵٪ از مطالعات، از وبسایت PhishTank استفاده کردند، در حالی که، برای مجموعه داده‌های قانونی، ۲۹ یا ۳۶٫۲۵٪ از مطالعات، از وبسایت الکسا استفاده کردند. بیشتر نویسندگان از طبقه‌بندی جنگل تصادفی استفاده کردند، که ۳۸٫۷۵٪ از ۸۰ مقاله است. اگرچه الگوریتم جنگل تصادفی در میان الگوریتم‌های سنتی یادگیری ماشین، بیشتر مورد استفاده قرار می‌گیرد، اما دقت الگوریتم تکاملی شبکه عصبی پیچشی، ۹۹٫۹۸٪ است و در میان تمام مطالعاتی که در این نظرسنجی گنجانده شده است بهترین است.

## References

- [1] Hannousse, A., & Yahiouche, S. (2021). Towards benchmark datasets for machine learning based website phishing detection: An experimental study. *Engineering Applications of Artificial Intelligence*, 104, 104347.
- [2] Anti-Phishing Working Group. (2006). Phishing Activity Trends Report-May, 2006. [http://www.anti-phishing.org/reports/apwg\\_report\\_May2006.pdf](http://www.anti-phishing.org/reports/apwg_report_May2006.pdf).
- [3] Wandera (2020). Mobile Threat Landscape 2020: Understanding the key trends in mobile enterprise security in 2020. Technical Report. <https://www.wandera.com/mobile-threat-landscape/>.
- [4] Safi, A., & Singh, S. (2023). A Systematic Literature Review on Phishing Website Detection Techniques. *Journal of King Saud University-Computer and Information Sciences*.
- [5] Yang, L., Zhang, J., Wang, X., Li, Z., Li, Z., & He, Y. (2021). An improved ELM-based and data preprocessing integrated approach for phishing detection considering comprehensive features. *Expert Systems with Applications*, 165, 113863.

باید در DNS ثبت شوند. عدم وجود رکورد DNS، به عنوان نشانه فیشینگ در نظر گرفته شده است (دودویی).

۶- نمایه گوگل: وبسایت‌های فیشینگ، عمر کوتاهی داشته و غالباً از طریق لینک‌های مستقیم در ایمیل به قربانی ارسال می‌شوند و نیازی به نمایه گوگل ندارند. از این رو، نمایه بودن یا نبودن در گوگل را به عنوان ویژگی در نظر گرفتیم (دودویی).

۷- Page rank: وبسایت‌های فیشینگ، محبوبیت چندانی ندارند، از این رو، فرض می‌شود که نسبت به سایت‌های قانونی، رتبه صفحه پائینی داشته باشند. بنابراین از Openpagerank برای به دست آوردن این ویژگی استفاده کردیم (عدد صحیح).

پیاده‌سازی کار به این صورت است که پس از آماده‌سازی داده‌ها (تمیز سازی داده‌ها، جداسازی برچسب‌ها از داده‌ها و نرمال‌سازی در صورت لزوم)، مدل‌های مختلف شبکه عمیق شامل، شبکه Dense، شبکه LSTM و شبکه RNN ساخته می‌شود. بر اساس آزمایش و خطا، تعداد بهینه لایه‌ها و نرون‌ها را برای هر مدل تعیین می‌کنیم. سپس با استفاده از نمودارهایی، میزان loss، و تعداد اپک بهینه را برای هر مدل تعیین می‌کنیم و مدل را آموزش می‌دهیم. سپس مدل را روی داده‌های تست (۲۰ درصد از کل داده‌ها) ارزیابی می‌کنیم. در صورت عدم کسب دقت رضایت‌بخش، با استفاده از cross-validation بازه‌های مختلف داده‌ها را به عنوان داده آموزش یا تست ارزیابی می‌کنیم. در نهایت، بهترین دقت به دست آمده به عنوان مبنای مقایسه با کارهای قبلی استفاده شده و نتیجه‌گیری مربوطه انجام می‌شود. به دلیل حجم بالای کار پیاده‌سازی روش پیشنهادی، ما قصد داریم در مقاله تحقیقی دیگری فرایند پیاده‌سازی ایده خود را به طور کامل شرح دهیم و در این مقاله به مرور تحقیقات گذشتگان به طور خلاصه و مفید پرداختیم و نیز ایده پیشنهادی خود را به طور خلاصه شرح دادیم.

## ۵- نتیجه‌گیری

کار انجام شده در این مقاله شامل، بررسی آن دسته از مطالعاتی است که عملکرد روش‌های شناسایی وبسایت فیشینگ را تحلیل کردند. این مقاله مجموعه داده‌های استفاده شده و الگوریتم‌های



- [6] Jain, A. K., & Gupta, B. B. (2018). PHISH-SAFE: URL features-based phishing detection system using machine learning. In *Cyber Security: Proceedings of CSI 2015* (pp. 467-474). Springer Singapore.
- [7] Jain, A. K., & Gupta, B. B. (2018). Two-level authentication approach to protect from phishing attacks in real time. *Journal of Ambient Intelligence and Humanized Computing*, 9, 1783-1796.
- [8] Sindhu, S., Patil, S. P., Sreevalsan, A., Rahman, F., & AN, M. S. (2020, October). Phishing detection using random forest, SVM and neural network with backpropagation. In *2020 International Conference on Smart Technologies in Computing, Electrical and Electronics (ICSTCEE)* (pp. 391-394). IEEE.
- [9] Zhu, E., Ju, Y., Chen, Z., Liu, F., Fang, X. (2020). DTOF-ANN: an artificial neural network phishing detection model based on decision tree and optimal features. *Appl. Soft Comput. J.* 95., <https://doi.org/10.1016/j.asoc.2020.106505>
- [10] Alkawaz, M. H., Steven, S. J., Hajamydeen, A. I., & Ramli, R. (2021, April). A comprehensive survey on identification and analysis of phishing website based on machine learning methods. In *2021 IEEE 11th IEEE Symposium on Computer Applications & Industrial Electronics (ISCAIE)* (pp. 82-87). IEEE.
- [11] Basit, A., Zafar, M., Javed, A. R., & Jalil, Z. (2020, November). A novel ensemble machine learning method to detect phishing attack. In *2020 IEEE 23rd International Multitopic Conference (INMIC)* (pp. 1-5). IEEE

## Copy – Move Detection using genetic algorithm

Farah Maki Alhijel<sup>1</sup>, Zahra Yaghoubi<sup>1\*</sup>

<sup>1</sup>Engineering Department, Imam Khomeini International University, Qazvin, Iran

### Article Information

#### Original Research Paper

#### Received:

2023 February 13

#### Accepted:

2023 June 20

#### Keywords:

Phishing website detection,  
Machine learning, Deep  
learning

#### Corresponding Author\*:

z.yaghoubi@eng.ikiu.ac.ir

### Abstract

The increasing popularity of the Internet has led to the dramatic growth of e-commerce. However, such activities have significant security challenges, mainly due to cyber fraud and identity theft. Therefore, checking the legitimacy of visited web pages is a very important task to secure the identity of customers and prevent phishing attacks. The use of machine learning methods, and deep learning is widely recognized as a promising solution. Research is full of studies that use machine learning and deep learning methods to detect website phishing. However, their findings depend on the data set and are far from generalizable. The two main reasons for the lack of generalization are impractical replication and lack of appropriate benchmark data sets for fair evaluation of systems. Furthermore, phishing methods are constantly evolving and the proposed models do not keep up with the rapid changes. In this article, we review the methods of identifying phishing sites from legal sites and finally reach the final conclusion.



:10.22034/ABMIR.2023.19703.1020:

E-ISSN: [2821-2037](https://doi.org/10.22034/ABMIR.2023.19703.1020) /© 2023. Published by Yazd University This is an open access article under the CC BY 4.0 License (<https://creativecommons.org/licenses/by/4.0/>).

