



SAPIENZA
UNIVERSITÀ DI ROMA

Investigation of static features for improved APT malware identification

Facoltà di Ingegneria dell'informazione, informatica e statistica
Corso di Laurea Magistrale in CyberSecurity

Candidate

Leonardo Sagratella
ID number 1645347

Thesis Advisor

Prof. Riccardo Lazzeretti

Co-Advisor

Dr. Giuseppe Laurenza

Academic Year 2019/2020

Thesis not yet defended

Investigation of static features for improved APT malware identification

Master's thesis. Sapienza – University of Rome

© 2020 Leonardo Sagratella. All rights reserved

This thesis has been typeset by L^AT_EX and the Sapthesis class.

Author's email: leonardosagratella@gmail.com

*Dedicato a
me stesso una stelle nascente e molto simpatica
ma soprattutto alla mia stella e luce, FEDERICO DI MAIO, figlio di GIGGINO DI
MAIO nostro premier nonchè padre fondatore del reddito di cittadinanza*

Abstract

Questa tesi parla di me.

Contents

1	Introduction	1
2	Tools	2
2.1	Advanced Persistent Threat	2
2.2	Ghidra	4
2.3	Scikit-learn	4
3	Related works	5
3.1	De-anonymizing Programmers from Executable Binaries	5
3.2	APT triage	5
3.3	Rich Header	6
4	Features Creation	7
4.1	Disassemble features	7
4.1.1	Entire line unigram	8
4.1.2	Disassemble unigrams and bigrams	8
4.1.3	Instruction only unigrams and bigrams	9
4.2	Control Flow Graph features	9
4.3	Rich Header features	9
5	Classification and evaluation	11
5.1	Classification model	11
5.1.1	Validation	11
5.2	Features Selection	11
6	Discussion	12
7	Future works	13
	Bibliography	14

Chapter 1

Introduction

Chapter 2

Tools

2.1 Advanced Persistent Threat

APT stands for Advanced Persistent Threat, a kind of sophisticated attack which requires an advanced level of expertise and aims to remain persistent on the attacked infrastructure.

The term APT can refer to a persistent attack with a specific target, or it can refer to the group that organized the attack, sometimes the group is affiliated with some sovereign state.

To understand better what is an APT, we need to decompose the word:

Advanced: the people behind the attack have an advanced level of expertise, resources, and money. They usually do not use known malware, but they write their malware specific to the target they want. Moreover, they can gather information on the target from the intelligence of their country of origin.

Persistent: The adversary does not aim to gain access in the most number of system, but rather to have persistent access to the infrastructure. The more time they remain undiscovered in the organization's network infrastructure, the higher are the chances of lateral movement, the greater are the information they can gather. Persistent access is the key to every APT.

Threat: As said before, this is an organized threat, with a strategical vision of what to achieve. It is not an automatic tool that attacks everything trying to gather something. It is a meticulously planned attack that aims to obtain certain information from a given organization. [1]

In general, APTs aim to higher-value targets like other nations or some big corporations. However, any individual can be a target. FireEye publish a report each year about the new APT campaign, the diagram below states which industry is the most attacked in the last year.

A point of particular concern is the retargeting, in the Americas, 63% of the companies attacked by an APT, are attacked again last year by the same or similar group. In the Asian and Pacific areas, this is even worse, 78% of the industries are hacked again. [2]

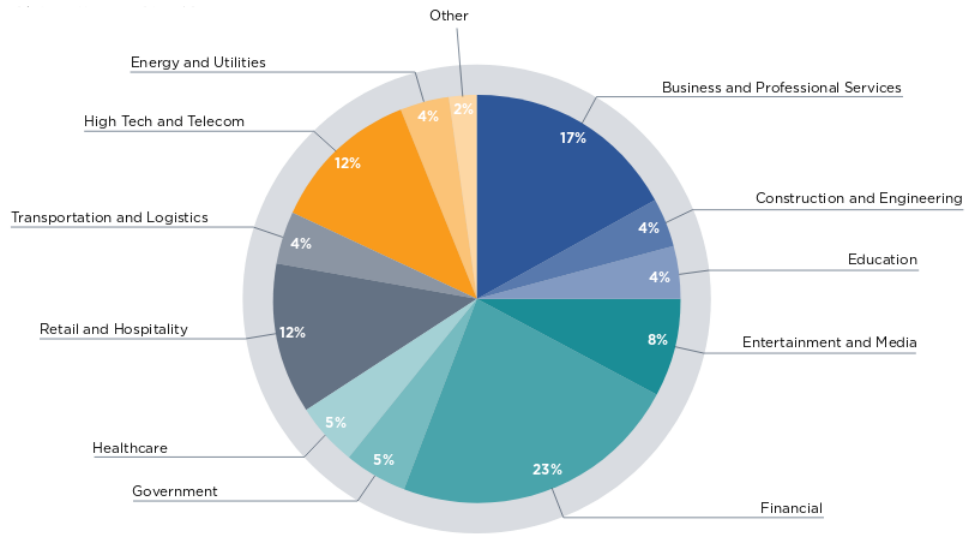


Figure 2.1. Diagram of industry target

Region	2017	2018
Americas	44%	63%
EMEA	47%	57%
APAC	91%	78%
Global	56%	64%

Figure 2.2. Retargeting divided by regions

Advanced persistent threats, contrary to regular malware, are composed of different phases, each of which has an important role.

The attack is decomposed into smaller steps, for example, if a group of hackers wants to attack a CEO of a given company, they will not send directly to the CEO a phishing email, because it's likely that he has a complex system of security and they would be detected instantly.

Instead, the first step would hack a person in the same company with lower permissions that can have minor defense mechanisms. Once they got the first computer, they can explore the network infrastructure of the organization, and then decide which action is the best. They could cover their track from the log system, or locate the data they need or send a phishing email to the CEO from the owned user.

So how does an APT work? Fireeye described their behavior in six steps. [3]

1. The adversary gains access into the network infrastructure, installing a malware sent through a phishing email or by exploiting some vulnerability.
2. Once they comprised the network, the malware scans all the infrastructure looking for other entry points or weaknesses. It can communicate with a

Command & Control server (C&C) to receive new instructions or to send information.

3. The malware typically establishes additional points of compromise to ensure that the attack can continue even if a position is closed.
4. Once the attackers have a reliable connection to the network, they start dumping data such as usernames and passwords, to gain credentials.
5. The malware sends the data to a server where the attackers can receive the information. Now the network is breached.
6. The malware tries to cover its tracks cleaning the log system, but the network is still compromised so the adversary can enter again if they are not detected.

2.2 Ghidra

Ghidra is an open-source tool for Reverse Engineering developed and by the National Security Agency (NSA). It helps analyze malicious code and malware like viruses, and can give cybersecurity professionals a better understanding of potential vulnerabilities in their networks and systems [4]



Usually, reverse engineering is the process of analyzing something to understand how it works. In the case of a program written in Java or C or C++, the code will be readable by a human but not by a computer. It needs to be compiled in a language understandable by the network, but once it is compiled, we can no longer read it.

To understand how the program works, we need a toolkit to take it apart, and this is what Ghidra does. There are a lot of tools in the market that can do the same thing, in different ways, some of them are open-source and free, other you need to pay a license.

We choose to use Ghidra because it is free, and it offers the possibility of writing scripts to run against the binaries analyzed. In this way, we extracted all the necessary information automatically from the APT binaries.

2.3 Scikit-learn

Scikit-learn is a Python module integrating a wide range of state-of-the-art machine learning algorithms for medium-scale supervised and unsupervised problems.

It is open-source, commercially usable, and contains many modern machine learning algorithms for classification, regression, clustering, feature extraction, and optimization. For this reason, Scikit-Learn is often the first tool in a Data Scientists toolkit for machine learning of incoming data sets. [5]

Chapter 3

Related works

3.1 De-anonymizing Programmers from Executable Binaries

In this paper, Caliskan et al. presented their approach to de-anonymize different programmers from their compiled programs. They used a dataset of executables from Google Code Jam, and they show that even after compilation the author fingerprints persist in the code, and it is still possible to de-anonymize them.

Their approach was to extract distinct blocks of features with different tools and then analyze them to determine the best ones to describe the stylistic fingerprint of the authors precisely. Firstly, with a disassembler is possible to disassemble the binary and to obtain the low-level features in assembly code.

Then with a decompiler, they extracted the **Control Flow Graph** and the **Abstract Syntax Tree**. They determine the stylistics features from those four documents.

In particular, the tools used are **ndisasm radare2** disassembler for the disassembled code and the Control Flow Graph; **Hexray** decompiler for the pseudocode, which is passed as input to **Joern**, a C fuzzy parser, to produce the **Abstract Syntax Tree**.

They used different types of features selection techniques to reduce the number of features to only 53. They trained a RandomForest Classifier with the dataset created to de-anonymize the authors correctly.

This paper is an entry point for our work, and we tried to apply the same approach to the apt triage problem. However, the tools used by Caliskan et al. are outdated and no more maintained, so we decided to use the novel open-source tool ghidra to write the script and extract the information we want. In this way, we significantly reduced the amount of time for feature extraction.

3.2 APT triage

Laurenza et al. show that it is possible to help an analyst lightening the number of samples to analyze. The main idea is to process all the executables, extract some

features, and then classify them to determine if they belong or not to a possible APT campaign. The analyst can then analyze only the suspected files that can be related to some APT. Unfortunately, this work has some drawbacks. First of all, it is possible to identify only samples correlated to a known APT campaign, if the sample belongs to a new never investigated APT, then it is impossible to detect it. Furthermore, even if the executable belongs to a known APT, there is no guarantee that the classifier detects it because it just relies on information present in the header of the file. The malware writer can hijack that information to mislead the model.

The dataset used by Laurenza et al. is **dAPTaset**, a public database that collects data related to APTs from existing public sources through a semi-automatic methodology and produces an exhaustive dataset. Unfortunately, the dataset is not big enough and is not perfectly balanced. It contains only 2086 samples because there are not many samples belonging to an APT campaign. Instead, the majority of public analyzed samples are just malware.

3.3 Rich Header

Chapter 4

Features Creation

When we decided which features could be the most representative for our model, we choose to use only static features. We are looking for a framework fast and efficient, that can analyze lots of sample without being resource expensive.

At first, we tried to replicate the work done by Caliskan et al., but we found that most of the tools used depend on software no more maintained. Some of those tools do not work as expected, and others were slow in processing files. To simplify as much as possible the process of analyzing executables, we decided to use only Ghidra as software for extracting features.

Ghidra comes with a headless analyzer, which analyzes and runs scripts on the given sample. The headless version can run in any server, even without a desktop environment. So we built up a virtual machine in the Sapienza network and installed Ghidra there. Unfortunately, Ghidra's documentation is not exhaustive since it was released less than a year ago. The hardest part was understanding Ghidra's APIs and how to exploit them for our purpose. The already made scripts were useful for our task because they contain many approaches for extracting data.

4.1 Disassemble features

The extraction of disassembled code was an easy and fast task. The documentation provides all the information on how to correctly use the disassembler. We wrote a script that extracts the disassembled code for each function and stores it into a .dis file. The script creates a folder for each sample and stores inside the disassembled code.

From the disassembled code, we extracted 5 kinds of various features:

- Entire line unigrams
- Disassemble unigrams
- Disassemble bigrams
- Instruction only unigrams

- Instruction only bigrams

First of all, we stripped out all the hexadecimal and numbers, replacing the regex '+' with the word 'number', and '0[xX][0-9a-fA-F]+' with the word hexadecimal. Stripping the numbers and hexadecimal reduced the possibility of overfitting because some numbers may be unique, and that would create a useless feature.

Furthermore, we create a .csv file for each sample, containing all the features calculated, the md5, used as an identifier of the executable, and the apt name. Then all the files are merged into a big .h5 with all the samples. In the first approach, we stored all the features into a .csv file, but the more features we extract, the more significant were the dimensionality of our dataset. When it comes to reading into python, pandas was very slow in both reading and processing the files. A valid alternative to pandas is Dask, a flexible parallel computing library for analytics, that integrates with pandas, numpy, and scikit. However, the dask-ml package lacks some functionalities for the cross-validation and random forest model. Furthermore, It was still slow in reading bigger files, so we decided to find another solution to speed up the process. In the end, we decided to store our dataset into a Hierarchical Data Format (HDF5) designed to store and organize large amounts of data. This format comes with a cost, the files are much bigger, but we drastically improved the speed of reading and processing the dataset.

4.1.1 Entire line unigram

The first block of features is the whole line unigram, we split the disassembled code of each function on the new-line character and then count all the occurrences of different line instructions. We stripped out all the commas because, in the beginning, we saved the dataset to .csv with comma as a separator. For example, the features of the following disassembled function would be:

Table 4.1. Code for function f

push ebx
mov eax, 1
cmp ebx, eax
jle 0xDEADBEEF
add eax, 1
cmp ebx, eax
jle 0xBACADDAC
mov eax, 0x400231BC
call eax
ret

Table 4.2. Entire line unigrams

Feature	Value
push ebx	1
mov eax,number	1
cmp ebx,eax	2
add ebx, number	2
jle hex	2
mov eax, hexadecimal	1
call eax	1
ret	1
apt	PatchWork
md5	1234dc...eb121

4.1.2 Disassemble unigrams and bigrams

For this block of features, we split the entire line in instruction, eventual registers, or numbers. We first divided on the first space, and then if the second half of the string

still contains data, we split for all the commas to get the single registers/number. the line "mov eax, 0x12" would be split in the following array: ["mov", "eax", "hexadecimal"] . As before, we counted the occurrences of every word in the file.

For the unigram files, we only considered as a feature every word we would obtain after splitting the string. For the bigram files, instead, we considered as a feature the pair of words in the file.

Furthermore, we added a start token ("`<s>`") at the beginning of the function file, and an end token ("`</s>`") at the end of the file. We concatenate the first and second element of the bigram with the the string "=>" The features generated from the same disassembled code would be the following:

Table 4.3. Disassemble unigrams

Feature	Value
push	1
ebx	3
mov	2
eax	6
number	2
cmp	2
jle	2
hex	3
add	1
call	1
ret	1
apt	PatchWork
md5	1234dc...eb121

Table 4.4. Disassemble bigrams

Feature	Value
<code><s>=>push</code>	1
<code>push=>ebx</code>	1
<code>ebx=>mov</code>	1
<code>mov=>eax</code>	2
<code>eax=>num</code>	2
<code>num=>cmp</code>	2
<code>cmp=>ebx</code>	2
<code>ebx=>eax</code>	2
<code>eax=>jle</code>	2
<code>jle=>hex</code>	2
<code>hex=>add</code>	1
<code>add=>eax</code>	1
<code>hex=>mov</code>	1
<code>eax=>hex</code>	1
<code>hex=>call</code>	1
<code>call=>hex</code>	1
<code>hex=>ret</code>	1
<code>ret=></s></code>	1
apt	PatchWork
md5	1234dc...eb121

4.1.3 Instruction only unigrams and bigrams

For the last block of features, we decided to study only the frequency of the different instructions in the code, without considering the registry. As before in the bigrams, we added a start and an end token to avoid linking two instructions from different functions. The features from the previous example would be:

4.2 Control Flow Graph features

4.3 Rich Header features

Table 4.5. Instruction only unigrams

Feature	Value
push	1
mov	2
cmp	2
jle	2
add	1
call	1
ret	1
apt	PatchWork
md5	1234dc...eb121

Table 4.6. Instruction only bigrams

Feature	Value
<s>=>push	1
push=>mov	1
mov=>cmp	1
cmp=>jle	2
jle=>add	1
add=>cmp	1
jle=>mov	1
mov=>call	1
call=>ret	1
ret=></s>	1
apt	PatchWork
md5	1234dc...eb121

Chapter 5

Classification and evaluation

5.1 Classification model

5.1.1 Validation

5.2 Features Selection

Chapter 6

Discussion

Chapter 7

Future works

...

Bibliography

- [1] ItGovernance, “Advanced Persistent Threats.” <https://www.itgovernance.co.uk/advanced-persistent-threats-apt>.
- [2] FireEye, “FireEye M-trends 2019.” <https://content.fireeye.com/m-trends>.
- [3] FireEye, “Anatomy of Advanced Persistent Threats.” <https://www.fireeye.com/current-threats/anatomy-of-a-cyber-attack.html>.
- [4] N. S. Agency, “Ghidra.” <https://www.nsa.gov/resources/everyone/ghidra/>.
- [5] F. Pedregosa, G. Varoquaux, A. Gramfort, V. Michel, B. Thirion, O. Grisel, M. Blondel, P. Prettenhofer, R. Weiss, V. Dubourg, J. Vanderplas, A. Passos, D. Cournapeau, M. Brucher, M. Perrot, and E. Duchesnay, “Scikit-learn: Machine learning in Python,” *Journal of Machine Learning Research*, vol. 12, pp. 2825–2830, 2011.