

A Survey on malicious detection technique using data mining and analyzing in web security

¹Komal Ahuja, ²Amandeep
Mtech student (CSE)
GRIMT Radaur, Haryana, India

Abstract - The security in email refers to the collective measures, used to secure the access and content of email account or service. It is allow an individual or organization to protect the overall access to one or more email address and accounts. Any program that is created to harm the computer system operations or data is termed as malicious programs. Malicious programs contain viruses, worms, adware, Trojans, backdoors, spywares, bots, root kits etc. All malwares are sometimes loosely termed as virus (worms, Trojans) Commercial anti-malware products are still called antivirus. The main purpose of this survey paper is based on study of various technique i.e. security enhancement, detection and protection from malwares or many technique is used for malware detection and analyzing the behavior of malwares.

Index Terms - Data mining, Worm detection, Security, Malware, Propagation, Detection, NB, KNN

I. INTRODUCTION

Emails can help initiate an attack in two ways—via attachments or URLs. Successful attacks have range from malware downloads to phishing incidents to targeted attack, which result in compliance concern issues, data breaches, financial loss and Phishing attempts to acquire data from a specific companies or individuals. Attackers collect the personal information and increase their probability of success. Normal documents difficult to differentiate from exploit documents. A solution have ability of uncovering known and zero-day used in attachments like, Adobe PDF, Microsoft® Office® or the other document formats can offer more advanced defenses. In order to effectively address the spam issue in the untrustworthy Internet, the author argue that receivers must gain more control over if and when a message should delivered to them. Asynchronous messages on the Internet are delivered using two different models: Receiver-pull and sender pull. They differ in who initiates the message delivery process.

II. MALICIOUS CODE

Malware includes computer viruses, Trojan horses, ransom ware, root kits, key loggers, dialers, worms, spyware, adware and other malicious programs. The main part of active malware threats are usually Trojans or worms rather than viruses fig 1. Shows Anti-virus vendors are facing huge quantities of suspicious files every day. These files are come together from variety of sources that contain third party providers, files reported by customers and dedicated honey pots either explicitly or automatically.

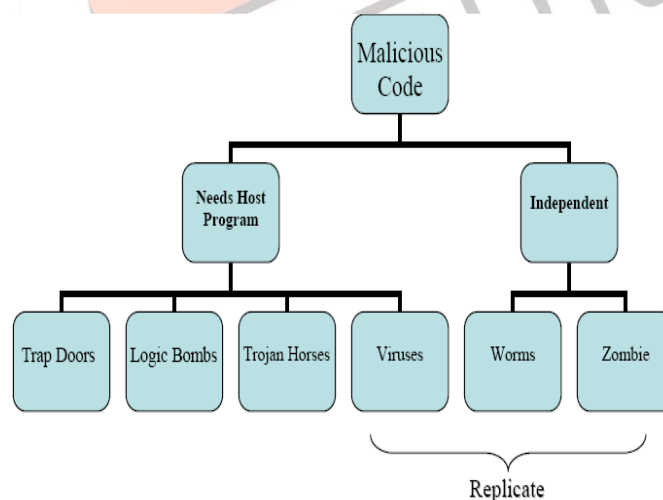


Figure 1:- Malicious Code infrastructure

A. Malware Detection

It is important to understand the behavioral feature of a malware that only possible by executing malware binary. Execution of malware binary has normal layouts or may show signs of some irregular termination in execution. Detection is about identifying whether code is malicious or genuinely benign. Robust malware detection depends on the capability of difficult to handling malware efficiently. The Code obfuscation changes malware syntax but not its intended behavior, which

has to be maintained. In both reverse engineering and deobfuscating techniques generally begin with some type of static program analysis, which can be represented as an idea of program semantics. For this purpose a component called as Malware detector, which can be explained as a system that attempts to identify malware using signatures and other heuristics parameters. Two common obfuscation techniques are Polymorphism and Metamorphism. The malware detector can be terminate as the safeguard interface of system whose functionality depends upon the behavior of executables. Two input components for malware detector.

- Signature or behavioral parameters of given code.
- Executable code under inspection.

a) *Polymorphic Malware :*

If a virus is in running order so that it look different every times it repeat again and again, but keeping the original code intact. This type of virus is referred as polymorphic virus. It contain of malicious code which is encrypted along with decrypted module. Polymorphic code is a method now commonly execute in malware that uses a polymorphic generator to change the code while protect the original algorithm intact.

b) *Metamorphic Malware:*

These are a type of body-polymorphic, where body of virus itself changes from one instance to another. Metamorphic malwares used the different types of obfuscation techniques to reprogram themselves into a new transformed code such as similar to original code. The metamorphic nature of the malware enables malicious code to mutate while spreading from one side to other side the network and produce signature based detection completely ineffective.

III. TECHNIQUES USES IN MALICIOUS CODE

A. Signature-Based Techniques

The most of the antivirus tools are based on detection with signature based techniques. These signatures are generated by examining the disassembled code of malware binary. Many disassemblers and debuggers are available which help in disassembled the portable executable. The Disassembled code is analyzed and features are extracted. These features are used for collect the signature of individual malware family.

B. Behavior-Based Techniques

The goal is to be analyzing the behavior of known or unknown malwares. Behavior parameters include various factors i.e. source/destination address of malwares, types of attachment and countable statistical features. Accordingly each of the above technique can be further applied to using static analysis, dynamic analysis or hybrid analysis.

a) *Data Mining:*

Data mining methods detected the patterns in large amounts of data, i.e. byte code, and used these patterns to detect future instances in the similar data. Their structure used classifiers to detect new malicious and execute. A classifier is a rule set and detection model generated by the data mining algorithm, which was trained over a given set of training data. One of the main problems faced by the virus community is to device methods for detecting new malicious programs have not been analyzed. In every day eight to ten malicious programs are created and most cannot be accurately detected until than signatures have been generating for them. During this time, systems are protected by signature-based algorithms are vulnerable to attacks.

b) *Naive Bayes classifier:*

This algorithm was importantly a collection of Naïve Bayes algorithms that supported on whole organization for an example. In Naive Bayes algorithm, it is able to classify the examples in the test set of malicious executables program. This method is used a machine with 1GB of RAM, and the size of the binary data was very big to get in to memory. Thus solve this problem we divided it in to smaller parts that could easily get into memory and hence training the naïve bayes (NB) algorithm. The Naive Bayes algorithm required a table chart of all strings or bytes to assess its possibilities. In every classifier, there is a rule set. In short it is used to calculate a group of cluster for ambiguity or malicious code.

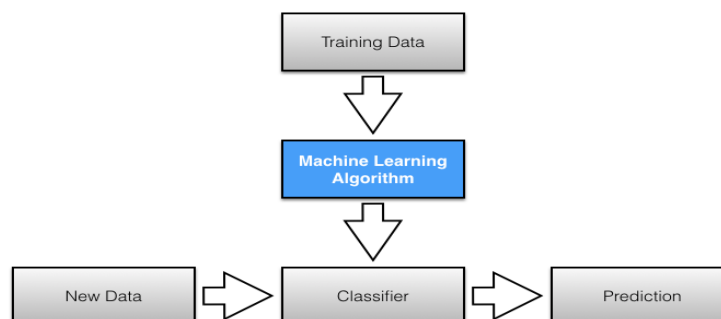


Figure 2: A simplified diagram of the general model building procedure for pattern classification.

IV. RELATED STUDY

Dipali Kharche¹, Anuradha Thakare (2015), Internet worm means malware computer programs that spread one computer to another computer and repeated itself. Malware includes, worms, computer viruses, root kits, dialers, adware, Trojan horse, malicious, spyware,. It is programmed by attackers to interrupt computer process, gather delicate Information, or gain entry to private computer systems. [1]

Ms. Milan Jain Dr. Bikram Pal (2014), We have three data mining algorithms to produce new classifiers with separate features: Naïve – Bayes, a Multi Classifier and RIPPER, and the comparison between three methods. It contain root kit data collection, performance and classification evaluation phases and data pre-processing.[2]

Rajkumar E.V.1 Aravindharamanan. (2014), Securing the web is a huge challenge that the modern era of computers have seen. The threat levels increase day by day there by making the network vulnerable to attacks To protect websites from attacks. many innovative strategies are brought into the field of cyber security. But still malware has remained a serious cause of concern to server administrators and web developers. [3]

Dharmesh Kumar Babubhai Patel, Sahjanand Harshadbhai Bhatt (2014), Discuss various data mining (DM) techniques that author have successfully applied for cyber security. This research investigates the use of data mining (DM) methods for malware detection and proposed a framework is possible for the traditional signature detection method. These applications are contain malicious code detection by mining binary executables by anomaly detection, and data stream mining [4]

Shubair Abdulla, Sureswaran Ramadass (2014), presents a worm detection system that leverages the effectiveness of learning machines and the reliability of IP-Flow. Typically, a host infected by an email worm or scanning initiates a significant amount of traffic that does not rely on Domain name system to translate names into numeric IP addresses [5]

Bhura Parull Rajiv Kumar Gurjwar (2014), Suggested four types of attack discussed in existing system but there are some problems mentioned is related to classifier algorithms: More than one algorithm used for each layer in layered approach. RFA gives good classification rate (i.e 99.16%, 99.25% for PROBE,R2L and DOS layer respectively).[6]

S.A. Joshi, Varsha S.Pimprale (2013), With the tremendous development in information technology, network security is one of the challenging matter and so as Intrusion Detection system. IDS are an essential component of the network to be secured. The traditional Intrusion Detection system (IDS) are difficult to manage various newly arising attacks. To deal with these new problems of networks, data mining based Intrusion Detection system (IDS) are opening new research avenues [7].

De Ocampo,Frances Bernadette c Del Castillo Trisha Mari L (2013), A Signature-based IDS helps in maintaining the integrity of data in a network controlled environment. Unfortunately, this type of Intrusion Detection system depends on predetermined intrusion patterns that are basically created. If the signature database of the Signature-based Intrusion Detection system is not updated, network attacks just pass through the intrusion detection system (IDS) without being noticed [8].

Raviraj Choudhary, Ravi Saharan (2011), present a approach that conducts an feature search on a set of computer viruses. These methods use mnemonics patterns to detect future instances in similar data and detect mnemonics patterns in large amounts of data, and. The author use apriori algorithm for select features to detect the malicious executables. Through those features author make detection model for trained over a given set of training data. [9]

Irena Koprinska, Josiah Poon (2006), study semi-supervised and supervised classification of e-mails. Author consider two tasks: spam e-mail filtering and filing e-mails into folders. Firstly, in a supervised learning setting, The author investigate the use of random forest for automatic spam e-mail filtering and e-mail filing into folders .The author show that random forest is good choice for these tasks as it runs fast on high dimensional databases, is highly accurate and is easy to tune, outperforming popular algorithms i.e. support vector machines and decision trees [10].

Sr. No	Techniques	Author	Year	Finding
1.	Data mining technique, Random Forest, Decision Tree, Bayesian Network	Dipali Kharche, Anuradha Thakare	2015	Data mining technique is efficient for detect the internet worms.
2.	Multiple Classification Algorithms, Clustering Techniques, nearest neighbor method, genetic algorithm	Ms. Milan Jain, Ms. Milan Jain	2014	Multiple classification algorithm is used for protected the computer system from malwares and it is consume less time. The main challenging issue is a malicious problem.
3.	Malware detection , analysis web security	Rajkumar E.V.1 Aravindharamanan	2014	The review many malware detection systems, analyzing them and also depicts various ways of detecting them. A new simulation must be designed to contain real system samples.
4.	Data mining technique and classification techniques	Dharmesh Kumar Babubhai Patel, Sahjanand Harshadbhai Bhatt	2014	Data-mining framework that detects new, previously unseen malicious executables accurately and automatically. This approach is costly and oftentimes ineffective
5.	Machine Learning Algorithms, k-nearest	Shubair Abdulla, Sureswaran Ramadass	2014	High curacy rates performance and good ac NB algorithm is much faster in terms of prediction

	neighbors, Naïve Bayes (NB)			time.
6.	NIDS technique, J48, SVM, Random Forest, Random Tree, Rotation Forest, KNN	Bhura Parul1, Rajiv Kumar Gurjwar	2014	The Layered approach is very effective system & gives good result. Classification rate very low.
7.	Network Intrusion Detection System, data mining	S. A. Joshi, Varsha S. Pimprale	2013	Efficiency and accuracy of intrusion detection system are increased. Network security is one of the issue in the IDS.
8.	Signature based intrusion detection system with network attack, Machine learning	De Ocampo, Frances Bernadette c Del Castillo Trisha Mari L	2013	Signature-based Intrusion Detection System (IDS) helps in maintaining the integrity of data in a network controlled environment. The goal is to correlate data between the logs of the anomaly-based IDS and the packet that has capturing in order to determine if a network traffic is really malicious or not.
9.	Apriori algorithm	Raviraj Choudhary, Ravi Saharan*2	2011	Using apriori algorithm we can select top L features. These L features can be used to recognize whether a file is virus or not.
10.	Supervised and semi-supervised classification of e-mails, RF approach, TFV	Irena Koprinska, Josiah Poon	2006	It runs fast on large and high dimensional databases, is easy to tune and is highly accurate. Majority of datasets consist of only a single set of features with no obvious way to divide them.

a. Technique to detect the malicious code

V. CONCLUSION

Various papers in literature have been studied. Data mining technique is very useful for detections of internet worms. Data mining is used such as Random forest, decision tree. Multiple classification algorithms are used i.e. nearest neighbor method, genetic algorithm. In this algorithm is efficient for computer systems can be protected from malwares and consume less time but main issue in this technique it has malicious problem. We have implemented the apriori algorithm and select the L features using the apriori algorithm. L feature is check the weather a file is virus or not.

VI. REFERENCES

- [1] Dipali Kharche, Anuradha Thakare, "Internet Worm Classification and Detection using Data Mining techniques" May-June 2015.
- [2] Ms. Milan Jain and Dr Bikram Pal, "Malicious Detection Using Multiple Classification Algorithms & Their Comparison Using Different Clustering Techniques" August 2014
- [3] Rajkumar E.V.1 Aravindharamanan, "A state of the art review on various malware detection and analysis in web security" June 2014.
- [4] Dharmesh Kumar Babubhai Patel1, Sahjanand Harshadbhai Bhatt, "Implementing Data Mining for Detection of Malware from Code" April 2014.
- [5] Shubair Abdulla1, Sureswaran Ramadass2, Altyeb Altaher2, and Amer Al-Nassiri, "Employing Machine Learning Algorithms to Detect Unknown Scanning and Email Worms" March 2014.
- [6] Bhura Parul and Rajiv Kumar Gurjwar, "A Review on attacks classification using decision tree algorithm" Feb 2014.
- [7] S.A. Joshi, Varsha S. Pimprale, "Network Intrusion Detection System (NIDS) based on Data Mining" Jan 2013.
- [8] De Ocampo, Frances Bernadette C. and Del Castillo, Trisha Mari L., "Automated signature creator for a signature based intrusion detection system with network attack detection capabilities" 2013.
- [9] Raviraj Choudhary and Ravi Saharan, "Feature detection approach from viruses through mining" June 2011.
- [10] Irena Koprinska, Josiah Poon, James Clark, Jason Chan, "Learning to classify e-mail" 2007.