

به نام خدا

عنوان پروژه: شناسایی حملات فیشینگ^۱ با استفاده از الگوریتم های یادگیری ماشین^۲

تاریخ: ۲۱-آبان-۱۴۰۰

نوع مسئله: یادگیری بانظارت (طبقه بندی)

کارهای گذشته:

(مهمت کارکمز و همکاران، ۲۰۲۱)، دسته بندی لینک های فیشینگ^۳ با استفاده از شبکه های عصبی عمیق رو دیتاست لینک های پر خطر

با توجه به روند رو به رشد اینترنت، تعداد رایانه های متصل به اینترنت روز به روز در حال افزایش است. تقریباً همه شرکت ها در حال انتقال اعمال اصلی خود از دنیای واقعی به دنیای سایبری هستند، اگر چه آن موضوع باعث افزایش فروش آنها میشود میتواند آسیب پذیری های زیادی مانند حملات سایبری برای این شرکت ها را در پی دارد به ویژه با ساختار ناشناس اینترنت. فیشینگ یکی از انواع محبوب حملات سایبری است که از ناآگاهی کاربر برای آسیب رساندن به آن ها سو استفاده می شود. در برخی از پژوهش ها از سیستم های تشخیص مبتنی بر قانون بعنوان یک مکانیسم پیشگیری استاتیک و سیستم های مبتنی بر یادگیری ماشین به عنوان یک مکانیسم پیشگیری پویا کمک میگیرد. در مقاله نام برده شده نویسندگان یک سیستم تشخیص فیشینگ مبتنی بر شبکه های عصبی عمیق^۴ را با تجزیه و تحلیل نشانی اینترنتی وب سایت های مشکوک پیاده سازی کرده اند. اگر چه در تمام تحقیقاتی قبلی که توسط محققان دیگر انجام شده مجموعه داده های مورد استفاده از منابع مختلفی جمع آوری شده بودند و که وب سایت های مجاز و فیشینگ در آنها مشخص است. اما نویسندگان مقاله نام برده ابتدا یک مجموعه داده پرخطر ایجاد می کنند که فقط شامل وب سایت

^۱phishing

^۲Machine learning

^۳Phishing

^۴Deep Neural Network (DNN)

های مشکوک است که به سایت فیشتانک^ه گزارش شده اند. تحقیقاتی تجربی نویسندگان مقاله نام برده نشان می دهد که سیستم پیشنهادی توسط نویسندگان مقاله نام برده کارایی بسیار خوبی از نظر دقت و زمان اجرا ارائه می دهد.

(او.اسماعیل و همکاران، ۲۰۱۰)، کاهش خطر حملات فیشینگ با استفاده از ابزار های هوش مصنوعی

در مقاله نام برده شده یک پژوهش جامع به همراه تجزیه و تحلیل در مورد فیشینگ، ویشینگ^و اسمیشینگ^ی برای بهره برداری از دانش در پیاده سازی یک ابزار هوشمند برای شناسایی و محافظت در مقابل موارد ذکر شده ارائه شده است. این یک مشکل جدید مهندسی اجتماعی است که زندگی روزمره ما را آسیب پذیر و دشوار می کند. تحقیق نام برده شده به ویژه بر روی فیشینگ از طریق ایمیل متمرکز است زیرا، به طور مستقیم پیامدهای جدی تری در مقایسه با سایر روش ها در رابطه با معاملات مالی دارد. شایان ذکر است که امن نمودن حجم بسیار زیاد معاملات آنلاین بسیار چالش برانگیز است زیرا روزانه چندین روش برای نقض حریم خصوصی افراد به منظور سرقت اعتبار آنها ابداع می شود. هزینه این نوع حملات بیش از میلیون ها دلار به صورت سالانه است. ابزارهای زیادی برای حل این مشکل ارائه شده است. اما متأسفانه، این معضل هنوز پا برجاست. در مقاله نام برده شده توسط نویسندگان یک روش برای توسعه یک ابزار هوشمند برای مقابله با این خطر پیشنهاد شده است.

(یزدان احمد السیرا و همکاران، ۲۰۲۰)، الگوریتم های فرایادگیری و درختان تصمیم اضافی

برای شناسایی صفحات فیشینگ

فیشینگ نوعی حمله مهندسی-وب اجتماعی^ا در فضای مجازی است که در آن مجرمان اطلاعات ارزشمند کاربران ناآگاه را به سرقت می برند. اقدامات متقابل موجود در قالب نرم افزار ضد فیشینگ

^هPhishTank

^وVishing

^یSmishing

^اsocial web-engineering

و روش های محاسباتی برای شناسایی فعالیت های فیشینگ اثبات شده است. با این حال، هکرها روش های جدیدی را برای خنثی سازی این اقدامات متقابل به کار می گیرند. با توجه به ماهیت تکامل یافته حملات فیشینگ، نیاز به اقدامات متقابل جدید و کارآمد بسیار مهم احساس می شود زیرا تأثیر حملات فیشینگ اغلب فاجعه بار است. طرح های هوش مصنوعی سنگ بنای روش های مقابل مدرنی است که برای کاهش حملات فیشینگ استفاده می شود. اقدامات متقابل یا روش های مقابله ای فیشینگ مبتنی بر هوش مصنوعی دارای نواقصی به ویژه بالا بودن میزان هشدار کاذب و عدم توانایی تفسیر نحوه عملکرد بیشتر روش های فیشینگ است. در مقاله نام برده شده نویسندگان چهار مدل فرا آموزش: آدا بوست اکسترا تری، بگینگ اکسترا تری، روتیشن فارست اکسترا تری، لجیت بوست اکسترا تری،^۳ بر مبنای الگوریتم طبقه بندی اکسترا تری توسعه داده اند. مدل ها توسط نویسندگان این مقاله روی مجموعه داده های وب سایت های فیشینگ آموزش داده شده و عملکرد آنها مورد ارزیابی قرار گرفته است. دقت تشخیص مدل ها بالای ۹۷ درصد با نرخ مثبت کاذب ناچیز ۰،۰۲۸ بوده است. علاو بر این مدل های پیشنهاد شده توسط نویسندگان مقاله نام برده شده در تشخیص حملات فیشینگار مدل های مبتنی برای یادگیری ماشینی برتری دارند. از این رو نویسندگان مقاله نام برده شده مدل های فرا آموزش را برای تشخیص حملات فیشینگ اتخاذ نموده اند.

(هانگ لی . همکاران، ۲۰۱۷)، یوارال نت^۴: یادگیری لینک های نماینده با یادگیری عمیق برای شناسایی لینک های مخرب

^۱meta-learner

^۲AdaBoost-Extra Tree (ABET)

^۳Bagging - Extra tree (BET)

^۴Rotation Forest - Extra Tree (RoFBET)

^۵LogitBoost-Extra Tree (LBET)

^۶URLNET

نشانی های مخرب محتوای ناخواسته را میزبانی می کنند و برای ارتکاب جرایم اینترنتی استفاده می شوند. تشخیص به موقع آنها ضروری است. به طور سنتی، این کار با استفاده از لیست های سیاه انجام می شود، که نمی تواند جامع باشد و نمی تواند نشانی های مخرب تازه ایجاد شده را شناسایی کند. برای حل این مسئله، در سال های اخیر شاهد تلاش های زیادی برای شناسایی لینک های مخرب با استفاده از یادگیری ماشین بوده ایم. محبوب ترین و مقیاس پذیرترین روش ها با استخراج ویژگی های بگ اف وردی^۵ از خصوصیات رشته واژه نشانی ها استفاده می کنند و به دنبال آن مدل های یادگیری ماشین مانند ماشین بردار پشتیبان^۶ استفاده می شود. همچنین ویژگی های دیگری نیز وجود دارد که توسط متخصصان برای بهبود عملکرد پیش بینی مدل طراحی شده است. این رویکردها از چندین محدودیت رنج می برند: شماره یک- عدم توانایی در دستیابی موثر معنایی و الگوهای پی در پی در رشته نشانیوبسایت ها -شماره دو- نیاز به انجام مهندسی ویژگی های به صورت دستی و شماره سه- عدم توانایی در مدیریت ویژگی های غیبی و تعمیم دادن آن ها برای آزمایش داده ها. برای رفع این چالش ها نویسندگان مقاله نام برده یو آر ال نت را پیشنهاد می دهند، یک چارچوب یادگیری عمیق برای یادگیری یک جاسازی نشانی های غیرخطی که نشانی های مخرب را مستقیماً از نشانی های معمولی تشخیص می دهد. به طور خاص نویسندگان مقاله نام برده شده شبکه های عصبی کانولوشنی^۷ را به هر دو گروه حروف و کلمات رشته نشانی ها اعمال می کنند تا نشانی تعبیه شده را در یک چارچوب مشترک بهینه سازی شده یاد بگیرد. این رویکرد به مدل اجازه می دهد تا انواع مختلفی از اطلاعات معنایی را بدست آورد، که توسط مدل های موجود امکان پذیر نبود. نویسندگان مقاله نام برده شده همچنین تعبیه پیشرفته کلمات را پیشنهاد می دهیم تا مشکل بسیاری از کلمات نادر مشاهده شده در این کار را برطرف کنیم. نویسندگان مقاله نام برده آزمایش های گسترده ای را روی یک مجموعه داده در مقیاس بزرگ انجام می دهند

^۵Bag-of-words

^۶SVM

^۷Convolutional

که این آزمایشات عملکرد قابل توجهی نسبت به روش های موجود نشان می ده. نویسندگان مقاله نام برده شده برای ارزیابی عملکرد اجزای مختلف یوارال نت مطالعات فرسایش را انجام می دهند. (دایساکی میامتو و همکاران، ۲۰۰۹)، ارزیابی روش های مبتنی بر یادگیری ماشینی برای شناسایی سایت های فیشینگ

در مقاله نام برده شده، نویسندگان این مقاله عملکرد روش های مبتنی بر یادگیری ماشینی را برای تشخیص سایت های فیشینگ ارائه می دهند. نویسندگان این مقاله از ۹ تکنیک یادگیری ماشینی استفاده می کنند که عبارتند از: آدابوست،^{۱۸} بگینگ،^{۱۹} ماشین های بردار پشتیبان،^{۲۰} درخت طبقه بندی و رگرسیون،^{۲۱} لجستیک رگرسیون،^{۲۲} جنگل های تصادفی،^{۲۳} شبکه های عصبی،^{۲۴} نایو بایز^{۲۵} و درخت های رگرسیون افزایشی بیزی.^{۲۶} نویسندگان این مقاله اجازه می دهند که این تکنیک های یادگیری ماشینی، شیوه های اکتشافی را با هم ترکیب کنند، و همچنین اجازه می دهند که روش های تشخیص مبتنی بر یادگیری ماشینی، سایت های فیشینگ را از سایرین متمایز کنند. نویسندگان این مقاله مجموعه داده های خود را که متشکل از ۱۵۰۰ سایت فیشینگ و ۱۵۰۰ سایت قانونی است تجزیه و تحلیل می کنند، نویسندگان مقاله نام برده شده داده های خود را با استفاده از روش های تشخیص مبتنی بر یادگیری ماشینی طبقه بندی می کنند، و عملکرد هر دسته بند را اندازه گیری می کنند. نویسندگان این مقاله از معیار اندازه گیری اف وان،^{۲۷} نرخ خطا، و مساحت تحت منحنی ارو سی به^{۲۸} عنوان معیارهای عملکرد هر الگوریتم دست بند استفاده کردند. بالاترین مقدار اف وان برابر با ۰,۸۵۸۱، کم ترین میزان خطا برابر با ۱۴,۱۵ درصد و بالاترین میزان ای یو سی برابر با ۰,۹۳۴۲

^۱AdaBoost

^۲Bagging

^۳Support Vector Machines

^۴Classification and Regression Trees

^۵Logistic Regression

^۶Random Forest

^۷Neural Network

^۸Naive Bayes

^۹Bayesian Additive Regression Trees

^{۱۰}F1 measure

^{۱۱}AUC

می باشد که همگی در مورد آدابوست مشاهده شده اند. همچنین نویسندگان این مقاله مشاهده کرده اند که ۷ روش از ۹ روش تشخیص مبتنی بر یادگیری ماشین عملکرد بهتری نسبت به روش تشخیص سنتی دارند.

(یوسف عراقی و همکاران، ۲۰۱۳)، شناسایی فیشینگ: بررسی ادبیات^{۲۹}

درمقاله نام برده نویسندگان به بررسی مقالات گوناگون در مورد روش های تشخیص حملات فیشینگ می پردازند. حملات فیشینگ، آسیب پذیری های موجود در سیستم ها را به واسطه عامل انسانی مورد هدف قرار می دهند. بسیاری از حملات سایبری از طریق مکانیزم هایی گسترش می یابند که از نقاط ضعف موجود در کاربران بهره برداری می کنند، که این مسئله کاربران را به ضعیف ترین عنصر در زنجیره امنیتی تبدیل می کند. مشکل فیشینگ گسترده است و هیچ راه حلی برای کاهش موثر تمام آسیب پذیری های ناشی از آن وجود ندارد، بنابراین تکنیک های متعددی برای کاهش حملات خاص اجرا می شوند. هدف مقاله نام برده شده بررسی بسیاری از تکنیک های کاهش حملات فیشینگ است که اخیرا پیشنهاد شده اند. در مقاله نام برده دسته های مختلفی از تکنیک های سطح بالا که جهت کاهش تعداد حملات فیشینگ وجود دارد ارائه شده است، مانند: تشخیص، دفاع تهاجمی، اصلاح، و پیش گیری، که نویسندگان این مقاله معتقداند شناسایی حملات فیشینگ حیاتی ترین روش برای کاهش تعداد حملات فیشینگ می باشد.

(جیان کاو و همکاران، ۲۰۱۶)، شناسایی لینک های مخرب بر اساس باز ارسال در شبکه های اجتماعی آنلاین

در سال های اخیر شبکه های اجتماعی آنلاین^{۳۰} مانند فیسبوک، توئیتر و... در میان کاربران اینترنت بسیار محبوب شده اند. متأسفانه مهاجمان از آنها برای مخفی کردن نشان های مخرب نیز استفاده میکنند. با توجه به اهمیت شناسایی نشانی مخرب در شبکه های اجتماعی آنلاین، چندین راه توسط

^{۲۹}A Literature Survey

^{۳۰}OSN

اپراتور های شبکه های اجتماعی آنلاین، شرکت های امنیتی و محققان دانشگاهی ارائه شده است، بیشتر این راه حل ها از روش های یادگیری ماشین برای آموزش مدل های طبقه بندی براساس انواع مختلف مجموعه ویژگی ها استفاده می کنند. با این حال بیشتر این مدل های آموزش دیده بی اثر هستند زیرا ویژگی های انتخاب شده آنها معمولی است. در مقاله نام برده شده نویسندگان به ویژگی های مبتنی بر باز ارسال تمرکز کرده اند زیرا ارتباط ویژه ای بین رفتار ارسال و انتشار نشانی های مخرب وجود دارد. در ابتدا نویسندگان مقاله یک تجزیه و تحلیل جامع از مجموعه ویژگی های نشانی های معمولی انجام داده اند. سپس برخی از ویژگی های مبتنی بر ارسال و باز ارسال را طراحی کرده اند و چندین ویژگی مبتنی بر نمودار را انتخاب کرده اند تا با ویژگی های قبلی ترکیب شوند تا یک مدل شناسایی لینک های مخرب را با آن ها آموزش دهند نویسندگان مقاله نام برده سیستم مورد نظر خود را با استفاده از حدود ۱۰۰۰۰۰ پیام اصلی از سیناویو^۱ جمع آوری شده ارزیابی میکنند. میزان دقت بالا و نرخ مثبت کاذب پایین نشان می دهد که ویژگی های مبتنی بر ارسال و باز ارسال در شناسایی نشانی های مخرب در شبکه های اجتماعی آنلاین بسیار موثرتر از سایر ویژگی های مرسوم هستند. از نظر دانش نویسندگان مقاله نام برده شده این مقاله اولین مقاله ای می باشد که ویژگی های مبتنی بر ارسال و باز ارسال را در شبکه های اجتماعی آنلاین تجزیه و تحلیل می کند و سهم ارزشمندی در تحقیقات این زمینه ارائه می دهد.

(انی ویزالی و همکاران، ۲۰۱۸)، پی ای دی-ام ال: شناسایی ایمیل های فیشینگ با استفاده از تکنیک های یادگیری ماشینی کلاسیک

در دوران مدرن، کلیه خدمات به صورت آنلاین صورت می پذیرد و همه از آن برای سرعت بخشیدن به فعالیت های روزمره خود استفاده می کنند. این شامل فعالیت های اجتماعی و همچنین مالی است که شامل استفاده از اطلاعات حساس برای انجام وظیفه مورد نظر است. با افزایش استفاده از چنین امکاناتی اهمیت ایمن سازی داده های مورد استفاده برای این اقدامات مطرح می شود. طی

^۱Sina Weibo

^۲PED-ML

دهه گذشته فیشینگ با سرقت اطلاعات حساس برای دستیابی به این امکانات، به تهدیدی جدی برای جامعه تبدیل شده است. این موضوع جزو سودآورترین جرایم اینترنتی محسوب می شود و طبق آمار محققان ای ام بی ایکس فورس^{۳۳} انجام داده اند، تعداد افرادی که قربانی چنین فعالیت هایی می شوند به طرز چشمگیری در حال افزایش می باشد. از آنجا که خطر ایمیل های فیشینگ به طور مداوم در حال افزایش است، نیاز به شناسایی و غلبه بر چنین شرایطی یکی از مهمترین وظایف موجود است. در مقاله نام برده شده، نویسندگان از مدل غیر متوالی مانند رویکرد ماتریس اسناد اصطلاحی به دنبال تجزیه ارزش واحد^{۳۴} و ضریب ماتریس غیر منفی^{۳۵} برای ساخت مدل شناسایی ایمیل فیشینگ به عنوان یک مشکل طبقه بندی نظارت شده برای شناسایی ایمیل های فیشینگ از ایمیل های قانونی استفاده خواهند کرد.

(عبدل حاتم سابس و همکاران، ۲۰۱۷)، سیستم هوشمند تشخیص وب سایت های فیشینگ با استفاده از دست بند جنگل تصادفی

فیشینگ به عنوان تقلیدی از وب سایت یک شرکت معتبر با هدف قرار دادن اطلاعات خصوصی کاربران آن شرکت تعریف می شود. به منظور پیش گیری از حملات فیشینگ، راه حل های مختلفی پیشنهاد شده است. با این حال، تنها یک گلوله جادویی نمی تواند این تهدید را به طور کامل از بین ببرد. داده کاوی یک تکنیک امیدوار کننده است که برای تشخیص حملات فیشینگ استفاده می شود. در مقاله نام برده شده، یک سیستم هوشمند برای تشخیص حملات فیشینگ توسط نویسندگان ارائه شده است. نویسندگان این مقاله از تکنیک های داده کاوی مختلف برای تصمیم گیری در مورد دسته های مختلف وب سایت ها استفاده کرده اند: قانونی یا فیشینگ. الگوریتم دسته بندی مختلفی به منظور ساخت سیستم های هوشمند با دقت بالا برای تشخیص وب سایت های فیشینگ مورد استفاده قرار گرفته اند و برای ارزیابی دقت عملکرد تکنیک های داده کاوی از معیار ای یو سی^{۳۶}

^{۳۳}IBMs X-Force

^{۳۴}Singular Value Decomposition (SVD)

^{۳۵}Nonnegative Matrix Factorization (NMF)

^{۳۶}AUC

که نشان دهنده مساخت زیر نمودار اروسی^{۳۷} است و اندازه گیری اف^{۳۸} استفاده می شود. نتایج نشان بدست آمده توسط نویسندگان مقاله می دهند که جنگل تصادفی^{۳۹} با دست یابی به دقت ۹۷,۳۶ بهترین دقت و عملکرد را در بین دسته بند های مختلف برای تشخیص دارد و زمان اجرای جنگل تصادفی بسیار سریع است و می تواند با وب سایت های مختلفی برای تشخیص فیشنگ برخورد کند.

(شاهن ماندال و همکاران، -)، مروری بر شناسایی لینک های فیشنگ با استفاده از الگوریتم های خوشه بندی فیشنگ نوعی حمله مهندسی اجتماعی می باشید. در این حملات مهاجم نقش یک نهاد قانونی را بازی میکند و از طریقی با قربانی ارتباط برقرار میکند و از کاربر میخواهد لینکی را باز کند که از لحاظ ظاهری به گونه طراحی شده که شبیه به یک وب سایت قانونی است مهاجم اطلاعات کاربران را برای سرقت هویت و سرقت حساب ها و غیره ذخیره میکند. در مقاله نام برده شده تمرکز نویسندگان بر روی حملات فیشنگ که مبتنی بر لینک هستند می باشد. راهکار های ارائه شده توسط نویسندگان این مقاله بیشتری بر روی الگوریتم های دسته بندی متمرکز شده است به نسبت الگوریتم های خوشه بندی. هدف نویسندگان آزمایش و مقایسه نتایج هر دو نوع الگوریتم است تا متوجه شوند کدام یک عملکرد بهتری در شناسایی لینک های مخرب دارد. فرض اصلی رویکرد نویسندگان این مقاله یک مدل یادگیری ماشینی ترکیبی است که از دو مرحله: بررسی با لیست سیاه و لیست سفید و شناسایی مبتنی بر ابتکار برای افزایش دقت الگوریتم پیشنهاد شده توسط نویسندگان مقاله.

(نوشین امیری و همکاران، ۱۳۸۹)، استفاده از شبکه های عصبی خودرمنزگار موازی برای تشخیص صفحات جعلی اینترنتی

^{۳۷}(ROC) curves

^{۳۸}F-measure

^{۳۹}Random Forest

^{۴۰}Classification

^{۴۱}Clustering

به کمک صفحات جعلی اینترنت تلاش می شود اطلاعات محرمانه یک کاربر مانند رمز حساب های بانکی و گذرواژه پست الکترونیکی به سرقت برده شود، این صفحات جعلی درواقع مشابه صفحات وب سایت های معتبر موجو مانند درگاه های پرداخت اینترنتی، یاهو و گوگل ساخته می شوند و به گونه ای کاربران به سمت این صفحات کشانده می شوند به این نوع حملات اینترنتی حمله فیشینگ گفته می شود تشخیص برخط صفحات فیشینگ به کمک نرم افزار های هوشمند میتواند از به سرقت رفتن اطلاعات کاربران جلوگیری کند و امنیت را در فضای وب افزایش دهد. در مقاله نام برده شده نویسندگان یک روش مبتنی بر شبکه های عصبی مصنوعی از نوع خود رمزنگار معرفی کرده اند در روش پیشنهاد شده توسط محققین مقاله نام برده شده از دو شبکه عصبی خود رمزنگار موازی استفاده کرده اند که یکی از آنها با صفحات معمولی و دیگری با صفحات جعلی آموزش دیده است، در زمان تشخیص بر اساس بردار های رمز شده به دست آمده از هر دو شبکه موازی و یک لایه شبکه عصبی معمولی مانند سافت مکس و نوع صفحه ورودی را تشخیص می دهند، در کار برد های عملی هرگاه صفحه ورودی جعلی تشخیص داده شود به سرعت از طریق مرور گر به کاربر اخطار داده می شود. یا دسترسی مسدود می شود. نتایج حاصل از آزمایش روش پیشنهادی به کمک مجموعه داده های فیشینگ وب سایت و معیار های صحت متوسط و دقت به خوبی نشان می دهند که شبکه های عصبی خود رمزنگار موازی عملکرد قوی تری نسبت به سایر روش های یادگیری ماشین در تشخیص صفحات جعلی اینترنتی دارد.

(مهدیه بهارلو و همکاران ، ۱۳۹۸)، بهبود روش شناسایی وب سایت فیشینگ با استفاده از داده کاوی روی صفحات وب

فیشینگ یک نوع حمله اینترنتی در سطح وب است که هدف آن سرقت مشخصات فردی کاربران برای دزدی آنلاین است. فیشینگ دارای اثر منفی در از بین بردن اعتماد بین کاربران در کسب و کارهای الکترونیکی است؛ بنابراین در مقاله نام برده شده سعی بر بررسی روشهای تشخیص وب سایت های فیشینگ با استفاده از داده کاوی شده است. شناسایی ویژگی های برجسته صفحات

فیشینگ یکی از پیش شرط‌های مهم در طراحی یک سیستم تشخیصی دقیق است؛ لذا در گام اول، برای شناسایی ویژگی‌های نفوذ فیشینگ یک لیست با ۳۰ ویژگی مطرح در وب‌سایت‌های فیشینگ آماده گردید. سپس برای افزایش کارایی سامانه‌های تشخیص فیشینگ روش جدیدی جهت کاهش ویژگی‌ها در دو مرحله مبتنی بر انتخاب ویژگی و استخراج ویژگی پیشنهاد شده است که موجب می‌شود تعداد ویژگی‌ها به طور قابل توجهی کاهش یابند. پس از آن عملکرد روش‌های درخت تصمیم جی ۴۸، جنگل تصادفی و بیزین ساده بر روی ویژگی‌های کاهش‌یافته مورد بررسی قرار گرفت شده است. نتایج حاصل از پژوهش محققان این مقاله نشان می‌دهد دقت مدل ایجاد شده برای تعیین صفحات فیشینگ با استفاده از کاهش ویژگی دو مرحله‌ای مبتنی بر پوششی و الگوریتم تحلیل مؤلفه اصلی در روش جنگل تصادفی ۹۶۰۵۸٪ می‌باشد که نسبت به سایر روش‌ها نتیجه مطلوبی است

(فاطمه صف آرا و همکاران، ۱۳۹۸)، افزایش دقت شناسایی صفحات جعلی وب با استفاده از

الگوریتم بهینه سازی گفتار و شبکه عصبی مصنوعی

ایجاد صفحات جعلی در محیط وب یا فیشینگ از جمله حملات سایبری است که نیازمند ملاحظات فرماندهی و کنترل می‌باشد. در حملات فیشینگ افراد به سمت صفحات جعلی که توسط فیشر یا سارق ساخته شده هدایت می‌شوند و اطلاعات مهم آنها توسط فیشر به سرقت می‌رود. الگوریتم‌های یادگیری ماشین و داده کاوی، الگوریتم‌های رایج برای طبقه بندی و تشخیص وبسایت‌های جعلی هستند. طبقه بندی وبسایت‌ها بر اساس ویژگی‌هایی که از آن سایت استخراج می‌شود صورت می‌گیرد. بنابراین انتخاب ویژگی تأثیر زیادی در نتایج طبقه بندی دارد. امروزه الگوریتم‌های فراابتکاری متعددی جهت انتخاب ویژگی و بهینه سازی عملکرد الگوریتم‌های طبقه بندی ارائه شده اند. در مقاله نام برده شده، الگوریتم فراابتکاری گفتار به منظور انتخاب ویژگی‌های مناسب برای طبقه بندی وبسایت‌های جعلی مورد استفاده قرار گرفته است. در این راستا، بهبودی بر الگوریتم فراابتکاری گفتار پیشنهاد شده و الگوریتم گفتار بهبودیافته، ویژگی‌های مناسب را از

میان کل ویژگی های موجود انتخاب کرده و به شبکه عصبی مصنوعی ارسال می کند تا در جهت طبقه بندی وبسایت ها مورد استفاده قرار گیرند. نتایج پیاده سازی الگوریتم پیشنهادی توسط نویسندگان مقاله نام برده شده نشان می دهد که این الگوریتم با دقت نهایی ۹۸/۶۴٪ نسبت به الگوریتم استاندارد بهینه سازی گفتار عملکرد بهتری داشته است. علاوه بر این، نتایج حاکی از برتری نسبت به سه الگوریتم فراابتکاری بهینه سازی ذرات، کرم شب تاب و خفاش است. همچنین در مقاله نام برده شده، الگوریتم پیشنهادی توسط محققین مقاله نامبرده شده با تعدادی از الگوریتم های طبقه بندی ارائه شده در پژوهش های پیشین روی مجموعه داده مشابه، مقایسه شده و برتری آن نشان داده شده است.

(مهدی دادخواه و همکاران، ۱۳۹۵)، ارائه رویکردی به منظور شناسایی و پیش بینی وب سایت

های فیشینگ به وسیله الگوریتم های کلاس بندی بر اساس مشخصه های صفحات وب

امروزه مهمترین ریسک و چالش مورد توجه در تجارت و بانکداری الکترونیک، خطر کلاهبرداری آنلاین و حملات فیشینگ است. حملات فیشینگ همواره به عنوان یکی از ابزارهای پرکاربرد برای مهاجمان، به منظور سرقت کلمه های عبور و رمزهای الکترونیک کاربران در مبادلات الکترونیک بوده است. در این نوع کلاهبرداری، مهاجمان نامه های الکترونیک با ادعاهای مختلف به قربانی ارسال می کند و با تکنیک های مختلفی قربانی را به صفحه های جعلی خود هدایت می کند سپس اقدام به سرقت اطلاعات حساس کاربران مانند رمزهای عبور می نماید. صفحات وب، نامه های الکترونیک و آدرسهای فیشینگ دارای ویژگی هایی هستند که از آن ها می توان برای شناسایی این حملات استفاده کرد. در مقاله نام برده شده رویکردی جهت شناسایی و پیش بینی وب سایت های فیشینگ با استفاده از الگوریتم های کلاس بندی بر اساس مشخصه های صفحات وب ارائه شده است که نرخ خطای کمتری نسبت به سایر تکنیک های مقابله با حملات فیشینگ، به خصوص تکنیک های مشابه مبتنی بر الگوریتم های داده کاوی دارد. در رویکرد ارائه شده توسط محققین مقاله نام برده شده، ویژگی های قابل استفاده در شناسایی صفحات فیشینگ بر اساس میزان تاثیر در

شناسایی این حملات وزن بندی شده سپس با اعمال الگوریتم های کلاس بندی بر روی مجموعه داده های مرتبط، الگویی به منظور شناسایی این حملات استخراج می گردد که قادر به شناسایی حملات فیشینگ بوده و نرخ خطای کمتری را نسبت به سایر روشهای مشابه پیشین نیز دارا می باشد.

(سلمان کمالی زاده و همکاران، ۱۳۹۴)، ارزیابی روش های شناسایی وب سایت فیشینگ

فیشینگ یکی از تکنیک های مهندسی اجتماعی برای فریب کاربران است که به معنای تلاش برای به دست آوردن اطلاعات محرمانه مانند نام کاربری، گذرواژه یا اطلاعات حساب بانکی است. امروزه از مهم ترین چالش های موجود در اینترنت، خطر حملات فیشینگ و کلاهبرداری های اینترنتی است. این حملات تنها در آمریکا، سالیانه چندین میلیارد دلار خسارت به بار می آورد. از این رو، پژوهشگران تلاش های زیادی در جهت شناسایی و مقابله با این گونه حملات داشته اند. در مقاله نام برده شده هدف محققین، ارزیابی روش های شناسایی وب سایت های فیشینگ است. مقاله نام برده شده از نظر هدف کاربردی و از نظر ماهیت از نوع توصیفی - تحلیلی است. در مقاله نام برده شده، ضمن معرفی حمله فیشینگ و روش های موجود، شناسایی وب سایت فیشینگ، بر اساس مطالعات انجام شده و تجارب محققان مقاله نام برده شده با پیشنهاد معیارهایی، روش های شناسایی وب سایت فیشینگ مورد ارزیابی قرار داده اند. نتایج به دست آمده از پژوهش های محققین مقاله نامبرده شده حاکی از آن است که روش هایی که از تکنیک های مختلف شناسایی در کنار هم استفاده می کنند و همچنین اکثر ویژگی های صفحات وب را بررسی می کنند، در شناسایی حمله از موفقیت بیشتری برخوردار می باشند.

در جدول زیر روش های تحقیق مقالات مختلف را مورد بررسی و مقایسه قرار داده ایم:

نام مقاله	الگوریتم مورد استفاده	دقت	نرخ مثبت کاذب	روش
الگوریتم های فرایادگیری و درختان تصمیم اضافی برای شناسایی صفحات فیشینگ	آدابوست اکسترا تری ، بگینگ اکسترا تری، روتیشن فارست اکسترا تری، لجیت بوست اکسترا تری	۹۷.۴۸۵ دقت میانگین مدل ها	۰,۰۳۸	الگوریتم های نامبرده شده را روی داده ها آموزش داده و نتایج حاصل از آنها را با هم مقایسه کرده اند
ارزیابی روش های مبتنی بر یادگیری ماشین برای شناسایی سایت های فیشینگ	آدابوست، بگینگ ، ماشین های بردار پشتیبان ، درخت طبقه بندی و رگرسیون، لجستیک رگرسیون، جنگل های تصادفی، شبکه های عصبی، نایو بایز و درخت های رگرسیون افزایشی بیزی	۹۱,۶۲ دقت میانگین مدل ها	۶,۲۱	داده های خود را با استفاده از روش های تشخیص مبتنی بر یادگیری ماشین طبقه بندی می کنند، و عملکرد هر دسته بند را اندازه گیری می کنند. برای ارزیابی مدل های توسعه داده شده از معیار های اف وان، نرخ خطا، و مساحت تحت منحنی اراو سی استفاده کردند
دسته بندی لینک های فیشینگ با	شبکه های عصبی مصنوعی	۸۶.۶۴	۱۳.۹۶	شبکه های عصبی مختلفی را به همراه

استفاده از شبکه های عصبی عمیق رو دیتاست لینک های پر خطر				الگوریتم های مختلف روی داده ها آموزش داده و عملکرد آنها را با هم مقایسه کرده اند
شناسایی لینک های مخرب بر اساس باز ارسال در شبکه های اجتماعی آنلاین	بیز نت، جنگل تصادفی، جی ۴۸	۸۴.۷۴ ۸۲.۲۲ ۷۹.۰۷	۹.۰۹ ۱۰.۲۳ ۱۰.۷۷	اینها از تکنیک های انتخاب ویژگی و استخراج ویژگی استفاده کرده اند سپس الگوریتم ها را آموزش داده و دقت و کارایی آنها را مورد ارزیابی قرار داده اند
سیستم هوشمند تشخیص وب سایت های فیشینگ با استفاده از دست بند جنگل تصادفی	جنگل تصادفی	۹۷.۳۶		الگوریتم را با استفاده از داده ها آموزش داده سپس با معیار ای یوسی دقت آنرا مورد ارزیابی قرار داده اند و نتایج را روی نمودار اوس نمایش داده اند و الگوریتم های مختلف را با هم مقایسه کرده اند

یوارال نت :یادگیری لینک های نماینده با یادگیری عمیق برای شناسایی لینک های مخرب	شبکه های عصبی مصنوعی	۹۹,۲۹	۰.۷۶۸۳	شبکه های عصبی مختلف را طراحی و آموزش داده اند و کارایی هر مدام را بدست آورده اند
افزایش دقت شناسایی صفحات جعلی وب با استفاده از الگوریتم بهینه سازی گفتار و شبکه عصبی مصنوعی	شبکه های عصبی ، الگوریتم بهینه سازی گفتار	۹۸,۶۴	_____	با استفاده از الگوریتم گفتار ویژگی های مناسب را به شبکه های عصبی داده و دقت آن را ارزیابی کرده اند.
استفاده از شبکه های عصبی خودرمزنگار موازی برای تشخیص صفحات جعلی اینترنتی	شبکه های عصبی خود رمزنگار	۹۳,۴	_____	از دو شبکه عصبی خود رمزنگار با معماری موازی استفاده کرده اند
بهبود روش شناسایی وب سایت فیشینگ با استفاده از داده کاوی روی صفحات وب	جنگل تصادفی ، بیز ساده، جی ۴۷	۹۷,۲۵۹۲ ۹۲,۹۸۰۶ ۹۵,۹۷۴۷	_____	از روش های کاهش ابعاد و انتخاب ویژگی استفاده کرده سپس الگوریتم های خود را آموزش داده و دقت آنها را مورد ارزیابی قرار داده اند

فیشینگ چیست؟

فیشینگ نوعی حمله مهندسی اجتماعی است که اغلب برای سرقت اطلاعات کاربر از جمله شماره کارت اعتباری استفاده می شود. این نوع حملات زمانی اتفاق می افتند که یک مهاجم، خود را به عنوان یک موجودیت قابل اعتماد نشان داده و قربانی را فریب می دهد تا یک ایمیل، پیام فوری یا پیام متنی را باز کند. سپس گیرنده فریب داده می شود تا روی یک پیوند مخرب کلیک کند، که می تواند منجر به نصب بدافزار، مسدود شدن سیستم به عنوان بخشی از حمله باج افزار یا افشای اطلاعات حساس شود. حملات فیشینگ از روزهای اولیه اینترنت وجود داشته است. مجرمان سایبری اولین حملات فیشینگ را در اواسط دهه ۱۹۹۰ با استفاده از سرویس آمریکا آنلین^۱ برای سرقت رمزهای عبور و اطلاعات کارت اعتباری مورد استفاده قرار دادند. در حالی که حملات مدرن از مدل های مهندسی اجتماعی مشابه ای استفاده می کنند، مجرمان سایبری از تاکتیک های تکامل یافته تری استفاده می کنند. فیشینگ در اصل یک روش حمله است که از تاکتیک های مهندسی اجتماعی استفاده می کند تا افراد را وادار به انجام اقدامی کند که برخلاف منافع شان است. با درک بهتر انواع حملات فیشینگ و نحوه شناسایی آنها، سازمان ها می توانند به طور موثرتری از کاربران و داده های آنها محافظت کنند. انواع حملات فیشینگ عبارتند از: فیشینگ ایمیل^۲، آج تی تی پی اس^۳ فیشینگ، اسپیر فیشینگ^۴ و ... می باشد. که هر کدام از این نوع حملات از روش های مختلفی برای آسیب رساندن به کاربران استفاده می کنند و روش های مختلفی می توان آنها را شناسایی کرد

^۱America Online(AOL)

^۲Email phishing

^۳HTTPS phishing

^۴Spear phishing

در این پروژه تمرکز ما بر روی شناسایی وب سایت های فیشینگ و استخراج ویژگی مهم برای شناسایی صفحات فیشینگ می باشد.

ضرورت شناسایی حملات فیشینگ چیست؟

با توجه به روند رو به رشد اینترنت، تعداد رایانه های متصل به اینترنت روز به روز در حال افزایش است. تقریباً همه شرکت ها در حال انتقال اعمال اصلی خود از دنیای واقعی به دنیای سایبری هستند، اگر چه این موضوع باعث افزایش فروش آنها میشود میتواند آسیب پذیری های زیادی مانند حملات سایبری برای این شرکت ها را در پی داشته باشد به ویژه با ساختار ناشناس اینترنت. فیشینگ یکی از انواع محبوب حملات سایبری است که از ناآگاهی کاربر برای آسیب رساندن به آن ها سو استفاده می شود. فیشینگ یکی از شدیدترین حملات سایبری است که محققان علاقه مند به یافتن راه حلی برای آن هستند. در فیشینگ، مهاجمان کاربران نهایی را فریب می دهند و اطلاعات شخصی آنها را می دزدند. برای به حداقل رساندن آسیب ناشی از فیشینگ باید در اسرع وقت شناسایی شود. حملات فیشینگ مختلفی مانند سمیشینگ^۶، ویشینگ^۷... برای آسیب رساندن و سرقت اطلاعات کاربران وجود دارد. تکنیک های مختلف تشخیص فیشینگ بر اساس لیست سفید^۸، لیست سیاه^۹، مبتنی بر محتوا^۵، یادگیری ماشین و ... وجود دارد. در این پروژه تمرکز ما بر روی روش های مبتنی بر یادگیری ماشین می باشد. فیشینگ نوعی حمله مهندسی-وب اجتماعی^{۱۰} در فضای مجازی است که در آن مجرمان اطلاعات ارزشمند کاربران ناآگاه را به سرقت می برند. کارآمد بودن اقدامات متقابل موجود در قالب نرم افزار ضد فیشینگ و روش های محاسباتی برای شناسایی فعالیت های فیشینگ اثبات شده است. با این حال، هکرها روش های جدیدی را برای خنثی سازی این اقدامات متقابل به کار می گیرند. با توجه به ماهیت تکامل یافته حملات فیشینگ، نیاز به اقدامات متقابل جدید و کارآمد بسیار احساس می شود زیرا تأثیر حملات فیشینگ اغلب فاجعه بار است. طرح

^۶Smishing

^۷Vishing

^۸Whitelisting

^۹Blacklisting

^{۱۰}content-based

^{۱۱}social web-engineering

های هوش مصنوعی سنگ بنای روش های متقابل مدرنی است که برای کاهش حملات فیشینگ استفاده می شود.

هوش مصنوعی^۲: هوش مصنوعی که گاهی اوقات هوش ماشینی نیز نامیده می شود شبیه سازی فرایندهای هوش طبیعی^۳ توسط ماشین ها به ویژه سیستم های رایانه ای است. به عبارت دیگر، هوش مصنوعی به سامانه هایی گفته می شود که می توانند واکنش هایی مشابه رفتارهای هوشمند انسانی از جمله، درک شرایط پیچیده، شبیه سازی فرایندهای تفکری و شیوه های استدلالی انسانی و پاسخ موفق به آن ها، یادگیری و توانایی کسب دانش و استدلال برای حل مسائل را داشته باشند، بطور خلاصه هوش مصنوعی را دانش ساخت و طراحی عامل هوشمند تعریف کرده اند. این علم کاربرد های فراوانی در علوم رایانه، علوم مهندسی، تجارت، پزشکی و بسیاری از علوم دیگر دارد بعنوان مثال: در پزشکی تجزیه و تحلیل صدا قلب، ربات های پرستار، ارائه مشاوره و پیش بینی احتمال مرگ بیمار برای هر روش جراحی....، در امور مالی و تجارت تجزیه و تحلیل بازار های مالی، پیش بینی قیمت سهام ها، معاملات الگوریتمی، مدیریت دارای و... از کاربرد های هوش مصنوعی در این علوم هستند. هوش مصنوعی، موضوعی بسیار گسترده است که شاخه های متعددی دارد. شاخه های هوش مصنوعی عبارتند از: یادگیری ماشینی، شبکه های عصبی^۵، سیستم های خبره^۶، پردازش زبان طبیعی^۷، تشخیص گفتار^۸ و بینایی ماشین^۹ و- رباتیک^{۱۰} و منطق فازی است.

یادگیری عمیق^۲: یادگیری عمیق که در زبان فارسی به یادگیری ژرف نیز ترجمه شده است بخشی از خانواده یادگیری ماشینی می باشد که بر روش های تمرکز دارد که مبتنی بر الگوریتم های شبکه عصبی مصنوعی^{۱۳} هستند. این الگوریتم ها د تلاش اند که مغز انسان را شبیه سازی کنند. به طور خلاصه در یادگیری عمیق شبکه های عصبی مصنوعی و الگوریتم های مشابه مغز بشر از مجموعه های عظیم داده مهارت های مورد نظر را فرا می

^۲Artificial Intelligence

^۳Natural Intelligence

^۴machine learning

^۵Neural network

Expert Systems

^۷Natural language processing

^۸speech recognition

^۹Machine vision

^{۱۰}robotic

^{۱۱}Fuzzy logic

^{۱۲}Deep learning

^{۱۳}Artificial neural network

گیرند. همانطور که ما از طریق تجربه چیزهای جدید یاد می گیریم الگوریتم یادگیری عمیق نیز با هر بار تکرار یک کار مهارت خود را نسبت به دفعات قبلی بهبود می بخشد. دلیل استفاده از عبارت یادگیری عمیق این است که شبکه های عصبی لایه های مختلف یا عمیقی دارند که یادگیری را ممکن می سازد.

داده کاوی: داده کاوی فرایندی برای تبدیل داده های خام به اطلاعات مفید می باشد، داده کاوی فرایند استخراج و کشف الگوها در مجموعه داده های بزرگ است که شامل روش هایی در محل تلاقی یادگیری ماشین ، آمار و سیستم های پایگاه داده است. به عبارت دیگر داده کاوی یک زیرشاخه بین رشته ای علوم کامپیوتر و آمار با هدف کلی استخراج اطلاعات (با روشهای هوشمند) از مجموعه داده و تبدیل اطلاعات به یک ساختار قابل درک برای استفاده بیشتر است.

یادگیری ماشین چیست؟

یادگیری ماشینی شاخه ای از هوش مصنوعی^۵ و علوم کامپیوتر^۶ است که بر استفاده از داده ها و الگوریتم ها برای تقلید از روشی که انسان ها یاد می گیرند تمرکز دارد و به تدریج دقت آن را بهبود می بخشد. یادگیری ماشین به عنوان بخشی از هوش مصنوعی در نظر گرفته می شود که به مطالعه الگوریتم های کامپیوتری می پردازد که می تواند به طور خودکار از طریق تجربه و با استفاده از داده ها بهبود یابد. الگوریتم های یادگیری ماشین مدلی را بر اساس داده های نمونه می سازند که به داده های آموزشی معروف است تا پیش بینی ها یا تصمیم گیری ها را بدون برنامه ریزی صریح انجام دهند. الگوریتم های یادگیری ماشین در کاربردهای متنوعی مانند پزشکی، فیلتر کردن ایمیل، تشخیص گفتار و بینایی کامپیوتری استفاده می شوند. انواع الگوریتم های یادگیری ماشین عبارت اند از:

^۵Data mining

^۶artificial intelligence

^۷computer science

درخت تصمیم^{۷۲}، لاجستیک رگرسیون^{۷۸} ماشین بردار پشتیبان^{۷۹}، دسته بندی بیز^{۸۰} نزدیک ترین همسایه^{۷۱}.... برخی از الگوریتم های استفاده شده در این پژوهش به شرح زیر می باشند:

درخت تصمیم:

الگوریتم درخت تصمیم به خانواده الگوریتم های یادگیری ماشین با نظارت^{۸۲} تعلق دارد. می توان از این الگوریتم برای حل مسائل طبقه بندی و رگرسیون استفاده کرد. هدف این الگوریتم ایجاد مدلی است که مقدار یک متغیر هدف را پیش بینی می کند، که برای این منظور درخت تصمیم از نمایش درختی برای حل مسئله استفاده می کند. گره برگ مربوط به یک برچسب کلاس است و ویژگی ها در گره داخلی درخت نشان داده می شوند.

جنگل تصادفی^{۷۳}:

جنگل های تصادفی یک روش یادگیری جمعی^{۸۴} برای طبقه بندی، رگرسیون و سایر وظایف است که با ساختن تعداد زیادی درخت تصمیم در زمان آموزش عمل می کند. برای کارهای طبقه بندی، خروجی جنگل تصادفی کلاسی است که توسط اکثر درختان انتخاب شده است. برای وظایف رگرسیون، میانگین یا میانگین پیش بینی درختان منفرد برگردانده می شود.

نایو بیز^{۷۵}:

الگوریتم ساده بیز یک الگوریتم یادگیری نظارت شده است که بر اساس قضیه بیز است و برای حل مسائل طبقه بندی استفاده می شود. طبقه بندی کننده ساده بیز یکی از ساده ترین و مؤثرترین

^{۷۲}Decision tree

^{۷۸}Logistic Regression

^{۷۹}Support vector machine

^{۸۰}Naive Bayes Classifiers

^{۷۱}KNN

^{۷۲}Supervised machine learning

^{۷۸}Random Forest

^{۷۹}ensemble learning

^{۷۵}Naïve Bayes Classifier

الگوریتم‌های طبقه‌بندی است که به ساخت مدل‌های یادگیری ماشین سریع کمک می‌کند که بتوانند پیش‌بینی‌های سریع انجام دهند.

انتخاب ویژگی چیست؟

هنگام ساخت یک مدل یادگیری ماشین در زندگی واقعی، تقریباً نادر است که همه متغیرهای مجموعه داده برای ساخت یک مدل مفید باشند. افزودن متغیرهای اضافی قابلیت تعمیم مدل را کاهش می‌دهد و همچنین ممکن است دقت کلی یک طبقه‌بندی کننده را کاهش دهد. علاوه بر این، افزودن متغیرهای بیشتر و بیشتر به یک مدل، پیچیدگی کلی مدل را افزایش می‌دهد. انتخاب ویژگی فرآیند کاهش تعداد متغیرهای ورودی هنگام توسعه یک مدل پیش‌بینی کننده است. کاهش تعداد متغیرهای ورودی برای کاهش هزینه محاسباتی مدل سازی و در برخی موارد برای بهبود عملکرد مدل مطلوب است. روش‌های انتخاب ویژگی مبتنی بر آمار شامل ارزیابی رابطه بین هر متغیر ورودی و متغیر هدف با استفاده از آمار و انتخاب آن دسته از متغیرهای ورودی است که قوی‌ترین رابطه را با متغیر هدف دارند. این روش‌ها می‌توانند سریع و موثر باشند، اگرچه انتخاب معیارهای آماری به نوع داده متغیرهای ورودی و خروجی بستگی دارد. هدف از انتخاب ویژگی در یادگیری ماشینی یافتن بهترین مجموعه از ویژگی‌ها است که به فرد امکان می‌دهد مدل‌های مفیدی از پدیده‌های مورد مطالعه بسازد. تکنیک‌های انتخاب ویژگی در یادگیری ماشینی را می‌توان به طور کلی به دسته‌های زیر طبقه‌بندی کرد:

نظارت شده^{۷۶}

این تکنیک‌ها را می‌توان برای داده‌های برچسب دار استفاده کرد و برای شناسایی ویژگی‌های مرتبط برای افزایش کارایی مدل‌های نظارت شده مانند طبقه‌بندی و رگرسیون استفاده می‌شود.

^{۷۶}Supervised

بدون نظارت^{۷۶}

این تکنیک ها را می توان برای داده های بدون برچسب استفاده کرد.

از آنجا که مسئله ما یک مسئله طبقه بندی می باشد به شرح چند مورد از روش ها انتخاب ویژگی مرتبط با طبقه بندی می پردازیم:

۱- فیلتر کردن^{۷۷}

روش های فیلتر، ویژگی های ذاتی ویژگی های اندازه گیری شده را از طریق آمار تک متغیره به جای عملکرد اعتبارسنجی متقابل، انتخاب می کنند. این روش ها سریع تر و از نظر محاسباتی هزینه کمتری نسبت به روش های واریپر^{۷۹} دارند. هنگام برخورد با داده های با ابعاد بالا، از نظر محاسباتی استفاده از روش های فیلتر ارزان تر است.

۲- واریپر:

واریپر ها به روشی برای جستجوی فضای همه زیرمجموعه های ممکن ویژگی ها، ارزیابی کیفیت آنها با یادگیری و ارزیابی طبقه بندی کننده با آن زیر مجموعه ویژگی نیاز دارند. فرآیند انتخاب ویژگی بر اساس یک الگوریتم یادگیری ماشین خاص است که ما سعی می کنیم آن را بر روی یک مجموعه داده معین قرار دهیم. این یک رویکرد جستجوی حریصانه را با ارزیابی همه ترکیب های ممکن از ویژگی ها در برابر معیار ارزیابی دنبال می کند. روش های واریپر معمولاً منجر به دقت پیش بینی بهتری نسبت به روش های فیلتر می شوند. ایده انتخاب زیرمجموعه ویژگی ها این است که بتوانیم بهترین ویژگی هایی را که برای کار طبقه بندی مناسب است پیدا کنیم. ما باید درک کنیم

^{۷۶}unsupervised

^{۷۷}filter

^{۷۸}wrapper

که همه ویژگی‌ها تاثیر یکسانی روی خروجی مدل نداشته و برخی ممکن است مرتبط‌تر از بقیه باشند.

الگوریتم بهینه ساز ازدحام ذرات:

در علوم محاسباتی، بهینه‌سازی ازدحام ذرات یک روش محاسباتی است که با تلاش مکرر برای بهبود راه‌حل کاندید با توجه به معیاری از کیفیت، یک مسئله را بهینه می‌کند. با داشتن جمعیتی از راه‌حل‌های کاندید، که در اینجا ذرات نامیده می‌شوند، و حرکت دادن این ذرات در فضای جستجو بر اساس فرمول ساده ریاضی بر روی موقعیت و سرعت ذره، مشکلی را حل می‌کند. بهینه‌سازی ازدحام ذرات یکی از الگوریتم‌های الهام‌گرفته از طبیعت است و روش کار آن بسیار ساده است به این صورت که راه‌حل بهینه را در فضای راه‌حل جستجو می‌کند. با سایر الگوریتم‌های بهینه‌سازی متفاوت است به گونه ای که فقط تابع هدف مورد نیاز است و وابسته به گرادیان یا هر شکل دیفرانسیل هدف نیست. همچنین دارای هایلپرپارامترهای بسیار کمی است. بهینه‌سازی ازدحام ذرات ابتدا توسط کندی و ابرهارت به‌عنوان یک الگوریتم تکاملی مبتنی بر جمعیت برای شبیه‌سازی رفتار مشارکتی پرندگان در یافتن غذا ابداع شد. مزیت اصلی این الگوریتم نسبت به سایر الگوریتم‌های بهینه‌سازی، توانایی آن در دستیابی به همگرایی سریع در بسیاری از مسائل بهینه‌سازی پیچیده است. علاوه بر این، این الگوریتم دارای چندین مزیت جذاب از جمله سادگی با معادلات ریاضی کمتر و داشتن پارامترهای کمتر در پیاده سازی است.

ابزارها:

با توجه به اینکه با پایتون به راحتی می‌توان فرآیندهای دشوار را مدیریت کرد و استفاده از آن ساده است، در این پژوهش ما از زبان برنامه نویسی پایتون برای توسعه مدل‌ها خود استفاده میکنیم همچنین این زبان برنامه نویسی کتابخانه‌های فراوانی برای کار با الگوریتم‌های هوش مصنوعی و پردازش داده‌ها را دارا می‌باشد. پایتون مجموعه وسیعی از کتابخانه‌ها را برای توسعه هوش

مصنوعی ارائه می‌دهد که شامل موارد پایه‌ای نیز هست که در زمان برنامه نویسی، صرفه جویی می‌کند. پایتون به دلیل کد جمع و جور و خواندنی‌اش مشهور است و از نظر قابلیت استفاده عملاً بی‌نظیر است. ساده و مختصر بودن پایتون دلیلی است که آن را با سایر زبان‌های برنامه نویسی متفاوت می‌کند و به زمان کدنویسی کمتری نیاز دارد. همچنین به توسعه دهنده اجازه می‌دهد تا الگوریتم‌ها را بدون اجرا کردن، سریع آزمایش کند.

داده ها:

به منظور آموزش و اعتبار سنجی روش پیشنهادی از مجموعه داده های فیشینگ وب سایت پایگاه داده یو سی ای استفاده شده است این مجموعه داده شامل ۱۱۰۵۵ نمونه وب سایت است که ۶۱۵۷ وب سایت از آنها فیشینگ و ۴۸۹۸ مورد دیگر وب سایت های مشروع هستند، از هر وب سایت ۳۰ ویژگی استخراج شده است این ویژگی ها عمدتاً صحیح و دوقدری هستند که به شرح زیر می باشند:

دسته	ویژگی	مقدار
ویژگی‌های مرتبط با آدرس سایت	Using the IP Address	-1, 1
	Long URL to Hide the Suspicious Part	1, 0, -1
	Using URL Shortening Services "TinyURL"	1, -1
	URL's having "@" Symbol	1, -1
	Redirecting using "//"	-1, 1
	Adding Prefix or Suffix Separated by (-) to the Domain	-1, 1
	Sub Domain and Multi Sub Domains	-1, 0, 1
	HTTPS (Hyper Text Transfer Protocol with Secure Sockets Layer)	-1, 1, 0
	Domain Registration Length	-1, 1
	Favicon	1, -1
	Using Non-Standard Port	1, -1
	The Existence of "HTTPS" Token in the Domain Part of the URL	-1, 1
ویژگی‌های غیرعادی مرتبط با صفحات و پیوندها	Request URL	1, -1
	URL of Anchor	-1, 0, 1
	Links in <Meta>, <Script> and <Link> tags	1, -1, 0
	Server Form Handler (SFH)	-1, 1, 0
	Submitting Information to Email	-1, 1
ویژگی‌های مرتبط با اطلاعات کد منبع سایت	Abnormal URL	-1, 1
	Website Forwarding	0, 1
	Status Bar Customization	1, -1
	Disabling Right Click	1, -1
	Using Pop-up Window	1, -1
ویژگی‌های مرتبط با دامنه	IFrame Redirection	1, -1
	Age of Domain	-1, 1
	DNS Record	-1, 1
	Website Traffic	-1, 0, 1
	PageRank	-1, 1
	Google Index	1, -1
	Number of Links Pointing to Page	1, 0, -1
	Statistical-Reports Based Feature	-1, 1

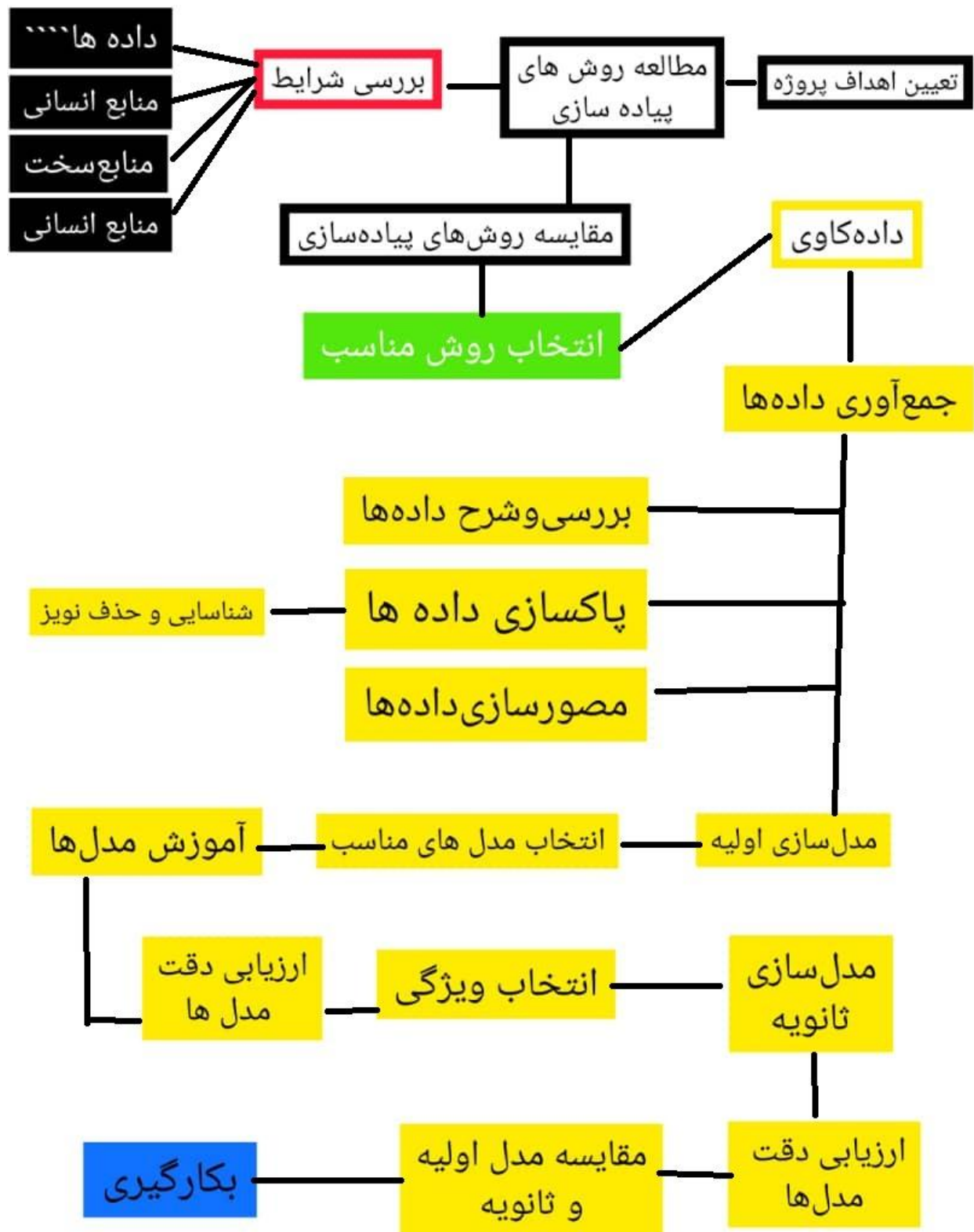
روش کار در این پژوهش:

اگر مجموعه ای از ویژگی ها به ما داده شود، چگونه می توانیم بهینه ترین زیر مجموعه را بشناسیم؟ از آنجا که ما میخواهیم انتخاب ویژگی ها برای هر الگوریتم دسته بندی به صورت انحصاری صورت گیرد در این پروژه از روش دوم یعنی وارپپر برای انتخاب ویژگی استفاده میکنیم. تا اینکار

باعث افزایش دقت مدل ها شود. برای وزن دهی به ویژگی ها ما از الگوریتم بهینه ساز ازدحام ذرات^۱ استفاده میکنیم، روش کار این الگوریتم به این صورت می باشد که برای هر مجموعه داده یک تابع هزینه، هزینه را محاسبه می کند و سپس الگوریتم ازدحام ذرات سعی می کند که مجموعه داده ای را پیدا کند که برای الگوریتم دسته بند مورد نظر کمترین تابع هزینه را به همراه داشته باشد بدین منظور ما یک تابع هزینه نیز نیاز داریم که الگوریتم بهینه ساز ازدحام ذرات آن را برای هر الگوریتم طبقه بندی بهینه کند. در این پژوهش علاوه بر شناسایی صفحات فیشینگ با الگوریتم های یادگیری ماشین ما می خواهیم با استفاده از روش های انتخاب ویژگی هزینه آموزش را کاهش و کیفیت مدل ها را افزایش داده برای این منظور ابتدا مدل ها انتخاب شده را با تمام داده ها آموزش داده و دقت آنها را سنجیده سپس با استفاده از روش های انتخاب ویژگی^۲ ویژگی های را که بیشترین تاثیر را روی خروجی مدل داشته شناسایی میکنیم و الگوریتم ماشین بردار پشتیبان با کرنل اربی اف را با کمک الگوریتم ازدحام ذرات بهینه میکنیم و دقت مدل توسعه داده شده را محاسبه کرده و با سایر الگوریتم ها مقایسه میکنیم.

^۱particle swarm optimization

^۲Feature selection



شرح پروژه:

ابتدا کتابخانه های پانداز: برای خواندن فایل و تجزیه و تحلیل دیتا فریم، نامپای: برای کار با آرایه های نامپای، مت پلات لیب: برای رسم نمودار، سیپورن: برای رسم نمودار، سایکیت لرن برای: استفاده از الگوریتم های یادگیری ماشینی، استاندارد کردن دیتا و اندازه گیری خطا فراخوانی شده اند. همچنین کتابخانه پایسوارم برای کار با الگوریتم ازدحام ذرات فراخوانی شده است. دیتاست مورد نظر ما با فرمت ای اراف اف^۳ ذخیره شده است که این فرمت مناسب کار در محیط پایتون نمی باشد به همین علت با قطعه کد زیر ما فرمت فایل را به سی اس وی^۴ تغییر می دهیم:

^۳.arff

^۴csv

```
[11] path_to_directory="/content/"
files = [arff for arff in os.listdir(path_to_directory) if arff.endswith(".arff")]

def toCsv(content):
    data = False
    header = ""
    newContent = []
    for line in content:
        if not data:
            if "@attribute" in line:
                attri = line.split()
                columnName = attri[attri.index("@attribute")+1]
                header = header + columnName + ","
            elif "@data" in line:
                data = True
                header = header[:-1]
                header += '\n'
                newContent.append(header)
        else:
            newContent.append(line)
    return newContent

for z,file in enumerate(files):
    with open(path_to_directory+'/'+file , "r") as inFile:
        content = inFile.readlines()
        name,ext = os.path.splitext(inFile.name)
        new = toCsv(content)
        with open(name+".csv", "w") as outFile:
            outFile.writelines(new)
```

مجموعه داده ما دارای ۳۱ ستون و ۱۰۵۵ سطر می باشد، تمام ستون ها شامل مقادیری کتگوریکال می باشند و در این مجموعه داده هیچ سطر با مقادی نامعلوم^۵ وجود ندارد.قطعه کد زیر پنج سطر ابتدای مجموعه داده ما را نشان می دهد:

```
[15] train.head()
```

	having_IP_Address	URL_Length	Shortening_Service	having_At_Symbol	double_qlash_redirecting	Prefix_Suffix	having_Sub_Domain	SSLfinal_State	Domain_registration_length	Feixicom	port	HTTPS_token	Request_URL	URL_of_Anchor	Links_in_tags	SFX
0	-1	1	1	1	-1	-1	-1	-1	-1	1	1	-1	1	-1	1	-1
1	1	1	1	1	1	-1	0	1	-1	1	1	-1	1	0	-1	-1
2	1	0	1	1	1	-1	-1	-1	-1	1	1	-1	1	0	-1	-1
3	1	0	1	1	1	-1	-1	-1	-1	1	1	-1	-1	0	0	-1
4	1	0	-1	1	1	-1	1	1	-1	1	1	1	1	0	0	-1

کدهای زیر توصیفی از ویژگی ها را به صورت دیتافریم نمایش می دهد:

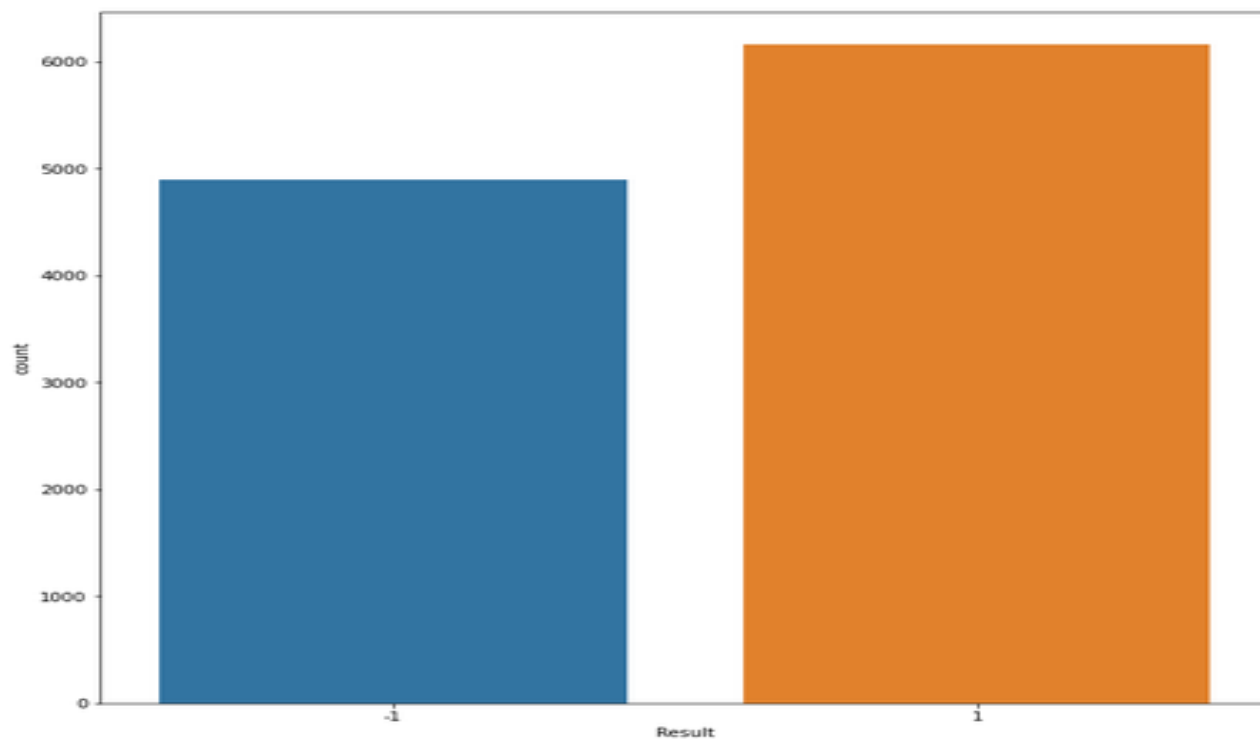
`train.describe().transpose()`

	count	mean	std	min	25%	50%	75%	max
having_IP_Address	11055.0	0.313795	0.949534	-1.0	-1.0	1.0	1.0	1.0
URL_Length	11055.0	-0.633198	0.766095	-1.0	-1.0	-1.0	-1.0	1.0
Shortlinking_Service	11055.0	0.738761	0.673998	-1.0	1.0	1.0	1.0	1.0
having_Ai_Symbol	11055.0	0.700588	0.713598	-1.0	1.0	1.0	1.0	1.0
double_slash_redireciting	11055.0	0.741474	0.671011	-1.0	1.0	1.0	1.0	1.0
Prefix_Suffix	11055.0	-0.734962	0.678139	-1.0	-1.0	-1.0	-1.0	1.0
having_Sub_Domain	11055.0	0.063953	0.817518	-1.0	-1.0	0.0	1.0	1.0
SSLFinal_State	11055.0	0.250927	0.911892	-1.0	-1.0	1.0	1.0	1.0
Domain_registration_length	11055.0	-0.336771	0.941629	-1.0	-1.0	-1.0	1.0	1.0
FavIcon	11055.0	0.628584	0.777777	-1.0	1.0	1.0	1.0	1.0
port	11055.0	0.728268	0.685324	-1.0	1.0	1.0	1.0	1.0
HTTP_&token	11055.0	0.675079	0.737779	-1.0	1.0	1.0	1.0	1.0
Request_URL	11055.0	0.186793	0.982444	-1.0	-1.0	1.0	1.0	1.0
URL_of_Anohor	11055.0	-0.076526	0.715138	-1.0	-1.0	0.0	0.0	1.0
Links_in_tags	11055.0	-0.118137	0.763973	-1.0	-1.0	0.0	0.0	1.0
BFH	11055.0	-0.595749	0.759143	-1.0	-1.0	-1.0	-1.0	1.0
submitting_to_email	11055.0	0.635640	0.772021	-1.0	1.0	1.0	1.0	1.0
Abnormal_URL	11055.0	0.705292	0.708949	-1.0	1.0	1.0	1.0	1.0
Redirect	11055.0	0.115694	0.319872	0.0	0.0	0.0	0.0	1.0
on_mouseover	11055.0	0.762099	0.647490	-1.0	1.0	1.0	1.0	1.0
RightClick	11055.0	0.913885	0.405991	-1.0	1.0	1.0	1.0	1.0
popupWindow	11055.0	0.613388	0.789818	-1.0	1.0	1.0	1.0	1.0
iframe	11055.0	0.816915	0.576784	-1.0	1.0	1.0	1.0	1.0
age_of_domain	11055.0	0.061239	0.998168	-1.0	-1.0	1.0	1.0	1.0
DNSSReoord	11055.0	0.377114	0.926209	-1.0	-1.0	1.0	1.0	1.0
web_traffic	11055.0	0.287291	0.827733	-1.0	0.0	1.0	1.0	1.0
Page_Rank	11055.0	-0.483673	0.875289	-1.0	-1.0	-1.0	1.0	1.0
Google_Index	11055.0	0.721574	0.692369	-1.0	1.0	1.0	1.0	1.0
Links_pointing_to_page	11055.0	0.344007	0.569944	-1.0	0.0	0.0	1.0	1.0
Statistical_report	11055.0	0.719584	0.694437	-1.0	1.0	1.0	1.0	1.0
Result	11055.0	0.113885	0.993539	-1.0	-1.0	1.0	1.0	1.0

قطعه کد زیر نمودار همبستگی بین ویژگی ها را رسم و با فرمت پی ان جی ذخیره می کند:

```
plt.figure(figsize=(30,40),dpi=100)
sb.heatmap(train.corr(),annot=True,linewidths=0.2)
plt.title("Relationship between features and label")
plt.savefig("/content/corr_n.png")
```

نمودار زیر نشان می‌دهد که چه میزان از لینک هایموجود در دیتاست مخرب و چه میزان از لینک ها سالم هستند:



قطعه کد زیر میزان داده های متعلق به هر مقدار عدد در هر ستون را نشان می‌دهد:

```
for column in train.columns:  
    print("value count of "+column+": " '\n',train[column].value_counts())  
    print(20*"=+=")
```


در عکس زیر مقادیر یکتای موجود در هر ستون را می بینیم:

```
for column in train.columns:
    print(column+":",train[column].unique())

having_IP_Address: [-1  1]
URL_Length: [ 1  0 -1]
Shortining_Service: [ 1 -1]
having_At_Symbol: [ 1 -1]
double_slash_redirecting: [-1  1]
Prefix_Suffix: [-1  1]
having_Sub_Domain: [-1  0  1]
SSLfinal_State: [-1  1  0]
Domain_registration_length: [-1  1]
Favicon: [ 1 -1]
port: [ 1 -1]
HTTPS_token: [-1  1]
Request_URL: [ 1 -1]
URL_of_Anchor: [-1  0  1]
Links_in_tags: [ 1 -1  0]
SFH: [-1  1  0]
Submitting_to_email: [-1  1]
Abnormal_URL: [-1  1]
Redirect: [0 1]
on_mouseover: [ 1 -1]
RightClick: [ 1 -1]
popUpWidnow: [ 1 -1]
Iframe: [ 1 -1]
age_of_domain: [-1  1]
DNSRecord: [-1  1]
web_traffic: [-1  0  1]
Page_Rank: [-1  1]
Google_Index: [ 1 -1]
Links_pointing_to_page: [ 1  0 -1]
Statistical_report: [-1  1]
Result: [-1  1]
```

با استفاده از کد زیر داده ها برای آموزش و اعتبارسنجی مدل ها به دو بخش ترین و تست تقسیم می کنیم:

```
[29] X=train.drop("Result",axis=1).values
y=train.Result
x_train,x_test,y_train,y_test=train_test_split(X,y,test_size=0.2,random_state=42,)
```

قطعه کد زیر تمام مدل ها را می سازد:

```
[30] ##models
dt_gini=DecisionTreeClassifier(criterion="entropy",max_depth=5)
dt_entropy=DecisionTreeClassifier(criterion="gini",max_depth=5)
knn=KNeighborsClassifier(n_neighbors=5)

rf=RandomForestClassifier()

lr=LogisticRegression(random_state=42)

mlp=MLPClassifier(hidden_layer_sizes=(40,),max_iter=500)

poly_svm=SVC(kernel='poly')
linear_svm=SVC(kernel='linear')
sigmoid_svm=SVC(kernel='sigmoid')
rbf_svm=SVC(kernel='rbf')

model_list=[dt_gini,dt_entropy,knn,rf,lr,mlp,poly_svm,linear_svm,sigmoid_svm,rbf_svm]
```

قطعه کد زیر مدل ها را با دیتا ها که از قبل آماده شده اند آموزش داده و سپس دقت آنها را با معیار های مختلفی می سنجد و نمودار های کانفیژن ماتریکس را برای هر الگوریتم رسم می کند:

```
train_list=list()
test_list=list()
for model in model_list:
    model.fit(X,y)
    train_score=model.score(x_train,y_train)
    test_score=model.score(x_test,y_test)
    print("train score: ",train_score,"and test score is: ",test_score)
    test_list.append(test_score)
    train_list.append(train_score)
    predict_test=model.predict(x_test)
    print('\n',classification_report(y_test,predict_test),'\n')
    cm=confusion_matrix(y_test,predict_test)
    ax= plt.subplot()
    sb.heatmap(cm, annot=True, fmt='g', ax=ax);
    ax.set_xlabel('Predicted labels');ax.set_ylabel('True labels');
    ax.set_title('Confusion Matrix');
    plt.show()
    print('\n'+10*" _ " +'\n')
```

تابع لاس پر پارتيکال^۶ که در زیر به نمایش گذاشته شده است مقدار تابع هزینه را برای هر ذره محاسبه کرده این تابع توسط تابع لاس^۷ فراخوانی میشود که خود تابع لاس مجموعه لاس ها محاسبه شده توسط لاس پر پارتيکال را به صورت یک آرایه نامپای برمیگرداند:

```
def loss_per_particle(m, alpha):  
  
    total_features = 30  
    if np.count_nonzero(m) == 0:  
        X_sub = X.value  
    else:  
        X_sub = X[:,m==1]  
  
    classifier.fit(X_sub, y)  
  
    P = (classifier.predict(X_sub) == y).mean()  
  
    j = (alpha * (1.0 - P) + (1.0 - alpha) * (1 - (X_sub.shape[1] / total_features)))  
    return j
```

```
def loss(x, alpha=0.88):  
  
    n_particles = x.shape[0]  
    j = [loss_per_particle(x[i], alpha) for i in range(n_particles)]  
    return np.array(j)
```

در ادامه ما بهینه ساز ازدحام ذرات را تعریف کرده ایم با کمک یک حلقه و دو تابعی که در بالا نام برده شده است سعی کردیم برای هر الگوریتم بهترین مجموعه داده را پیدا کنیم و با استفاده از کد زیر هر داده روی مجموعه داده انحصاری خود آموزش داده شده و سپس دقت آن با معیار های مختلفی سنجیده شده است و در پایان دیتافریمی ایجاد شده است که دقت هر الگوریتم در شرایط مختلف نشان میدهد:

^۶loss_per_particle
^۷loss

pso

```
In [ ]: class PSO(object):
    def __init__(self,particle_num,particle_dim,iter_num,c1,c2,w,max_value,min_value):

        self.particle_num = particle_num
        self.particle_dim = particle_dim
        self.iter_num = iter_num
        self.c1 = c1
        self.c2 = c2
        self.w = w
        self.max_value = max_value
        self.min_value = min_value

    def swarm_origin(self):
        particle_loc = []
        particle_dir = []
        for i in range(self.particle_num):
            tmp1 = []
            tmp2 = []
            for j in range(self.particle_dim):
                a = random.random()
                b = random.random()
                tmp1.append(a * (self.max_value - self.min_value) + self.min_value)
                tmp2.append(b)
            particle_loc.append(tmp1)
            particle_dir.append(tmp2)

        return particle_loc,particle_dir

    def fitness(self,particle_loc):

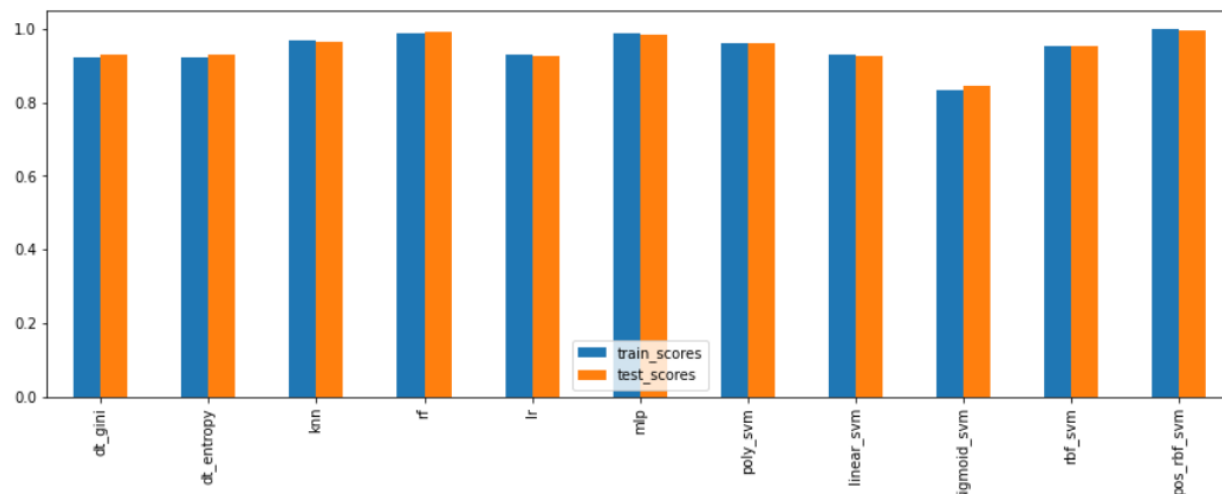
        fitness_value = []
        for i in range(self.particle_num):
            rbf_svm = SVC(kernel = 'rbf', C = particle_loc[i][0], gamma = particle_loc[i][1])
            cv_scores = cross_val_score(rbf_svm,trainX,trainY,cv =3,scoring = 'accuracy')
```

```
[ ] score_df["train_score_pos"]=pos_train_list
score_df["test_score_pos"]=pos_test_list
score_df
```

تذکر: این بخش از کد به دلیل سنگین بودن محاسبات آن در زمان اجرا نیازمند سخت افزار قوی برای اجرا می باشد ،متاسفانه به همین علت خروجی این قطعه کد برای بنده در دسترس نمی باشد!!!

نتایج :

تعداد حملات فیشینگ نسبت به گذشته رشد چشم گیری داشته است جلوگیری و شناسایی به موقع این حملات تشخیص صفحات فیشینگ علاوه بر اینکه امنیت کاربران را افزایش دهد می تواند از وارد شدن لطمه به اعتبار کسب و کار ها جلوگیری کند و از ضرر و زیان های بسیاری جلوگیری کند. یکی از راهکارای پیشنهاد شده برای شناسایی حملات فیشینگ استفاده از هوش مصنوعی برای تشخیص صفحات فیشینگ می باشد، الگوریتم های هوش مصنوعی می توانند داده های بسیار زیادی را در لحظه تحلیل کرد و به سرعت نتیجه را ارائه دهند. به همین علت هوش مصنوعی میتواند سرعت سیستم های تشخیص صفحات فیشینگ را افزایش دهد. ما در این پژوهش تلاش کردیم تا الگوریتم های مختلف هوش مصنوعی را جهت شناسایی صفحات فیشینگ مورد ارزیابی قرار دهیم که این الگوریتم ها عملکرد خوبی را داشته اند ما میتوانیم به کمک تکنیک های انتخاب ویژگی، ویژگی های که در شناسایی صفحات فیشینگ نقش بسزایی دراند را شناسایی کرده و الگوریتم ها را فقط با این ویژگی های استخراج شده آموزش دهیم این کار سرعت آموزش الگوریتم ها را افزایش داده و همچنین سرعت شناسایی آنها توسط مدل ها را افزایش میدهد، الگوریتم ازدحام ذرات میتواند ویژگی های مهم متناسب با هر الگوریتم طبقه بندی را شناسایی کند و سبب افزایش دقت آنها شود ، نتایج حاصل از مقایسه نشان می دهد که الگوریتم ها داری سرعت بالاتری نسبت به حالت اولیه هستن و در زمان آموزش آنها کاهش میابد و از دادن اطلاعات اضافی به مدل جلوگیری می شود. نمودار زیر به خوبی ن نشان میدهد که دقت الگوریتم ماشین بردار پشتیبان که با کمک الگوریتم بهینه سازی ازدحام ذرات بالاتر از سایر الگوریتم های می باشد و عملکرد بهتری را دارد:



در جدول زیر دقت الگوریتم های مختلف در هر دو بخش تست و ترین نشان داده شده که الگوریتم پیشنهادی ما دقت در هر دو بخش نسبت به سایر الگوریتم ها با دقت نزدیک به ۹۹ دو هر دو بخش عملکرد بهتری نسبت به سایر الگوریتم ها داشته:

	dt_gini	dt_entropy	knn	rf	lr	mlp	poly_svm	linear_svm	sigmoid_svm	rbf_svm	pos_rbf_svm
train_scores	0.921303	0.921303	0.967209	0.989258	0.928991	0.982135	0.960199	0.928087	0.832542	0.954545	0.999869
test_scores	0.928087	0.928087	0.964722	0.991407	0.925825	0.977838	0.959294	0.927635	0.843962	0.953415	0.995499

پیشنهادهای کاربردی:

با استفاده از نتایج حاصل از این پژوهش می توان الگوریتم های مناسب و همچنین ویژگی های که بیشترین تاثیر بر سالم بودن یا نبودن یک وب سایت را دارند شناسایی کرده و از آنها برای توسعه سامانه جهت افزایش امنیت کاربران در مقابل حملات فیشینگ استفاده کرد.

پیشنهاد آتی:

با استفاده از هوش مصنوعی در فضای وب می توان امنیت کاربران را به طور قابل توجه ای افزایش داد که موضوع سبب اعتماد بیشتر کاربران به کسب کار ها و ... می شود از آنجا که تعداد قربانیان حملات فیشینگ بسیار زیاد می باشد توسعه سیستم های برای شناسایی حملات فیشینگ بسیار ضروری می باشد، پیشنهاداتی برای توسعه سیستم های جهت شناسایی حملات فیشینگ به شرح زیر می باشند:

۱- سامانه متشکل از دو بخش توسعه داده شود بخش اول وظیفه استخراج و بهینه سازی ویژگی های وب سایت را دارد و ویژگی های که استخراج می کند را در اختیار بخش دوم قرار می دهد بخش دوم الگوریتم یادگیری جمعی بوستینگ^۸ است که مجموعه ای از الگوریتم های مختلف طبقه بند شامل می شود که سعی می کنند تشخیص دهند که آیا صفحه مورد نظر یک صفحه سالم است یا خیر و در صورت ناسالم بودن صفحه از طریق مرورگر به کاربر اخطار داده شود.

۲- مدلی جهت تشخیص صفحات فیشینگ با استفاده از ترکیب شبکه های عصبی با الگوریتم ژنتیک توسعه داده شود.

محدودیت ها:

- ۱- نداشتن سخت افزار مناسب برای اجرای سریع برنامه ها
- ۲- دسترسی محدود به سرویس های مانند گوگل کولب
- ۳- زمان بر بودن انجام محاسبات

REFERENCES:

^۸Boosting

- ۱-Korkmaz, M., Kocyigit, E., Sahingoz, O. K., & Diri, B. (2020). Deep Neural Network Based Phishing Classification on a High-Risk URL Dataset. In *SoCPaR* (pp. 648-657).
- ۲-Salem, O., Hossain, A., & Kamala, M. (2010, June). Awareness program and ai based tool to reduce risk of phishing attacks. In *2010 10th IEEE International Conference on Computer and Information Technology* (pp. 1418-1423). IEEE.
- ۳-Alsariera, Y. A., Adeyemo, V. E., Balogun, A. O., & Alazzawi, A. K. (2020). Ai meta-learners and extra-trees algorithm for the detection of phishing websites. *IEEE Access*, 8, 142532-142542.
- ۴-Le, H., Pham, Q., Sahoo, D., & Hoi, S. C. (2018). URLNet: Learning a URL representation with deep learning for malicious URL detection. *arXiv preprint arXiv:1802.03162*.
- ۵-Miyamoto, D., Hazeyama, H., & Kadobayashi, Y. (2008, November). An evaluation of machine learning-based methods for detection of phishing sites. In *International Conference on Neural Information Processing* (pp. 539-546). Springer, Berlin, Heidelberg.
- 6-Khonji, M., Iraqi, Y., & Jones, A. (2013). Phishing detection: a literature survey. *IEEE Communications Surveys & Tutorials*, 15(4), 2091-2121.
- 7-Cao, J., Li, Q., Ji, Y., He, Y., & Guo, D. (2016). Detection of forwarding-based malicious URLs in online social networks. *International Journal of Parallel Programming*, 44(1), 163-180.
- 8-Vazhayil, A., Harikrishnan, N. B., Vinayakumar, R., Soman, K. P., & Verma, A. D. R. (2018). PED-ML: Phishing email detection using classical machine learning techniques. In *Proc. 1st AntiPhishing Shared Pilot 4th ACM Int. Workshop Secur. Privacy Anal.(IWSPA)* (pp. 1-8). Tempe, AZ, USA.
- 9-Subasi, A., Molah, E., Almkallawi, F., & Chaudhery, T. J. (2017, November). Intelligent phishing website detection using random forest classifier. In *2017 International conference on electrical and computing technologies and applications (ICECTA)* (pp. 1-5). IEEE.
- 10-Mondal, S., Maheshwari, D., Pai, N., & Biwalkar, A. (2019, December). A Review on Detecting Phishing URLs using Clustering Algorithms. In *2019 International Conference on Advances in Computing, Communication and Control (ICAC3)* (pp. 1-6). IEEE.
- ۱۱- امیری, & نادری. (۲۰۲۰). استفاده از شبکه‌های خودرمن‌نگار موازی به منظور تشخیص صفحات جعلی اینترنتی. *مجله نوآوری های فناوری اطلاعات و ارتباطات کاربردی*, ۲(۲), ۴۱-۵۰.
- ۱۲- بهارلو, & یاری. بهبود روش شناسایی وب سایت فیشینگ با استفاده از داده‌کاوی روی صفحات وب. *دوفصلنامه فناوری اطلاعات و ارتباطات ایران*, ۴۳(۴۳), ۲۷.
- ۱۳- صف آرا, & صباح نو. (۲۰۲۰). افزایش دقت شناسایی صفحات جعلی وب با استفاده از الگوریتم بهینه سازی گفتار و شبکه عصبی مصنوعی. *فصلنامه علمی-پژوهشی فرماندهی و کنترل*, ۳(۴), ۷۲-۹۱.
- ۱۴- دادخواه مهدی, داورپناه جزی محمد, & سعیدی مبارکه مجید. ارائه رویکردی به منظور شناسایی و پیش بینی وب سایت‌های فیشینگ به وسیله الگوریتم‌های کلاس بندی بر اساس مشخصه های صفحات وب.
- ۱۵- کمالی زاده, سلمان, & شاه‌محمدی. ارزیابی روش های شناسایی وب‌سایت فیشینگ. *فصلنامه علمی پژوهش‌های اطلاعاتی و جنایی*, ۱۱(۴۱), ۳۸-۹.

