

بِسْمِ اللَّهِ الرَّحْمَنِ الرَّحِيمِ



دانشکده مهندسی کامپیوتر و فناوری اطلاعات

ارایه روشی بر مبنای پردازش سیگنال جهت تحلیل

آلودگی فایل‌های اجرایی

ابوالفضل سرکردۀ بی

استاد راهنما:

دکتر علی‌اکبر پویان

اساتید مشاور:

دکتر حمید حسن‌پور

مهندس جواد کیا

پایان نامه جهت اخذ درجه کارشناسی ارشد

۱۳۹۱ بهمن ماه

ب

دانشگاه صنعتی شاهرود

دانشکده: مهندسی کامپیوتر و فناوری اطلاعات

پایان نامه کارشناسی ارشد آقای ابوالفضل سرکردی

تحت عنوان: ارایه روشی بر مبنای پردازش سیگنال جهت تحلیل آلودگی فایل‌های اجرایی

در تاریخ ۹۱/۱۱/۲۴ توسط کمیته تخصصی زیر جهت اخذ مدرک کارشناسی ارشد مورد ارزیابی و با درجه مورد پذیرش قرار گرفت.

امضاء	اساتید مشاور	امضاء	اساتید راهنمای
	دکتر حمید حسن پور		دکتر علی اکبر پویان
	مهندس جواد کیا		

امضاء	نماينده تحصيلات تمكيلي	امضاء	اساتيد داور
	مهندس محسن فرهادي		دکتر مرتضی زاهدی
			دکتر حسین مروی

نقد پیم به

م در و مادر عزیزتر از جانم

و

برادر فدا کارم

و

خواهران هم بانم

مشکر و قدردانی

پس از سپاسگزاری از درگاه خداوند منان، برخود لازم می‌دانم از استادگر اتقدر جناب آقای

دکتر پویان که انجام این پایان نامه بدون گمک های ایشان امکان پذیر نبود و چنین بپاس زحماتی که

ایشان در طول بیش از ۶ سال برای اینجانب کشیدند، کمال مشکر و قدردانی را داشته باشم؛

چنین از اساتید محترم جناب آقای دکتر حسن پور که بارا هنایی هایشان در طول انجام پایان نامه

گمک شایانی به من کردند و جناب آقای دکتر زاده‌ی سپاسگزارم؛

در پایان از جناب آقای مهندس کیا که مراد انجام این پایان نامه بسیار یاری رسادند مشکر می

کنم.

تعهد نامه

اینجانب ابوالفضل سرکرده بی دانشجوی دوره کارشناسی ارشد رشته مهندسی کامپیوتر- گرایش

هوش مصنوعی دانشکده مهندسی کامپیوتر دانشگاه صنعتی شاهرود نویسنده پایان نامه ارایه روشنی بر

مبناي پردازش سیگنال جهت تحلیل آلودگی فایل های اجرایی تحت راهنمائی آقای دکتر علی-

اکبر پویان متعهد می شوم:

- تحقیقات در این پایان نامه توسط اینجانب انجام شده است و از صحت و اصالت برخوردار است .
- در استفاده از نتایج پژوهش های محققان دیگر به مرجع مورد استفاده استناد شده است .
- مطالب مندرج در پایان نامه تاکنون توسط خود یا فرد دیگری برای دریافت هیچ نوع مدرک یا امتیازی در هیچ جا ارائه نشده است .
- کلیه حقوق معنوی این اثر متعلق به دانشگاه صنعتی شاهرود می باشد و مقالات مستخرج با نام «دانشگاه صنعتی شاهرود» و یا «Shahrood University of Technology» به چاپ خواهد رسید.
- حقوق معنوی تمام افرادی که در به دست آمدن نتایج اصلی پایان نامه تأثیرگذار بوده اند در مقالات مستخرج از پایان نامه رعایت می گردد.
- در کلیه مراحل انجام این پایان نامه، در مواردی که از موجود زنده (یا بافت های آنها) استفاده شده است، ضوابط و اصول اخلاقی رعایت شده است.
- در کلیه مراحل انجام این پایان نامه، در مواردی که به حوزه اطلاعات شخصی افراد دسترسی یافته یا استفاده شده است اصل رازداری، ضوابط و اصول اخلاق انسانی رعایت شده است.

تاریخ

امضای دانشجو

مالکیت نتایج و حق نشر

- کلیه حقوق معنوی این اثر و محصولات آن (مقالات مستخرج ، کتاب ، برنامه های رایانه ای ، نرم افزار ها و تجهیزات ساخته شده است) متعلق به دانشگاه صنعتی شاهرود می باشد . این مطلب باید به نحو مقتضی در تولیدات علمی مربوطه ذکر شود .
- استفاده از اطلاعات و نتایج موجود در پایان نامه بدون ذکر مرجع مجاز نمی باشد.

چکیده

با افزایش روزافرون بدافزارهای جدید و ناشناخته و همچنین استفاده از تکنیک‌های پنهان‌سازی در بدافزارها ارایه روش‌هایی جدید جهت شناسایی آن‌ها بسیار احساس می‌شود. در آنتی‌ویروس‌های تجاری از ترکیب روش‌های مبتنی بر علامت‌های هویتی و همچنین روش‌های مبتنی بر یادگیری جهت شناسایی بدافزارها استفاده می‌شود. با این وجود به دلیل مناسب نبودن روش‌های مبتنی بر یادگیری موجود در آنتی‌ویروس‌ها جهت تشخیص بدافزارهای جدید، این نرم افزارها برای تشخیص بدافزار بسیار به بروز بودن پایگاه داده وابسته هستند. در نتیجه امروزه بهبود روش‌های مبتنی بر یادگیری و ارایه روش‌های جدید در این زمینه بخش اصلی تحقیقات در تشخیص بدافزار را تشکیل می‌دهد.

شناسایی فایل‌های اجرایی آلوده شده با کدهای مخرب یکی از اصلی‌ترین کارهای یک آنتی‌ویروس می‌باشد. که برای تشخیص کدهای مخرب ناشناخته و جدید از روش‌های مبتنی بر یادگیری استفاده می‌شود. هدف از انجام این پایان‌نامه تلاش جهت بدست آوردن یک الگوی سیگنالی از فایل‌های اجرایی می‌باشد تا بتوان با استفاده از آن الگو و تکنیک‌های پردازش سیگنال یک روش تشخیص فایل‌های اجرایی آلوده ارایه کرد. برای این کار ابتدا به بررسی روش‌های مختلف تحلیل فایل اجرایی پرداخته و سپس روش مناسبی که از تحلیل قسمت منبع فایل اجرایی بدست می‌آید معرفی می‌شود. جهت اثبات کارا بودن روش معرفی شده، براساس آن یک بردار ویژگی ۲۵۶ تایی از تکرار ترکیب بایت‌ها در قسمت آخر فایل اجرایی بدست می‌آید و با استفاده از تکنیک‌های دسته‌بندی متن روشهای تشخیص فایل اجرایی آلوده جدید، با دقت ۹۹.۱۰ درصد ارایه می‌شود.

کلمات کلیدی: فایل آلوده، کد مخرب ناشناخته، آنتی‌ویروس، پردازش سیگنال، تشخیص فایل آلوده.

لیست مقالات مستخرج از پایان نامه

1. ابوالفضل سرکردہ یی، علی اکبر پویان، حمید حسنپور، "ارایه روشی هوشمند بر مبنای دسته بندی متن جهت تشخیص کدهای مخرب ناشناخته". یازدهمین کنفرانس سیستم های هوشمند ایران، دانشگاه خوارزمی اسفند ۱۳۹۱.

فهرست مطالب

فصل اول	۱
۱	مقدمه
۲	۱-۱ مقدمه
۳	۲-۱ دسته بندی انواع بدافزارها
۳	۱-۲-۱ ویروس
۵	۲-۲-۱ کرم شبکه
۶	۳-۲-۱ اسپ تروا(تروجان)
۹	۳-۱ تعریف مسئله
۱۰	۴-۱ مقایسه فایل اجرایی سالم و آلوده
۱۰	۱-۴-۱ اطلاعات سرآیند فایل اجرایی
۱۱	۲-۴-۱ تغییرات اشاره گر دستورالعمل
۱۲	۳-۴-۱ تفاوت در قسمت منبع
۱۳	فصل دوم
۱۳	۲ بررسی روش های تشخیص آلودگی فایل های اجرایی
۱۴	۱-۲ مقدمه
۱۵	۲-۲ روش‌های مبتنی بر علامتهای هویتی
۱۵	۱-۲-۲ نحوه تحلیل بدافزارها
۱۷	۲-۲-۲ نقاط ضعف

۱۸..... روشهای مبتنی بر یادگیری و تخمین ۳-۲

۱۸..... تحلیل براساس n بایت استخراج شده از فایل اجرایی ۱-۳-۲

۱۹..... براساس اطلاعات سرآیند فایل اجرایی ۲-۳-۲

۲۰..... براساس کد عملگر فایل اجرایی ۳-۳-۲

فصل سوم

۲۲..... روش های تحلیل فایل های اجرایی آلوود ۳

۲۳..... مقدمه ۱-۳

۲۴..... براساس تحلیل کد اسمبلی ۲-۳

۲۵..... براساس کد عملگر ۱-۲-۳

۲۶..... براساس تغییرات اشاره گر دستورالعمل ۲-۲-۳

۲۷..... براساس بایت های فایل اجرایی ۳-۳

۲۸..... براساس بایت های تشکیل دهنده ۱-۳-۳

۳۰..... براساس ترکیب بایت های تشکیل دهنده ۲-۳-۳

فصل چهارم

۳۴..... پیاده سازی و ارزیابی نتایج ۴

۳۵..... مقدمه ۱-۴

۳۵..... استخراج ویژگی ۲-۴

۳۶..... وزن دهی ویژگی ها ۳-۴

۳۶..... روش TF ۱-۳-۴

۳۷..... روش normTF ۲-۳-۴

۳۷..... روش logTF ۳-۳-۴

۳۷	روش ITF	۴-۳-۴
۳۸	روش Sparck	۵-۳-۴
۳۸	روش IDF	۶-۳-۴
۳۸	روش TFIDF	۷-۳-۴
۳۸	روش normTFIDF	۸-۳-۴
۳۹	روش TFRF	۹-۳-۴
۳۹	انتخاب روش دسته بندی	۴-۴
۴۰	پایگاه داده	۵-۴
۴۱	معیارهای ارزیابی	۶-۴
۴۱	نرخ تشخیص درست مثبت	۱-۶-۴
۴۲	نرخ تشخیص غلط مثبت	۲-۶-۴
۴۲	دقت	۳-۶-۴
	نتایج	۷-۴
	فصل پنجم	۴۷
۴۷	نتیجه گیری و کارهای آینده	۵
۴۸	نتیجه گیری	۱-۵
۴۹	کارهای آینده	۲-۵
	منابع	۵۰

فهرست اشکال

شكل (۱-۳) نمودار فرکانس بایت های تشکیل دهنده قسمت آخر فایل اجرایی. (الف-۱)، (ب-۱)، (ج-۱)،

(پ-۱) (ی-۱). برای فایل های اجرایی سالم. (الف-۲)، (ب-۲)، (ج-۲)، (پ-۲) (ی-۲)، برای فایل

های اجرایی آلوده ۳۰

شكل (۲-۳) نمودار تعداد تکرار ۲ بایت ها در قسمت آخر فایل اجرایی. (الف-۱)، (ب-۱)، (ج-۱)، (پ-۱)

(ی-۱). برای فایل های اجرایی سالم. (الف-۲)، (ب-۲)، (ج-۲)، (پ-۲) (ی-۲)، برای فایل های

اجرایی آلوده ۳۳

شكل ۴-۱ نمودار کارایی روش های مختلف وزن دهی براساس معیارهای معرفی شده ۴۳

فهرست جداول

جدول ۱-۱ بررسی کارایی روش های مختلف وزن دهی بازای تغییر مقادیر K. (الف) روش TF. (ب)

روش .TFIDF. (ج) روش .logTF. (د) روش .ITF. (ه) روش .Sparck. (و) روش .normTF

..... (ز) روش .TFRF. (ح) روش .normTFIDF

فصل اول

ا مقدمہ

۱-۱ مقدمه

با افزایش روز افرون بدافزارها و همچنین استفاده از تکنیک های مختلف جهت پیچیده تر کردن آن ها نیاز به روش های هوشمند برای شناسایی بدافزارها بسیار احساس می شود به طور کلی می توان نرم افزار های کامپیوتری را به دو دسته بدافزار^۱ و خوب افزار^۲ تقسیم کرد. منظور از خوب افزار همه نرم افزارهایی می باشد که توسط کاربر جهت کاربردی خاص بروی کامپیوتر مورد استفاده قرار گرفته و بدون اجازه کاربر توانایی اجرا نداشته باشد و در صورت اجرا هیچ گونه فعالیت خرابکارانه ای را انجام ندهند. از سوی دیگر بدافزار به نرم افزاری گفته می شود که بدون اجازه کاربر اجرا شده و فعالیت های خرابکارانه ای را انجام می دهد که این فعالیت ها می توانند شامل سرقت اطلاعات شخصی افراد، از بین بردن اطلاعات، تحت کنترل گرفتن سیستم کاربر، اخلال در عملکرد سیستم کاربر و... باشد.

به صورت کلی می توان گفت یک بدافزار مجموعه ای از کدها است که جهت انجام فعالیت مخرب در سیستم کاربر اجرا می شود. یکی از روش هایی که منجر به اجرای کدهای مخرب^۳ در سیستم کاربر می شود اجرای این کدها در قالب یک فایل اجرایی آلوده^۴ می باشد، به عبارت دیگر آن فایل اجرایی با آن کدهای مخرب آلوده شده و تبدیل به یک فایل اجرایی آلوده شده است.

در نرم افزارهای آنتی ویروس^۵ شناسایی فایل های اجرایی آلوده یا همان بدافزارهایی که از نوع یک فایل اجرایی می باشند به عنوان بخش مهمی از فعالیت های یک آنتی ویروس در نظر گرفته می شود . آنتی ویروس های موجود در شناسایی بدافزارهای جدیدی که قبل از این فعالیت آن ها ارایه نشده است

Malicious Software(Malware)^۱

Good Software^۲

Malicious Code^۳

Infected File^۴

Anti Viruss^۵

دچار مشکل هستند. همچنین شناسایی نمونه^۱های جدید بدافزار های قدیمی با دقت خوبی انجام نمی شود. در ادامه فصل به دسته بندی بدافزارها می پردازیم.

۲-۱ دسته بندی انواع بدافزارها

بدافزارها را می توان براساس نحوه عملکرد آن ها به دسته^۲ های مختلفی تقسیم کرد که هر کدام از دسته ها شامل نمونه^۳ های مختلفی می شود که ممکن است شامل زیر نمونه^۴ هایی نیز باشد. در ادامه مهمترین دسته ها و نمونه ها و زیر نمونه ها معرفی می شود.

۱-۲-۱ ویروس^۵

ویروس بدافزاری است که با تغییر دادن کد یک فایل اجرایی و اضافه کردن کدهای مخرب ، آن فایل را آلوده می کند و زمانی که آن برنامه اجرا شود، مجموعه کدهای مخرب نیز اجرا خواهد شد. این دسته از بدافزارها به نمونه های مختلفی تقسیم می شوند که در ادامه به توضیح آن ها می پردازیم.

۱-۲-۱-۱ ویروس فایل^۶

در این نمونه، از فایل های سیستمی سیستم عامل جهت منتشر و توزیع کردن کدهای مخرب استفاده می شود. این نمونه شامل این مواردمی شود: ویروس هایی که از فایل های اجرایی به عنوان میزبان^۷

Variant^۸

Category^۹

Type^{۱۰}

Subtype^{۱۱}

Virus^{۱۲}

File Virus^{۱۳}

استفاده کرده و آن ها را آلوده می کنند و یا یک کپی از یک فایل می سازند که این کپی نسخه ای از آن ویروس می باشد و یا ویروس هایی که خود را در مسیرهای مختلف کپی کرده و به ویروس هایی که فایل های سیستمی را آلوده کرده اند لینک می شوند. ویروس اسکریپت^۱ یک زیر مجموعه از این نوع ویروس ها می باشد که با یکی از زبان های اسکریپتی نوشته می شود.

۲-۱-۲-۱ ویروس بوت سکتور^۲

این نمونه، سکتور بوت یا یک رکورد بوت را آلوده می کند و یا سکتور بوت فعال را تغییر می دهد و زمانی که کامپیوتر بوت می شود ویروس نیز فعال می شود. تعداد زیادی از این نمونه ویروس زمانی که یکبار اجرا شوند مانع بوت شدن کامپیوتر می شوند.

۳-۱-۲-۱ ویروس ماکرو^۳

این نمونه با زبان ماکرو اسکریپتی نرم افزار های کاربردی مانند پردازش متن، حسابداری و... نوشته می شود که با تغییر دادن ویژگی های این نوع زبان ها از یک فایل آلوده به یک فایل سالم انتقال پیدا می کنند. اغلب این نوع ویروس ها برای نرم افزار های مجموعه آفیس^۴ وجود دارند و به دلیل اینکه این نوع ویروس ها با زبان یک نرم افزار کاربردی نوشته شده اند بروی هر بستر^۵ که آن نرم افزار اجرا شوند آن ویروس می تواند فعالیت کند به همین دلیل این نمونه از ویروس ها به بستر وابسته نیستند.

Host^۱

Script Virus^۲

Boot Sector Virus^۳

Macro Virus^۴

Microsoft Office^۵

Platform^۶

۴-۱-۲-۱ ویروس ایمیل^۱

این نمونه برخلاف موارد قبلی که براساس نحوه آلوده سازی تقسیم بندی شده اند براساس نحوه توزیع ویروس نام گذاری شده است. ایمیل می تواند برای انتقال همه نمونه های قبلی استفاده شود و زمانی که یکی از افرادی که این ایمیل را دریافت کرده اند فایل ضمیمه را اجرا کنند این ویروس از طریق ایمیل آن فرد به آدرس هایی که آن فرد ذخیره کرده است ارسال می شود.

۴-۱-۲-۱ ویروس چند نوعی^۲

در این نمونه هسته ویروس ها یکسان بوده ولی تغییراتی کمی در پیاده سازی اعمال شده است تا براحتی توسط آنتی ویروس ها شناسایی نشوند.

۲-۲-۱ کرم شبکه

یک بدافزار خود پخش^۳ است که بروی شبکه توزیع می شود و اغلب از طریق اینترنت گسترش پیدا می کند. برخلاف ویروس برای اجرا شدن، تکرار شدن و توزیع شدن به برنامه های دیگر و انجام کارهایی مانند باز کردن ایمیل آلوده وابسته نیست. کرم ها جهت یافتن حفره های امنیتی و نقاط آسیب پذیر در سیستم شبکه را اسکن کرده و با استفاده از آن ها پخش می شوند و خود را در سیستم های آسیب پذیر کپی می کنند. کرم ها اغلب براساس اینکه از چه محیطی برای پخش استفاده می کنند دسته بندی می شوند در برخی موارد نیز براساس سرعت پخش کرم ها نیز این دسته بندی انجام می شود که در ادامه نمونه های مختلف کرم ها را بررسی می کنیم:

Email Virus^۱

Multi-Variant Virus^۲

Self-Propagate^۳

کرم ایمیل که با آلوده کردن ضمیمه ایمیل پخش می شود و همچنین کرم هایی که با استفاده از کanal^۱ IM و یا IRC^۲ خود را پخش می کنند و به همین نام شناخته می شوند. کرم اینترنتی نیز که با استفاده پروتکل های وب مانند FTP^۳ و یا استفاده از صفحات وب و آلوده کردن سیستم هایی افرادی که به آن صفحه مراجعه می کنند، خود را پخش می کند همچنین کرم هایی نیز براساس سرعت پخش شدن دسته بندی می شوند به عنوان مثال کرم وارول^۴ که در کمتر از ۱۵ دقیقه سرورهای آسیب پذیر موجود در اینترنت را آلوده می کند. همچنین ممکن است که یک کرم با بکارگیری تعداد زیاد کرم دیگر فعالیت خود را گسترش دهد مانند کرم سوارم^۵ که به صورت هوشمند با تعداد زیادی کرم همکاری می کند.

یکی از تفاوت هایی که بین ویروس و کرم وجود دارد این است که ویروس های ابتدایی خود را به صورت فایل به فایل و با تکنیک تکرار فایل منتشر می کردند اما کرم ها از طریق اینترنت و به صورت میزبان به میزبان منتشر یک کرم می توانند از طریق یک لینک مهندسی اجتماعی شده که در یک ایمیل قرار دارد باشد که می توانند نقاط آسیب پذیر میزبان را شناسایی کند.

۱-۲-۳ اسب ترووا(تروجان)

یک بدافزار است که در قالب یک نرم افزار خوب و بدون خطر قرار گرفته است و پس از آن که تروجان برای اولین بار اجرا شد کنترل نرم افزار به مهاجم داده می شود و مهاجم می توانند عملیات خود را انجام

Instant Message^۱

Internet Relay Chat^۲

File Transfer Protocol^۳

Warhol Worm^۴

Swarm Worm^۵

Trojan Horse^۶

دهد. برای جلوگیری از شناسایی تروجان ها از تکنیک Process injection استفاده می شود که در آن پروسس تروجان با فعال شدن یک پروسس به ظاهر سالم فعال می شود و بدین ترتیب پروسس تروجان براحتی شناسایی نمی شود. در ادامه به بررسی نمونه های مختلف تروجان می پردازیم.

۱-۳-۲-۱ تروجان در پشتی^۱

نام دیگر این نوع تروجان دسترسی از راه دور نیز می باشد این نمونه یک دسترسی از راه دور بروی سیستم آلوده شده^۲ برای مهاجم فراهم می کند. برای مثال از این نمونه می توان تروجان DoS مطرح کرد که اگر گسترش و پراکندگی این تروجان مناسب باشد می توان از آن به عنوان بستری برای حمله های DDoS^۳ استفاده کرد.

۱-۳-۲-۲ تروجان جمع آوری اطلاعات^۴

در این نمونه اطلاعات به صورت پنهانی از سیستم آلوده شده جمع آوری و برای مهاجم ارسال می شود به این نوع تروجان Stealth ware گفته می شود. Spyware هم از این نمونه است که به آن Spybot نیز گفته می شود در این نمونه، تروجان به صورت مخفی جهت جمع آوری اطلاعات بر روی یک سیستم نصب می شود. Key logger نیز از این نمونه است که خود را در یک نرم افزار مرورگر و یا به عنوان درایور یک دستگاه خود را نصب می کند و داده هایی که با استفاده از کیبورد وارد می شود را کنترل می کند و اطلاعات آن را برای مهاجم ارسال می کند. Notifier نیز از این نمونه است که وظیفه‌ی

Backdoor Trojan^۱

Infected System^۲

Distributed Denial of Service^۳

Data-collection Trojan^۴

تایید کردن آلود شدن موفق سیستم را دارد و اطلاعات مانند IP سیستم و پورت های باز آن را ارسال می کند.

Rootkit ۳-۳-۲-۱

مجموعه ای از برنامه هاست که مهاجم جهت دسترسی به یک سیستم و پنهان ماندن و شناسایی نشدن استفاده می کند. این نمونه جهت پنهان کردن پروسس ها، فایل ها و اطلاعات رجیستری طراحی شده اند. به عبارت دیگر مهاجم جهت پنهان کردن خود و جلوگیری از ردیابی و همچنین پنهان کردن فعالیت هایی که انجام می دهد از این نوع تروجان استفاده می کند. دو نوع Rootkit وجود دارد که براساس نحوه نصب rootkit دسته بندی شده اند: ۱- روتکیت مد هسته ۲- روتکیت مد کاربر. روتکیت مد هسته به این دلیل که براحتی پنهان می شود و قوی تر است بیشتر مورد استفاده قرار می گیرد.

Botnet و Bot ۴-۳-۲-۱

لغت Bot از کلمه Robot مشتق شده است به معنی کارگر و کسی که کار می کند در کامپیوتر برای یک پروسس اتوماتیک نیز از این عبارت استفاده می شود Botها با اهداف زیادی در اینترنت فعالیت می کنند گستره کاربرد Botها از موتورهای جستجو گرفته تا Game Bots و IRC Channel Bot می باشد می توان گفت Botها یک فایل اجرایی هستند که در معرض خطر قرار می دهند، کنترل می کنند و میزبان را عضو Botnet می کنند. اینگونه نرم افزارها بدون اطلاع کاربر بروی سیستم آنها نصب می شود Botها می توانند مانند ویروس ها مخفی و مانند کرمها در شبکه پخش شوند.

Bot ها سیستم ها را آلوده می کنند و کنترل آن سیستم را به مهاجم^۱ می دهند تا از راه دور بتواند سیستم را کنترل کند و با ارسال دستورات برای Bot موجود در سیستم آلوده طبق دستورات به اهداف خویش برسد به این سیستم های آلوده Zombie نیز گفته می شود.

۱-۳ تعریف مسئله

در این پایان نامه هدف ارایه الگوی سیگنالی از فایل اجرایی جهت تشخیص فایل اجرایی آلوده از فایل اجرایی سالم است که منظور ما -در این پایان نامه- از فایل اجرایی، فایل اجرایی قابل حمل^۲ بیتی ویندوز می باشد. به عبارت دیگر می خواهیم برای هر فایل اجرایی یک الگوی سیگنالی ارایه کنیم تا براساس آن آلودگی فایل اجرایی به کد مخرب را تشخیص دهیم. این کد مخرب می تواند فعالیت های مخرب مختلفی مانند حذف یا خراب کردن فایل های سیستمی، دزدیدن اطلاعات کاربر، در اختیار مهاجم قرار دادن سیستم کاربر و ... باشد. لازم به ذکر است که فرض ما بر این است که فایل های اجرایی به صورت غیر فشرده^۳ می باشند و همچنین فایل اجرایی باید به صورت ایستا تحلیل شود و آلوده بودن آن به صورت ایستا تشخیص داده شود. تاکنون روشی بر مبنای پردازش سیگنال جهت تشخیص فایل اجرایی آلوده ارایه نشده است و این پایان نامه می تواند آغازی در ارایه روش های تشخیص آلودگی فایل اجرایی بر مبنای پردازش سیگنال باشد.

Bot Master^۱

Portable Executable (PE32)^۲

Unpack^۳

۴-۱ مقایسه فایل اجرایی سالم و آلوده

همانطور که گفته شد منظور از آلودگی یک فایل اجرایی این است که آن فایل شامل کدهای مخرب بوده که بدون ایجاد خلل در فعالیت های فایل میزبان آن کدها اجرا شده و به سیستم آسیب می زند. جهت شناسایی فایل سالم از فایل آلوده می باشد تفاوت های موجود بین این دو نوع فایل اجرایی بررسی شده و با توجه به این تفاوت ها راهکارهایی جهت شناسایی فایل های آلوده از فایل سالم ارایه کرد.

۴-۱-۱ اطلاعات سرآیند^۱ فایل اجرایی

در سرآیند فایل اجرایی اطلاعاتی مانند اندازه فایل، تعداد قسمت ها ، اطلاعات کنترلی و... وجود دارد که این اطلاعات محتوای کلی آن فایل را بیان می کند. پس از آنکه کد مخرب یک فایل سالم را آلوده کرد اطلاعاتی مانند اندازه آن فایل نیز تغییر می کند ولی اطلاعاتی که در سرآیند فایل آلوده وجود دارد همان اطلاعات فایل سالم است. به عبارت دیگر در فایل آلوده اطلاعات سرآیند با محتوی فایل مطابقت ندارد.

به عنوان مثال در فایل سالم اندازه فایل با مقداری که در سرآیند به عنوان اندازه فایل آمده است یکسان است اما در فایل های آلوده در بدافزارهایی که پس از ایجاد تغییرات در محتوای فایل سالم و آلوده کردن آن به کدهای مخرب اطلاعات سرآیند را تغییر نمی دهند بین مقادیر موجود در سرآیند فایل و مقادیر واقعی فایل تفاوت وجود دارد که براساس این تفاوت می توان فایل آلوده را شناسایی کرد یعنی اگر در فایلی اطلاعات سرآیند با محتوی فایل متفاوت باشد آن را به عنوان فایل اجرایی آلوده در نظر می گیریم.

لازم به ذکر است این تفاوت بین فایل سالم و آلوده بسیار رایج بوده و یکی از روش های شناسایی فایل های آلوده وجود این مقادیر متفاوت بوده است اما می توان گفت در بدافزارهای امروزی این عدم تطابق وجود ندارد. به عبارت دیگر پس از آلوده کردن فایل، اطلاعات سرآیند را با توجه به تغییراتی که فایل رخ داده است تغییر می دهد.

۲-۴-۱ تغییرات اشاره گر دستورالعمل^۱

تفاوت مهم دیگر بین فایل سالم و آلوده تغییرات اشاره گر دستورالعمل می باشد. در فایل سالم تغییرات این اشاره گر نرمال بوده و تغییرات شدیدی ندارد و از ابتدا تا انتهای یک روند دارد. اما در فایل های آلوده به این دلیل که کدهای مخرب به فایل سالم اضافه می شود تغییرات این اشاره گر نرمال نیست. به عنوان مثال معمولاً این اشاره گر در فایل آلوده ابتدا به انتهای فایل یعنی نقطه شروع کدهای مخرب پرس^۲ کرده و پس از اجرای آن دستورالعمل ها دوباره به شروع برنامه اصلی برگشته و آن ها را اجرا می کند به عبارت دیگر در فایل های آلوده در ابتدا اشاره گر دستورالعمل به انتهای فایل پرس کرده و سپس به شروع برنامه اصلی بر می گردد که این مورد اصلی ترین تفاوت در تغییرات اشاره گر بین فایل اجرایی آلوده و سالم می باشد که می توان از آن جهت تشخیص فایل سالم از آلوده استفاده کرد.

Instruction Pointer(IP)^۱

Jump^۲

۳-۴-۱ تفاوت در قسمت^۱ منبع^۲

فایل اجرایی از قسمت های مختلفی تشکیل شده است که پس از آلوده شدن یک فایل سالم تغییرات در قسمت منبع محسوس تر است. این قسمت که در انتهای فایل اجرایی قرار دارد در فایل سالم شامل اطلاعات متنی از قبیل توضیحات فایل اجرایی، آیکون، ورژن ، توضیحات کامپایلر و... می باشد و بخش اصلی آن را کاراکترهای متنی تشکیل می دهد.

برای آنکه یک فایل آلوده شود باید کدهای مخرب به آن فایل اضافه شوند بدین منظور کدهای مخرب باید به گونه ای به فایل سالم اضافه شوند که فایل میزبان خراب نشود و پس از آلوده شدن نیز به درستی اجرا شود. بدین منظور کدهای مخرب معمولا در انتهای فایل اجرایی یعنی قسمت آخر آن قرار می گیرند و با تغییر در شروع اجرای فایل ابتدا کد مخرب اجرا شده و پس از کد اصلی برنامه اجرا می شود. به طور کلی می توان گفت در فایل آلوده همواره بخشی از کدهای مخرب در این قسمت از فایل آلوده قرار دارند و براساس آن می توان با تحلیل این قسمت به روشهای جهت شناسایی فایل آلوده رسید.

Section^۱

Resource (.rsrc)^۲

فصل دوم

۲ برسی روش‌های تشخیص آلودگی فایل‌های اجرایی

۱-۲ مقدمه

جهت شناسایی یک فایل اجرایی آلوده در بین فایل های اجرایی سالم روش های مختلفی ارایه شده

است که می توان این روش ها را به دو دسته کلی تقسیم کرد:

۱-روش های مبتنی بر علامت های هویتی^۱ ۲-روش های مبتنی بر یادگیری و تخمین^۲

در دسته اول همانطور که از نام آن مشخص است روش ها براساس علامت^۳ ها و نشان هایی که در فایل ها آلوده وجود دارد کار می کنند و اگر در فایلی آن علامت وجود داشته باشد آن فایل به عنوان بدافزار در نظر گرفته می شود. به عبارت دیگر در این دسته برای هر نوع بدافزار یک علامت یا نشان در نظر گرفته شده و هر فایلی که آن علامت را داشته باشد یک بدافزار از آن نوع شناخته می شود و اگر فایلی هیچ کدام از علامت های انواع مختلف بدافزارها را نداشته باشد به عنوان یک فایل سالم در نظر گفته می شود. نقطه ضعف این روش این است که قدرت شناسایی بدافزارهای جدید را ندارد.

اما در دسته دوم، روش ها براساس ویژگی های کلی که در اکثر بدافزارها وجود دارد فایل های آلوده را شناسایی می کنند به عبارت دیگر در این دسته، روش ها به دنبال الگوهایی می گردند که در اکثر بدافزارها وجود داشته باشد و سپس هر فایلی که آن الگو را داشته باشد را به عنوان فایل آلوده در نظر می گیرند همچنین در این روش ها اگر فایلی الگویی شبیه الگوی فایل آلوده داشته باشد به عنوان یک فایل مشکوک به آلوده بودن شناخته می شود.

در آنچه ویروس های امروزی برای شناسایی فایل های آلوده از هر دو دسته روش مذکور به صورت ترکیبی استفاده می شود و روش های دسته دوم تلاش می کنند ضعف موجود در دسته اول یعنی همان

Signature Base^۱

Heuristic^۲

Tag^۳

عدم شناسایی بدافزارهای جدید را پوشش دهنده. در ادامه به بررسی این دو دسته روش می پردازیم و سپس بطور کلی این دو دسته روش را باهم مقایسه می کنیم.

۲-۲ روش‌های مبتنی بر علامت‌های هویتی

همانطور که در ابتداء گفته شد نکته اصلی در این روش‌ها بdst آوردن علامت‌های هویتی^۱ بازای هر نوع بدافزار می باشد که این علامت‌ها در پایگاه داده ای ذخیره شده و برای بررسی هر فایل به این پایگاه داده مراجعه می شود و در فایل ورودی علامت‌های موجود در پایگاه داده جست و جو می گردد. اگر هر یک از علامت‌ها در فایل پیدا شود آن فایل به عنوان فایل آلوده در نظر گرفته می شود. می توان گفت اصلی ترین چیزی که قدرت تشخیص این دسته از روش‌ها را تحت تاثیر قرار می دهد بروز بودن پایگاه داده است. نکته دیگر درمورد این روش‌ها نحوه بdst آوردن علامت‌های هویتی است که به آن نحوه تحلیل بدافزارها نیز گفته می شود. در واقع با این روش به دنبال نشان‌ها و علامت‌هایی می گردیم که در صورت مشاهده آن علامت‌ها در فایل‌های دیگر آن فایل‌ها را به عنوان بدافزار بشناسیم که بعد از بdst آوردن این علامت‌ها، می بایست آن‌ها را به پایگاه داده اضافه کرد تا سیستم بتواند بدافزارهایی از آن نوع را شناسایی کند.

۱-۲-۲ نحوه تحلیل بدافزارها

برای آنکه روش‌های این دسته بتوانند بدافزارهای بیشتری را شناسایی کنند می بایست پس از شروع فعالیت بدافزارهای جدید، آن بدافزارها را تحلیل کرد. همانطور که گفته شد منظور از تحلیل بدافزارها

Signature^۱

نحوه بدست آوردن علامت های هویتی برای بدافزارها می باشد به دو روش می توان بدافزارها را تحلیل کرد که در ادامه به بررسی این روش ها می پردازیم.

۱-۲-۱ تحلیل ایستا^۱

در این روش جهت بدست آوردن علامت های بدافزارها مجموعه دستورالعمل های تشکیل دهنده یک بدافزار مورد تحلیل قرار می گیرد. در واقع برای تحلیل یک بدافزار به بررسی دستورالعمل های آن پرداخته و وجود مجموعه دستورالعمل ها در یک فایل را به عنوان یک علامت برای آلوده بودن آن فایل در نظر گرفت. که این کار برای انواع بدافزارها انجام شده و برای هر نوع از بدافزارها مجموعه ای از دستورالعمل ها به عنوان علامت در نظر گرفته شده و وجود آن دستورالعمل ها در جای مشخص در یک فایل، نشان دهنده آلوده بودن فایل می باشد.

۲-۱-۲ تحلیل پویا^۲

در این روش هدف بررسی رفتاری بدافزارها می باشد. به عبارت دیگر برای بدست آوردن علامت مشخص که نشان دهنده آلوده بودن یک فایل باشد فعالیت هایی که یک فایل در زمان اجرا انجام می دهد مورد بررسی قرار می گیرد. برای بدست آوردن علامت های هویتی ابتدا بدافزار در یک محیط امن اجرا می شود سپس فعالیت های بدافزار در زمان اجرا زیرنظر گرفته می شود حال در بین فعالیت های حین اجرای یک بدافزار به دنبال رفتاری می گردیم که می توان آن را به عنوان یک علامت درنظر گرفت و گفت که اگر فایل دیگری چنین رفتار مشابه ای را داشته باشد آن فایل نیز یک بدافزار است این علامت ها نیز به پایگاه داده اضافه خواهد شد.

Static Analysis ^۱

Dynamic Analysis ^۲

۲-۲-۲ نقاط ضعف

در روش های مبتنی بر علامت های هویتی روش کار به این شکل است که با مراجعه به پایگاه داده یک علامت انتخاب شده و در فایل مورد نظر در جای مشخص شده آن علامت جست و جو می‌گردد. در صورت وجود هر یک از علامت های موجود در پایگاه داده در فایل مورد بررسی، آن فایل به عنوان یک بدافزار شناخته می‌شود. می‌توان نتیجه گرفت اگر فایل مورد بررسی یک بدافزار جدید باشد که قبل از نیامده است این فایل سالم در نظر گرفته نمی‌شود. همچنین اگر علامت های بدافزار جدید بدست آمده باشد اما در پایگاه داده وجود نداشته باشد باز هم این فایل سالم در نظر گرفته نمی‌شود.

می‌توان گفت این دسته از روش ها برای تشخیص بدافزارهایی که علامت آن ها بدست آمده است کاملاً به بروز بودن پایگاه داده وابسته است و حتی در مورد بدافزارهایی که نوع^۱ جدیدی از بدافزارهای قدیمی - که علامت آن ها در پایگاه وجود دارد - هستند این دسته از روش ها با مشکل مواجه می‌شوند. همچنین توانایی تشخیص بدافزارهای جدید در این دسته از روش ها وجود ندارد پس نتیجه می‌گیریم که با استفاده از روش های این دسته هرگز نمی‌توان به قدرت تشخیص ۱۰۰ درصد رسید به همین دلیل در آنتی ویروس ها همراه با روش های این دسته از روش های مبتنی بر یادگیری و تخمین استفاده می‌شود تا بتواند قدرت تشخیص را در آنتی ویروس ها افزایش دهد. اما امروزه با وجود استفاده از روش های مبتنی بر یادگیری و تخمین هنوز دقت آنتی ویروس ها به نقطه مطلوب نرسیده است به عنوان نمونه براساس آزمایشی که سایت^۲ مقایسه نرم افزارهای آنتی ویروس در سال ۲۰۰۸ انجام داده است نشان داد که قدرت تشخیص ۱۶ آنتی ویروس مطرح امروزی که پایگاه داده آن ها بروز رسانی نشده بود حداقل به

Variant^۱

www.av-comparatives.org^۲

۷۰ درصد رسیده است. پس می توان نتیجه گرفت جهت افزایش دقت آنتی ویروس ها می بایست بر روی بهبود روش های مبتنی بر یادگیری و تخمین کار کرد. در ادامه به بررسی این دسته از روش ها می پردازیم.

۳-۲ روش های مبتنی بر یادگیری و تخمین

در این دسته از روش ها تلاش می شود با استفاده از روش های اتوماتیک، بدافزارهایی که توسط روش های مبتنی بر علامت های هویتی قابل شناسایی نیستند را تشخیص داد به عبارت دیگر این دسته از روش ها به عنوان مکمل دسته قبلی جهت رسیدن قدرت تشخیص به ۱۰۰ درصد مورد استفاده قرار می گیرند. برای این کار می بایست در بین بدافزارها به دنبال الگو^۱ هایی بگردیم که در آن ها مشترک باشد. می توان روش های موجود در این دسته را براساس نحوه تحلیل و بررسی فایل اجرایی به چند زیر دسته کلی تقسیم کرد که در ادامه به بررسی آن ها می پردازیم.

۳-۳-۱ تحلیل براساس n بایت استخراج شده از فایل اجرایی^۲

در این دسته، فایل اجرایی به صورت یک فایل متنی که از ۰ و ۱ تشکیل شده است مورد بررسی قرار می گیرد به عبارت دیگر ما مجموعه ای از فایل های اجرایی را به صورت فایل های متنی مت Shankل از ۰ و ۱ تبدیل کرده و سپس صورت مسئله مانند دسته بندی فایل های متنی می شود. فایل های متنی ما به دو دسته سالم و آلوده تقسیم می شود. در دسته بندی متون، کلمه به عنوان ویژگی اصلی مورد استفاده قرار می گیرد که در اینجا دنباله n بایت ها مانند کلمه در دسته بندی متون در نظر گرفته می شود.

Pattern^۱

Byte N-gram^۲

سپس می توان الگوریتم های دسته بندی مانند K-NN و روش های وزن دهی که در دسته بندی متون استفاده می شود را بروی این مسئله اعمال کرد. نکته ای که در اینجا حائز اهمیت است این است که در این روش ها فایل اجرایی به صورت معنی دار بررسی نمی شود. به عبارت دیگر محتوای فایل اجرایی و اینکه مثلا از چه دستورالعمل هایی در آن فایل اجرایی استفاده شده است بررسی نمی شود و فقط از دنباله ای از ۰ و ۱ ها استفاده می شود. می توان انتظار داشت که این روش به دلیل اینکه به صورت معنی داری با فایل اجرایی برخورد نمی کند دارای سرعت خوبی باشد. به عنوان نمونه در [۱] از روش n بایت به عنوان روش استخراج ویژگی استفاده شده و برای دسته بندی نیز روش نزدیک ترین همسایه^۱ مورد استفاده قرار گرفته است؛ پایگاه داده ای که این روش بر روی آزمایش شد مشتمل بر ۶۵ بدافزار و فایل سالم بود که با دقت ۹۸ درصد فایل های آلوده مشخص شدند. در [۲] نیز از همین روش برای استخراج ویژگی استفاده شده است.

۲-۳-۲ براساس اطلاعات سرآیند فایل اجرایی^۳

در این دسته همانطور که از نام آن پیدا است شناسایی فایل های آلوده براساس اطلاعات موجود در سرآیند فایل انجام می شود. اطلاعات موجود در سرآیند از قبیل اطلاعات فیزیکی مثل زمان ایجاد فایل، سایز فایل و... و همچنین اطلاعات ساختار منطقی مثل اطلاعات مربوط به قسمت ها و... و اطلاعات مربوط به DLL ها توابع استفاده شده در برنامه و موارد دیگر جهت شناسایی فایل های اجرایی آلوده می توان استفاده کرد. به عنوان مثال زمانی که فایل اجرایی سالم آلوده می شود کد های مخرب به انتهای فایل سالم اضافه شده و اندازه فایل زیاد می شود. نکته ای که باعث تشخیص فایل آلوده می شود این

K-NN^۱

Accuracy^۲

Portable Executable Header Information^۳

است که اندازه فایل زیاد شده ولی اطلاعات درون سرآیند که شامل اندازه فایل می شود تغییری نکرده است که این تفاوت منجر به شناسایی فایل اجرایی آلوده می شود.

۳-۲-۳ براساس کد عملگر فایل اجرایی^۱

در زبان ماشین دستورالعمل ها شامل کد عملگر^۲ و یک یا بیشتر عملوند^۳ می شوند که براساس کد عملگر های موجود در فایل اجرایی می توان به آلودگی آن پی برد. در این دسته از روش ها مانند دسته قبل فایل اجرایی به صورت معنی دار مورد بررسی قرار می گیرد. در این دسته کد عملگر استفاده شده در فایل اجرایی بررسی می شود. به عبارت دیگر از فایل اجرایی کد عملگرهای استفاده شده در آن را استخراج می کنیم. در ادامه می توانیم بر اساس روش های مختلف آن فایل را تحلیل کنیم. مثلا مجموعه کدهای عملگر فایل اجرایی را به عنوان یک فایل متنی در نظر گرفته و در ادامه مانند دسته بندی متون عمل کنیم. یکی از کارهایی که می توان انجام داد این است که کد عملگر را به عنوان کلمه در نظر گرفته و دنباله ای از n کلمه(کد عملگر) را به عنوان ویژگی در نظر بگیریم و سپس الگوریتم های دسته بندی را روی آن ها اعمال کنیم.

به عنوان نمونه در [۳] با استفاده از دنباله دستورالعمل های مفید برای تشخیص که از روش های اتوماتیک بدست آمده بودند با دقت ۹۸.۴ درصد فایل های آلوده شناسایی شدند. همچنین در [۴] نحوه توزیع تکرار کد عملگر ها در فایل های آلوده و سالم مورد تحلیل قرار گرفت و این نتیجه حاصل شد که توزیع کد عملگر ها (خصوصیت کد عملگرهای کم کاربرد) در فایل آلوده نسبت فایل سالم بسیار متفاوت

OpCode Information ^۱

Opcode(Operational Code) ^۲

Operand ^۳

است و در [۵] نیز با استفاده از n کد عملگر^۱ به عنوان ویژگی و بررسی ۵ روش دسته بندی و آزمایش بروی مجموعه ای با بیش از ۳۰۰۰۰ فایل سالم و آلوده دقت ۹۹ درصد بدست آمده است. نکته ای در این دسته حائز اهمیت است این است که در برخی موارد بخصوص در فایل های آلوده نمی توان به طور کامل به دستورالعمل های تشکیل دهنده آن فایل اجرایی را بدست آورد و به عبارت دیگر نمی توان گفت همواره دستورالعمل های بدست آمده از یک فایل اجرایی به طور ۱۰۰ درصد درست است که این موضوع نقطه ضعف بزرگی برای این دسته به حساب می آید.

Opcode n-gram ^۱

فصل سوم

۳ روشهای تخلیل فایل های اجرایی آلوده

۱-۳ مقدمه

پس از بررسی تفاوت های موجود بین فایل اجرایی آلوده و فایل اجرایی سالم و همچنین بررسی کارهای انجام شده در تشخیص فایل های آلوده، در این فصل می خواهیم تحلیل فایل های اجرایی آلوده را از جنبه های مختلف بررسی کنیم. در واقع در این فصل تلاش می کنیم بازی هر فایل اجرایی یک الگوی سیگنالی داشته باشیم که بتوان براساس آن فایل های آلوده را شناسایی کرد. برای این کار باید با توجه به تفاوت های فایل های سالم و آلوده به تحلیل فایل های اجرایی پرداخت. بطور کلی دو تحلیل برای ارایه یک الگوی سیگنالی بازی هر فایل اجرایی در این فصل مورد بررسی قرار گرفته است: ۱- براساس کد اسمنبلی ۲- براساس بایت های فایل اجرایی

در دسته اول ابتدا فایل اجرایی به کد اسمنبلی تبدیل می شود و سپس بروی آن تحلیل های مختلفی انجام می شود پس در این دسته ابتدا باید بروی فایل اجرایی یک عملیات پیش پردازش^۱ انجام داد که از نظر زمانی برای سیستم ما یک نقطه ضعف محسوب می شود. علاوه بر این کد اسمنبلی بدست آمده به طور کامل دقیق نیست زیرا نرم افزار های معروف تبدیل فایل اجرایی به اسمنبلی^۲ مانند آیدا^۳ نیز در برخی موارد اشتباه می کنند.

در دسته دوم فایل اجرایی براساس بایت های تشکیل دهنده فایل که در مبنای ۱۶ می باشند تحلیل می شوند و تقریباً می توان گفت نسبت به دسته اول پیش پردازش ندارد و در نتیجه دسته دوم بسیار سریعتر از دسته اول می باشد و با توجه به اهمیت زمان در سیستم تشخیص فایل آلوده این امر نقطه قوت بزرگی برای این دسته به حساب می آید.

Preprocessing^۱

Disassembler^۲

IDA^۳

۲-۳ براساس تحلیل کد اسمنبلی

در این دسته محتوای با معنی فایل اجرایی مورد تحلیل قرار می‌گیرد و می‌خواهیم براساس اینکه چه کدهای اسمنبلی در فایل اجرایی جود دارد آلودگی فایل را تشخیص دهیم و براساس کد اسمنبلی بازی هر فایل اجرایی یک الگوی سیگنالی داشته باشیم. برای این کار اولین مرحله تبدیل فایل اجرایی به کد اسمنبلی می‌باشد تا در ادامه تصمیم بگیریم چگونه کد اسمنبلی را تحلیل کنیم. در تبدیل فایل اجرایی به کد اسمنبلی با دو مشکل روبرو می‌شویم:

۱- زمانبر بودن ۲- دارای خطابودن

در تبدیل فایل اجرایی به کد اسمنبلی زمان تبدیل، به اندازه فایل اجرایی بستگی دارد و هرچه فایل بزرگتر باشد در پیش پردازش زمان بیشتری صرف می‌شود علاوه بر این ممکن است کد اسمنبلی بدست آمده از یک فایل اجرایی دارای اشتباه باشد که این امر در فایل‌های آلوده بیشتر دیده می‌شود. علت مشکل دوم این است که در بدافزارها بروی قسمت کد مخرب تکنیک‌های پنهان‌سازی^۱ اعمال می‌شود که این امر موجب می‌شود کد اسمنبلی بدست آمده از یک فایل اجرایی دقیق نباشد و برخی دستورات به اشتباه بدست آمده باشند.

نکته دیگری که در این دسته حائز اهمیت است وجود کدهای اضافی و بدون استفاده در فایل‌های آلوده می‌باشد. در فایل‌های آلوده جهت جلوگیری از تحلیل درست فایل، کدهای بی‌فایده و اضافی زیادی به کد مخرب اضافه می‌شود مثلاً حلقه با تعداد تکرار بالا که در پایان حلقه هیچ مقداری تغییر نکرده است و یا جمع و تفریق یک مقدار با مقادیر مختلف و جایجایی یک مقدار که در نهایت آن مقدار نیز تغییر نمی‌کند.

¹Obfuscation

پس از بدست آوردن کد اسambilی رویکردهای متفاوتی را می توان جهت بدست آوردن الگوی سیگنالی از آن اتخاذ کرد که در اینجا دو رویکرد مورد بررسی قرار می گیرد: ۱- براساس کد عملگر ۲- براساس تغییرات اشاره‌گر دستورالعمل که در ادامه به بررسی آن ها می پردازیم.

۱-۲-۳ براساس کد عملگر

همانطور که گفته شد پس از تبدیل فایل اجرایی به کد اسambilی با جدا کردن کد عملگرها از عملوند ها ادامه تحلیل بروی کد عملگرها انجام می شود. چند روش می توان برای تحلیل کد عملگرها ارایه کرد. در اولین روش تعداد تکرار کد عملگر به عنوان الگوی سیگنالی در نظر گرفته می شود که در آن تعداد تکرار هر کدام از کد عملگرها مانند Move, Add, Sub,... براساس آن نماینده ای برای آن فایل معرفی می شود. در این روش می توان فقط تعدادی از کد عملگرها را در نظر گرفت. به عبارت دیگر پس از آنکه تعداد تکرار همه کد عملگرها بدست آمد از بین آن ها کد عملگر هایی که در تشخیص فایل سالم از آلوده نقشی ندارد را حذف کرد و سپس براساس تعداد تکرار کد عملگرهای باقی مانده الگوی سیگنالی فایل را بدست آورد.

رویکردنی دیگری که می توان در این دسته ارایه کرد این است که تعداد تکرار کلاس های مختلف عملگرها را در نظر گرفت. به عبارت دیگر ابتدا باید کد عملگرها را کلاس بندی کرد، که برای این کار می توان از دسته بندی ارایه شده توسط شرکت سازنده پردازنده استفاده کرد و پس از آن بازای دیده شدن هر کد عملگر یک واحد به کلاسی که آن کد عملگر عضو آن می باشد اضافه کرد و سپس تعداد تکرار هر کلاس مدنظر قرار بگیرد.

با توجه به هزینه بر بودن تبدیل فایل اجرایی به کد اسambilی و همچنین عدم دقت ۱۰۰ درصد در این تبدیل، این دسته از روش ها در پایان نامه مورد استفاده قرار نگرفت.

۲-۲-۳ براساس تغييرات اشاره گر دستورالعمل

همانطور که گفته شد در فایل آلوده تغييرات اشاره گر دستورالعمل نسبت به تغييرات اين اشاره گر در فایل سالم متفاوت می باشد که دليل آن اضافه شدن کد مخرب به فایل سالم است. بر همين اساس می توان تغييرات مقادير اشاره گر دستورالعمل را با توجه به زمان به عنوان الگوي سيگنالي هر فایل اجرايی در نظر گرفت. به عبارت ديگر از زمان آغاز اجرا تا پایان اجرای دستورالعمل هاي يك فایل اجرائي مقادير اين اشاره گر ذخیره می شود و نمودار اين تغييرات بر حسب زمان رسم می شود.

جهت بدست آوردن مقادير اشاره گر دستورالعمل دو رویکرد می توان ارایه کرد. در رویکرد اول فایل اجرايی در يك محبيط شببيه سازی شده اجرا شده و در هر لحظه مقادير اين اشاره گر ذخیره شود و سپس نمودار مقادير مختلف بر حسب زمان رسم می شود. چند نکته در باب اين رویکرد قابل ذكر است اول اينكه ايجاد محبيط شببيه سازی شده و امن که بتوان در آن فایل را اجرا کرد هزينه بر است. دوم اين که هر فایل يکبار به طور كامل اجرا شده و سپس مقادير اشاره گر آن فایل تحليل شود اين کار بسيار زمان بر است و با توجه به اينكه زمان در تشخيص آلودگی فایل بسيار مهم است اين روش مفيد نیست. مازاد بر موارد ذكر شده فرض ما بر اين است که باید فایل به صورت ايستا تحليل شود.

در رویکرد دوم ابتدا باید فایل اجرائي به کد اسمبلي تبديل شود و سپس برای بدست آوردن تغييرات مقادير اشاره گر دستورالعمل از ابتداي برنامه شروع به بررسی کد عملگر کرده و با توجه به آن مقدار اين اشاره گر را بدست آوريم و اين کار را تا پایان برنامه ادامه دهيم که در اين روش نيز باید ازا ابتدا تا انتهای فایل اجرائي بررسی شود و بسيار زمان بر است و علاوه بر اينكه ماهيت اين روش زمان بر است در فایل هاي آلوده به دليل وجود بسيار زياد کدهاي زايد و حلقه هاي با تكرار بالا و تو در تو اين روش را از نظر زمانی با مشكل مواجه می کند. همچنين با توجه به اينكه صحيح بدست آوردن مقادير اشاره گر بسيار وابسته به تبديل درست و كامل فایل اجرائي به کد اسمبلي می باشد اما همانطور که گفته شد به دليل پنهان سازی در فایل هاي آلوده اين امر امكان پذير نیست.

۳-۳ براساسس بايت های فايل اجرائي

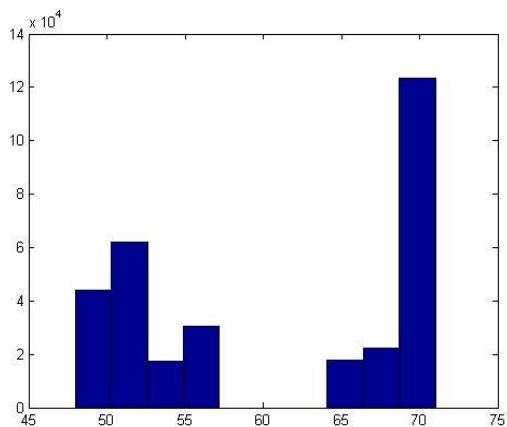
در اين دسته محتواي بدون معنى فايل اجرائي تحليل مى شود به عبارت ديگر مجموعه بايت هاي تشکيل دهنده فايل اجرائي که در ظاهر نمی توان از آن ها به آلودگی فايل رسيد معياري برای تشخيص آلودگی فايل مى شود. کارهايي که تاکنون در اين دسته انجام شده است بروي کل فايل اجرائي اعمال شده است که منجر به در بر گرفتن بخش زيادي از بايت ها مى شود که در تشخيص آلودگی فايل نقشی ندارند.

بر همين اساس ما به دنبال ارایه روشی بوديم تا براساسس آن بخشی از فايل اجرائي که بيشترین تغييرات را در هنگام آلودگی دارد؛ مورد تحليل قرار گرفته و براساسس آن الگوي سيگنانالي فايل را بدست آوريم. با توجه به اينکه همواره بخشی از کد مخرب در قسمت آخر فايل اجرائي يعني قسمت منبع قرار دارد به جای آن که کل فايل اجرائي تحليل شود فقط قسمت منبع مورد تحليل قرار مى گيرد. در اين روش، مرحله پيش پردازش شامل جداسازی قسمت آخر فايل اجرائي و ذخیره آن در يك فايل متنی است که شامل بايت هاي اين قسمت در مبناي ۱۶ است. همچنین جهت تحليل راحت تر فايل هاي متنی بين هر دو بايت از حرف G به عنوان جدا کننده استفاده كردیم. در ادامه کليه تحليل ها بروي قسمت آخر فايل های اجرائي سالم و آلوده که در قالب فايل متنی ذخیره شده اند انجام مى شود.

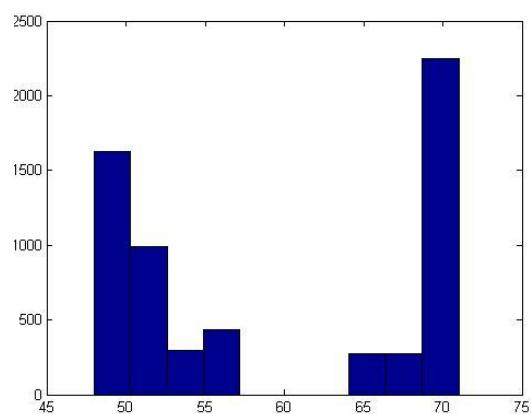
پس از جداسازی قسمت آخر فايل اجرائي، با دو رویکرد به تحليل اين قسمت پرداختيم: ۱- براساسس بايت های تشکيل دهنده ۲- براساسس ترکيب بايت ها. در ادامه هر کدام از روش ها به طور كامل ارایه مى شود.

۱-۳-۳ براسas بايت هاي تشکيل دهنده

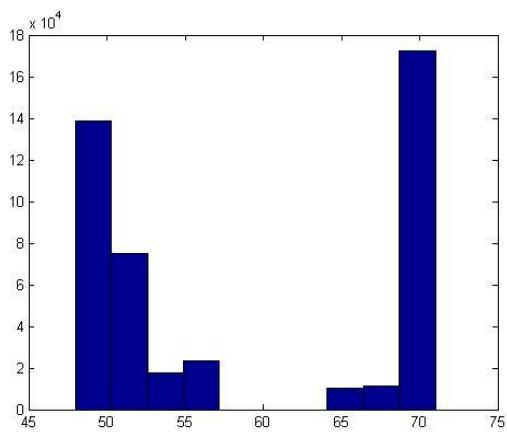
در اين روش جهت بدست آوردن الگوي سيگنالي، تعداد تكرار هر کدام از بايت ها مد نظر قرار می گيرد همانطور که گفته شد مجموعه بايت ها در مبنای ۱۶ شامل ارقام ۰-۹ و حروف A-G که G نقش جدا کننده برای ۲ بايت را دارد می شوند. برای راحتر شدن کار از کد اسکی آن ها که بين ۵۷تا۴۸ برای ارقام و ۷۱تا۶۵ برای حروف می باشد؛ استفاده می کنيم. جهت بررسی روش ارایه شده چند نمونه فایل سالم و آلوده در نظر گرفته شد و سپس با جدا کردن قسمت آخر آن فایل ها به بررسی نمودار فرکانس بايت های تشکیل دهنده قسمت آخر پرداختیم.



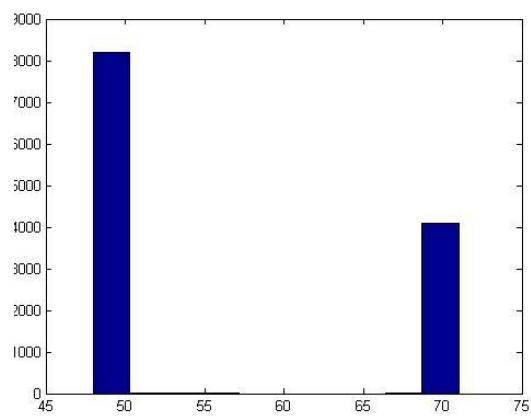
الف-۲



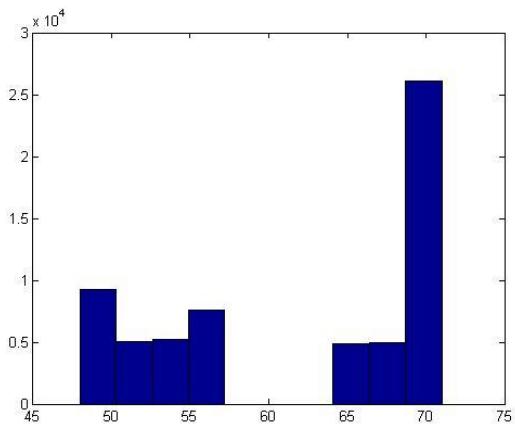
الف-۱



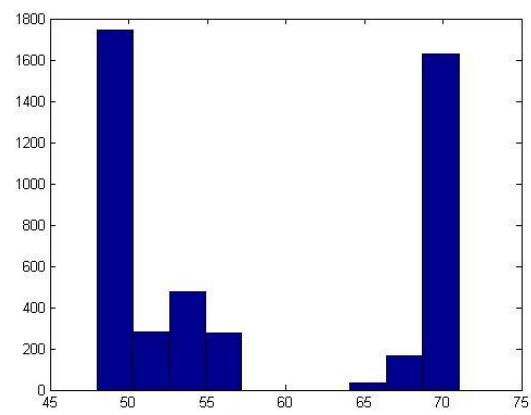
۲- β



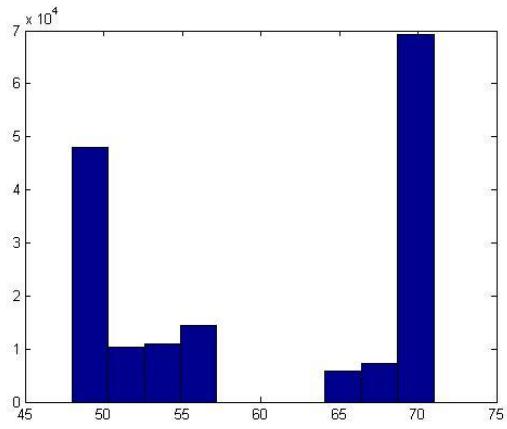
۱- β



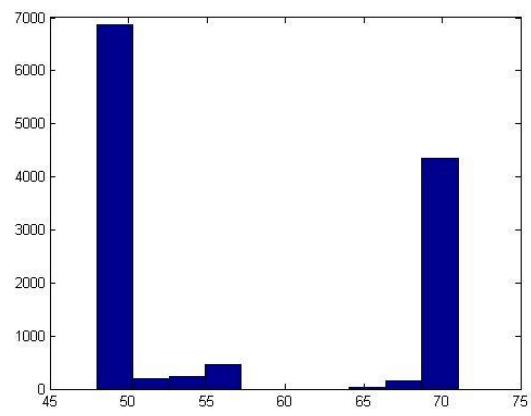
۲- γ



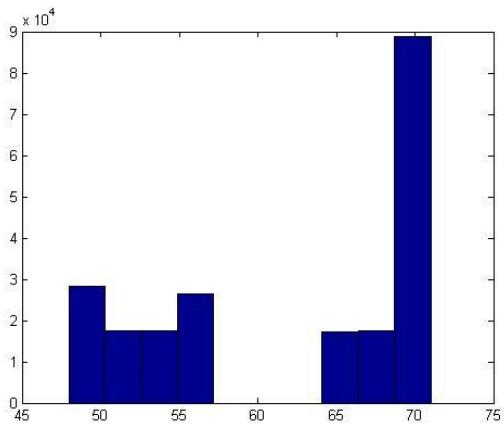
۱- γ



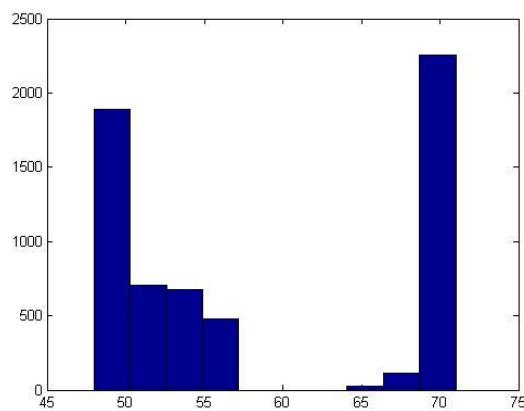
۲- ϕ



۱- ϕ



۲-ی



۱-ی

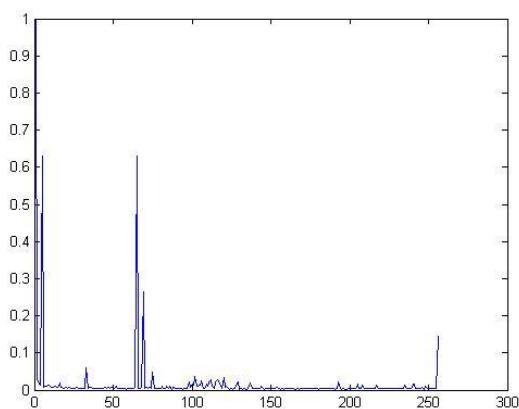
شکل (۱-۳) نمودار فرکانس بایت های تشکیل دهنده قسمت آخر فایل اجرایی. (الف-۱)، (ب-۱)، (ج-۱)، (پ-۱) (ی-۱). برای فایل های اجرایی سالم. (الف-۲)، (ب-۲)، (ج-۲)، (پ-۲) (ی-۲)، برای فایل های اجرایی آلوه

در شکل (۱-۳) نمودار تعداد تکرار بایت های تشکیل دهنده قسمت آخر فایل اجرایی آمده است در این نمودارها محور افقی شامل کد اسکی بایت های تشکیل دهنده این قسمت می باشد و محور عمودی نیز تعداد تکرار آن بایت ها را نشان می دهد. همانطور که گفته شد به این دلیل تحلیل بروی قسمت آخر فایل اجرایی انجام شد که همواره بخشی از کد مخرب در این قسمت از فایل قرار دارد پس باید ارتباطی بین تکرار بایت ها و وجود کد مخرب در فایل آلوه پیدا کرد و یا اینکه ارتباطی بین وجود کاراکتر متنی در فایل سالم و تعداد تکرار بایت ها در فایل سالم دیده شود. با توجه به دو ارتباط مطرح شده به این نتیجه رسیدیم تا به جای بررسی تعداد تکرار ۱ بایت، به بررسی تکرار ترکیب بایت ها بپردازیم.

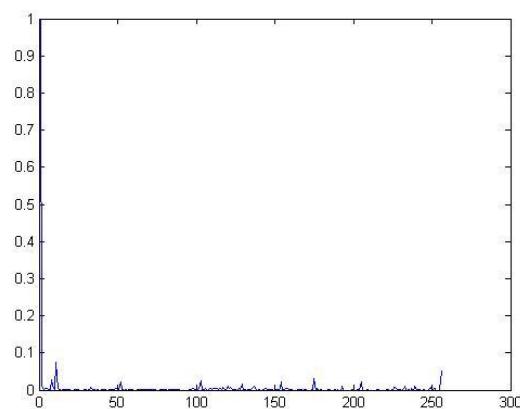
۲-۳-۳ براساس ترکیب بایت های تشکیل دهنده

پس از بررسی نمودار تکرار بایت ها در قسمت آخر فایل اجرایی، مشخص شد براساس آن نمی توان فایل آلوه را تشخیص داد. در این قسمت هدف بررسی ترکیب های مختلف بایت ها جهت رسیدن به بهترین ترکیب می باشد. به عبارت دیگر با توجه به اینکه مبنای تحلیل قسمت آخر فایل اجرایی بر این

اساس است که در قسمت منبع فایل سالم بیشتر کاراکترهای متنی وجود دارد ولی این قسمت در فایل آلوده شامل کدهای مخرب می باشد؛ تعداد تکرار ترکیبی از بایت ها را بررسی کنیم که این اختلاف بین فایل سالم و آلوده را بهتر نشان دهد. با توجه به اینکه کاراکترهای متنی به صورت ۲ بایت در کد اسکی^۱ نشان داده می شوند و همچنین ۲ بایت نماینده بهتری برای کدهای مخرب می باشد؛ اگر تعداد تکرار ۲ بایت ها الگوی سیگنالی یک فایل اجرایی را تشکیل دهد جهت تشخیص آلودگی یک فایل مناسب تر است. در مبنای ۱۶ دو بایت ۲۵۶ حالت مختلف می توانند کنار هم قرار بگیرند و در نتیجه الگوی سیگنالی یک فایل اجرایی نشان دهنده فرکانس این ۲۵۶ حالت مختلف در قسمت آخر فایل اجرایی می باشد.

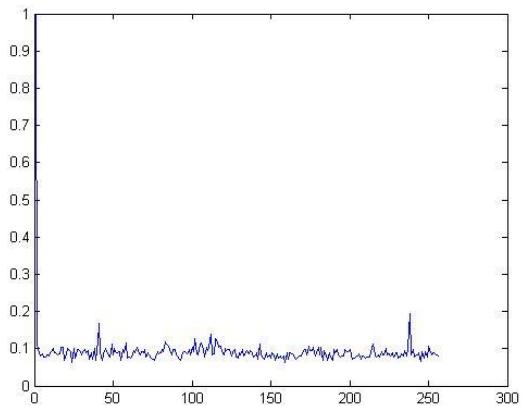


الف-۲

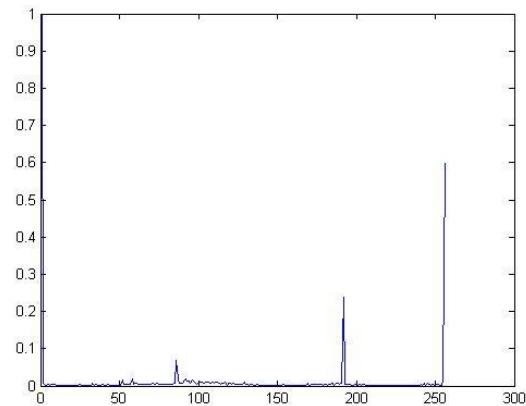


الف-۱

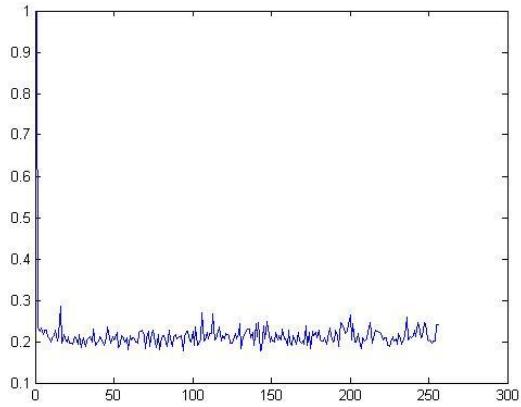
ASCII^۱



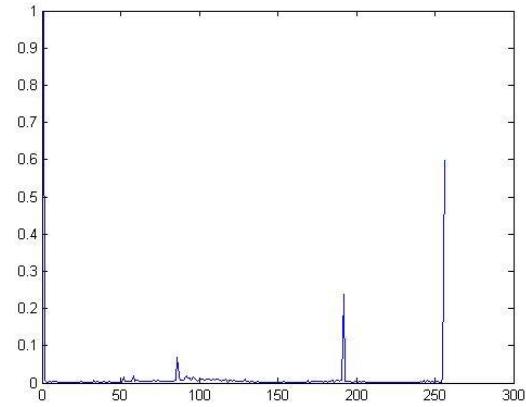
۲- β



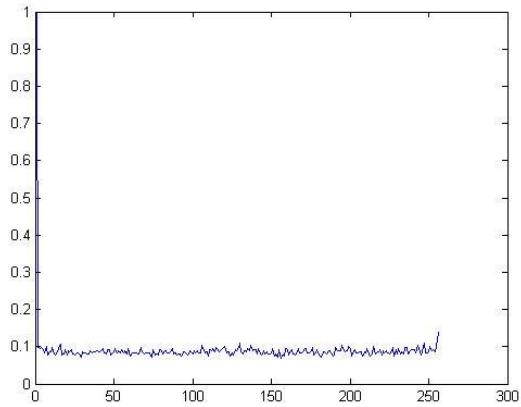
۱- β



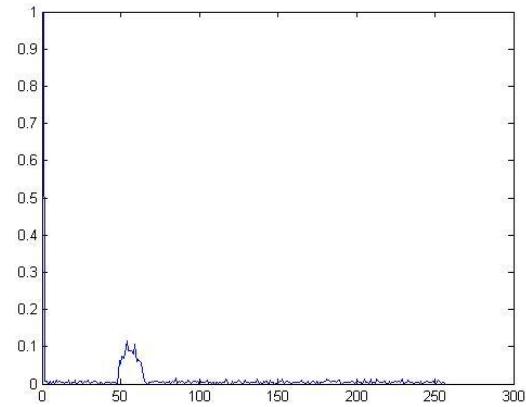
۲- γ



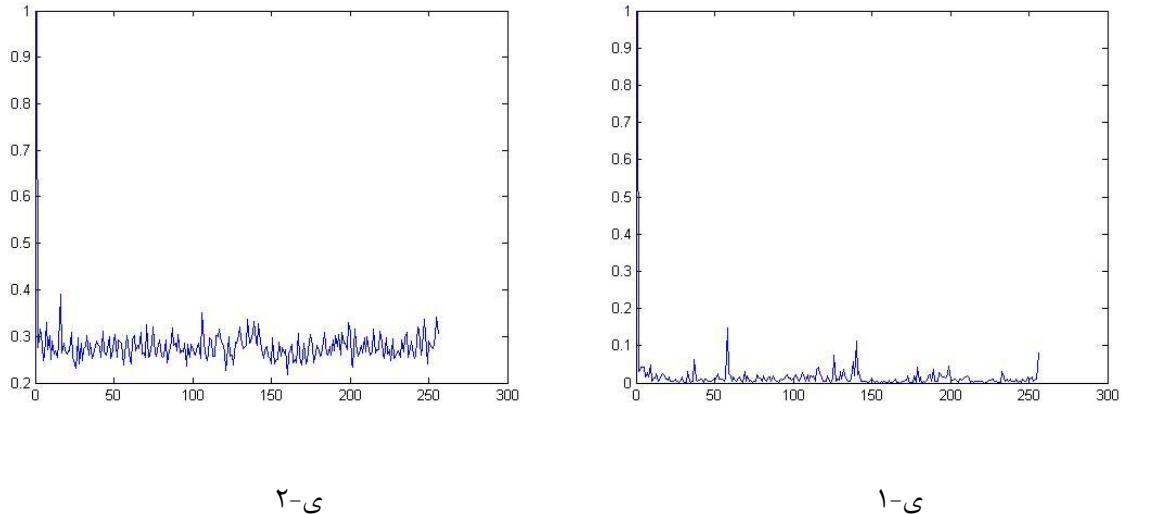
۱- γ



۲- ζ



۱- ζ



شکل (۲-۳) نمودار تعداد تکرار ۲ بایت ها در قسمت آخر فایل اجرایی. (الف-۱)، (ب-۱)، (ج-۱)، (پ-۱) (ی-۱). برای فایل های اجرایی سالم. (الف-۲)، (ب-۲)، (ج-۲)، (پ-۲) (ی-۲)، برای فایل های اجرایی آلوده

در شکل (۲-۳) نمودار تعداد تکرار ۲ بایت ها در قسمت آخر فایل اجرایی آمده است. در این نمودارها محور افقی نشان دهنده ۲۵۶ حالت مختلف ۲ بایت ها و محور عمودی تعداد تکرار آن ها می باشد که برای مقایسه بهتر نمودارها مقادیر تعداد تکرار بایت ها نرمال شده اند و تعداد تکرار هر حالت مقداری بین ۰ و ۱ دارد. در نمودار فایل های سالم مقادیر تکرار حالت های مختلف معمولاً کمتر از مقادیر تعداد تکرار در فایل آلوده می باشد. می توان گفت میانگین تعداد تکرار در فایل سالم کمتر از میانگین تعداد تکرار در فایل آلوده می باشد و با تحلیل نمودارهای بدست آمده می توان به این نتیجه رسید که در روش دوم تفاوت بیشتری بین نمودار فایل های سالم و آلوده قابل مشاهده است و می توان براساس آن و اعمال روش های دسته بندی^۱ سیستم تشخیص آلودگی فایل اجرایی را تولید کرد.

Classification^۱

فصل چهارم

۴ پیاده‌سازی و ارزیابی نتایج

۱-۴ مقدمه

در ادامه کار می بایست براساس نتایجی که در فصل سوم بدست آمد یک روش مناسب جهت تشخیص فایل آلوده از سالم ارایه کنیم. برای این کار ابتدا باید ویژگی مناسبی را انتخاب و با بررسی روش های مختلف وزن دهی یک روش وزن دهی کارا ارایه کنیم.

در این فصل براساس بهترین روش ارایه شده در فصل سوم یعنی بدست آوردن فرکانس ۲ بایت ها در قسمت آخر فایل اجرایی استخراج ویژگی را انجام می دهیم و از این ویژگی ها به عنوان نماینده هر فایل اجرایی استفاده می کنیم. سپس با اعمال روش های مختلف وزن دهی و همچنین استفاده از روش نزدکترین همسایه به عنوان روش دسته بندی، کارایی هر کدام از روش های وزن دهی را مورد بررسی قرار می دهیم و همچنین تغییرات تعداد همسایگی و رابطه آن با تشخیص فایل های اجرایی آلوده را بررسی می کنیم.

۲-۴ استخراج ویژگی

با توجه به اینکه از فایل اجرایی فقط قسمت آخر آن در قالب یک فایل متنی، معیاری جهت تشخیص آلودگی در نظر گرفته شد می توان برای یافتن ویژگی مناسب، صورت مسئله را به این شکل تغییر داد که می خواهیم فایل های متنی را به دو دسته آلوده و سالم تقسیم و در ادامه کار براساس روش های دسته بندی متن رفتار کنیم.

ویژگی ای که در دسته بندی متن استفاده می شود کلمه است. پس باید در صورت مسئله جدید منظورمان از کلمه را در یک فایل متنی تشکیل شده از بایت هایی در مبنای ۱۶ بیان کنیم. با توجه به اینکه در کد اسکریپت ها با ۲ بایت نشان داده می شوند و همچنین در کد اسمنبلی نیز ۲ بایت ارتباط بهتری با دستورالعمل ها دارد هر ۲ بایت را به عنوان یک کلمه در نظر می گیریم و با توجه به اینکه که

بایت ها در مبنای ۱۶ می باشند هر متن از ۲۵۶ کلمه یا ویژگی تشکیل شده است. به عبارت دیگر بازی هر فایل اجرایی یک بردار ویژگی ۲۵۶ تایی داریم.

۳-۴ وزن دهی ویژگی ها

در این قسمت به بررسی روش های مختلف وزن دهی ویژگی ها که برای بدست آوردن بهترین روش وزن دهی مورد بررسی قرار گرفته اند می پردازیم. منظور از وزن دهی به ویژگی این است که مشخص کنیم هر ویژگی به چه میزان می تواند نشانگر محتوای آن متن باشد. اگر فرض کنیم مجموعه کل متن ها و $C = \{c_1, \dots, c_m\}$ مجموعه دسته ها باشد، هر متن $D = \{d_1, \dots, d_n\}$ ویژگی (t_k, d_i) نشان داده می شود که در آن t_k تعداد ویژگی های متمایز در مجموعه D است که تعداد آن در روش ما ۲۵۶ است و وزن ویژگی t_k می باشد.

۱-۳-۴ روش^۱ TF

در این روش وزن ویژگی t_k در متن d_i برابر با تعداد تکرار ویژگی در آن متن می باشد. ایده اصلی این روش بر این اساس است که هرچه تعداد تکرار یک ویژگی بیشتر باشد ارزش آن ویژگی بیشتر است. رابطه (۱-۴) بیانگر روش بدست آوردن این نوع وزن دهی می باشد

$$w_{ki} = tf(t_k, d_i) = \begin{cases} \#(t_k, d_i) & t_k \in \text{vector of } d_i \\ 0 & t_k \notin \text{vector of } d_i \end{cases} \quad (1-4)$$

که در آن $\#(t_k, d_i)$ برابر با تعداد تکرار ویژگی t_k در متن d_i می باشد.

^۱ Term Frequency

۲-۳-۴ روش normTF

معمولاً طول متن های موجود در مجموعه D برابر نمی باشد برای آنکه طول متن بروی وزن ویژگی تاثیری نداشته باشد از این روش که نرمال شده روش TF می باشد استفاده می کنیم که این امر باعث می شود مقدار وزن ویژگی ها بین (0,1) باشد. برای این کار تعداد تکرار هر ویژگی در هر متن را بر طول آن متن تقسیم می کنیم رابطه (۲-۴) نشان دهنده این نوع وزن دهی می باشد:

$$w_{ki} = \text{normTF}(t_k, d_i) = \frac{tf(t_k, d_i)}{\sqrt{\sum_k (tf(t_k, d_i))^2}} \quad (2-4)$$

که در آن (t_k, d_i) از رابطه (۱-۴) بدست می آید.

۳-۳-۴ روش logTF

ممکن است در یک متن تعداد تکرار یک ویژگی خیلی بزرگتر از ویژگی های دیگر باشد و این امر بروی دسته بندی تاثیر منفی بگذارد برای رفع این اثر نامطلوب و یکسان کردن محدوده اختصاص داده شده از این روش استفاده می کنیم:

$$w_{ki} = \text{logTF}(t_k, d_i) = \log(tf(t_k, d_i)) \quad (3-4)$$

۴-۳-۴ روش^۱ ITF

این روش با استفاده از قوانین آماری ارایه شده است [۶] که در آن معمولاً مقدار r یک در نظر گرفته می شود.

$$w_{ki} = \text{ITF}(t_k, d_i) = 1 - \frac{r}{r + tf(t_k, d_i)} \quad (4-4)$$

^۱Inverse Term Frequency

۵-۳-۴ روش Sparck

این روش با استفاده از تیوری های آماری ارایه شده است [۷] و که در آن k تعداد کل واژگی های

متمايز در مجموعه D و $p_k = \sum_D tf(t_k, d_i)$ می باشد.

$$w_{ki} = \text{Sparck}(t_k, d_i) = tf(t_k, d_i) * (k - \log(p_k)) \quad (5-4)$$

۶-۳-۴ روش IDF

در این روش وزن واژگی برای کل مجموعه تعیین می شود [۸] که در آن $|D|$ تعداد کل متن های مجموعه D و $|D(t_k)|$ تعداد متن هایی است که واژگی t_k در آن ها وجود دارد و در نتیجه هر چه یک واژگی در تعداد متن بیشتری وجود داشته باشد وزن کمتری می گیرد.

$$w_{ki} = \text{idf}(t_k) = \log \frac{|D|}{|D(t_k)|} \quad (6-4)$$

۷-۳-۴ روش TFIDF

همانطور که گفته شد روش IDF به هر واژگی در کل مجموعه یک وزن می دهد به همین دلیل از ترکیب این روش با روش TF می توان به هر واژگی به هر متن یک وزن اختصاص داد که در آن $\text{idf}(t_k)$ از رابطه (۶-۴) بدست می آید.

$$w_{ki} = \text{tfidf}(t_k, d_i) = tf(t_k, d_i) * idf(t_k) \quad (7-4)$$

۸-۳-۴ روش normTFIDF

برای از بین بردن اثر طول بر روی وزن واژگی ها در روش TFIDF، از این روش استفاده می شود [۹]:

$$w_{ki} = \text{normTFIDF}(t_k, d_i) = \frac{\text{tfidf}(t_k, d_i)}{\sqrt{\sum_k (\text{tfidf}(t_k, d_i))^2}} \quad (8-4)$$

TFRF ۹-۳-۴ روش

در این روش توزیع ویژگی t_k در دسته از پیش تعریف شده c_j استفاده می شود. [۱۰] برای این منظور فاکتور ارتباط^۱ بازی هر ویژگی t_k در دسته c_j به صورت رابطه (۹-۴) تعریف می شود:

$$rf(t_k, c_j) = \log \left(2 + \frac{|D(t_k, c_j)|}{\sum_{m=1, m \neq j}^{|C|} |D(t_k, c_m)|} \right) \quad (9-4)$$

که در آن $c_j \in C$ و $|D(t_k, c_j)|$ تعداد متن هایی عضو مجموعه D و از دسته c_j که ویژگی t_k را داشته باشند است. $\sum_{m=1, m \neq j}^{|C|} |D(t_k, c_m)|$ مجموع تعداد مستنداتی است که عضو مجموعه D باشند و از دسته ای به جز دسته c_j قرار داشته باشند که دارای ویژگی t_k هستند. همانطور که مشاهده می شود بازی هر ویژگی این فاکتور رابطه‌ی مستقیم دارد با تعداد متن هایی از همان دسته متن اصلی هستند و دارای آن ویژگی باشند و همچنین رابطه معکوس دارد با تعداد متن هایی از دسته های دیگر که دارای این ویژگی باشند. در نهایت وزن ویژگی t_k از مستند d_i بعد از نرمال سازی از رابطه (۱۰-۴) بدست می آید که در آن c_{d_i} طبقه متن d_i می باشد.

$$w_{ki} = TFRF(t_k, d_i) = \frac{tf(t_k, d_i) * rf(t_k, c_{d_i})}{\sqrt{\sum_k (tf(t_k, d_i))^2 * (rf(t_k, c_{d_i}))^2}} \quad (10-4)$$

۴-۴ انتخاب روش دسته بندی

پس از اعمال روش های مختلف وزن دهی در مرحله بعدی باید یک روش دسته بندی مناسب انتخاب کرد. با توجه به اینکه روش نزدیک ترین همسایه یکی از موثرترین و سریعترین روش ها در دسته بندی متن می باشد ما نیز این روش را به عنوان روش دسته بندی انتخاب کردیم. در این روش مجموعه ای از

Relevancy Factor(rf)^۱

نمونه ها به عنوان مجموعه آموزش^۱ در نظر گرفته می شود. در زمان تست بازای هر نمونه از مجموعه تست فاصله آن با تمام نمونه های مجموعه آموزش براساس معیار فاصله در نظر گرفته شده بدست می آید سپس براساس تعداد همسایگی در نظر گرفته شده K فاصله‌ی کمتر بررسی می شود. دسته‌ای که بیشترین تعداد نمونه در بین K همسایه‌ی انتخاب شده را داشته باشد به عنوان دسته نمونه تست مورد نظر انتخاب می شود. بر این اساس مقدار K معمولاً عددی فرد در نظر گرفته می شود. در این پایان نامه برای تعیین فاصله بین دو بردار از فاصله اقلیدسی^۲ که در رابطه (۱۱-۴) آمده است استفاده می کنیم:

$$\sqrt{\sum_{j=1}^k (d_{1j} - d_{2j})^2} \quad (11-4)$$

۵-۴ پایگاه داده

در روش های ارایه شده در زمینه تشخیص فایل های اجرایی آلوده، پایگاه داده استاندارد وجود ندارد و چون کارهایی که در این زمینه انجام می شوند اغلب با هدف تجاری سازی و استفاده در آنتی ویروس ها صورت می گیرد پایگاه داده خود را در اختیار دیگران نمی گذارند و معمولاً در هر روشی یک پایگاه داده معرفی می شود. لازم به ذکر است در برخی روش ها از پایگاه داده یکسان استفاده شده بود ولی آن پایگاه داده نیز در دسترس نبود. بر همین اساس در این پایان نامه سعی کردیم پایگاه داده مناسبی ارایه کنیم. برای آزمایش روش پیشنهادی، ۱۶۷ فایل اجرایی آلوده به همراه ۱۶۷ فایل اجرایی سالم مجموعه کل داده های ما را که ۳۳۴ فایل اجرایی می باشد را تشکیل می دهند. که فایل های آلوده شامل دو دسته از کدهای مخرب رایج به نام های^۳ Sality و^۴ Virut می شوند و فایل های سالم نیز از مجموعه

Train^۱

Euclidian Distance^۲

W32.Sality^۳

W32.Virut^۴

فایل های اجرایی ویندوز هفت^۱ بدست آمده است. که برای آزمایش به روش نزدیکترین همسایه از ۱/۳ داده ها برای تست و ۲/۳ آن ها برای آموزش استفاده می کنیم . بنا بر این داده تست ما از ۱۱۱ فایل شامل ۵۶ فایل آلوده و ۵۵ فایل سالم تشکیل شده است.

۶-۴ معیارهای ارزیابی

جهت ارزیابی روش های دسته بندی معیار های مختلفی مطرح شده است ولی در دسته بندی فایلهای آلوده و سالم سه معیار اصلی وجود دارد.

۶-۴-۱ نرخ تشخیص درست مثبت^۲

اولین معیار نرخ تشخیص درست مثبت می باشد که معادل است با تقسیم تعداد فایل هایی که به درستی آلوده تشخیص داده شده اند بر کل فایل های آلوده که در رابطه(۱۲-۴) آمده است:

$$TPR = \frac{TP}{TP+FN} \quad (12-4)$$

که در آن ^۳TP برابر است با تعداد فایل های که به درستی آلوده تشخیص داده شده اند و ^۴FN تعداد فایل های آلوده ای که اشتباهات سالم در نظر گرفته شده اند. هر چه این عدد به ۱ نزدیکتر باشد یعنی از بین فایل های اجرایی آلوده فایل های بیشتری شناسایی شده اند و در نتیجه کارایی آن روش بیشتر است.

Windows 7 ^۱

TPR(True Positive Ratio)^۲

True Positive ^۳

False Negative ^۴

۲-۶-۴ نرخ تشخیص غلط مثبت^۱

معیار دوم نرخ تشخیص غلط مثبت می باشد که معادل است با نسبت تعداد فایل های سالم که اشتباه آلوده تشخیص داده شده اند به کل فایل های سالم که در رابطه (۱۳-۴) آمده است:

$$FPR = \frac{FP}{FP + TN} \quad (13-4)$$

که در آن FP^{\ddagger} نشان دهنده تعداد فایل های سالمی است که اشتباه آلوده تشخیص داده شده اند و TN^{\ddagger} برابر است با تعداد فایل هایی که به درستی سالم تشخیص داده شده اند. در این معیار نرخ فایل های سالمی که اشتباه آلوده تشخیص داده می شود بیان می شود. در یک سیستم تشخیص فایل آلوده هر قدر این معیار به صفر نزدیک شود کارایی سیستم بیشتر است. به این دلیل که اگر این سیستم فایل های سالم زیادی را آلوده تشخیص دهد اشتباهها اجازه فعالیت به آن فایل های سالم نمی دهد کارایی کاربر پایین می آید.

۳-۶-۴ دقت

معیار سوم دقت است که معادل با تقسیم تعداد کل فایل هایی که به درستی دسته بندی شده اند تقسیم بر تعداد فایل های مجموعه تست می باشد که در رابطه (۱۴-۴) آمده است:

$$Accuracy = \frac{TP + TN}{TP + TN + FP + FN} \quad (14-4)$$

در واقع در این معیار دقت سیستم در تشخیص درست فایل های سالم و آلوده ارزیابی می شود و در آن هم تشخیص درست فایل های آلوده نقش دارد و هم تشخیص درست فایل های سالم. هر اندازه مقدار

FPR(False Positive Ratio)^۱

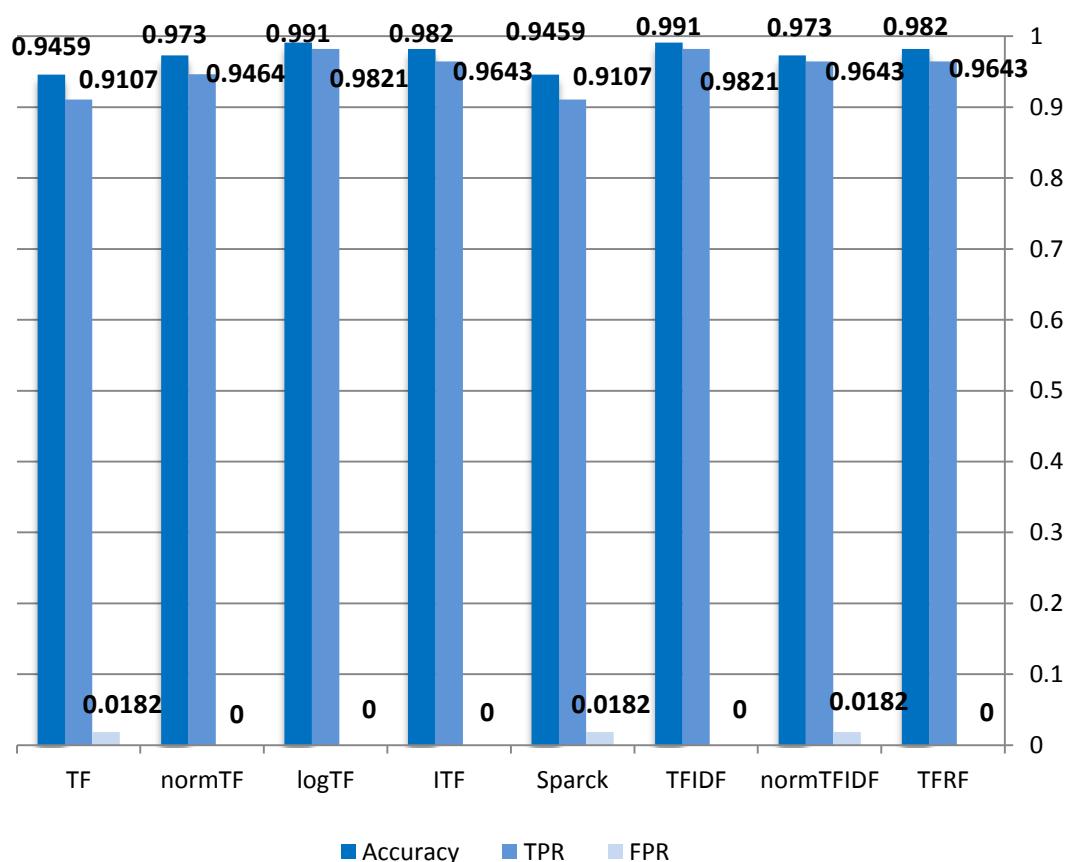
False Positive ^۲

True Negative ^۳

این معیار به ۱ نزدیکتر باشد یعنی فایل های بیشتری درست تشخیص داده شده اند و در نتیجه کارایی سیستم بیشتر است.

۷-۴ نتایج

پس از پیاده سازی روش پیشنهادی و استفاده از مجموعه پایگاه داده معرفی شده نتایج زیر حاصل شد.



شکل ۱-۴ نمودار کارایی روش های مختلف وزن دهنده براساس معیارهای معرفی شده

در شکل (۴) نمودار کارایی روش های مختلف وزن دهنده زمانی که در روش نزدکترین همسایه مقدار k برابر با یک باشد آمده است. همانطور که در شکل مشخص است روش های logTF و TFIDF دقت ۰.۹۹۱۰ و ۰.۹۸۲۱ TPR و FPR صفر بهترین کارایی را دارند و با توجه به اینکه در روش

تعداد متن هایی که یک ویژگی در آن ها وجود دارد در ارزش گذاری برای یک ویژگی دخیل می کند پیش بینی می شود در مجموعه هایی با تعداد بیشتر فایل اجرایی کارایی بهتری نسبت به روش logTF داشته باشد.

همچنین برای بررسی تاثیر تغییرات مقدار K در روش دسته بندی برای روش های مختلف وزن دهی مقادیر ۱ و ۳ و ۵ و ۷ و ۹ برای K در نظر گرفته شد و معیارهای معرفی شده بررسی گردید.

جدول ۴-۱ بررسی کارایی روش های مختلف وزن دهی بازی تغییر مقادیر K.
 (الف) روش TF. (ب) روش TFRF
 (ج) روش normTF. (د) روش logTF. (ه) روش ITF. (و) روش Sparck. (ز) روش TFIDF. (ح) روش normTFIDF

TF	k=1	k=3	k=5	k=7	k=9
Accuracy	0.9459	0.964	0.964	0.955	0.955
TPR	0.9107	0.9464	0.9464	0.9464	0.9464
FPR	0.0182	0.0182	0.0182	0.0364	0.0364

الف

normTF	k=1	k=3	k=5	k=7	k=9
Accuracy	0.973	0.964	0.964	0.964	0.964
TPR	0.9464	0.9286	0.9286	0.9286	0.9286
FPR	0	0	0	0	0

ب

logTF	k=1	k=3	k=5	k=7	k=9
Accuracy	0.991	0.982	0.982	0.991	0.991
TPR	0.9821	0.9643	0.9643	1	1
FPR	0	0	0	0.0182	0.0182

ج

ITF	k=1	k=3	k=5	k=7	k=9
Accuracy	0.982	0.982	0.991	0.991	0.982
TPR	0.9643	0.9821	1	1	1
FPR	0	0.0182	0.0182	0.0182	0.0364

۵

Sparck	k=1	k=3	k=5	k=7	k=9
Accuracy	0.9459	0.964	0.964	0.955	0.955
TPR	0.9107	0.9464	0.9464	0.9464	0.9464
FPR	0.0182	0.0182	0.0182	0.0364	0.0364

۶

TFIDF	k=1	k=3	k=5	k=7	k=9
Accuracy	0.991	0.982	0.982	0.991	0.991
TPR	0.9821	0.9643	0.9643	0.9821	0.9821
FPR	0	0	0	0	0

۷

normTFIDF	k=1	k=3	k=5	k=7	k=9
Accuracy	0.973	0.982	0.973	0.973	0.973
TPR	0.9643	0.9821	0.9643	0.9643	0.9643
FPR	0.0182	0.0182	0.0182	0.0182	0.0182

ز

TFRF	k=1	k=3	k=5	k=7	k=9
Accuracy	0.982	0.982	0.982	0.982	0.982
TPR	0.9643	0.9643	0.9643	0.9643	0.9643
FPR	0	0	0	0	0

ح

در جدول (۱-۴) با تغییر مقادیر K در روش نزدیک ترین همسایه کارایی روش های مختلف وزن دهنده آمده است. در روش TF با افزایش مقدار K کارایی افزایش و در روش های TFIDF و logTF با افزایش همسایگی کارایی کاهش پیدا می کند در نتیجه نمی توان گفت که هر چقدر تعداد همسایگی های در

نظر گرفته شده بیشتر باشد کارایی افزایش پیدا می کند و باید با بررسی جوانب مختلف همسایگی مناسب را پیدا کرد.

فصل پنجم

۵ نتیجه کسری و کارهای آینده

۱-۵ نتیجه گیری

در این پایان نامه تلاش شد روشی ارایه شود تا بازای هر فایل اجرایی یک الگوی سیگنالی داشته باشیم و براساس آن بتوانیم فایل های اجرایی آلوده را شناسایی کنیم. به عبارت دیگر رویکردی جدید جهت تشخیص فایل اجرایی آلوده معرفی شد که براساس آن بتوان با استفاده از تکنیک های پردازش سیگنال فایل آلوده را تشخیص داد. تاکنون روشی که برمبانی پردازش سیگنال فایل آلوده را شناسایی کند مطرح نشده است. بدین منظور روش های مختلف تحلیل ایستای فایل اجرایی مورد بررسی قرار گرفت این روش ها در دو دسته کلی الف) براساس کد اسمبلی و ب) براساس بایت های تشکیل دهنده ارایه شد

در تشخیص فایل اجرایی آلوده از سالم دو شاخص جدید در این پایان نامه معرفی شد: الف) تغییرات اشاره گر دستور العمل و ب) تفاوت در قسمت منبع که می توان از آن ها در سیستم های تشخیص فایل آلوده استفاده کرد. با در نظر گرفتن تفاوت در قسمت منبع و بدست آوردن تعداد تکرار ۲ بایت ها در آن، به عنوان نماینده یک فایل اجرایی و استفاده از روش های دسته بندی متون اثبات کردیم که این روش می تواند در جهت تشخیص فایل اجرایی آلوده موثر باشد.

برای اثبات اینکه تحلیل قسمت آخر فایل اجرایی و بدست آوردن الگوی سیگنالی از آن قسمت می تواند شاخص خوبی برای تشخیص آلودگی فایل اجرایی باشد. پس از بدست آوردن تعداد تکرار ۲ بایت ها در قسمت آخر فایل اجرایی، هر ۲ بایت را به عنوان یک کلمه در نظر گرفتیم و روش های مختلف وزن دهی در دسته بندی متون را بروی آن ها اعمال کردیم. برای دسته بندی از روش نزدیک ترین همسایه و از معیارهای دقت و TPR و FPR برای ارزیابی کارایی روش های وزن دهی استفاده کردیم که در نهایت

روش TFIDF به همراه روش نزدیکترین همسایه با مقدار $K=1$ با دقت ۰.۹۹۱۰ و TPR ۰.۹۸۲۱ و FPR صفر بهترین کارایی را داشت.

۲-۵ کارهای آینده

برای ادامه کار می توان در دو جنبه روش های این پایان نامه را دنبال کرد. جنبه اول در بهبود روشی که بر مبنای دسته بندی متن ارایه شده بود. که در این زمینه می توان روش ارایه شده را بروی مجموعه داده هایی با تعداد بالا فایل سالم و آلوده آزمایش کرد و سپس براساس نتایج بدست آمده به بهبود روش مذکور پرداخت. همچنین بررسی روش های دیگر دسته بندی مانند درخت تصمیم^۱ ، روش بیزین^۲ و... می تواند بروی این روش های وزن دهی اعمال شود. کار دیگری که می توان در این زمینه ارایه کرد یافتن روش وزن دهی جدیدی است بر مبنای صورت مسئله ای که تعریف شده است. به عبارت دیگر با توجه به اینکه روش وزن دهی موجود برای دسته بندی متن ارایه شده است تلاش شود روش وزن دهی جدیدی با توجه بردار ویژگی ارایه شده مطرح کرد.

در جنبه دوم برای ادامه کار می توان بروی ارایه روش دسته بندی با استفاده از پردازش سیگنال کار کرد. به عبارت دیگر با توجه به اینکه تحلیل قسمت منبع می تواند ما را به الگوی سیگنالی برساند که براساس آن بتوان فایل آلوده را شناسایی کرد؛ به دنبال بدست آوردن الگوی سیگنالی مناسب از قسمت آخر فایل اجرایی باشیم به عنوان مثال می توانیم اختلاف ۲ بایت های مجاور در این قسمت را مورد بررسی قرار دهیم.

Decision Tree ^۱

Navie Bayes ^۲

منابع

- [1] Abou-Assaleh T, Cercone N, Keselj V, Sweidan R. N-gram based detection of new malicious code. In: Proceedings of the 28th annual international computer software and applications conference; 2004
- [2] Moskovich R, Stopel D, Feher C, Nissim N, Elovici Y. Unknown malcode detection via text categorization and the imbalance problem. In: IEEE Intelligence and Security Informatics, Taiwan; 2008.
- [3] Siddiqui M, Wang MC, Lee J. Data mining methods for malware detection using instruction sequences. Artificial Intelligence and Applications 2008.
- [4] Bilar D. Opcodes as predictor for malware. International Journal of Electronic Security and Digital Forensic;1:2, 2007
- [5] Moskovich R, Feher, Tzachar N, Berger E, Gitelman M, Dolev S, et al. Unknown malcode detection using OP CODE representation. In: European conference on intelligence and security informatics 2008 (EuroISI08), Esbjerg, Denmark; 2008
- [6]E. Leopold and J. Kindermann. Text categorization with support vector machines: How to represent texts in input space Journal of Machine learning 46 423-444- 2002
- [7]K .sparck Jones, Indexing term weighting. Information Storage and Retrieval vol .9pp619-663 ,1995
- [8]G. Salton , c.s.Yang, On the specification of termvalue in automatic indexing,journal of documentation, vol29 no4 pp351-357 ,1973
- [9]G . Saltpn, J Allan, A. Singhal, Automatic text decomposition and structuring, Inforamtionprocessing and mangment, vol 32 no 2 pp127-138 ,1996
- [10]M .LAn, S.Y. Sung, H.B LOW, C.L.Tan, A comparative Study on term Weighting Schemes for Text Categorization, IEEE International on Neural Networks pp 546-551, 2005
- [11] Ding Yuxin, Yuan Xuebing, Zhou Di, Dong Li, An Zhancha Feature representation and selection in malicious code detection methods based on static system calls, computers & security 30514-524,2011
- [12] Zhang B, Yin J, Hao J, Zhang D, Wang S. Malicious codes detection based on ensemble learning. Autonomic and trusted computing;.. p. 468–27. 2007

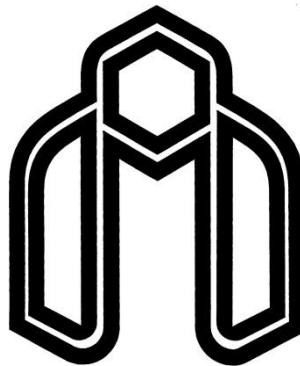
- [13] Menahem E, Shabtai, Rokach L, Elovici Y. Improving malware detection by applying multi-inducer ensemble. Computational Statistics and Data Analysis 2008
- [14] I. Santos et al., Opcode sequences as representation of executables for data-mining-based unknown malware detection, Inform. Sci. 2011
- [15] Karim E, Walenstein A, Lakhotia A, Parida L. Malware phylogeny generation using permutations of code. Journal in Computer Virology;1(1–2):13–23. 2005
- [16] Dolev S, Tzachar N. Malware signature builder and detection for executable code. Patent application; 2008
- [17] Schultz M, Eskin E, Zadok E, Stolfo S. Data mining methods for detection of new malicious executables. In: Proc. of the IEEE symposium on security and privacy. p. 178–84. ; 2001
- [18] Henchiri O, Japkowicz N. A feature selection and evaluation scheme for computer virus detection. In: Proceedings of ICDM-2006, Hong Kong;. p. 891–95. 2006
- [19] M. Christodorescu and S. Jha, “Static Analysis of Executables to Detect Malicious Patterns”, in Twelfth Security Symposium, pp. 169-186, 2003
- [20] Dolev S, Tzachar N. Malware signature builder and detection for executable code. Patent application; 2008.
- [21] Kolter J, Maloof M. Learning to detect and classify malicious executables in the wild. Journal of Machine Learning Research;7:2721–44. 2006
- [22] D. K. S. Reddy, A. K. Pujari. N-gram analysis for computer virus detection. Journal of Compute Virology. France;2: 231239, 2006.
- [23] D. Reddy, S. Dash, A. Pujari. New malicious code detection using variable length n-grams. International Conference on Information Systems Security. Germany;276-288; 2006.
- [24] Bergeron J, Debbabi M, Desharnais J, Erhioui MM, Lavoie Y, Tawbi N. Static detection of malicious code in executable programs. In: Symposium on requirements engineering for information security (SREIS’01), Indianapolis, Indiana, USA; March. p. 157e166; 2001.
- [25] M. Christodorescu, S. Jha, Static analysis of executables to detect malicious patterns, in: Proceedings of the 12th USENIX Security Symposium, , pp. 169–186; 2003.

Abstract

Due to proliferate of new malware and using obfuscation techniques in malware, research in unknown malware detection is essential . In the commercial antivirus, complimentary combination of signature based method and heuristic is used; nevertheless heuristic feeble performance results in reliable detection rate depending on database updating. Therefore today the most research in malware detection conducted to enhance and improve heuristic method.

The detection of files infected by unknown malicious code is important task in antivirus using heuristic methods for unknown malicious code detection. In this thesis has studied various some kinds of methods analyzing portable executable (PE) file and then introducing a signal for each PE file. The most reliable method is concluded based on Resource section analysis of PE files . To show reliable performance of proposed method, an unknown malicious code detection method based on text classification technique with using feature vector extracted from proposed method achieve 99.10 accuracy.

Keywords:Infected file, Unknown malicious code, Signal processing, unknown malware detection



**Shahrood University of Technology
Faculty of Computer Engineering**

Infected Files Analysis Based on Signal Processing

Thesis

Submitted in Partial Fulfillment of the Requirements for the Degree of
Master of Science (M.Sc)

Aboulfazl Sarkardei

Supervisor

Dr. Ali Akbar Pouyan

Associate supervisors

Dr. Hamid Hassanpour
Javad Kia

Date: February 2013