

تاریخ دریافت مقاله: ۹۵/۰۷/۰۵  
تاریخ پذیرش مقاله: ۹۵/۱۰/۲۴

## روشی جدید جهت شناسایی بدافزارهای فراریخت با تحلیل ایستای ثباتها و کدهای عملیاتی

هادی گلباگی

دانشجوی کارشناسی ارشد دانشکده مهندسی و علوم رایانه - دانشگاه شهید بهشتی - تهران - ایران  
پست الکترونیکی: h.golbaghi@mail.sbu.ac.ir

مجتبی وحیدی اصل\*

استادیار دانشکده مهندسی و علوم رایانه - دانشگاه شهید بهشتی - تهران - ایران  
پست الکترونیکی: mo\_vahidi@sbu.ac.ir

علیرضا خلیلیان

دانشجوی دکتری دانشکده مهندسی رایانه - دانشگاه اصفهان - اصفهان - ایران  
پست الکترونیکی: khalilian@eng.ui.ac.ir

### چکیده

بین خانواده‌های مختلف بدافزارها باشد. به منظور ارزیابی روش پیشنهادی، آزمایش‌هایی بر روی ۴۵۰ فایل متشکل از فایل‌های سالم و بدافزار فراریخت از ویروس‌های خانواده‌های G2, MP-2, CGEN, NGVCK, VLC و کرم‌های خانواده MWOR انجام شد. نتایج آزمایش‌ها نشان می‌دهد که روش پیشنهادی در چهار الگوریتم دسته‌بندی و در تمامی خانواده‌های بدافزارها به دقت شناسایی ۱۰۰ درصدی دست یافته است که با معیار ROC اندازه‌گیری شده است.

واژه‌های کلیدی: بدافزار، فراریختی، مبهم‌سازی، کد، موتورهای تکثیر، تحلیل ایستا

شناسایی بدافزارهای فراریخت مسئله دشواری است، زیرا در هر انتشار ساختار کد تغییر می‌یابد درحالی‌که عملکرد و رفتار بدافزار ثابت می‌ماند. تاکنون روش‌های مختلفی برای شناسایی بدافزارهای فراریخت پیشنهاد شده‌اند. گاه روش‌ها با جایگزین کردن دستورات مشابه یا درج کد زائد در بدنه بدافزار دچار شکست می‌شوند و در برخی موارد سربار محاسباتی بالا، دقت شناسایی کم و کارایی ضعیف روش‌ها، آن‌ها را دچار چالش می‌کند. در این مقاله روشی جدید برای شناسایی بدافزارهای فراریخت پیشنهاد می‌شود که با تحلیل ایستای کدهای عملیاتی و ثبات‌های مورد استفاده کار می‌کند. این ایده می‌تواند اساس تمایز

\* نویسنده مسئول

یک امضای<sup>۸</sup> متفاوت ایجاد می‌شود که از شناسایی توسط ضد بدافزارها فرار کند [۴]. علاوه بر این، انواع دیگری از بدافزارها مانند اسب تروا<sup>۹</sup> و کرم‌هایی<sup>۱۰</sup> مبتنی بر همین فنون فراریختی وجود دارند [۴]. در یک دسته‌بندی کلی می‌توان روش‌های کشف بدافزار را به دو دسته ایستا<sup>۱۱</sup> و پویا<sup>۱۲</sup> تقسیم کرد. روش‌های ایستا بسیار محبوب و رایجند زیرا بدون اجرای بدافزار و از تحلیل ساختار کد آن، قادرند آن‌ها را کشف نمایند. به‌طور معمول، در روش ایستا برای شناسایی بدافزارها، روش‌ها و ابزار تشخیص با استفاده از روش‌های مبتنی بر امضا به‌طور گسترده‌ای استفاده می‌شود. ساختارهای رشته‌ای کوتاه از بایت‌ها که منحصربه‌فرد بوده، الگویی از آن‌ها استخراج می‌شود، سپس امضاها در بانک‌های اطلاعاتی نگهداری می‌شوند و عملیات شناسایی مبتنی بر آن‌ها صورت می‌گیرد [۳]. در مقابل، روش‌های پویا با اجرای بدافزار می‌توانند آن‌ها را کشف کنند که این کار معمولاً نیاز به تقلید<sup>۱۳</sup> دارد [۶].

روش‌های متعددی بر مبنای تحلیل ایستای بدافزار کار می‌کنند. یکی از این روش‌ها شناسایی بر اساس خانواده بدافزارها [۶] است که این روش فقط محدود به فایل‌های اجرایی قابل حمل می‌باشد. همچنین در مقاله [۷] بر روی ترتیب دستورات اجرایی و گراف جریان کنترلی در کدهای اسمبلی کار شده که با استفاده از فنون مبهم‌سازی کد که جریان کنترلی دستورات را تغییر می‌دهند، روش دچار چالش خواهد شد. روشی دیگر استفاده از هیستوگرام<sup>۱۴</sup> فراوانی کدهای عملیاتی [۸] است که از طریق درج کد زائد و جانشینی دستورات مشابه دچار شکست می‌شود. روش شباهت گراف کدهای عملیاتی [۹] روش موثری است اما اگر از روش جانشینی دستورات مشابه استفاده شود دچار چالش خواهد شد. همچنین در مقاله‌های [۴، ۱] کارهایی بر مبنای کدهای عملیاتی در زبان اسمبلی صورت

بر مبنای گزارش‌های سالیان اخیر بدافزارها به‌طور روزافزونی تهاجمی‌تر شده‌اند. تحلیلگران شرکت سمانتک<sup>۱</sup> نشان داده‌اند که بدافزاری به اسم Reveton<sup>۲</sup> حدود ۵۰۰ هزار رایانه را در عرض حدود ۱۸ روز آلوده کرده است [۱]. به‌طور متوسط در هر حادثه جرایم سایبری ۱۹۷ دلار از دست می‌رود [۱]. در سال ۲۰۱۳ طبق تخمین‌ها حدود ۵۵۶ میلیون کاربر در سرتاسر جهان یک نمونه از جرایم سایبری را تجربه کرده‌اند. در سال ۲۰۱۱ مشتری‌ها بیشتر از ۴۹ هزار گزارش مختلف از تهدیدات خانواده‌های مختلف بدافزارها را به مرکز حفاظت از بدافزارهای مایکروسافت<sup>۳</sup> داده‌اند. تعدادی از این گزارش‌ها تهدیداتی از جنس خانواده‌های بدافزارهای چندریختی<sup>۴</sup> و فراریختی<sup>۵</sup> بوده است [۱]. در وضعیت کنونی بدافزارهای مختلف رشد بسیار زیادی کرده‌اند و باعث اختلالات اساسی در سامانه‌های حساس<sup>۶</sup> و حیاتی مانند اقتصادی، مالی، نظامی، پزشکی، سیاسی و غیره شده‌اند [۲]. مزایای اقتصادی بسیاری که در این حوزه وجود دارد، صنعت بدافزارهای مخرب را به بازاری مستعد برای تولیدکنندگان و نویسندگان بدافزارها تبدیل کرده است. از این رو روش‌های جدید و هوشمند شناسایی هم در بسیاری از مواقع در حرکت همگام با تولید بدافزارهای جدید و شناسایی آن‌ها، عاجز بوده‌اند [۳]. یکی از انواع بدافزارها که هنوز هم چالش‌های عمده‌ای در کشف آن وجود دارد، بدافزارهای فراریخت هستند. ساختار این بدافزارها در هر تکثیر با استفاده از فنون مختلف مبهم‌سازی کد<sup>۷</sup>، تغییر می‌یابد اما عملکرد اصلی آن‌ها حفظ می‌شود و این مسئله، شناسایی آن‌ها را بسیار دشوار و پیچیده می‌کند [۲]. مولد یک بدافزار فراریخت همیشه در تکثیرهای بعدی کد، جهش ایجاد می‌کند به‌طوری که برای هر ویروس در هر نفوذ،

8-Signature

9-Trojan Horse

10-Worm

11-Static Analysis

12-Dynamic Analysis

13-Emulation

14-Histogram

1-Semantec

2-Trojan.Ransomlock.G

3-<https://www.microsoft.com/en-us/security/portal/mmpc/default.aspx>

4-Polymorphic Malware

5-Metamorphic

6-Critical Systems

7-Code Obfuscation

گرفته است که با درج کد برنامه‌های سالم و استفاده از جانشینی دستورات مشابه، روش‌ها با افت دقت شناسایی روبرو خواهند شد.

روش‌های شناسایی موجود دارای مشکلات دیگری نیز هستند. برخی از روش‌ها [۱۰،۳] به دلیل سربار محاسباتی بالا، زمانی در مقابل ویروس‌ها موثر عمل می‌کنند که دارای زمان کافی برای شناسایی باشند و ضد بدافزار و بانک اطلاعاتی آن، به نرخ سریعی به‌روزرسانی شوند در غیر این صورت تأثیر مثبت آن‌ها کاهش پیدا می‌کند [۱۲]. به دلیل وجود چالش‌های مطرح شده، کشف کامل و قطعی همه انواع و گونه‌های بدافزارهای چندریختی و فراریختی با اعمال روش‌های مختلف و پیچیده جدید مبهم‌سازی، هنوز نیاز به پژوهش دارد و مسئله تحقیقاتی ارزشمندی است.

در این مقاله با تکیه بر تفاوت نرخ تکرار کدهای عملیاتی در خانواده‌های مختلف بدافزار فراریخت و فایل‌های سالم، روش کشف جدیدی پیشنهاد شده است. این روش با استفاده از چهار معیار استخراج شده حاصل از تحلیل ایستا، هر فایل ورودی مشکوک را به دو دسته بدافزار یا سالم دسته‌بندی می‌کند. روش پیشنهادی با تأکید روی پالایش و انتخاب صحیح کدهای عملیاتی در صد افزایش کارایی است. همچنین در شمارش کدهای عملیاتی دستورات مشابه را نیز همسان در نظر گرفته که در مقابل روش مبهم‌سازی جانشینی دستورات مشابه مقاوم است. شمارش ثبات‌ها به عنوان یک معیار نیز در این روش مدنظر قرار گرفته شده است. به منظور ارزیابی روش پیشنهادی، آزمایش‌هایی بر روی ۴۵۰ فایل متشکل از فایل‌های سالم و بدافزار فراریخت از ویروس‌های خانواده‌های G2, MPCGEN, NGVCK, VLC و کرم‌های خانواده MWOR صورت گرفته است. نتایج آزمایش‌ها نشان می‌دهند در کنار دقت شناسایی ۱۰۰ درصدی، عدم افت دقت شناسایی در جانشینی دستورات مشابه، سادگی و سربار کم محاسباتی از مزایای روش پیشنهادی هستند. ساختار ادامه مقاله به این صورت است: در بخش دوم

کارهای مرتبط مورد بررسی قرار می‌گیرد. در بخش سوم روش پیشنهادی تشریح می‌شود. بخش چهارم به ارزیابی و تفسیر نتایج اختصاص دارد. در نهایت بخش پنجم نیز به نتیجه‌گیری و کارهای آینده می‌پردازد.

## ۲- کارهای مرتبط

کانفورما و همکاران [۱] روشی برای شناسایی ویروس‌های فراریخت ارائه کرده‌اند که مبتنی بر شناسایی ویژگی‌های خاص از کدهای اسمبلی مانند دستوراتی که محتوای ثبات‌ها را تغییر می‌دهند یا دستوراتی که کنترل جریان را تغییر می‌دهند، کار می‌کند. آزمایش‌های صورت گرفته در این مقاله موفقیت ۹۷ درصدی در شناسایی بدافزارهای فراریخت را نشان می‌دهد. از مزایای این روش می‌توان به سادگی آن اشاره کرد. همچنین پیاده‌سازی آن برای ضدبدافزارها کار آسانی می‌باشد و بار زیادی برای سیستم ندارد.

رویکرد مهرا و همکاران در مقاله [۶] بر روی شناسایی و طبقه‌بندی بدافزارهای فراریخت بر اساس خانواده‌های آن‌ها تمرکز دارد. روش پیشنهادی به این صورت است که گراف جریان کنترلی رسم شده و گراف فراخوانی API<sup>۱</sup> ها ایجاد می‌شود. این رویکرد هر بدافزار فراریخت را بر اساس ویژگی‌های خانواده‌شان که از هیستوگرام و فرمول اندازه‌گیری کای دو، که بر اساس تحلیل پویا است، طبقه‌بندی می‌کنند. در این مقاله، دقت در الگوریتم‌های مختلف طبقه‌بندی از ۸۹ تا ۹۹/۱۰ درصد به دست آمده است.

روش پیشنهادی کانفورما و همکاران [۴] به معرفی فنون شناسایی مبتنی بر این فرض که یک اثر جانبی مشترک بین بسیاری از موتورهای فراریخت وجود دارد، تکیه کرده است. این روش برای حدود ۱۰۰۰ برنامه مورد آزمایش و بررسی قرار داده شده و بر اساس نتایج آن به‌طور دقیق ویروس‌های فراریخت و غیرفراریخت را دسته‌بندی کرده است. این روش دارای دقت نزدیک به ۹۹

15-Application programming interface

درصدی در شناسایی است و از معایب آن نیز می‌توان به این مورد اشاره کرد که اگر از روش جانشینی دستورات مشابه یا اضافه کردن کد زائد به کد بدافزار استفاده شود، توزیع دستورات تکراری تغییر می‌کند و روش پیشنهادی دچار شکست در شناسایی خواهد شد.

محمد بن خمس و همکاران در مقاله [۱۱] ویژگی‌های تغییر داده نشده در ساختار بدافزارهای فراریخت را برای استفاده در فرایند شناسایی بدافزار با استفاده از ماشین بردار پشتیبانی<sup>۱۶</sup> استخراج می‌کند. خصوصیات n-gram به‌طور مستقیم از ساختار دودویی بدافزار استخراج شده که این خصوصیات به عنوان امضا در نظر گرفته می‌شوند. این خصوصیات مقادیر قابل توجهی از تعداد خصوصیات انتخاب n-gram در حالت اصلی را کاهش می‌دهد. این روش ترکیبی از استخراج امضا n-gram و ماشین بردار پشتیبانی است. نتایج ارزیابی‌ها دقتی در حدود ۹۹ درصد و نرخ مثبت کاذب پایینی را برای روش پیشنهادی ثبت کرده است. تمرکز اصلی کومار و همکاران در روش پیشنهادی [۱۳]، روی تجزیه و تحلیل و مقایسه برخی از محبوب‌ترین ابزار تشخیص بدافزارهای مخرب است. در این مقاله، نویسندگان برای تشخیص بدافزارهای مخرب یک عملکرد مقایسه‌ای با ابزار و فنون موجود ارائه داده‌اند. در مقاله مذکور<sup>۱۷</sup> ابزار شناخته شده تشخیص بدافزارهای مخرب و ۲۹ بدافزار به عنوان یک مجموعه داده برای مقایسه در نظر گرفته شده است.

چن و همکاران در مقاله [۳]، جدا از این که بر اساس محتویات فایل‌های استخراج شده از نمونه فایل‌ها، بدافزار مخرب را شناسایی می‌کنند، مطالعه و بررسی چگونگی استفاده از گراف مربوط به ارتباطات فایل‌ها را برای تشخیص بدافزارهای مخرب مورد توجه قرار داده‌اند و یک الگوریتم انتشار جدید بر اساس گراف ساخته شده برای تشخیص بدافزارهای مخرب تازه و ناشناخته را پیشنهاد کرده‌اند. آزمایش‌ها روی ۶۷۵ فایل متشکل از فایل‌های بدافزار فراریخت و سالم انجام گرفته است. بر طبق نتایج

16- Support Vector Machine

به‌دست آمده، این روش موثرتر و کاراتر از روش‌های قبلی شناسایی بدافزارها مبتنی بر داده کاوی بوده است. گرچه روش دارای سربار محاسباتی بالا و کارایی کمی می‌باشد.

### ۳- روش پیشنهادی

تحلیل و بررسی روش‌های پیشنهادی در ادبیات و نقاط قوت و ضعف آن‌ها نشان می‌دهد که موتورهای فراریختی همه چیز را در ساختار کد بدافزار تغییر نمی‌دهند. عملیات روی تعداد بسیاری از ثبات‌ها عوض نمی‌شود؛ بعضی از ثبات‌ها اصلاً تغییر نمی‌کنند و بسیاری از کدهای عملیاتی در بدنه کد نیز بی‌تغییر می‌مانند. این بینش در ساختار بدافزارهای فراریخت انگیزه طراحی روش جدیدی برای کشف بدافزارهای فراریخت است.

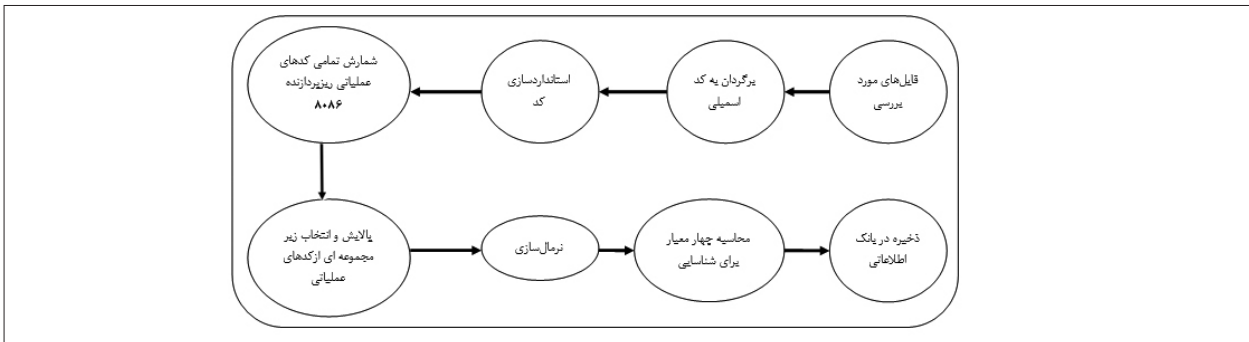
روش پیشنهادی این مقاله، شناسایی بدافزار فراریخت مبتنی بر تحلیل ایستا فرض می‌کند که معیار مشترکی بین بسیاری از موتورهای فراریخت وجود دارد که در بدنه و ویروس‌ها تعداد زیادی از ثبات‌ها، دستورات عملیاتی یا کدهای عملیاتی تکرار می‌شوند. این معیار می‌تواند پایه و اساسی برای تمایز بین خانواده‌های مختلف بدافزارها باشد و با استفاده از این تمایز می‌توان طبقه‌بندی صحیحی از آن‌ها انجام داد. روش پیشنهادی شامل تحلیل نرخ تکرار دستورات برنامه یا همان کدهای عملیاتی و ثبات‌ها است. سپس طبق دسته‌بندی نرخ تکرارها، فایل‌های سالم و خانواده‌های مختلف بدافزارهای فراریخت را طبقه‌بندی خواهد کرد. مراحل مختلف پیش‌پردازش اولیه روش پیشنهادی به این صورت است:

پس از پیش‌پردازش انجام شده و محاسبه معیارها،  $n$  و  $k$  روند کلی روش پیشنهادی برای شناسایی بدافزارهای فراریخت به صورت نشان داده شده در شکل ۱ و شکل ۲ خواهد بود.

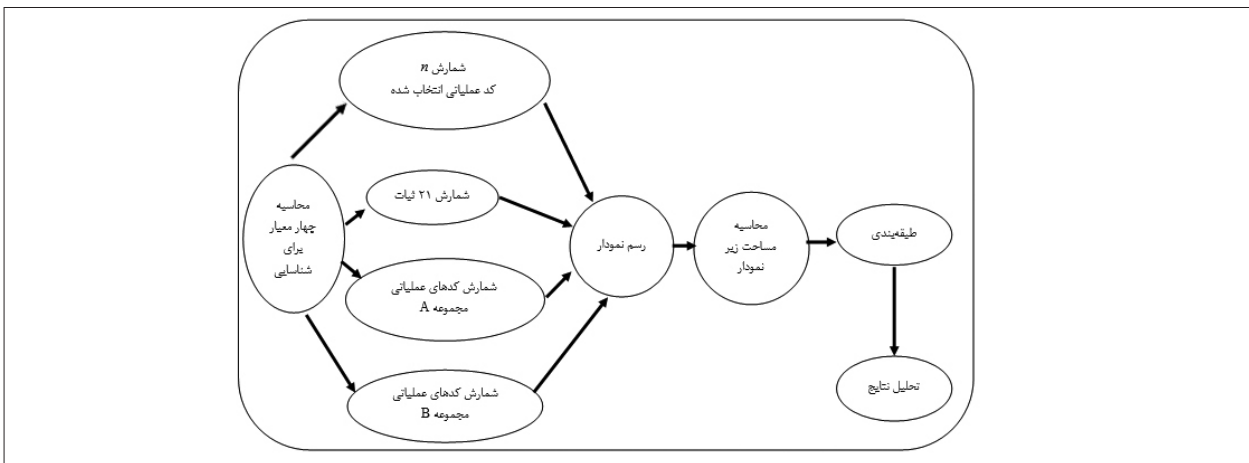
در مقاله [۴] از مجموع تمامی کدهای عملیاتی، ۱۵۴ مورد را بر مبنای تعیین سطح آستانه نرخ تکرار دو، برای فرایند شناسایی نهایی انتخاب کرده است اما امکان بهبود

الگوریتم: مرحله پیش‌پردازش روش پیشنهادی  
 ورودی: فایل مورد بررسی  
 خروجی: کدهای عملیاتی و ثبات‌های شمارش شده  
 شروع مراحل:  
 ۱- برگردان فایل مورد بررسی به کد اسمبلی  
 ۲- استانداردسازی کد اسمبلی فایل مورد بررسی  
 ۳- شمارش تمامی کدهای عملیاتی ریزپردازنده ۸۰۸۶  
 ۴- پالایش و انتخاب  $n$  کد عملیاتی و حذف  $m$  مورد که نرخ تکرار آن‌ها به حد آستانه تعیین شده  $k$  نرسیده باشد  
 ۵- نرمال‌سازی نرخ تکرارها  
 ۶- محاسبه چهار معیار روش پیشنهادی برای شناسایی و طبقه‌بندی فایل‌های سالم و خانواده‌های بدافزارهای فراریخت  
 ۷- ذخیره در بانک اطلاعاتی مربوط به شمارش ثبات‌ها و کدهای عملیاتی خانواده‌های مختلف  
 اتمام مرحله پیش‌پردازش

الگوریتم: روند کلی شناسایی بدافزار در روش پیشنهادی  
 ورودی: فایل مورد بررسی  
 خروجی: اعمال معیارها و طبقه‌بندی فایل مورد بررسی به فایل سالم یا بدافزار  
 شروع مراحل:  
 ۱- تا مرحله ششم مرحله پیش‌پردازش، مجدداً تکرار خواهد شد.  
 ۲- محاسبه چهار معیار به صورت زیر:  
 معیار  $(\alpha)$  شمارش  $n$  کد عملیاتی انتخاب شده طبق حد آستانه  $k$   
 معیار  $(\beta)$  شمارش تمامی ثبات‌های ریزپردازنده ۸۰۸۶ بجز ثبات فلگ که ۲۱ مورد هستند  
 معیار  $(\gamma)$  شمارش کدهای عملیاتی مجموعه  $A$  که در جدول (۱) نشان داده شده‌اند  
 معیار  $(\delta)$  شمارش کدهای عملیاتی مجموعه  $B$  که در جدول (۱) نشان داده شده‌اند  
 ۳- رسم نمودار معیارهای محاسبه شده به صورت جداگانه  
 ۴- محاسبه مساحت زیر نمودار هر معیار به صورت جداگانه  
 ۵- طبقه‌بندی بر اساس نتایج به دست آمده از معیارها  
 ۶- تحلیل نتایج شناسایی  
 اتمام مراحل شناسایی



شکل ۱: گام اول نمای کلی روش پیشنهادی مربوط به پیش‌پردازش اولیه



شکل ۲: گام دوم نمای کلی روش پیشنهادی

روش پیشنهادی خواهد شد. پس از شمارش کل فایل‌های خانواده‌های فراریخت و سالم فقط کدهای عملیاتی پالایش

در انتخاب حد آستانه و حذف تعداد بیشتری از کدهای عملیاتی که در روند شناسایی تبعات افزایش کارایی

و انتخاب خواهند شد که میزان تکرار آن‌ها در کل فایل‌ها بیشتر از حد آستانه  $k$  خواهد بود و انتخاب حد آستانه صحیح بر روی کارایی بسیار تاثیرگذار است. برای انتخاب حد آستانه و حذف  $m$  مورد از کدهای عملیاتی، یکی از دلایل این است که کدهای عملیاتی که در ساختار کد یک فایل، دارای اهمیت بسیار کمی هستند و نرخ تکرار آن‌ها بسیار پایین است و اصولاً تمایزی بین خانواده‌های مختلف بدافزارها و فایل‌های سالم ایجاد نمی‌کنند باید حذف شوند. پس از فرایند انتخاب و پالایش و حذف  $m$  مورد که تعداد تکرار کمتر از حد آستانه  $k$  داشتند،  $n$  کد عملیاتی از تمامی کدهای عملیاتی برای فرایند شناسایی نهایی انتخاب شدند که در جدول (۱) نمایش داده شده‌اند. پس از پالایش کدهای عملیاتی، چون نرخ تکرار آن‌ها در فایل‌های مختلف بسیار متفاوت است و این که نرخ تکرارها را در یک محدوده قرار داد باید آن‌ها را نرمال‌سازی کرد. داده‌ها به‌طور مستقیم قابل مقایسه نیستند به این دلیل که تحلیل هر فایل یک بعد از تعداد کل دستورات را بیان می‌کند در حالی که از یک برنامه با برنامه دیگر متفاوت خواهد بود. به‌منظور از بین بردن این اختلاف در برنامه‌های نمونه و ایجاد داده قابل مقایسه شدن، نیاز به نرمال‌سازی هر جزء از جمعیت تحلیل شده از جنبه‌های از تعداد کل دستورات است. برای نرمال‌سازی از فرمول (۱) استفاده شده است:

نرخ تکرار هر کد عملیاتی

(۱) نرمال شده نرخ تکرار هر کد عملیاتی =  $\frac{\text{مجموع تکرار تمامی کدهای عملیاتی همان فایل مورد بررسی}}{\text{تعداد کل دستورات در آن فایل}}$

اصلی برای شناسایی و طبقه‌بندی کردن فایل‌ها محاسبه شوند. همان‌طور که در جدول (۱) نشان داده شده است برای هر یک از چهار معیار، کدهای عملیاتی مختلف و ثبات‌هایی در نظر گرفته شده‌اند که در زیر به تفکیک مورد بحث قرار می‌گیرند:

معیار  $\alpha$ : برای این معیار کل  $n$  کد عملیاتی انتخاب شده مدنظر قرار خواهند گرفت.

معیار  $\beta$ : برای این معیار ۲۱ ثبات در نظر گرفته شده‌اند و

هر یک در تمامی فایل‌ها شمارش شده و محاسبه می‌شوند. معیار  $\gamma$ : این معیار به کدهای عملیاتی (کدهای عملیاتی مجموعه A) می‌پردازد که با توجه به نرخ تکرار آن‌ها و تحلیل‌هایی که بعد از شمارش تمامی فایل‌ها صورت گرفته است نرخ تکرار بالایی در فایل‌های سالم دارند و در فایل‌های بدافزار یا اصلاً تکرار نشده‌اند و یا تعداد زیر سه بار تکرار را دارند.

معیار  $\delta$ : این معیار به کدهای عملیاتی (کدهای عملیاتی مجموعه B) می‌پردازد که با توجه به نرخ تکرار آن‌ها و تحلیل‌هایی که بعد از شمارش تمامی فایل‌ها صورت گرفته است نرخ تکرار بالایی در فایل‌های بدافزار دارند و در فایل‌های سالم یا اصلاً تکرار نشده‌اند و یا تعداد زیر سه بار تکرار را دارند.

سپس نمودار خطی تکرار کدهای عملیاتی که نمودار  $X$ ها مربوط به ثبات‌ها یا کدهای عملیاتی و نمودار  $Y$ ها به تعداد تکرار هر ثبات یا کد عملیاتی اختصاص دارد، رسم خواهد شد. برای محاسبه چهار معیار، مساحت زیر نمودار مربوط به تکرار همان کدهای عملیاتی و ثبات‌های جدول (۱) محاسبه می‌شود. سپس برای هر فایل این اعداد در بانک اطلاعاتی ذخیره شده تا این که همه فایل‌ها مورد شمارش قرار گیرند. پس از شمارش کلیه فایل‌ها و محاسبه تمامی معیارها، کلیه اطلاعات در بانک اطلاعاتی ذخیره می‌شود تا مورد تحلیل و طبقه‌بندی قرار گیرد.

در ارتباط با گسترش روش پیشنهادی به پردازنده‌های امروزی نیز می‌توان در مرحله پیش‌پردازش کدهای عملیاتی پردازنده‌های جدیدتر را نیز به لیست شمارش اضافه کرد و به این شکل روش به پردازنده‌های جدید نیز بسط داده خواهد شد. همچنین با توجه به این که برای شناسایی بدافزار، برگردان کد وجود دارد و روال همیشگی همه انواع بدافزارها است و چون مرحله پیش‌پردازش یک بار انجام می‌گیرد هزینه محاسباتی آن سرشکن خواهد شد.



جدول ۱: معیارهای شناسایی و ارزیابی در روش پیشنهادی همراه با کدهای عملیاتی و ثبات‌های انتخاب شده

معیارها کدها	معیار $\alpha$ : n کد عملیاتی انتخاب شده	معیار $\beta$ : ۲۱ ثبات	معیار $\gamma$ : کدهای عملیاتی مجموعه A	معیار $\delta$ : کدهای عملیاتی مجموعه B
کدهای عملیاتی یا ثبات‌ها	ADC,ADD,AND,ARPL,BT,CALL,CDQ,CLC,CLD CLI,CMP,CWD,DEC,DIV,HLT,IDIV,IMUL,IN INC,INS,INT,JA,JNBE,IAE,JNB,JB,JNAE JBE,JNA,JC,JCXZ,JECXZ,JE,JZ,JG,JNLE JGE,JNL,JL,JNGE,JLE,JNG,JMP,JNC,JNE JNZ,JNS,JO,JS,LDS,LEA,LEAVE,LES,LOCK LOOP,LOOPNZ,LOOPNE,MOV,MOVSB,MOVSD,MOVZX NEG,NOP,NOT,OR,OUT,POP,POPA,POPAD PUSH,PUSHA,PUSHAD,RCL,RCR,REP,REPE REPZ,REPNE,REPNZ,RET,RETF,ROL,ROR SAHF,SAL,SHL,SAR,SBB,SCAS,SETB SETNAE,SETE,SETZ,SETNE,SETNZ,SHR STC, STD,STI,STOS,SUB,TEST,XCHG,XOR	AX,BX CX,DX AH,AL BH,BL CH,CL DH,DL CS,DS ES,SS SI,DI SP,BP IP	NOP,MOVZX JLE - JNG,LEAVE JBE - JNA JG - JNLE SETNZ - SETNE CLD,INT SETZ - SETE JGE - JNL JL - JNGE,ROL,JS REPE - REPZ JNB - JAE MOVSB,SAHF SAR,JNS,CDQ IDIV,ARPL,INS,JO	STC,STI ROR RET - RETF PUSHA - PUSHAD LOOP,JNC JC,CWD CLI,BT XCHG,NEG ADC POPA - POPAD SBB

#### ۴- ارزیابی

ویروس VLC و ۵۰ کرم MWOR بوده‌اند. ابتدا همه ۴۵۰ فایل مورد بررسی، با استفاده از ابزار IDA PRO<sup>۱۷</sup> به کد اسمبلی برگردانده می‌شوند. سپس کد هر یک از فایل‌ها به منظور شمارش ابتدا باید استانداردسازی روی آن انجام شود. در این مرحله نویسه‌هایی که عملیات شمارش را با اختلال مواجه می‌سازند و در اصطلاح هیچ نقشی در ساختار کد ندارند از کد اصلی حذف می‌شوند و فقط ساختار اصلی کد باقی می‌ماند.

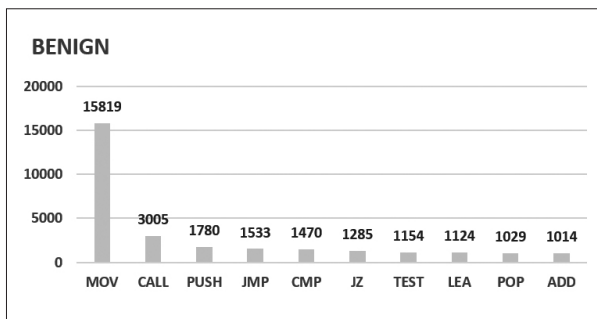
پایه و اساس روش پیشنهادی در شمارش کدهای عملیاتی است که در ریزپردازنده ۸۰۸۶ در زبان اسمبلی ۱۹۱ کد عملیاتی وجود دارد. تمامی ۱۹۱ مورد برای تمامی ۴۵۰ فایل مورد بررسی، در مرحله پیش‌پردازش شمارش می‌شوند. با آزمایش‌ها و ارزیابی‌های صورت گرفته و مشاهده تاثیر انتخاب حد آستانه در کارایی روش پیشنهادی، حد آستانه مقدار  $k=5$  تعیین شد که دلیل آن

در روش پیشنهادی این مقاله توجه اصلی به بهبود دقت شناسایی نسبت به مقاله‌های پیشین، افزایش کارایی سیستم از جمله میزان حافظه مصرفی و سرعت در حد قابل قبولی نسبت به روش‌های پیشین، پوشش نقاط ضعف روش‌های قبلی مانند جایگزینی دستورات مشابه که موجب کاهش درصد شناسایی می‌شده است، خواهد بود. حال با مدنظر قرار دادن نکات مطرح شده و فرضیات روش پیشنهادی آزمایش‌ها و ارزیابی‌ها انجام خواهد شد.

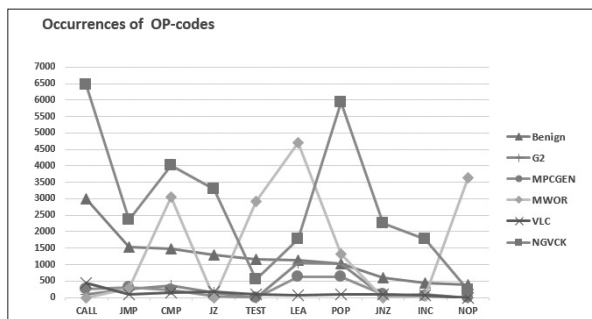
#### ۴-۱- آماده‌سازی مجموعه داده

به‌منظور ارزیابی روش پیشنهادی، ۴۵۰ فایل متشکل از ۴۰ فایل سالم از مجموعه Cygwin و ۴۱۰ فایل بدافزار فراریخت مورد بررسی قرار گرفته است. بدافزارها از پنج خانواده بدافزارهای فراریخت متشکل از ۲۵۰ ویروس NGVCK، ۵۰ ویروس G2، ۵۰ ویروس MPCGEN، ۱۰

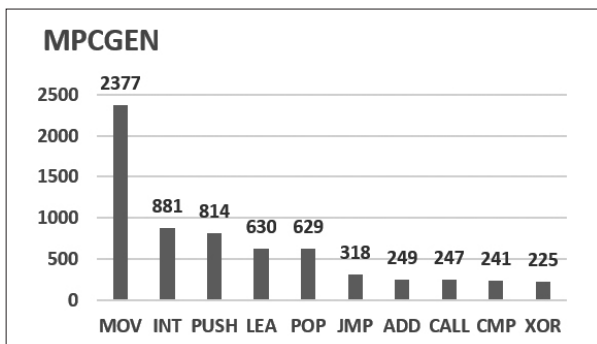
17-<https://www.hex-rays.com/products/ida/support/download.shtml>



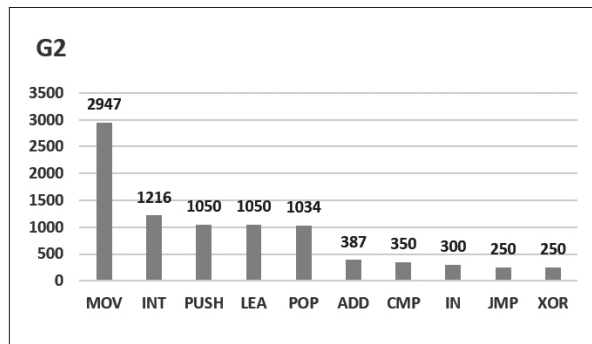
شکل ۴: نرخ تکرار ۱۰ کد عملیاتی پر تکرار در فایل‌های سالم



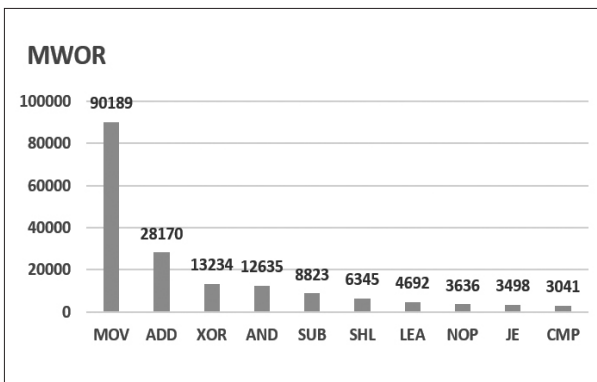
شکل ۳: تفاوت در رخداد برخی از کدهای عملیاتی مختلف در خانواده‌های متفاوت بدافزارهای فراریخت و فایل‌های سالم



شکل ۶: نرخ تکرار ۱۰ کد عملیاتی پر تکرار در خانواده بدافزار MPCGEN



شکل ۵: نرخ تکرار ۱۰ کد عملیاتی پر تکرار در خانواده بدافزار G2



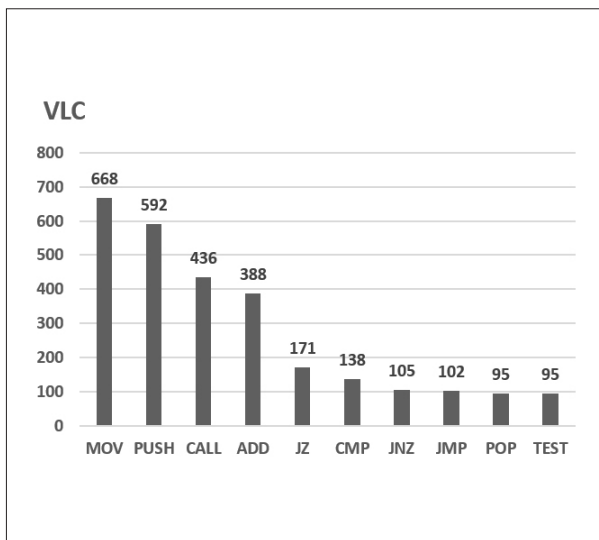
شکل ۷: نرخ تکرار ۱۰ کد عملیاتی پر تکرار در خانواده بدافزار MWOR

با نرخ تکرار متفاوتی پرتکرار هستند و این می‌تواند تمایزی را برای طبقه‌بندی کردن خانواده‌های مختلف ایجاد کند. برخی از کدهای عملیاتی فقط در برنامه‌های سالم هستند و هرگز در بدافزارها وجود ندارند. بعضی کدهای عملیاتی ممکن است در بدافزارها وجود داشته باشند اما هرگز در برنامه‌های سالم وجود ندارند. برخی از کدهای عملیاتی نه در بدافزارها وجود دارند و نه در برنامه‌های سالم که این موارد در ادامه معیارهایی برای طبقه‌بندی‌ها در روش پیشنهادی شده است. به‌عنوان مثال برخی از کدهای

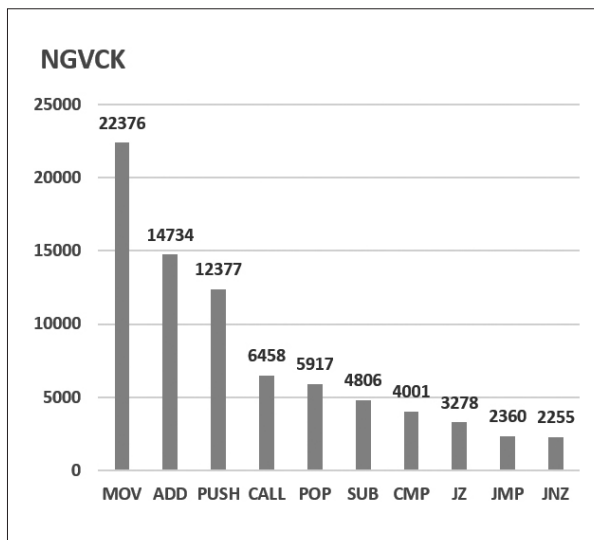
این است که مقدار کمتر از پنج باعث حذف کدهای عملیاتی کمتری می‌شود که در واقع چون نرخ تکرار پایینی داشتند نقشی در روند شناسایی نداشتند و انتخاب مقداری بالاتر از پنج باعث افت دقت شناسایی شده بود. با در نظر گرفتن  $n=83$ ,  $k=5$ , تعداد کدهای عملیاتی که پالایش و انتخاب شدند  $m=107$  می‌باشد و تعداد کدهای عملیاتی حذف شده نیز  $m=107$  خواهد بود. با مشاهده نرخ تکرار کدهای عملیاتی و ثبات‌ها در فایل‌های سالم و خانواده‌های مختلف از بدافزارها، کاملاً مشهود بود که تمایزی بین تعداد تکرارها برای کدهای عملیاتی مختلف وجود دارد. این تمایز در تعداد تکرارها پایه و اساس روش پیشنهادی این مقاله است. با توجه به شکل (۳) این تمایز به خوبی نشان داده شده است:

مقایسه تحلیل‌ها نشان می‌دهد که توزیع کدهای عملیاتی در برنامه‌های سالم با توزیع آن‌ها در بدافزارها متفاوت است. همچنین میزان نرخ تکرار ۱۰ کد عملیاتی پرتکرار در فایل‌های سالم و پنج خانواده بدافزار فراریخت که در شکل‌های (۴، ۵، ۶، ۷، ۸، ۹) نشان داده شده‌اند بیانگر این مسئله است که در هر یک از خانواده‌ها، کدهای عملیاتی متفاوتی

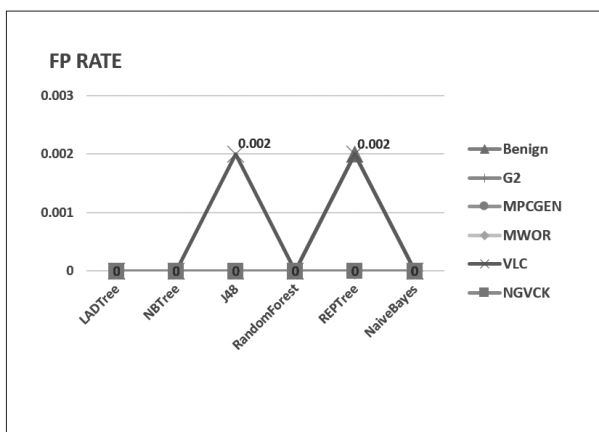




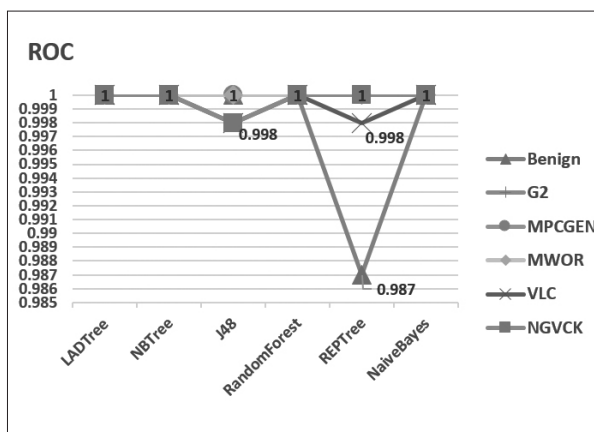
شکل ۹: نرخ تکرار ۱۰ کد عملیاتی پر تکرار در خانواده بدافزار VLC



شکل ۸: نرخ تکرار ۱۰ کد عملیاتی پر تکرار در خانواده بدافزار NGVCK



شکل ۱۱: مقایسه میزان نرخ مثبت کاذب برای پنج خانواده بدافزارهای فراریخت و فایل‌های سالم



شکل ۱۰: مقایسه میزان ROC برای پنج خانواده بدافزارهای فراریخت و فایل‌های سالم

است. تحلیل نتایج طبقه‌بندی‌ها در جدول (۲) نشان می‌دهد که روش پیشنهادی قادر است که به‌طور صحیحی طبقه‌بندی بین بدافزارهای فراریخت و فایل سالم را با دقت بسیار بالا انجام دهد. علاوه بر این نرخ مثبت کاذب برای برنامه سالم و بدافزارهای تولید شده با همه موتورهای تکثیر بسیار پایین و تقریباً صفر است. نتیجه شایان توجه این است که روش پیشنهادی در چهار الگوریتم دسته‌بندی دارای دقت شناسایی ۱۰۰ درصدی است و نرخ مثبت کاذب در این چهار الگوریتم مقدار صفر را نشان می‌دهد.

نکات حائز اهمیت در مورد نتایج به‌دست آمده این است که در بسیاری از حالات شناسایی ۱۰۰ درصدی را نشان

عملیاتی مانند MOV, ADD, PUSH, POP, AND و مواردی که در شکل‌ها مشهود است دارای نرخ بالای رخداده در برنامه‌های سالم و همچنین دارای تنوع وسیعی هستند درحالی‌که کدهای عملیاتی دیگر مقادیر رخداده و تنوع آن‌ها اختلاف نزدیکی دارد.

## ۴-۲- تحلیل نتایج

با توجه به تحلیل‌های صورت گرفته و استفاده از الگوریتم‌های دسته‌بندی مختلف، شکل‌های (۱۰، ۱۱، ۱۲، ۱۳) به دست آمده‌اند و نتایج کلی در جدول (۲) نشان داده شده

جدول ۲: نتایج مرحله طبقه‌بندی برای شش الگوریتم طبقه‌بندی

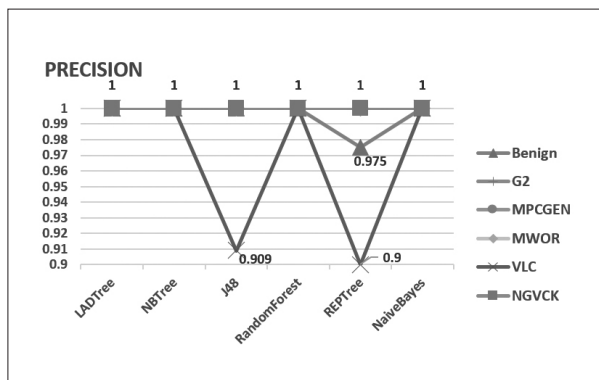
الگوریتم‌ها	مثبت واقعی	مثبت کاذب	دقت	بازیابی	معیار F	ناحیه ROC	فایل‌ها
LADTree	۱	۰	۱	۱	۱	۱	BENIGN
	۱	۰	۱	۱	۱	۱	G2
	۱	۰	۱	۱	۱	۱	MPCGEN
	۱	۰	۱	۱	۱	۱	MWOR
	۱	۰	۱	۱	۱	۱	VLC
	۱	۰	۱	۱	۱	۱	NGVCK
NBTree	۱	۰	۱	۱	۱	۱	BENIGN
	۱	۰	۱	۱	۱	۱	G2
	۱	۰	۱	۱	۱	۱	MPCGEN
	۱	۰	۱	۱	۱	۱	MWOR
	۱	۰	۱	۱	۱	۱	VLC
	۱	۰	۱	۱	۱	۱	NGVCK
J48	۱	۰	۱	۱	۱	۱	BENIGN
	۱	۰	۱	۱	۱	۱	G2
	۱	۰	۱	۱	۱	۱	MPCGEN
	۱	۰	۱	۱	۱	۱	MWOR
	۱	۰/۰۰۲	۰/۹۰۹	۱	۰/۹۵۲	۰/۹۹۸	VLC
	۰/۹۹۶	۰	۱	۰/۹۰۶	۰/۹۹۸	۰/۹۹۸	NGVCK
RandomForest	۱	۰	۱	۱	۱	۱	BENIGN
	۱	۰	۱	۱	۱	۱	G2
	۱	۰	۱	۱	۱	۱	MPCGEN
	۱	۰	۱	۱	۱	۱	MWOR
	۱	۰	۱	۱	۱	۱	VLC
	۱	۰	۱	۱	۱	۱	NGVCK
REPTree	۰/۹۷۵	۰/۰۰۲	۰/۹۷۵	۰/۹۷۵	۰/۹۷۵	۰/۹۸۷	BENIGN
	۱	۰	۱	۱	۱	۱	G2
	۱	۰	۱	۱	۱	۱	MPCGEN
	۱	۰	۱	۱	۱	۱	MWOR
	۰/۹	۰/۰۰۲	۰/۹	۰/۹	۰/۹	۰/۹۹۸	VLC
	۱	۰	۱	۱	۱	۱	NGVCK
NaiveBayes	۱	۰	۱	۱	۱	۱	BENIGN
	۱	۰	۱	۱	۱	۱	G2
	۱	۰	۱	۱	۱	۱	MPCGEN
	۱	۰	۱	۱	۱	۱	MWOR
	۱	۰	۱	۱	۱	۱	VLC
	۱	۰	۱	۱	۱	۱	NGVCK

جدول ۳: مقایسه روش پیشنهادی این مقاله با روش پیشنهادی مقاله [۴]

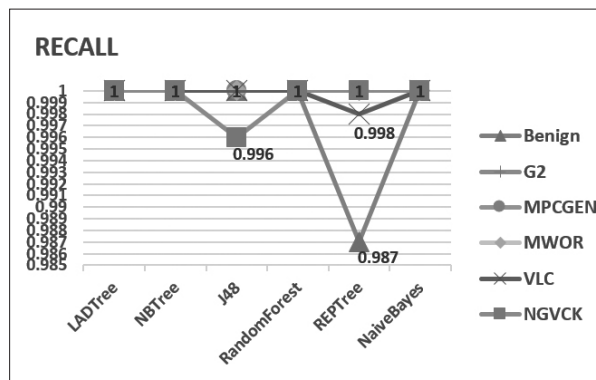
الگوریتم‌ها	نرخ مثبت کاذب		ناحیه ROC		بدازارها
	مقاله [۴]	روش پیشنهادی	مقاله [۴]	روش پیشنهادی	
LADTree	۰/۰۱۲	*	۰/۹۹	۱	BENIGN G2 MPCGEN NGVCK
	۰/۰۰۵	*	۰/۹۹	۱	
	۰/۰۰۵	*	۰/۹۹	۱	
	۰/۰۰۸	*	۱	۱	
NBTree	۰/۰۰۳	*	۰/۹۹	۱	BENIGN G2 MPCGEN NGVCK
	۰/۰۰۴	*	۰/۹۹	۱	
	۰/۰۰۷	*	۰/۹۹	۱	
	.	*	۱	۱	
J48	۰/۰۱۳	*	۰/۹۹	۱	BENIGN G2 MPCGEN NGVCK
	۰/۰۰۱	*	۰/۹۹	۱	
	۰/۰۰۴	*	۰/۹۹	۱	
	.	*	۱	۰/۹۹۸	
RandomForest	.	*	۰/۹۹	۱	BENIGN G2 MPCGEN NGVCK
	۰/۰۰۳	*	۰/۹۹	۱	
	۰/۰۰۱	*	۱	۱	
	۰/۰۰۳	*	۱	۱	
REPTree	۰/۰۰۵	۰/۰۰۲	۰/۹۹	۰/۹۸۷	BENIGN G2 MPCGEN NGVCK
	۰/۰۰۵	*	۰/۹۹	۱	
	۰/۰۰۴	*	۰/۹۹	۱	
	۰/۰۰۳	*	۰/۹۹	۱	

تعداد کدهای عملیاتی کمتر که در حدود نصف مقاله [۴] بوده است و با توجه به نصف شدن تعداد کدهای عملیاتی شمارش شده، سرعت شناسایی نیز افزایش پیدا کرده است که این بهبود نیز بسیار مطلوب است. در مقاله‌هایی که در بخش کارهای مرتبط نیز بحث شد یکی از نقاط ضعف این بود که با جایگزین کردن دستورات مشابه مانند JE و JZ و یا REPE و REPZ که دقیقاً یک عملیات را انجام می‌دهند اما ساختار لغوی متفاوت دارند روش پیشنهادی آن‌ها دچار شکست می‌شده است اما برای روش پیشنهادی این مقاله، تمامی دستورات مشابه در شمارش‌ها یکسان در نظر گرفته شده‌اند و با جایگزین کردن آن‌ها هیچ تاثیری در

می‌دهد که نشان‌دهنده این مسئله است که نفوذ بدافزارهای فراریخت تقریباً صفر خواهد بود. همچنین میزان نرخ مثبت کاذب در بسیاری از حالات برابر با صفر بوده است که باز هم بسیار مطلوب است. با توجه به نتایج به دست آمده، در چهار الگوریتم طبقه‌بندی - LADTree, NBTree, Random-Forest, Naivebayes مقدار ROC=1 و FP=0 می‌باشد که نشان‌دهنده عملکرد بسیار مطلوب روش پیشنهادی است. در دو الگوریتم طبقه‌بندی دیگر نیز نتایج مطلوبی کسب شده است که جای بهبود نیز دارد. همچنین کارایی روش پیشنهادی که بسیار دارای اهمیت است نسبت به مقالات مشابه افزایش پیدا کرده است که با در نظر گرفتن میزان



شکل ۱۳: مقایسه میزان دقت برای پنج خانواده بدافزارهای فراریخت و فایل‌های سالم



شکل ۱۴: مقایسه میزان بازیابی برای پنج خانواده بدافزارهای فراریخت و فایل‌های سالم

ایجاد کند. نتایج نشان می‌دهد که این روش یک معیار موثر برای نشان دادن تمایز بین خانواده‌های مختلف بدافزارها و فایل‌های سالم است. یافته‌های آزمایش‌ها نمایانگر این است که با وجود سادگی، روش پیشنهادی بسیار دقیق می‌باشد. در روش پیشنهادی سربار محاسباتی بسیار کم است، ایده اساسی روش به سادگی قابل فهم و در ضد بدافزارها به آسانی قابل پیاده‌سازی است، در مقابل روش‌های جانشینی دستورات مشابه و فنون مختلف مبهم سازی با توجه به شناسایی بر مبنای معیارهای مختلف، مقاوم می‌باشد و دارای کارایی مطلوبی از نظر حافظه مصرفی و سرعت شناسایی می‌باشد. همچنین در مقایسه با روش‌های مشابه موجود دارای بهبودهایی در ابعاد مختلف است.

برای کارهای آینده می‌توان تحلیل پویا را نیز به روش پیشنهادی اضافه کرد تا قدرت این روش افزایش یابد چون در مواردی بدافزارهای فراریخت از روش‌هایی از مبهم سازی استفاده می‌کنند که بدنه کد بدافزار فراریخت رمزنگاری خواهد شد که باعث می‌شود روش‌های تحلیل ایستا به طور کامل دچار شکست شوند که ایده پیشنهادی برای کار آینده این است که با اضافه کردن بخش تحلیل پویا و مشاهده کردن<sup>۱۸</sup> اجرای آن‌ها در محیط‌های امن<sup>۱۹</sup> و مجازی<sup>۲۰</sup> می‌توان این نقص را هم پوشش داد که در کارهای آینده مدنظر خواهد بود.

دقت شناسایی روش پیشنهادی مشاهده نمی‌شود که بهبود مناسبی نسبت به مقالات پیشین است.

نکته دیگر هم این‌که در مقالات مختلف در حوزه شناسایی بدافزارهای فراریخت معمولاً یک معیار یا شاخص برای شناسایی انتخاب می‌شود که اگر نویسندگان بدافزارهای فراریخت بتوانند آن معیار را دچار نقص و چالش کنند روش پیشنهادی با شکست روبرو خواهد شد. اما در روش پیشنهادی این مقاله چهار معیار متفاوت برای شناسایی و طبقه‌بندی خانواده‌های مختلف در نظر گرفته شده‌اند که حتی با شکست یکی از معیارها باز هم تاثیر ناچیزی در روند شناسایی کلی انجام می‌گیرد و روش با شکست مواجه نمی‌شود که نکته‌ای حائز اهمیت است. در جدول (۳) بهبودهای حاصل شده نسبت به مقاله [۴] که نتایجی بسیار مطلوب را کسب کرده بود نشان داده شده است. جدول (۳) نشان می‌دهد که تقریباً در تمامی حالات روش پیشنهادی این مقاله دارای بهبود بوده است.

##### ۵- نتیجه‌گیری و کارهای آینده

روش پیشنهادی این مقاله با شمارش و تحلیل ثبات‌ها و کدهای عملیاتی به شناسایی بدافزارهای فراریخت می‌پردازد و نشان می‌دهد که نرخ تکرار کدهای عملیاتی مختلف در خانواده‌های متفاوتی از بدافزارهای فراریخت و فایل‌های سالم، ناهمسان است که این مسئله می‌تواند وجه تمایزی را

18-Monitoring  
19-Sandbox  
20-Virtual Machine

fence Science Journal, 66(2), 138-145, 2016.

8- Rad, B. B., & Masrom, M. Metamorphic virus variants classification using opcode frequency histogram. arXiv preprint arXiv:1104.3228. 2011.

9- Runwal, N., Low, R. M., & Stamp, M. Opcode graph similarity and metamorphic detection. Journal in Computer Virology, 8(1-2), 37-52, 2012.

10- Toderici, A. H., & Stamp, M. Chi-squared distance and metamorphic virus detection. Journal of Computer Virology and Hacking Techniques, 9(1), 1-14, 2013.

11- Khammas, B. M., Monemi, A., Ismail, I., Nor, S. M., & Marsono, M. N. Metamorphic Malware Detection Based on Support Vector Machine Classification of Malware Sub-Signatures. TELKOMNIKA (Telecommunication Computing Electronics and Control), 14(3), 2016.

12- Al Daoud, E., Jebri, I. H., & Zaqibeh, B. Computer virus strategies and detection methods. Int. J. Open Problems Compt. Math, 1(2), 12-20, 2008.

13- Pandey, S. K., & Mehtre, B. M. Performance of malware detection tools: A comparison. In Advanced Communication Control and Computing Technologies (ICACCCT), 2014 International Conference on (pp. 1811-1817). IEEE, 2014.

## مراجع

1- Canfora, G., Mercaldo, F., Visaggio, C. A., & Di Notte, P. Metamorphic malware detection using code metrics. Information Security Journal: A Global Perspective, 23(3), 57-67, 2014.

2- Aycock, J. Computer Viruses and Malware (Advances in Information Security). Secaucus, 2006.

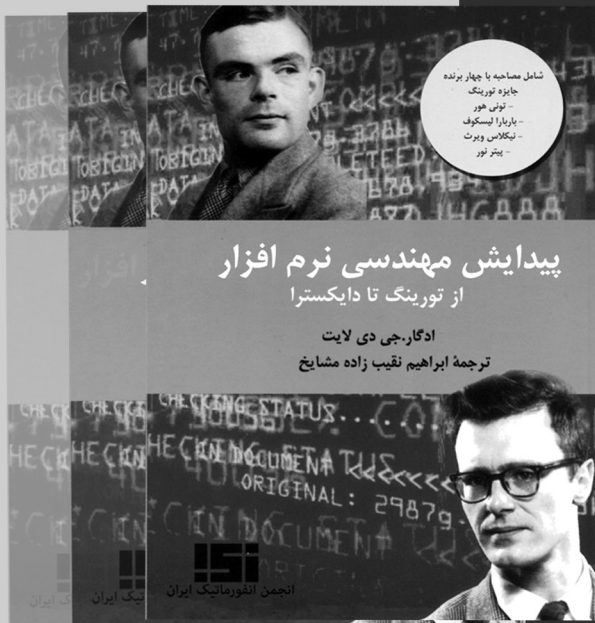
3- Chen, L., Li, T., Abdulhayoglu, M., & Ye, Y. Intelligent malware detection based on file relation graphs. In Semantic Computing (ICSC), 2015 IEEE International Conference on (pp. 85-92). IEEE, 2015.

4- Canfora, G., Iannaccone, A. N., & Visaggio, C. A. Static analysis for the detection of metamorphic computer viruses using repeated-instructions counting heuristics. Journal of Computer Virology and Hacking Techniques, 10(1), 11-27, 2014.

5- Szor, P. The art of computer virus research and defense. Pearson Education, 2005.

6- Mehra, V., Jain, V., & Uppal, D. DaCoMM: Detection and Classification of Metamorphic Malware. In Communication Systems and Network Technologies (CSNT), 2015 Fifth International Conference on (pp. 668-673). IEEE, 2015.

7- Kapoor, A., & Dhavale, S. Control Flow Graph Based Multiclass Malware Detection Using Bi-normal Separation. De-



# منتشر شد!

پیدایش مهندسی نرم افزار

ترجمه: ابراهیم نقیب زاده مشایخ

برای تهیه کتاب با دفتر انجمن انفورماتیک ایران

تماس بگیرید ۶۶۴۱۲۸۶۱