

論文要旨

「未知マルウェア対策技術に関する研究」

Countermeasure technologies against unknown malware

2018年3月

情報セキュリティ研究科

情報セキュリティ専攻

5655102 田中 恭之

Yasuyuki Tanaka

指導教員 後藤 厚宏

本論文は、「未知マルウェア対策技術に関する研究」と題し、5章と付録からなる。近年、マルウェアは増加の一途をたどり、アンチウイルスソフトの検出を逃れるマルウェアも多く出現し問題となっている。本論文では、3つの異なるアプローチを提案・検証し、これらのマルウェアに対し、感染前にできるだけ軽量に、検出できることを目指したものである。

第1章は、「序論」で、本研究の背景として、マルウェアの爆発的な増加を受けたアンチウイルスソフトの検出の限界、標的型攻撃で用いられるようなゼロデイ脆弱性を含むようなマルウェアの存在、マルウェアが配布されるインターネット上の悪性サイトの多様化の問題を挙げている。次に、従来からのマルウェア対策技術分野を、マルウェア自体の検出分野について3つ、マルウェアの通信の検出分野について2つ、計5個に分類している。本研究で取り組んだ3つのアプローチをこの5個の分野にマッピングし、マルウェア感染前かつ軽量に検出する本研究の目的との位置付けを明確にしている。さらに3つのアプローチを間接的か直接的かの観点で特徴付け、各アプローチのメリットを整理している。

第2章は、「統計的手法を用いたマルウェア検出」と題して、ファイルの静的な特徴に対して統計的手法を用いてマルウェア判定を行っている。ファイルの多数のヘッダ情報から、判定に有効な情報に絞りモデルを構築し、そのモデルを利用することで軽量にマルウェア判定する方式を提案し評価を行っている。本手法は特に、インターネット上のマルウェアダウンロードサイトから配布されるような一般に広く流通するマルウェアに等に有効であると結論づけている。本手法により、アンチウイルスソフトが検出できなかったマルウェアを未知マルウェアと定義し、それらを検出可能なことを示している。評価実験では、パッキングの有無、アンチウイルスソフトの検知有無を考慮した複数の検証条件で、マルウェアダウンロードサイトから収集した未知マルウェアを含むデータセットで、マルウェアか正常ファイルかの識別精度評価を行っている。結果、従来手法に比べ識別性能が高いことを示している。さらに選定した変数モデルを、代表的な機械学習手法であるサポートベクターマシンに適用し、さらに高精度で識別できることを示している。

第3章は、「攻撃コード構成の特徴を用いたROPコード検出」と題して、ROP(Return-Oriented Programming)と呼ばれる攻撃コードが、シェルコードの複合ルーチンの外側に配置されると言う攻撃コードを構成する上での特徴を

利用して検出を行う手法を検討している。本手法は、第 2 章での検出対象である一般に広く用いられるマルウェアでなく、標的型攻撃や未知であるゼロデイ脆弱性とともに用いられる、ある特定のマルウェアに有効であることを示している。ROP コードは、ここ数年の脆弱性においてホスト側の防御機構を突破することを目的として、多くの攻撃コードに付加されるものである。特に、昨今の標的型攻撃で用いられる手口で、攻撃コードを埋め込んだ悪性文書ファイルを送付し、被害者がファイルを開くことでマルウェア感染等が引き起こされ、情報の搾取等がなされるケースがある。この ROP コードを静的に検出することで、悪性文書ファイル判定を行う方式を提案し、実際の検体にて評価を行っている。さらに、各検体で共通に使われる ROP コードを分析し、一部の ROP コードは異なる脆弱性とともに汎用的に用いられることを示し、これらの ROP コードを検出することでゼロデイ対策とできる可能性を示している。

第 4 章は、「効果的なブラックリスティングの検討」と題し、第 2 章や第 3 章で示すマルウェア自体の特徴を捉えて検出するのではなく、マルウェアのダウンロードサイトを、効果的にブラックリスティングし未知マルウェア対策につなげることを目的としている。これにより、マルウェア自体のハッシュ値等の特徴が変化して、未知マルウェアがダウンロードされた場合でも、検出可能である。実際には、約 43,000 個のマルウェアダウンロードサイトを、Web クローラを用い、約 1.5 年にわたり長期観測を行い、幾つかのサイトはとても長い期間にわたりマルウェアダウンロードサイトとして生存を続けること、また幾つかのサイトは消滅と復活を繰り返しながらマルウェアダウンロードサイトとして活動を続けること明らかにしている。マルウェアの変化に応じて、マルウェアダウンロードサイトを 3 つのカテゴリに分けて、その生存と挙動の観点から、長期分析を行っている。生存観点では、10 個の特徴量を定義し、生存期間や生き返りについて分析している。また挙動観点では、11 個の特徴量を定義し、IP アドレスや配布されるマルウェアの変化の特徴を分析している。それぞれの特徴から効果的なブラックリスティングを検討し、その課題を整理している。

第 5 章は、「結論」で、本論文の提案についてまとめている。

付録は、参考文献と、この研究を外部発表したリストをまとめたものである。

This thesis is entitled "Countermeasures technologies against unknown malware", which consists of Chapter 5 and an appendix. In recent years, malware has increased rapidly, and many malware that evade detection of antivirus software also appear and become a problem. In this paper, we propose and verify three different approaches and aim to be able to detect these malware as light as possible before malware infection.

Chapter 1, "Introduction", describes the background of this research. The first is the limits of detection of anti-virus software because of explosive growth of malware. The second is the limits of detection of malware which includes zero-day vulnerability as used in targeted attacks. The third is the issue of diversification of malicious sites on the Internet where malware is distributed. Next, the conventional anti-malware technology category is classified into 3 malware detection areas and 2 malware communication detection areas, totaling 5 malware detection areas. We mapped the three approaches we engaged in this research into these five fields and clearly indicated the purpose of this research that to lightly detect malware before infection. We further characterized the three approaches either indirectly or directly, and summarize the merits of each approach.

Chapter 2 entitled "Malware Detection Using Statistical Methods" performs a malware detection using statistical methods for characteristics of files. We propose a lightweight malware judging method by constructing a suitable model from a large number of header information of a file. We concluded that this method is particularly effective for widely distributed malware, such as distributed from a malware download site on the Internet. By this method, malware that could not be detected by antivirus software is defined as unknown malware, and it shows that they can be detected. In the evaluation experiments, we evaluated the identification accuracy of malware or benign in terms of known or not, packed or not. As a result, we showed that the discrimination performance is higher than the conventional method. Furthermore, we showed the higher discrimination performance by applying our model into support vector machine which is a representative machine learning method.

Chapter 3 entitled "ROP Code Detection Using Characteristics of Attack

Code Structure" detects ROP code by using the characteristics that ROP code is arranged outside the shellcode encoding routine. This method is not a widely used malware that is the detection target in Chapter 2, but it indicates that it is effective for a specific malware used together with a target type attack and an unknown zero-day vulnerability. The ROP attack code is added to many attack codes for the purpose of breaking the host side defense mechanism in the vulnerability of the past few years. In particular, there are the following way that are used in recent targeted attacks. The attacker sends a malicious document file in which the attack code is embedded. Since the malicious document file is difficult to distinguish from a normal file at first glance, the victim opens the file and malware infections occur. After that, information is stolen by attacker. By statically detecting this ROP attack code, we propose a malicious document file judgment method and evaluate it with actual malicious document samples. Furthermore, we analyzed the ROP code commonly used in each sample, showed that some ROP codes can be used by different vulnerabilities. Therefore, by detecting these ROP codes, the possibility of countermeasures against zero-days vulnerability.

Chapter 4 entitled "Consideration of Effective Blacklisting", we do not detect the malware by using characteristics of file as shown in Chapter 2 and Chapter 3, but we will effectively take measures against unknown malware by blacklist the malware download site. As a result, even if unknown malware is downloaded due to a change in features such as a hash value of the malware itself, it can be detected. In this study, we analyzed approximately 43,000 malware download URLs to investigate malware distribution and the behavior of malware download sites over an extended period, i.e., over 1.5 years. We found that some sites survive for a very long time and are revived frequently, a finding not revealed in previous research. By focusing on the malware variation, we have identified three categories of malware download sites. We also analyzed sites in terms of IP address changes, anti-virus application results, URL features. We found that each category had different attacker operational and resource characteristics. Finally, based on our findings, we discuss effective countermeasures for each

category.

Chapter 5 "Conclusion" summarizes proposal and contribution of this thesis. The appendix summarizes the references and a list of papers about this thesis.