

MALICIOUS URL DETECTION BASED ON MACHINE LEARNING

S.SHANMUGAPRIYA , PREMKUMAR .S

Abstract - The project “MALICIOUS URL DETECTION BASED ON MACHINE LEARNING” is developed to Currently, the risk of network information in security is increasing rapidly in number and level of danger. The methods mostly used by hackers today is to attack end-to end technology and exploit human vulnerabilities. These techniques include social engineering, phishing, pharming, etc. One of the steps in conducting these attacks is to deceive users with malicious Uniform Resource Locators (URLs). As a results, malicious URL detection is of great interest nowadays. There have been several scientific studies showing a number of methods to detect malicious URLs based on machine learning and deep learning techniques. In this paper, we propose a malicious URL detection method using machine learning techniques based on our proposed URL behaviours and attributes. Moreover, bigdata technology is also exploited to improve the capability of detection malicious URLs based on abnormal behaviours. In short, the proposed detection system consists of a new set of URLs features and behaviours, a machine learning algorithm, and a bigdata technology. The experimental results show that the proposed URL attributes and behaviour can help improve the ability to detect malicious URL significantly. This is suggested that the proposed system may be considered as an optimized and friendly used solution for malicious URL detection.

I. INTRODUCTION

The project “MALICIOUS URL DETECTION BASED ON MACHINE LEARNING” is Uniform Resource Locator (URL) is used to refer to resources on the Internet. In, Sahoo et al. presented about the characteristics and two basic components of the URL as: protocol identifier, which indicates what protocol to use, and resource name, which specifies the IP address or the domain name where the resource is located. It can be seen that each URL has a specific structure and format. Attackers often try to change one or more components of the URL's structure to deceive users for spreading their malicious URL. Malicious URLs are known as links that adversely affect users. These URLs will

redirect users to resources or pages on which attackers can execute codes on users' computers, redirect users to unwanted sites, malicious website, or another phishing site, or malware download. Malicious URLs can also be hidden in download links that are deemed safe and can spread quickly through file and message sharing in shared networks. Some attack techniques that use malicious URLs include Drive-by Download, Phishing and Social Engineering, and Spam. According to statistics presented in, in 2019, the attacks using spreading malicious URL technique are ranked first among the 10 most common attack techniques. Especially, according to this statistic, the three main URL spreading techniques, which are malicious URLs, botnet URLs, and phishing URLs, increase in number of attacks as well as danger level. From the statistics of the increase in the number of malicious URL distributions over the consecutive years, it is clear that there is a need to study and apply techniques or methods to detect and prevent these malicious URLs.

Regarding the problem of detecting malicious URLs, there are two main trends at present as malicious URL detection based on signs or sets of rules, and malicious URL detection based on behaviour analysis techniques. The method of detecting malicious URLs based on a set of markers or rules can quickly and accurately detect malicious URLs. However, this method is not capable of detecting new malicious URLs that are not in the set of predefined signs or rules. The method of detecting malicious URLs based on behaviour analysis techniques adopt machine learning or deep learning algorithms to classify URLs based on their behaviours.

In our research, machine learning algorithms are used to classify URLs based on the features and behaviours of URLs. The features are extracted

S.Shanmugapriya, Assistant Professor , Department of Computer Applications , Erode Sengunthar Engineering College (Autonomous), Perundurai , Erode.

(Email : riyashanmu@gmail.com)

Premkumar.S, PG Scholar , Department of Computer Applications, Erode Sengunthar Engineering College (Autonomous), Perundurai , Erode. (Email : kumarprem25498@gmail.com).

from static and dynamic behaviours of URLs and are new to the literature. Those newly proposed features are the main contribution of the research. Machine learning algorithms are a part of the whole malicious URL detection system. Two supervised machine learning algorithms are used, Support vector machine (SVM) and Random Forest (RF). The paper is organized as follows. Section II reviews some recent works in the literature on malicious URL detection. The proposed malicious URLs detection system using machine learning is presented in Section III. In this section, the new features for URLs detection

process is also described in details. Experimental results and discussions are provided in Section IV. The paper is concluded by Section V.

II. OBJECT DETECTION- AN OVERVIEW

The experimental results show that the proposed URL attributes and behaviour can help improve the ability to detect malicious URL significantly. This is suggested that the proposed system may be considered as an optimized and friendly used solution for malicious URL detection. Regarding the problem of detecting malicious URLs, there are two main trends at present as malicious URL detection based on signs or sets of rules, and malicious URL detection based on behaviour analysis techniques. The method of detecting malicious URLs based on a set of markers or rules can quickly and accurately detect malicious URLs. However, this method is not capable of detecting new malicious URLs that are not in the set of predefined signs or rules. The method of detecting malicious URLs based on behaviour analysis techniques adopt machine learning or deep learning algorithms to classify URLs based on their behaviours. They propose a malicious URL detection method using machine learning techniques based on our proposed .

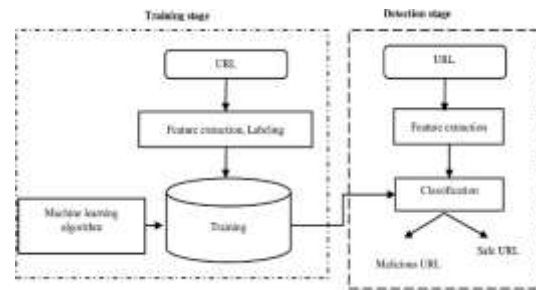


Figure: malicious URL detection.

URL behaviours and attributes. Moreover, bigdata technology is also exploited to improve the capability of detection malicious URLs based on abnormal behaviours. In short, the proposed detection system consists of a new set of URLs features and behaviours, a machine learning algorithm, and a bigdata technology. Malicious URLs are known as links that adversely affect users. These URLs will redirect users to resources or pages on which attackers can execute codes on users' computers, redirect users to unwanted sites, malicious website, or another phishing site, or malware download. Malicious URLs can also be hidden in download links that are deemed safe and can spread quickly through file and message sharing in shared networks. Some attack techniques that use malicious URLs include Drive-by Download, Phishing and Social Engineering, and Spam.

III. LITERATURE SURVEY

Malicious URL, a.k.a. malicious website, is a common and serious threat to cybersecurity. Malicious URLs host unsolicited content (spam, phishing, drive-by exploits, etc.) and lure unsuspecting users to become victims of scams (monetary loss, theft of private information, and malware installation), and cause losses of billions of dollars every year. It is imperative to detect and act on such threats in a timely manner. Traditionally, this detection is done mostly through the usage of blacklists. However, blacklists cannot be exhaustive, and lack the ability to detect newly generated malicious URLs. To improve the generality of malicious URL detectors, machine learning techniques have been explored with increasing attention in recent years. This article aims to provide a comprehensive survey and a structural understanding of Malicious URL

Detection techniques using machine learning. We present the formal formulation of Malicious URL Detection as a machine learning task, and categorize and review the contributions of literature studies that addresses different dimensions of this problem (feature representation, algorithm design, etc.).

A. SYSTEM ANALYSIS

A. Signature based Malicious URL Detection

Studies on malicious URL detection using the signature sets had been investigated and applied long time ago.

Most of these studies often use lists of known malicious URLs. Whenever a new URL is accessed, a database query is executed.

IV. EXISTING SYSTEM

If the URL is blacklisted, it is considered as malicious, and then, a warning will be generated; otherwise, URLs will be considered as safe. The main disadvantage of this approach is that it will be very difficult to detect new malicious URLs that are not in the given list. Machine Learning based Malicious URL Detection There are three types of machine learning algorithms that can be applied on malicious URL detection methods, including supervised learning, unsupervised learning, and semi supervised learning. And the detection methods are based on URL behaviours. In [1], a number of malicious URL system based on machine learning algorithms have been investigated. Those machine learning algorithms include SVM, Logistic Regression, Nave Bayes, Decision Trees, Ensembles, Online Learning, act. In this paper, the two algorithms, RF and SVM, are used. The accuracy of these two algorithms with different parameters setups will be presented in the experimental results.

THE DRAWBACKS OF EXISTING SYSTEM

- The system is not implemented Machine Learning Algorithm Selection.
- The system is not implemented URL Attribute Extraction and Selection.

V. PROPOSED SYSTEM

In the proposed system, machine learning algorithms are used to classify URLs based on the

features and behaviours of URLs. The features are extracted from static and dynamic behaviours of URLs and are new to the literature. Those newly proposed features are the main contribution of the research. Machine learning algorithms are a part of the whole malicious URL detection system. Two supervised machine learning algorithms are used, Support vector machine (SVM) and Random Forest (RF).

ADVANTAGES OF THE PROPOSED SYSTEM

- The proposed algorithms are suitable to utilized the usefulness of our new features selected for malicious URL detection.
- In the proposed work, SVM and RF are selected as an example to illustrate the good performance of the whole detection system, and are not our main focus. Readers are encouraged to implement some other algorithms such as Naïve Bayes, Decision trees, k-nearest neighbours, neural networks, etc.

VI. SYSTEM IMPLEMENTATION

In this module, the Service Provider has to login by using valid user name and password. After login successful he can do some operations such as

Login, Browse URLs Datasets and Train & Test Data Sets, View URLs Datasets Trained and Tested Accuracy in Bar Chart, View URLs Datasets Trained and Tested Accuracy Results, View Prediction of URLs Type, View URLs Type Ratio, Download Predicted Data Sets, View URLs Type Ratio Results, View All Remote Users

VII. MODULES

- Service Provider
- View and Authorize Users
- Remote User

1) VIEW AND AUTHORIZE USERS

In this module, the admin can view the list of users who all registered. In this, the admin can view the user's details such as, user name, email, address and admin authorize the users.

2) REMOTE USER

In this module, there are n numbers of users are present. User should register before doing any operations. Once user registers, their details will be stored to the database. After registration successful, he has to login by using authorized user name and password. Once Login is successful user will do some operations like register and login, predict URLs type, view your profile.

VIII. CONCLUSION

In this paper, a method for malicious URL detection using machine learning is presented. The empirical results in Tables V and VI have shown the effectiveness of the proposed extracted attributes. In this study, we do not use special attributes, nor do we seek to create huge datasets to improve the accuracy of the system as many other traditional publications. Here, the combination between easy-to-calculate attributes and big data processing technologies to ensure the balance of the two factors is the processing time and accuracy of the system. The results of this research can be applied and implemented in information security technologies in information security systems. The results of this article have been used to build a free tool [20] to detect malicious URLs on web browser

IX. REFERENCES

- [1] D. Sahoo, C. Liu, S.C.H. Hoi, "Malicious URL Detection using Machine Learning: A Survey". Corer, abs/1701.07179, 2017.
- [2] M. Kohji, Y. Iraqi, and A. Jones, "Phishing detection: a literature survey," IEEE Communications Surveys & Tutorials, vol. 15, no. 4, pp 2091–2121, 2013.
- [3] M. Cova, C. Kregel, and G. Vigna, "Detection and analysis of drivebydownload attacks and malicious JavaScript code," in Proceedings of the 19th international conference on World wide web. ACM, 2010, pp. 281– 290.
- [4] R. Hartfield and G. Loukas, "A taxonomy of attacks and a survey of defines mechanisms for semantic social engineering attacks," ACM Computing Surveys (CSUR), vol. 48, no. 3, p. 37, 2015.
- [5] Internet Security Threat Report (ISTR) 2019–Symantec. <https://www.symantec.com/content/dam/symantec/docs/reports/istr-24-2019-en.pdf> [Last accessed 10/2019].
- [6] S. Sheng, B. Wardman, G. Warner, L. F. Cramer, J. Hong, and C. Zhang, "An empirical analysis of phishing blacklists," in Proceedings of Sixth Conference on Email and Anti-Spam (CEAS), 2009.
- [7] C. Seifert, I. Welch, and P. Kisarazu, "Identification of malicious web pages with static heuristics," in Telecommunication Networks and Applications Conference, 2008. ATNAC 2008. Australasian. IEEE, 2008, pp. 91–96.
- [8] S. Sinha, M. Bailey, and F. Jahani an, "Shades of grey: On the effectiveness of reputation-based "blacklists"," in Malicious and Unwanted Software, 2008. MALWARE 2008. 3rd International Conference on. IEEE, 2008, pp. 57–64.
- [9] J. Ma, L. K. Saul, S. Savage, and G. M. Voelker, "Identifying suspicious URLs: an application of large-scale online learning," in Proceedings of the 26th Annual International Conference on Machine Learning. ACM, 2009, pp. 681–688.
- [10] B. Eshete, A. Villafuerte, and K. Teklemariam, "Inspect: Holistic analysis and detection of malicious web pages," in Security and Privacy in Communication Networks. Springer, 2013, pp. 149–166.
- [11] S. Parfait, "Phishing counter measures and their effectiveness— literature review," Information Management & Computer Security, vol. 20, no. 5, pp. 382–420, 2012.
- [12] Y. Tao, "Suspicious urn and device detection by log mining," Ph.D. dissertation, Applied Sciences: School of Computing Science, 2014.
- [13] G. Canfora, E. Med vet, F. Mercado, and C. A. Visagie, "Detection of malicious web pages using system calls sequences," in Availability, Reliability, and Security in Information Systems. Springer, 2014, pp. 226–238.
- [14] Leo Bierman.: Random Forests. Machine Learning 45 (1), pp. 5- 32, (2001).
- [15] Thomas G. Dieterich. Ensemble Methods in Machine Learning International Workshop on Multiple Classifier Systems, pp 1- 15, Cagliari, Italy, 2000.
- [16] Developer Information. https://www.phishtank.com/developer_info.php. [Last accessed 11/2019].
- [17] Uraeus Database Dump. <https://urlhaus.abuse.ch/downloads/csv/>. [Nagy troy nap 11/2019].
- [18] Dataset URL. http://downloads.majestic.com/majestic_million.csv. [Last accessed 10/2019].
- [19] Malicious_n_non-MaliciousURL. <https://www.kaggle.com/antonyj453/urldataset#data.csv>. [Last accessed 11/2019].
- [20] chrome.zip. https://drive.google.com/file/d/13G_Ndr4hMFx_qWyTEjHuOyJmHFWd0Gud/viewbicle=IwAR0SLVCrvjHHGmoHZH97nXN3BmDMY7jG4S0sKZYLAZjTFgeoJADfli64-g. [Last accessed 12/2019].