



آشکارسازی بدافزارها با استفاده از دسته‌بندی دنباله‌های با طول متغیر

فاطمه حسینی^۱، میترا میرزارضایی^{۲*} و آرش شریفی^۳

دانشگاه آزاد اسلامی، واحد علوم و تحقیقات، گروه مهندسی کامپیوتر، تهران، ایران

چکیده

در این مقاله روشی مبتنی بر گراف به عنوان استخراج ویژگی برای دنباله‌های با طول متغیر پیشنهاد می‌شود. روش پیشنهادی بدون ثابت کردن طول دنباله‌ها، با تعیین پر تکرارترین دستورها و گذاشتن باقی دستورها در مجموعه 'other' از لحاظ سرعت و حافظه صرفه‌جویی می‌کند. با توجه به میزان شباهت ویژگی‌ها، هر نمونه امتیازی می‌گیرد و از امتیازات جهت دسته‌بندی استفاده می‌شود. برای بهبود نتایج، دو رویکرد پیشنهاد می‌شود. در رویکرد نخست، ویژگی‌های استخراج شده از روش‌های امتیازدهی بر روی آپکد، هگزادسیمال و فراخوانی سیستمی در ورودی دسته‌بندها ترکیب می‌شوند. در رویکرد دوم، خروجی دسته‌بندهای مختلف ترکیب شده و از رأی اکثریت استفاده می‌شود. رویکرد پیشنهادی با دقت ۹۷٪ بدافزارهای دگرگون شده رایانه‌ای از مجموعه vxheaven را نه تنها شناسایی، بلکه دسته بدافزارها را نیز تعیین می‌کند؛ در حالی که روش‌های SSD و HMM تحت شرایط یکسان با دقت ۸۴٪ و ۸۰٪ توانستند بدافزارها را شناسایی کنند.

واژگان کلیدی: آشکارسازی بدافزارها، روش‌های مبتنی بر گراف، ترکیب دسته‌بندها، دسته‌بندی با طول متغیر، ماشین بردار پشتیبان

Malware Detection using Classification of Variable-Length Sequences

Fatemeh Hosseini¹, Mitra Mirzarezaee² & Arash Sharifi³

Islamic Azad University, Science and Research Branch, Department of Computer Eng.,
Tehran, Iran

Abstract

In this paper, a novel method based on the graph is proposed to classify the sequence of variable length as feature extraction. The proposed method overcomes the problems of the traditional graph with variable length of data, without fixing length of sequences, by determining the most frequent instructions and insertion the rest of instructions on the set of "other", save speed and memory. According to features and the similarities of them, a score is given to each sample and that is used for classification. To improve the results, the method is not used alone, but in the two approaches, this method is combined with other existing Technique to get better results. In the first approach, which can be considered as a feature extraction, extracted features from scoring techniques (Hidden Markov Model, simple substitution distance and similarity graph) on op-code sequences, hexadecimal sequences and system calls are combined at classifier input. The second approach consists of two steps, in the first step; the scores which obtained from each of the scoring Technique are given to the three support vector machine. The outcomes are combined according to the weight of each Technique and the final decision is taken based on the majority vote. Among the components of the support vector machine, when given a higher weight in the similarity graph method (the proposed method), the result is better, Because the similarity graph method is more accurate than the other two methods. Then, in the second section, considering the strengths and benefits of each classifier, classifier outputs are combined and the majority voting is used. Three methods have been tested for group combinations, including Ensemble

* Corresponding author

* نویسنده عهده‌دار مکاتبات

Averaging, Bagging, and Boosting. Ensemble Averaging consisting of the combination of four classifiers of random forests, a support vector machine (as obtained in the previous section), K nearest neighbors and naive Bayes, and the final decision is taken based on the majority vote; therefore, it is used as the proposed method. The proposed approach could detect metamorphic malware from Vxheaven set and also determines categories of malware with accuracy of 97%, while the SSD and HMM methods under the same conditions could detect malware with an accuracy of 84% and 80% respectively.

Keywords: Malware Detection, Graph Techniques, Combining Classifiers, Variable Length Classification, Support vector machine

بدافزار ساختار برنامه‌نویسی استاندارد ندارند، لذا دارای گراف‌های متفاوت با طول‌های متفاوتی هستند و نمی‌توان آن‌ها را با هم مقایسه کرد، مگر با استفاده از روش‌های نرمال‌سازی، طول تمام گراف‌ها را به یک اندازه در آورد. مشکل دیگر این است که به‌طور معمول بدافزارهای امروزی، شیوه‌نامه‌های بیشتری دارند؛ لذا حجم گراف تولید شده نیز بالا می‌رود. در این صورت گراف‌های مورد آزمایش بسیار بزرگ و پرتعداد می‌شوند و الگوریتم‌هایی که از حافظه استفاده می‌کنند، پس از چندین ساعت، با خطای حافظه متوقف می‌شوند؛ از این‌رو در این پژوهش، الگوریتمی طراحی شده که بر این مشکلات غلبه کند.

روش پیشنهادی مبتنی بر گراف است و بر این ایده بنا شده که دنباله‌های آپکد، هگزایسیمال و فراخوانی‌های سامانه‌ای در کنار یکدیگر در فایل‌های بدافزار می‌توانند دانش خوبی از میزان شباهت فایل‌های اجرایی نمایش دهند. دنباله‌های موجود، حاوی تعداد زیادی آپکد، دستورهای هگزایسیمال و تابع فراخوانی سامانه‌ای مجزا هستند؛ از این دستورها، تنها تعداد خاصی از آنها به‌طور مداوم ظاهر می‌شوند و در طبیعت ذاتی بدافزار یا فایل خوش‌خیم سهمیم هستند؛ لذا روش پیشنهادی در مواجهه با داده‌های با طول متغیر، بدون ثابت‌کردن طول دنباله‌ها، با تعیین پرتکرارترین دستورها و گذاشتن باقی دستورها در مجموعه 'other' از لحاظ سرعت و حافظه صرفه‌جویی می‌کند. برای بهبود نتایج و پیدا کردن یک الگوی مناسب شباهت بین دنباله‌های با طول متغیر، از روش ارائه‌شده به‌تنهایی استفاده نمی‌شود؛ بلکه از دو رویکرد ترکیب ویژگی‌ها در ورودی و ترکیب دسته‌بندها در خروجی استفاده می‌شود. در رویکرد نخست، از ترکیب روش‌های مدل مخفی مارکوف، فاصله جایگزینی ساده و روش مبتنی بر گراف در ورودی دسته‌بندها استفاده می‌شود. همچنین با بهره‌گیری از ترکیب نتایج دسته‌بندها می‌توان کارایی سامانه آشکارسازی بدافزارها را بهبود بخشید و نرخ خطای تشخیص را کاهش داد. در رویکرد دوم، ترکیب دسته‌بندها در خروجی مورد آزمایش قرار گرفته است که شامل دو مرحله است. در گام نخست، امتیازات حاصل از هر یک از روش‌های امتیازدهی به سه دسته‌بند ماشین بردار پشتیبان داده شده است. خروجی‌ها بر

۱- مقدمه

با توجه به اهمیت و جایگاه رایانه‌ها در عصر فناوری و این موضوع که تهدید و در معرض خطر قرار گرفتن رایانه‌ها می‌تواند تهدید جدی برای جامعه باشد؛ عرضه روشی هوشمند جهت آشکارسازی بدافزارها امری ضروری است. در حوزه آشکارسازی بدافزارها، اغلب با داده‌ها یا کدهای دارای طول متغیر سروکار داریم. ماهیت این نوع داده‌ها باعث می‌شود مقایسه، تجزیه و تحلیل و در نهایت دسته‌بندی آنها با چالشی بزرگ مواجه باشد. از آنجا که روش‌های دسته‌بندی موجود، اغلب با مجموعه ویژگی‌های با طول ثابت کار می‌کنند؛ باید به‌دنبال راه‌کاری بود تا داده‌های با طول نابرابر به قالبی مناسب برای این دسته‌بندها تبدیل شوند. ساده‌ترین راه برای رسیدن به این منظور، نرمال‌سازی طول بردار ویژگی است. اگر چه به‌طور معمول به‌دلیل سادگی این روش از آن زیاد استفاده می‌شود، اما نرمال‌سازی طول بردار ویژگی باعث تخریب داده‌ها و افت شدید کارایی تشخیص می‌شود. اگر بتوان یک خصوصیت رایج برای همه اعضای یک خانواده دگرگون‌شده تعیین کرد، پس می‌توان از این ویژگی برای تشخیص و دسته‌بندی آنها استفاده کرد. از آنجا که هیچ الگوریتمی وجود ندارد که برای تمامی شرایط و در تمامی زمان‌ها بهترین یادگیر را به‌وجود آورد و هیچ دسته‌بندی قادر به تشخیص صحیح همه الگوها در تمام شرایط نیست، استفاده از نتایج چند دسته‌بند با عنوان یادگیری دسته جمعی، یک روش مؤثر است که در آن به‌منظور بهبود دقت یادگیری، نتایج دسته‌بندها با یکدیگر ترکیب شده و یک سامانه مرکب شکل می‌گیرد. ترکیب نتایج دسته‌بندهایی که از ویژگی‌های متفاوتی استفاده می‌کنند می‌تواند کارایی سامانه تشخیص را بهبود بخشد [1]. به‌منظور جلوگیری از دست‌رفتن اطلاعات داده‌های ورودی، هدف، یافتن رویکردی است که الگوی مناسب‌تری برای داده‌های با طول متغیر و دسته‌بندی آن‌ها نمایش دهد؛ با استفاده از روش‌های جستجوی شباهت و امتیازدهی، روشی برای تجزیه و تحلیل، دسته‌بندی و آشکارسازی بدافزارها پیشنهاد می‌شود که براساس ماهیت متغیر بردارهای ویژگی عمل می‌کند. از آنجایی که کدهای

ارائه داده‌اند. آنها برای نخستین بار با استفاده از تجزیه و تحلیل ایستا، فراخوانی‌های سیستمی را از فایل‌های بدافزار استخراج کرده و از الگوریتم نایو بیز به عنوان تعیین‌کننده تقریبی ویروس استفاده کرده‌اند. در این پژوهش چهار گونه مختلف بدافزار مورد آزمایش قرار گرفته است. نتایج حاصل از روش پیشنهادی در مقایسه با هشت آنتی ویروس معرفی شده بهبود قابل توجهی را در شناسایی ویروس Win32 نشان می‌دهد که نشان از کارایی مناسب روش پیشنهادی دارد.

لندیچ و همکاران [5]، از سه روش استخراج ویژگی برای استخراج آپکد استفاده کرده‌اند. در روش انتخاب ویژگی، داده‌های بی‌ربط حذف شده و از الگوریتم‌های دسته‌بندی C4.5، Ibk، Ripper و K-نزدیک‌ترین همسایه استفاده شده است. خروجی دسته‌بندی‌های متعدد ترکیب شده و پیش‌بینی نهایی بر اساس رأی‌گیری حق و تو است. استراتژی تصمیم‌گیری حق و تو با معرفی رأی‌گیری حق و تو مبتنی بر اعتماد بهبود یافته است.

وانگ و همکاران [6]، نشان داده‌اند که ابزارهای بر اساس مدل پنهان مارکوف در آشکارسازی بدافزارهای دگرگون‌شده مؤثر هستند. آنها دویست بدافزار از خانواده دگرگون‌شده را جمع‌آوری کرده، نرخ تشخیص ۱۰۰٪ بوده است. در آشکارسازی بدافزارها استراتژی جدیدی از مدل مخفی مارکوف، معرفی شده که نتایج برتری نسبت به روش‌های مبتنی بر آستانه، به دست آورده است. روش مبتنی بر آستانه با استفاده از یک مدل پنهان مارکوف، یک فایل را به عنوان فایل آلوده معرفی می‌کند، اگر احتمال به دست آمده بیش از آستانه تعیین شده باشد. کلهر و همکاران [7]، نشان داده‌اند که چگونه رویکرد آستانه می‌تواند با این استراتژی جدید، برای کاهش سربار عملکرد ترکیب شود و در دقت و صحت نتایج هیچ کاهشی نباشد. آنها مدل پنهان مارکوف را برای چهار کامپایلر مختلف کد اسمبلی دست‌نوشته، سه کیت ساخت و ساز ویروس، و دو خانواده بدافزار دگرگون‌شده مورد بررسی قرار دادند. کاربرد پروفایل مدل پنهان مارکوف، که اطلاعات مربوط به موقعیت و حالات را در نظر می‌گیرد، می‌تواند برای آشکارسازی انواع خاصی از ویروس‌های دگرگون‌شده مؤثر باشد؛ اما برای ویروس‌هایی که بلاک‌های کد را دور از هم و با فاصله تغییر می‌دهند، به خوبی انجام نمی‌گیرد [8]. تحلیل مدل پنهان مارکوف برای دسته‌بندی بدافزارها نیز می‌تواند مؤثر واقع شود. *آناچاترا* و همکاران [9]، آزمایش خود را بر روی چهار کامپایلر مختلف پیاده‌سازی کرده و ۹۴۴۲ بدافزار گردآوری کردند. تعداد خوشه‌ها بین دو تا

اساس وزن‌دهی به هر یک از روش‌ها ترکیب شده و تصمیم نهایی بر اساس رأی اکثریت گرفته شده است. از میان ترکیبات ماشین بردار پشتیبان با وزن‌دهی متفاوت، زمانی که به روش گراف تشابه (روش ارائه‌شده) وزن بیشتری داده می‌شود، نتیجه بهتری می‌دهد؛ چون روش گراف تشابه، نسبت به دو روش دیگر معیارهای ارزیابی بالاتری دارد؛ سپس در بخش دوم با توجه به قدرت و مزایای هر دسته‌بند، از ترکیب خروجی دسته‌بندها استفاده شده است. این تنوع باعث شده است که از قدرت و مزایای سامانه‌های مختلف استفاده کند؛ بنابراین به عنوان روش پیشنهادی از آن استفاده شده است. روش پیشنهادی می‌تواند در دیگر حوزه‌هایی که با ویژگی‌های با طول متغیر سر و کار دارد به کار گرفته شود.

در ادامه در بخش ۲ به کارهای انجام‌شده اشاره می‌شود. معماری کامل روش پیشنهادی در بخش ۳ بیان شده است. در بخش ۴ شبیه‌سازی و نتایج آزمایش‌های انجام‌شده برای ارزیابی روش پیشنهادی مورد بررسی قرار می‌گیرد. در نهایت در بخش ۵ نتیجه‌گیری به عمل می‌آید.

۲- پیشینه پژوهش

در طول این سال‌ها، روش‌های زیادی برای آشکارسازی بدافزارها پیشنهاد شده است. *الذئاب* و همکاران [2]، از تجزیه و تحلیل n-گرام از محتوای دودویی، برای دسته‌بندی استفاده کرده‌اند. تجزیه و تحلیل اولیه، با استفاده از ماشین بردار پشتیبان^۱ توسط ارزش‌های مختلف n از ۱ تا ۵، صورت گرفته است. پایگاه داده آنها شامل ۲۴۲ بدافزار و ۷۲ فایل سالم بوده است و با دقت ۰/۹۶۵ بدافزارها را شناسایی کرده‌اند. *لین* و همکاران [3]، روشی عرضه کرده‌اند که ترکیبی از انتخاب و استخراج ویژگی است و به طور قابل توجهی ابعاد ویژگی‌های آموزش و دسته‌بندی را کاهش داده است. بر اساس رفتارهای بدافزارها که در محیط سندباکس جمع‌آوری شده، روش آنها از پنج مرحله تشکیل شده است: (۱) استخراج n-گرام ویژگی از فضای داده؛ (۲) ایجاد یک دسته‌بند ماشین بردار پشتیبان برای دسته‌بندی بدافزارها؛ (۳) انتخاب یک زیرمجموعه از ویژگی‌ها؛ (۴) تبدیل بردارهای ویژگی با بُعد بالا به بردار ویژگی با بعد کوچک‌تر و (۵) انتخاب مدل. آنها ۴۲۸۸ نمونه بدافزار را گردآوری کردند. دقت، حساسیت و F-measure روش ارائه‌شده به ترتیب با اندازه ۰/۷۹۶، ۰/۷۷۸۵ و ۰/۷۸۴۴ بوده است.

یو و همکاران [4]، یک روش آشکارسازی ویروس، بر اساس شناسایی دنباله فراخوانی سامانه‌ای تحت محیط ویندوز

^۱ Support Vector Machines (SVM)

قسمت‌های خطرناک در مدل‌های بدافزار دیده نشود، بالا می‌رود. اگرچه روش‌های هوشمند زیادی معرفی شده، اما توانایی‌های بدافزارها روزبه‌روز پیچیده‌تر شده است؛ به‌صورتی‌که هر روشی که برای آشکارسازی و شناسایی بدافزارها ابداع می‌شود، بدافزار با استفاده از استراتژی‌هایی جهت مخفی‌سازی کدهای خود، نوعی هوشمند ایجاد می‌کنند تا این روش‌ها قادر به آشکارسازی نباشند و بتوانند راه‌های امنیتی اخیر را کنار زنند و آنها را خنثی کنند. از آنجا که رویکردهای شناسایی و آشکارسازی بدافزارها نیز باید به‌روز شده و از هوشمندی برخوردار باشند، عرضه و توسعه روش‌های جدید و هوشمند امری ضروری است؛ لذا این مقاله، پژوهشی است بر روی یک الگوریتم جدید و هوشمند تا به‌کمک آن بتوان یک الگوی مناسب‌تری برای داده‌های با طول متغیر و دسته‌بندی آن‌ها ضمن افزایش نرخ تشخیص پیشنهاد دهد.

۳- روش پیشنهادی

در حوزه آشکارسازی بدافزارها، ویژگی‌های مختلف، بازنمایی‌های متفاوتی از فایل‌های اجرایی هستند که هر کدام حاوی یک نوع اطلاعات مفید در مورد آن فایل هستند. برای تشخیص یک الگوی مناسب از بدافزارها به‌طورمعمول به استخراج ویژگی‌های متفاوتی نیاز است. در این مقاله روشی مبتنی بر گراف به‌عنوان استخراج ویژگی برای دنباله‌های با طول متغیر پیشنهاد می‌شود. این روش با توجه به میزان شباهت ویژگی‌ها، بدون نیاز به ثابت‌کردن طول دنباله‌ها، یک امتیاز به هر نمونه می‌دهد و از این امتیازات جهت دسته‌بندی استفاده می‌شود. سه روش مدل مخفی مارکوف^۷، فاصله جایگزینی ساده^۸ و روش پیشنهادی گراف تشابه به‌عنوان استخراج ویژگی بررسی می‌شوند و از ترکیب نتایج آنها در ورودی استفاده خواهد شد.

بر خلاف کارهای پیشین که به‌طورمعمول از یک دسته ویژگی استفاده شده و فقط دو طبقه بدافزار و فایل‌های سالم وجود داشته است، در این مقاله از سه ویژگی دنباله‌های آپکد، هگزادسیمال و فراخوانی‌های سامانه‌ای در کنار یکدیگر استفاده خواهد شد و چهار طبقه مختلف داریم که سه طبقه آن مربوط به گونه‌های مختلف بدافزار است و سعی در تشخیص و تعیین دسته مطلوب برای هر یک از آنها هستیم. از آنجا که هر دسته‌بند تا حد خاصی قادر به تشخیص صحیح الگوها است و هیچ تک‌الگوریتم آموزشی مشخصی وجود ندارد که بتواند برای تمامی کاربردها و در تمامی شرایط بهترین و دقیق‌ترین باشد، لذا در رویکرد دوم، از ترکیب

پانزده مورد آزمایش قرار گرفت و با نه خوشه نتیجه بهتری به‌دست آمده است و با انحراف معیار ۰/۹۵ بدافزارها را خوشه‌بندی کرده‌اند. جاس و همکاران [10]، کاربرد روش‌های بی‌زین را در کنار مدل پنهان مارکوف برای آشکارسازی بدافزارها شرح داده‌اند. سینگ [11]، از مزایای روش‌های امتیازدهی با استفاده از مدل پنهان مارکوف استفاده و آن را با ماشین بردار پشتیبان ترکیب کرده است. مجموعه بدافزارها شامل درب پشتی از خانواده هربوت^۱، تروجان از خانواده پوشش امنیتی^۲، ویروس‌های لوح سخت^۳، خانواده متامورفیس^۴، ضدبات^۵، تروجان از زیرو اکسس^۶ بوده است. هر خانواده به‌صورت جداگانه در برابر فایل‌های سالم آموزش داده شده است. AUC برای هر پنج خانواده بدون استفاده از ویروس‌های لوح سخت، پوشش امنیتی، هربوت، ضدبات، زیرو اکسس با نسبت دگرگون‌کردن کد ۵۰٪ به ترتیب ۰/۷۹، ۰/۷۰، ۰/۹۵، ۱ و ۰/۹۳ بوده است.

استفاده از نتایج چند دسته‌بند با عنوان یادگیری دسته‌جمعی یک رویکرد مؤثر در یادگیری ماشین است که در آن به‌منظور بهبود دقت یادگیری نتایج دسته‌بندها با یکدیگر ترکیب شده و یک سامانه مرکب شکل می‌گیرد. قائمی و همکاران [12]، یک معماری جدید برای دسته‌بندی ترکیبی پرسش‌ها ارائه کرده‌اند. نتایج هر یک از دسته‌بندها توسط پنج روش رأی‌گیری وزن‌دار، فضای دانش رفتار، بی‌ز ساده، کلیشه تصمیم و دمپستر شفر ترکیب شده و خروجی نهایی را شکل می‌دهد. این روش ترکیبی متشکل از دو دسته‌بند مبتنی بر یادگیری ماشین (ماشین بردار پشتیبان و نمایش پراکنده) و یک دسته‌بند مبتنی بر قانون استفاده شده است. در پایان نتایج حاصل از دسته‌بندها با روش معمول در ترکیب دسته‌بندهای تک‌طبقه ترکیب شده‌اند و نتایج حاصل بیان‌کننده بهبود عملیات دسته‌بندی نسبت به روش‌های موجود است.

در نگاه نخست ویژگی‌های با طول ثابت، به‌دلیل سادگی و سرعت بالا و کاهش زمان، روشی کارا به‌نظر می‌آید، اما هیچ بدافزاری فقط از کدهای بدخواه ساخته نشده است. بدافزارها همیشه ترکیبی از قسمت‌های بی‌خطر و خطرناک هستند؛ حتی در بیش‌تر موارد قسمت‌های خطرناک نسبت به قسمت‌های بی‌خطر در اقلیت قرار دارند؛ درنتیجه، استفاده‌کردن از یک زیرمجموعه از کدها امکان این که

¹ Harebot

² Security Shield

³ Smart HDD

⁴ NGVCK

⁵ Zbot

⁶ ZeroAccess

⁷ Hidden Markov Model (HMM)

⁸ Simple Substitution Distance (SSD)

هر سطر در ماتریس، ماتریس احتمالی ایجاد می‌شود. برای هریک از دنباله‌های مجموعه آموزش این ماتریس ساخته شده و با استفاده از رابطه (۲) ماتریس E ایجاد می‌شود (m تعداد داده‌های مجموعه آموزش):

$$E = (\text{matrix}_{0ij} + \text{matrix}_{1ij} + \dots + \text{matrix}_{(m-1)ij}) / m \quad (1)$$

گراف مجموعه آموزش از روی این ماتریس ساخته می‌شود. برای هر یک از فایل‌های مجموعه آزمون، ماتریس D از روی ماتریس احتمالی دنباله مورد نظر ایجاد می‌شود و گراف تشابه D از روی این ماتریس برای هر دنباله ساخته می‌شود. در مرحله چهارم، امتیاز دنباله مورد نظر، با توجه به رابطه (۳) تعیین می‌شود:

$$\text{Score}(k) = d(D, E) = \frac{1}{N_2} \sum |E - D| \quad (2)$$

امتیاز حاصل بیان گر میزان شباهت دنباله‌ها به مجموعه آموزش است.

فرض کنید برای ویژگی آپکد، کلید مورد نظر به صورت زیر باشد:

MOV, CALL, ADD, XOR

حال به عنوان نمونه، اگر دنباله آپکد برای امتیازدهی

به صورت زیر باشد:

JMP, MOV, MOV, ADD, INC, INC, INC

با مقایسه کلید و دنباله مورد نظر، کد میانی به صورت

زیر بازنویسی می‌شود:

Other, MOV, MOV, ADD, Other, Other, Other

از روی این کد، ماتریس توزیع گراف که حاوی تعداد

وقوع هر آپکد به دنبال آپکد دیگر در فایل داده شده است،

(جدول ۱) ایجاد می‌شود؛ سپس با تقسیم هر مقدار به مجموع

هر سطر در ماتریس، ماتریس احتمالی (جدول ۲) ایجاد شده

و گراف تشابه آپکد از روی این جدول ساخته می‌شود

(شکل ۱).

(جدول ۱-۱): ماتریس توزیع گراف

(Table-1): digraph distribution matrix

Other	XOR	ADD	CALL	MOV	
2	0	0	0	1	MOV
0	0	0	0	0	CALL
0	0	0	0	1	ADD
0	0	0	0	0	XOR
2	0	3	0	0	Other

دسته‌بندی در خروجی استفاده می‌شود. انگیزه اصلی بر توسعه چنین روشی، کاهش نرخ خطا است. انواع ترکیب دسته‌بندی کننده‌ها شامل ترکیب ایستا^۱ و ترکیب پویا هستند. ترکیب ایستا شامل سه روش میانگین گروه^۲، بگینگ^۳ و بوستینگ^۴ است [1].

۱-۳- روش پیشنهادی مبتنی بر گراف

در این بخش، روش پیشنهادی مبتنی بر گراف (شبه‌کد) برای دسته‌بندی داده‌های با طول متغیر شرح داده می‌شود. این روش شامل چهار مرحله تعیین کلید، ایجاد کد میانی، ایجاد گراف تشابه از روی ماتریس توزیع و تعیین امتیاز است.

Procedure Graph Similarity

Begin

1: Create key

2: for each sample do

3: Create middle code based on key

4: End for

5: for each sample in train set do

6: Built e matrix on their middle code

7: End for

8: $E = (\text{matrix } e_{0ij} + \text{matrix } e_{1ij} + \dots + \text{matrix } e_{(m-1)ij}) / \text{number of samples on train set}$

9: Built E graph based on E matrix

10: for each file in test set do

11: Built D matrix for it

12: Built D graph based on D matrix

13: Calculate score for it

14: End for

15: Return scores

16: End procedure

(شبه‌کد ۱): روش پیشنهادی مبتنی بر گراف تشابه

(Pseudo Code 1): the proposed method based on the graph

مرحله نخست، در بین داده‌های آموزشی پرتکرارترین

دستورها را پیدا کرده و n تا، به عنوان کلید انتخاب می‌شوند.

در مرحله دوم، تمام داده‌ها بر اساس کلید بازنویسی

می‌شوند؛ هر یک از دستورها، اگر متعلق به کلید بودند آن را

نوشته در غیر این صورت از واژه "other" استفاده خواهد شد.

بدین ترتیب یک کد میانی ایجاد خواهد شد که از لحاظ

سرعت و حافظه صرفه جویی خواهد شد.

در مرحله سوم، از روی کد میانی ماتریسی ایجاد

می‌شود که حاوی تعداد وقوع هر دستور به دنبال دستور دیگر

در دنباله داده شده است؛ سپس با تقسیم هر مقدار به مجموع

¹ Static Structure

² Ensemble Averaging

³ Bagging

⁴ Boosting

۲-۴- معیارهای ارزیابی

از آنجا که سامانه تشخیص، چهار طبقه است، معیارهای ارزیابی سامانه پیشنهادی بر اساس سامانه چندطبقه است [13,14].

تعاریف مورد نیاز برای معیارهای ارزیابی به صورت زیر هستند:
 True Positive) Tp - تعداد بدافزارهایی که مشکوک شناخته شده‌اند.

True Negative) Tn - تعداد برنامه‌های سالمی که سالم شناخته شده‌اند.

False Positive) Fp - تعداد برنامه‌های سالمی که مشکوک شناخته شده‌اند.

False Negative) Fn - تعداد بدافزارهایی که سالم شناخته شده‌اند.

مقدار صحت سامانه^۲ میزان نزدیکی آن معیار به اندازه حقیقی سامانه گفته می‌شود.

$$\text{Average Accuracy} = \frac{\sum_{i=1}^l \frac{tp_i + tn_i}{tp_i + fp_i + tn_i + fn_i}}{l} \quad (3)$$

معیار دقت^۳ میزانی است که نشان می‌دهد اگر شرایط عوض نشود سامانه تا چه حدی همان نتیجه را به دست خواهد آورد و براساس فرمول (۵) محاسبه می‌شود:

$$\text{Precision } \mu = \frac{\sum_{i=1}^l TP_i}{\sum_{i=1}^l TP_i + FP_i} \quad (4)$$

معیار شفافیت یا اختصاصی بودن^۴ نشان‌دهنده نسبتی از حالت منفی است که درست تشخیص داده شده‌اند.

$$\text{Specificity } \mu = \frac{\sum_{i=1}^l TN_i}{\sum_{i=1}^l TN_i + FP_i} \quad (5)$$

معیار حساسیت^۵ نشان‌دهنده نسبتی از حالت مثبت است که درست شناخته شده‌اند.

$$\text{Sensitivity } \mu (\text{Recall}) = \frac{\sum_{i=1}^l TP_i}{\sum_{i=1}^l TP_i + FN_i} \quad (6)$$

F-measure دو معیار دقت و حساسیت را هم‌زمان مورد بررسی قرار می‌دهد. از این رو این معیار را به عنوان مجموع وزن دار این دو کمیت نیز می‌توان در نظر گرفت. درحقیقت به این معیار می‌توان به عنوان میانگین هارمونیک بین این دو معیار دقت و حساسیت نگاه کرد. بهترین پاسخ این معیار یک و بدترین پاسخ، صفر است.

$$F1_{\mu} = 2 \cdot \frac{\text{Precision}_{\mu} * \text{Recall}_{\mu}}{\text{Precision}_{\mu} + \text{Recall}_{\mu}} \quad (7)$$

² Accuracy

³ Precision (Positive Predictive Value)

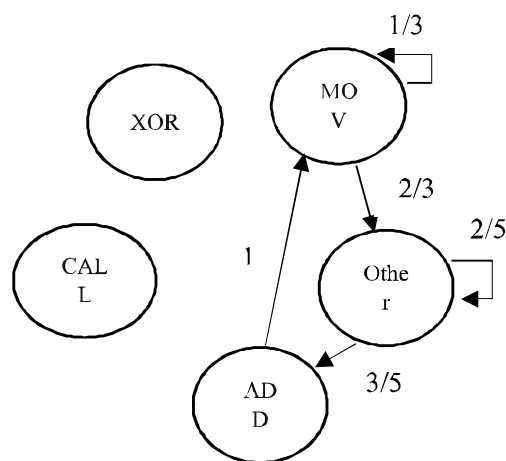
⁴ Specificity

⁵ Sensitivity

(جدول-۲): ماتریس احتمالی جدول (۱)

(Table-2): Probabilities matrix from Table 1

Other	XOR	ADD	CALL	MOV	
2/3	0	0	0	1/3	MOV
0	0	0	0	0	CALL
0	0	0	0	1	ADD
0	0	0	0	0	XOR
2/5	0	3/5	0	0	Other



(شکل-۱): گراف پیشنهادی بر روی آپکد

(Figure-1): The proposed graph for opcode

۴- شبیه‌سازی

در این بخش مجموعه داده، معیارهای ارزیابی و نتایج آزمایش‌های انجام شده برای ارزیابی کارایی روش پیشنهادی شرح داده می‌شود.

۴-۱- مجموعه داده

مجموعه بدافزارهای آزمایش‌ها از مجموعه VX Heaven Virus Collection^۱ گردآوری شده و شامل سه گونه مختلف ویروس‌ها، تروجان و درب پستی است. برنامه پاک و بی‌خطر از فایل‌های اجرایی ویندوز از پوشه Program Files یا System32 جمع آوری شده است. بدین ترتیب چهار طبقه مختلف شامل فایل‌های سالم، ویروس، تروجان و درب پستی (از هر طبقه دویست نمونه) وجود دارد.

استخراج ویژگی مبتنی بر محتوای اسمبلی فایل‌ها صورت گرفته است. سه ویژگی آپکد، هگزا دسیمال و فراخوانی سامانه‌ای برای هر نمونه از روی کد اسمبلی آنها استخراج خواهد شد.

¹ <http://vxheaven.org>

(جدول-۳): تأثیر مقادیر مختلف کلید بر کارایی روش گراف

تشابه آپکد

(Table-3): The effect of different values of key on the performance of proposed Opcode graph

F1	شفافیت	حساسیت	دقت	صحت	
0.57	0.86	0.57	0.57	0.78	N=10
0.68	0.89	0.68	0.68	0.84	N=25
0.52	0.84	0.52	0.52	0.76	N=50
0.68	0.89	0.68	0.68	0.84	N=100
0.66	0.89	0.66	0.66	0.83	N=120
0.65	0.88	0.65	0.65	0.83	N=150
0.56	0.85	0.56	0.56	0.78	N=200

(جدول-۴): تأثیر مقادیر مختلف کلید بر کارایی روش گراف

پیشنهادی هگزا دسیمال

(Table-4): The effect of different values of key on the performance of proposed Hexadecimal graph

F1	شفافیت	حساسیت	دقت	صحت	
0.61	0.87	0.61	0.61	0.80	N=10
0.48	0.83	0.48	0.48	0.74	N=25
0.50	0.83	0.50	0.50	0.75	N=50
0.46	0.82	0.46	0.46	0.73	N=100
0.41	0.80	0.41	0.41	0.71	N=120
0.48	0.83	0.48	0.48	0.74	N=150
0.59	0.86	0.59	0.59	0.79	N=200

(جدول-۵): تأثیر مقادیر مختلف کلید بر کارایی روش گراف

پیشنهادی فراخوانی سیستمی

(Table-5): The effect of different values of key on the performance of proposed system call graph

F1	شفافیت	حساسیت	دقت	صحت	
0.59	0.86	0.59	0.59	0.79	N=10
0.61	0.87	0.61	0.61	0.81	N=25
0.58	0.86	0.58	0.58	0.79	N=50
0.63	0.88	0.63	0.63	0.81	N=100
0.61	0.87	0.61	0.61	0.80	N=120
0.62	0.87	0.62	0.62	0.81	N=150
0.53	0.84	0.53	0.53	0.77	N=200

با توجه به جدول (۵)، مقادیر ۲۵ و ۱۰۰ برای ویژگی فراخوانی سامانه‌ای نتایج بهتری داشته است؛ اما ارزش ۱۰۰ از صحت و حساسیت بهتری برخوردار است. با توجه به این که فراخوانی سامانه‌ای می‌تواند رفتار یک قطعه کد بدافزار را نشان دهد، برای تمایز بهتر بین گونه‌های مختلف بدافزارها،

برای هر آزمایش، از اعتبارسنجی ضرب‌دوری با پنج فولد استفاده می‌شود. در این نوع اعتبارسنجی داده‌ها به پنج زیرمجموعه افزای می‌شوند. از این پنج زیرمجموعه، هر بار یکی برای اعتبارسنجی و چهارتای دیگر برای آموزش به کار می‌روند. این روال پنج بار تکرار می‌شود و همه داده‌ها به‌طور دقیق یکبار برای آموزش و یکبار برای اعتبارسنجی به کار می‌روند. در نهایت میانگین نتیجه این پنج بار اعتبارسنجی به‌عنوان یک تخمین نهایی برگزیده می‌شود.

۳-۴- نتایج

تابع امتیازدهی و اندازه کلید در روش پیشنهادی گراف تشابه، پارامترهایی هستند که باید مقدار بهینه آنها تعیین شود. تغییرات در هر یک از این پارامترها اثر قابل توجهی بر عملکرد الگوریتم خواهد داشت.

مقادیر مختلف اندازه کلید برای هر یک از ویژگی‌های فراخوانی سامانه‌ای، آپکد و هگزا دسیمال به‌صورت جداگانه مورد آزمایش قرار گرفته است. دنباله‌های موجود، حاوی تعداد زیادی آپکد، دستورهای هگزا دسیمال و تابع فراخوانی سامانه‌ای مجزا هستند؛ از این دستورها، تنها تعداد خاصی از آنها به‌طور مداوم ظاهر می‌شوند و در طبیعت ذاتی بدافزار یا فابل خوش‌خیم سهیم هستند. در جداول (۳ و ۴) تأثیر مقادیر مختلف N بر کارایی روش گراف تشابه نمایش داده شده است.

با توجه به جدول (۳)، مقادیر ۲۵ و ۱۰۰ برای آپکد نتایج بهتری داشته است. برای اینکه اندازه ماتریس بزرگ نشود و مقایسه آسان‌تر شود، مقدار ارزش ۲۵ برای کلید در نظر گرفته می‌شود. از این‌رو، ماتریس D و E برای آپکد در اندازه ۲۶*۲۶ ساخته شده است. این تعجب‌آور نیست که تعداد به‌نسبه کمی از دستورها نتایج بهتری از مقادیر بزرگ‌تر دارند، زیرا پرتکرارترین آپکدها می‌توانند رفتار کل آپکدها در نمونه‌های بدافزار را توصیف کنند.

با توجه به نتایج جدول (۴)، مقدار ده برای ویژگی هگزا دسیمال نتایج بهتری داشته است. از این‌رو، ماتریس D و E برای ویژگی هگزا دسیمال در اندازه ۱۱*۱۱ ساخته شده است. با توجه به ماهیت دستورهای هگزا دسیمال که اطلاعات خام برای شیوه‌نامه مورد نظر هستند، تعداد کمی از آنها به‌عنوان پرتکرارترین دستورها می‌تواند کل دستورهای هگزا دسیمال را توصیف کند.

(جدول-۷): ترکیب سه روش امتیازدهی بر روی هر سه ویژگی
(Table-7): The combination of three scoring functions on all three features

صحت	دقت	حساسیت	شفافیت	F1	خطا
0.93	0.86	0.86	0.95	0.86	0.07

پس از ترکیب سه روش امتیازدهی بر روی هر سه ویژگی، نسبت به هر یک از روش‌های امتیازدهی به صورت فردی نتایج پیشرفت چشم‌گیری داشته است؛ لذا در سامانه پیشنهادی از تمام ترکیبات ویژگی استفاده شده است. بدین صورت می‌توان از مزایای هر روش و هر ویژگی استفاده کرد و به تشخیص بدافزارهای ناشناخته کمک کرد.

در رویکرد دوم، ترکیب دسته‌بندها در خروجی مورد آزمایش قرار گرفته است. در بخش نخست، امتیازات حاصل از هر یک از روش‌های مدل مخفی مارکوف، فاصله جایگزینی ساده و گراف تشابه به سه دسته‌بند ماشین بردار پشتیبان داده شده است. خروجی‌ها بر اساس وزن‌دهی به هر یک از روش‌ها ترکیب شده و تصمیم نهایی بر اساس رأی اکثریت گرفته شده است. نتایج حاصل از این ترکیبات در جدول (۸) آورده شده است. سطر نخست این جدول مربوط به حالتی است که به روش گراف تشابه وزن سه و به دو رویکرد دیگر وزن یک داده شده است. سطر دوم، به روش فاصله جایگزینی ساده وزن سه و به دو روش دیگر وزن یک و سطر سوم به مدل مخفی مارکوف وزن سه و به دو رویکرد دیگر وزن یک داده شده است. سطر چهارم حالتی است که به هر سه روش وزن یکسان داده است.

(جدول-۸): ترکیب نتایج ماشین بردار پشتیبان
(Table-8): Combine results support vector machine

	صحت	دقت	حساسیت	شفافیت	F1
1	0.94	0.88	0.88	0.96	0.88
2	0.91	0.81	0.81	0.94	0.81
3	0.89	0.78	0.78	0.93	0.78
4	0.92	0.84	0.84	0.95	0.84

از میان ترکیبات ماشین بردار پشتیبان با وزن‌دهی متفاوت، زمانی که به روش گراف تشابه وزن بیشتری داده می‌شود، نتیجه بهتری می‌دهد؛ چون روش گراف تشابه نسبت به دو روش دیگر معیارهای ارزیابی بالاتری دارد. در بخش دوم با توجه به قدرت و مزایای هر دسته‌بند، از ترکیب خروجی دسته‌بندها استفاده شده است. سه روش ترکیب دسته‌بندها شامل میانگین گروه، روش بگینگ و آدابوست مورد آزمایش قرار گرفته است (جدول ۹).

ماتریس D و E برای ویژگی فراخوانی سامانه‌ای در اندازه ۱۰۰*۱۰۰ ساخته شده است.

مقدار بهینه تابع امتیازدهی در روش گراف تشابه نیز باید تعیین شود. به منظور تعیین تابع امتیازدهی مؤثر برای ارزیابی روش گراف تشابه، چهار تابع به صورت جداگانه در روش گراف تشابه اعمال شده است.

$$\text{Score1: } \sum_{i,j} |dij - eij| \quad (۸)$$

$$\text{Score2: } \frac{1}{n^2} \sum_{i,j} |dij - eij|^2 \quad (۹)$$

$$\text{Score3: } \frac{1}{n^2} \sum_{i,j} |dij - eij| \quad (۱۰)$$

$$\text{Score4: } \sum_{i,j} |dij - eij|^2 \quad (۱۱)$$

نتایج حاصل از اعمال هر یک از این توابع در جدول (۶) آمده است.

(جدول-۶): مقایسه توابع امتیازدهی در روش گراف پیشنهادی
(Table-6): Compare Scoring functions in the proposed graph method

تابع	صحت	دقت	حساسیت	شفافیت	F1
Score1	0.87	0.74	0.74	0.91	0.74
Score2	0.89	0.79	0.79	0.93	0.79
Score3	0.87	0.74	0.74	0.91	0.74
Score4	0.89	0.78	0.78	0.93	0.78

از آنجا که فاصله اقلیدسی فاصله مستقیم بین دو نقطه است که در حقیقت اندازه کوتاه‌ترین خط بین دو نقطه است، با توجه به نتایج به دست آمده، مربع فاصله اقلیدسی دو گراف تشابه (تابع ۲ و تابع ۴) نتایج بهتری نسبت به دیگر توابع داشته‌اند. تابع ۲ با توجه به اینکه تعداد گره‌ها در آن دخیل است با تفاوت اندکی نتایج بهتری نسبت به تابع ۴ در بر دارد. بنابراین از این تابع در روش گراف تشابه استفاده شده است. با توجه به امتیازات به دست آمده، یک الگوی جدید از مجموعه داده ایجاد شده است که به جای استفاده مستقیم از داده‌ها و ویژگی‌های استخراج شده با طول متغیر و حجم زیاد، می‌توان از این الگو برای ورودی دسته‌بندها استفاده کرد. در این مقاله از دسته‌بند ماشین بردار پشتیبان با روش یکی در مقابل یکی برای تعیین دسته مطلوب نمونه‌ها استفاده شده است.

در رویکرد نخست، هر سه روش‌های امتیازدهی شامل مدل مخفی مارکوف، فاصله جایگزینی ساده و گراف تشابه بر روی هر سه ویژگی اعمال و از ترکیب نتایج آنها در ورودی دسته‌بندها استفاده شده است (جدول ۷).

رویکرد پیشنهادی با صحت ۰/۹۷ و نرخ خطای تشخیص ۰/۰۳ بدافزارهای ناشناخته را شناسایی کرد، درحالی‌که روش فاصله جانشینی ساده، مدل مخفی مارکوف و آزمون خی دوی تحت شرایط یکسان با دقت ۰/۸۴، ۰/۸۷ و ۰/۸۷ و نرخ خطای تشخیص ۰/۱۸، ۰/۲۰ و ۰/۱۴ توانستند بدافزارها را شناسایی کنند. رویکرد پیشنهادی دارای عملکرد بهتری نسبت به دیگر روش‌های پایه دارد.

۵- نتیجه‌گیری

در این مقاله، یک روش پیشنهادی مبتنی بر گراف برای دسته‌بندی دنباله‌های با طول متغیر بدون ثابت کردن طول دنباله‌ها در حوزه آشکارسازی بدافزارها پیشنهاد شد. روش پیشنهادی شامل چهار مرحله تعیین کلید، ایجاد کد میانی، ایجاد گراف تشابه E و D از روی ماتریس توزیع و تعیین امتیاز است. روش پیشنهادی با تعیین پرتکرارترین دستورها و گذاشتن باقی دستورها در مجموعه 'other' از لحاظ سرعت و حافظه صرفه‌جویی کرده است. آزمایش‌ها نشان می‌دهند گراف تشابه پیشنهادی، ضمن غلبه بر مشکلات روش‌های گراف سنتی در مواجهه با داده‌های با طول متغیر، بهتر از روش‌های پایه همچون مدل مخفی مارکوف، روش فاصله جایگزینی ساده و آزمون خی دوی است. همچنین در این مقاله، از ترکیب دسته‌بندی در آشکارسازی بدافزارها با استفاده از دسته‌بندی دنباله‌های با طول متغیر استفاده شده است. روش پیشنهادی شامل دو رویکرد بود؛ در رویکرد نخست، روش‌های استخراج ویژگی شامل سه روش مدل مخفی مارکوف، فاصله جایگزینی ساده و روش گراف تشابه در ورودی ترکیب شدند؛ درحالی‌که در رویکرد دوم، از ترکیب دسته‌بندی در خروجی استفاده شد. در ترکیب روش‌های استخراج ویژگی، نسبت به هر یک از روش‌های امتیازدهی به صورت فردی نتایج پیشرفت چشم‌گیری داشته است. از میان ترکیبات دسته‌بندی در خروجی، میانگین گروه با ترکیب خروجی چهار دسته‌بند K-نزدیک‌ترین همسایه، جنگل تصادفی، ماشین بردار پشتیبان و نایو بیز با دقت ۰/۹۷ بالاترین نرخ تشخیص و کم‌ترین نرخ خطای تشخیص را داشت. در رویکرد پیشنهادی گراف می‌توان اشاره کرد که تغییرات جزئی در توابع امتیازدهی اثرات بسیار زیادی در نتایج در بر خواهد داشت. برای کشف دقت بیشتر، می‌توان به دنبال یک تابع امتیازدهی قوی‌تری بود. همچنین با استفاده از گراف روند کنترلی، در کنار گراف فراخوانی سامانه‌ای می‌توان گراف قدرتمندتری به دست آورد. روش پیشنهادی می‌تواند در دیگر حوزه‌هایی که با ویژگی‌های با طول متغیر سر و کار دارد به کار گرفته شود.

(جدول-۹): ترکیب دسته‌بندی‌ها

(Table-9): Combine Classifications

خطا	F1	شفافیت	حساسیت	دقت	صحت	
0.03	0.94	0.98	0.94	0.94	0.97	میانگین گروه
0.07	0.86	0.95	0.86	0.86	0.93	بگینگ
0.16	0.68	0.89	0.68	0.68	0.84	آدابوست

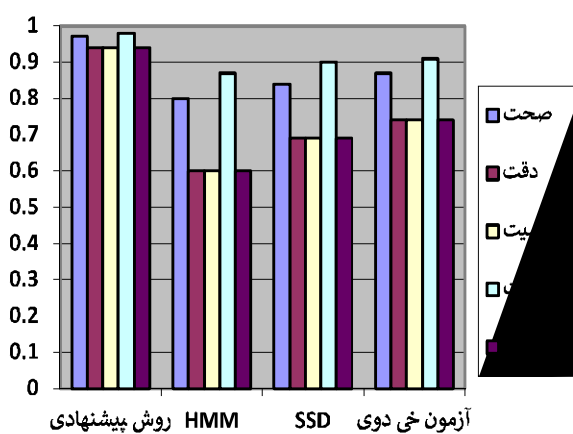
میانگین گروه از ترکیب چهار دسته‌بند جنگل تصادفی، ماشین بردار پشتیبان، K-نزدیک‌ترین همسایه و نایو بیز ایجاد شده و تصمیم بر اساس رأی اکثریت گرفته شده است؛ این تنوع باعث شده که از قدرت و مزایای سامانه‌های مختلف استفاده کند. بنابراین به عنوان روش پیشنهادی از آن استفاده شده است.

در جدول (۱۰) و شکل (۲) کارایی روش پیشنهادی با روش‌های مرسوم که در مقالات دیگر به دفعات مورد استفاده قرار گرفته، مقایسه شده است.

(جدول-۱۰): مقایسه روش پیشنهادی با روش‌های استفاده‌شده

(Table-10): Comparison of proposed methods with the methods used

خطا	F1	شفافیت	حساسیت	دقت	صحت	
0.03	0.94	0.98	0.94	0.94	0.97	روش پیشنهادی
0.18	0.60	0.87	0.60	0.60	0.80	HMM
0.20	0.69	0.90	0.69	0.69	0.84	SSD
0.14	0.74	0.91	0.74	0.74	0.87	آزمون خی دوی



(شکل-۲): مقایسه روش پیشنهادی با روش‌های استفاده‌شده

(Figure-2): Comparison of proposed methods with the methods used

[13] V. Asch, "Macro- and micro-averaged evaluation measures [[BASIC DRAFT]]," university of Antwerp, 2013.

[14] M. Sokolova and G. Lapalme, "A systematic analysis of performance measures for classification tasks," *Information Processing and Management*, vol. 45, pp. 427-437, 2009.



فاطمه حسینی مدرک کارشناسی خود

را در رشته مهندسی کامپیوتر-نرم افزار از دانشگاه آزاد واحد ارسنجان در سال ۱۳۸۸ دریافت کرد. در حال حاضر ایشان به عنوان دانشجوی کارشناسی ارشد رشته

هوش مصنوعی در واحد علوم و تحقیقات تهران مشغول به تحصیل هستند. پژوهش های ایشان بر روی امنیت اطلاعات و شبکه، داده کاوی و یادگیری ماشین متمرکز است.

نشانی رایانامه ایشان عبارت است از:

fatima.hosseini@srbiau.ac.ir



میترا میرزازاده مدرک کارشناسی خود

را در رشته مهندسی کامپیوتر-نرم افزار از دانشگاه آزاد واحد تهران مرکزی و همچنین مدارک کارشناسی ارشد و دکتری خود را در رشته مهندسی کامپیوتر-هوش مصنوعی از

دانشگاه آزاد اسلامی واحد علوم و تحقیقات دریافت کرد. در حال حاضر ایشان عضو هیأت علمی تمام وقت واحد علوم و تحقیقات هستند. زمینه های پژوهشی مورد علاقه ایشان شناسایی الگو و کاربردهای آن است.

نشانی رایانامه ایشان عبارت است از:

mirzarezaee@srbiau.ac.ir



آرش شریفی مدرک کارشناسی خود را در

رشته مهندسی کامپیوتر-سخت افزار از دانشگاه آزاد واحد تهران جنوب در سال ۱۳۸۳ دریافت و همچنین مدارک کارشناسی ارشد و دکتری خود را در رشته مهندسی کامپیوتر-هوش مصنوعی در

سال های ۱۳۸۶ و ۱۳۹۰ از دانشگاه آزاد واحد علوم و تحقیقات دریافت کرد. در حال حاضر ایشان عضو هیأت علمی تمام وقت واحد علوم و تحقیقات هستند. زمینه های پژوهشی مورد علاقه ایشان عبارتند از: محاسبات نرم، یادگیری عمیق و پردازش تصویر.

نشانی رایانامه ایشان عبارت است از:

a.sharifi@srbiau.ac.ir

- [1] J. Quinlan, "Bagging, Boosting and C4.5," 2006.
- [2] M. Alazab, R. Layton, S. Venkataraman and P. Watters, "Malware Detection Based on Structural and Behavioural Features of API Calls," Perth, WA, 2010.
- [3] C. T. Lin, N.-J. Wang, H. Xiao and C. Eckert, "Feature Selection and Extraction for Malware Classification," *Journal of Information Science and Engineering* 31, vol. 31, no. 3, pp. 965-992, 2015.
- [4] J. Xu, A. H. Sung, S. Mukkamala, and Q. Liu, "Obfuscated Malicious Executable Scanner," *Journal of Research and Practice in Information Technology*, vol. 39, pp. 181-197, 2007.
- [5] M. J. Landage and P. M.P. Wankhade, "Malware Detection with Different Voting Schemes," *COMPUSOFT, An international journal of advanced computer technology*, vol. 3, no. 1, pp. 450-456, 2014.
- [6] W. Wong and M. Stamp, "Hunting for metamorphic engines," *Journal in Computer Virology*, vol. 2, no. 3, pp. 211-229, 2006.
- [7] T. Kalbhor, "Dueling hidden Markov models for virus analysis," *Journal of Computer Virology and Hacking Techniques*, vol. 11, no. 2, pp. 103-118, 2015.
- [8] S. Attaluri, S. McGhee and M. Stamp, "Profile hidden Markov models and metamorphic virus detection," *Journal in Computer Virology*, vol. 5, no. 2, pp. 151-169, 2009.
- [9] C. Annachhatre and M. Stamp, "Hidden Markov models for malware classification," *Journal of Computer Virology and Hacking Techniques*, vol. 11, no. 2, pp. 59-73, 2015.
- [10] S. Josse and E. Filiol, "New Trends in Security Evaluation of Bayesian Network-Based Malware Detection Models," Maui, Hawaii USA, 2012.
- [11] T. Singh, F. D. Troia, V. A. Corrado, T. H. Austin and M. Stamp, "Support vector machines and malware detection," *Journal of Computer Virology and Hacking Techniques*, vol. 12, no. 4, pp. 203-212, 2016.
- [12] H. Ghaemi and M. Kahani, "Question classification using ensemble classifiers," *Quarterly Journal Signal and Data Processing*, vol. 29, number.3, pp.99, 1395.

[۱۲] هادی قائمی و محسن کاهانی، "دسته بندی

پرسش ها با استفاده از ترکیب دسته بندی ها،

فصل نامه علمی-پژوهشی پردازش علائم و داده ها،

شماره ۳، پیاپی ۲۹، صفحه ۹۹، سال ۹۵.