

An Advanced Approach to Polymorphic/Metamorphic Malware Detection using Hybrid Clustering Approach

¹Harshit Kumar, ²Parvinder Kaur
¹M.E Scholar, ²Assistant professor
^{1,2}SUSCET, Tangori, Mohali, India

Abstract - Malware Classification has been a challenging problem in the recent past and several researchers have attempted to solve this problem using various tools. It is security threat which can break machine operation while not knowing user's data and it's tough to spot its behavior. This paper proposes a novel technique using DBSCAN (Density based Kmeans) algorithmic rule to spot the behavior of malware. After classification from DBSCAN, pattern matching is applied using the instructions pattern in the generated reports. Among of these techniques a pattern based mostly technique is well famed for the detection of malware. For the moderation and improvement of the present system the signature based mostly technique is most popular. The results are found to be quite accurate and better than the existing ones in terms of accuracy.

Keywords - Malware, its sorts and detection based mostly techniques.

I. INTRODUCTION

Nowadays, there are numbers of computer security threat that cause by the malware attack have extraordinarily enlarged. It includes seventeen categories like worms, viruses, spyware, bug and together several malicious packages that cause a billion of losses to the operation worldwide. However, all types of malware have their specific objective, the foremost purpose is to interrupt the computer operation.

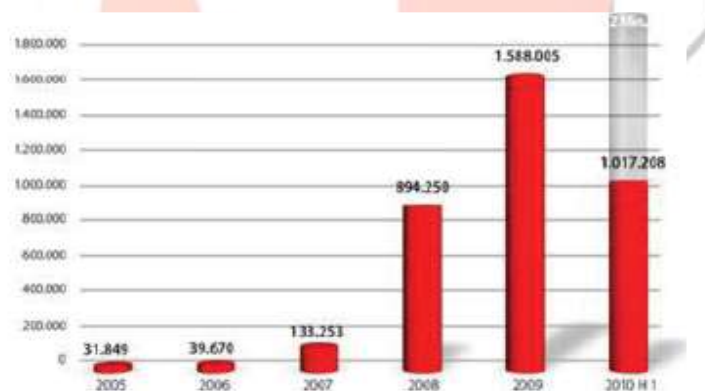


Figure. 1: Malware growth since 2005 to 2010

Malware is a one type of software which can harm the computer's operating system and may also can steal the personal information from the computer, malware can be made by using any programming language by the programmer. It is very difficult to define a malware with a single term or a single name. A malware can be consider as a **malicious software** or **malcode** or as a malicious code. Malware do the bulk of the intrusive activities on a system and that spreads itself across the hosts in a network. Malware is defined as software performing actions intended by an attacker, mostly with malicious intentions of stealing information, identity or other resources in the computing systems. There are different types of malware like adware, bots, Trojan horses, viruses, bugs, rootkits, spyware and worms. However, attributable to the technology advancement many malware writers try to use higher concealment techniques to avoid detection. The concealment technique is created with the combination of previous behaviour therefore on attack and at identical time to avoid the signature-based detection. In this, several common techniques that are commonly used like as polymorphic and metamorphic. Additionally analysis ought to be tired order to hunt out possible resolution on avoiding current technique of malware attacks.

There are four main generations in the gradual development of the stealing methodologies: secret writing, Polymorphism, Oligomorphism, and natural action. It displays the evolution timeline of camouflage techniques in malware.

2. Primitive Malware

When the story began, virus writing has a kind of programming fun for portable computer specialists to point off their technical skills, but it step by step became as a tool for several functions, like swiping the people’s information, like mastercard numbers, passwords, or bank account numbers, or for avenges functions, and so on.

3. Stealth Malware

Malware creators’ initial tries therefore on turn tail from redounded to appear of stealing techniques. Stealth virus is prepared to cover its signs and traces. Virus normally changes and modifies info resources on the system. For example, a file-hosted virus would possibly append its own code to the tip of Associate in possible file. If Associate in application examines the infected file, it'll discover the being code inside the file and catch the virus.

4. Camouflage Evolution

4.1 Encryption

Malware authors invariably try to improve the program to escape from the code analyzer technicians. Consequently, they will get longer for his or her created malware to live inside the wild and blow their own horns extra. The simplest and earliest methodology utilized by the malware programmers to achieve this goal was secret writing. Encrypted virus contains two basic sections: a secret writing loop and main body. Decrypt, or secret writing loop, might be a brief piece of code, that's responsible to cipher and rewrite code of main body. The foremost body is that particular code of malware, encrypted, and is not purposeful before its being decrypted by secret writing loop. Once the virus starts to run on host portable computer, initial the decrypt or loop ought to decipher the foremost body into machine possible code and purposeful info.

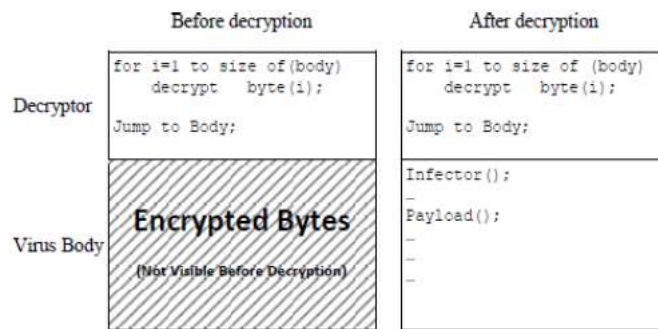


Figure.2: Structure of encrypted virus

4.2 Oligomorphism

It contains a collection of assorted decryptors, that are haphazardly selected for the replacement victim. Oligomorphism is not a heavy drawback for antivirus package as a result of it entirely makes a malware slightly harder to seem at. In distinction to encrypted virus, antivirus engine ought to check all possible decryptor instances instead of looking for only 1 decryptor, and it wishes extended time decryptors.

4.3 Polymorphic

Polymorphic techniques try to create analysis of the virus additional sturdy by ever-changing its look. The principal rule is to change appearance of code constantly, from a reproduction to a unique [7]. It ought to be administered in such the best method so no permanent common string keeps among variants of a deadly illness to be exploited by the antivirus scanner engine for detection purpose. Polymorphic techniques are rather hard to implement and manage.

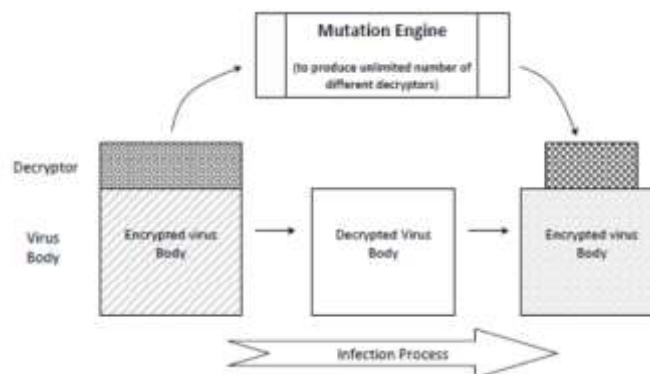


Figure 3: Polymorphic virus structure

4.4 Metamorphism

Unlike the three camouflage generations, metamorphic virus has no encrypted. Therefore, it does not would really like decryptor, but like polymorphic virus, it works a mutation engine, as well, instead of modifying the decryptor loop entirely, it mutates all its body. The foremost sententious definition of metamorphic viruses is introduced by Igor Muttik “Metamorphics are body-

polymorphics". Each new copy may need completely totally different structure, code sequence, size and descriptive linguistics properties, but behaviour of the virus does not modification.

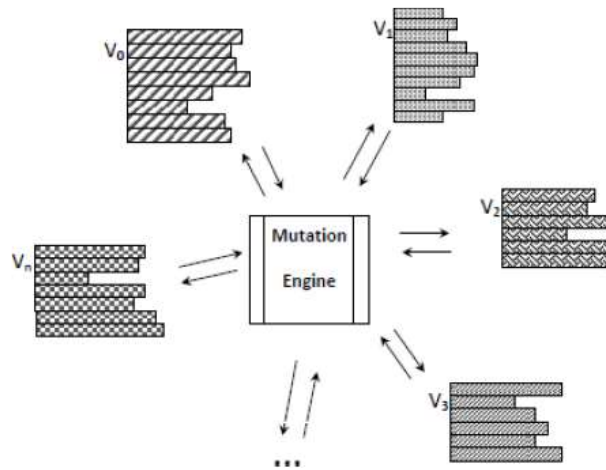


Figure.4: Metamorphic virus propagation scheme

Components of mutation engine square measure displayed. Once the virus finds things of its own code, it's to convert code into assembly instruction, that's finished by an inside dis-assembler. The code analyzer is responsible to give information for code device module. The code device could have some information like structure and flow chart of the program, subroutines, life quantity of variables and registers, and so on.

II. LITERATURE SURVEY

Malware could o be a program with the malicious intent designed to interrupt computer on it executes the network on that it communicates [1]. However all kinds of malware have specific objective, the foremost purpose is to interrupt computer operation. As results, the protection management should be enforced therefore on safeguard all code and knowledge against the replacement, modification or sub versioned. Signature-based matching technique is told the modish approaches to malware detection [2]. This scheme was commercially applied by the anti-virus or anti-spyware product among market. The foremost limitation of signature based technique is detection new malware. This technique is used distinctive bytes string forever fails to take unseen malware that typically known as zero day attacks.

In this, a flexible and automatic approach to extract the malware behavior by the perceptive systems that operate calls performed during virtualized execution in surroundings. Similarities and distances between the malware behaviors unit of the measurement computed that lets classifying malware [3] behaviours. The foremost choices of approach reside in coupling a sequence alignment technique to figure out similarities and leverage the distance to figure out associated distances. [4] The classification technique projected by work is pattern tree. However, this method still encompasses the limitation because of incorrectly classified number of malware behaviour. It applied machine learning rule to their malware detection method. [5] That technique has pattern adaptive info compression thus on counter the limitation of signature-based technique in present industrial anti-malware tool. The two limitations are signature-based technique. First, not all malicious programs have bit patterns that unit of the measurement proof of their malicious nature or collectively not recorded among the virus book of facts. Second, obfuscated malware that take many types of bit patterns will not engage on signature-based technique. The projected [6] technique used adaptive info compression model and prediction by the partial matching as learning engine to form two compression models. This method works on unstructured input, that is, raw executables, with Associate in underlying math compression model.

In this, consider the structural analysis and comparison of computer code supported by management flows. They used graph to enhance the results of comparison. The difficulties [7] of disassembling the code placed at specific locations within network traffic aren't absolutely self-addressed; however, the methodology is computationally valuable. There is a tendency to found that geological process and obfuscation modification the management flow in such some way that easy comparison on the structural characteristics is probably going to fail. The malware detection and classification tools area unit didn't address variants that may be created mechanically victimization newest polymorphic and metamorphic transformation engines that may manufacture variants that area unit shut to every different. Chouchane and Lakhotia[8] proposed using "engine signatures" to assist in detecting metamorphic malware. Basically, this technique evaluates collected forensic evidence from x86 code segments through a code scoring function. This score is a measure of how likely it is that the code has been generated by a known instruction substituting metamorphic engine.

In lately the pc systems could be a part of our life and have various applications in our society the Signature-based strategies are accustomed discover the Malware [9]. These strategies are used wide by industrial detectors. However assault techniques of malwares progressed reciprocally. New malware uses several techniques like polymorphism and geological process to run off from detection strategies. These techniques modify the look of malicious codes by encrypting the codes or commutation them with the different codes or mix of them in such way that the malicious origin of codes remains unchanging. Bayer et al. was used an (able to be made bigger or smaller) clustering approach to identify and group harmful programs or apps samples that show

almost the same behavior [10]. This approach also (does/completes) energetic/changing analysis to get the execution traces of harmful programs or apps programs using automated tools.

This approach is used to boost the efficiency of dynamic malware analysis systems [11]. It is a large sort of latest malicious files presently appears. It's because of mutations of only variety of malware programs. The projected system avoids analyzing malware binaries that just represent mutated instances of already analyzed polymorphic malware. It drastically decrease the quantity of some time required for analyzing a set of malware programs. Seyedet al. proposed harmful programs or apps detection that applied data mining which is based on the analysis of byte level file content [12]. This way of doing things also designed to provide protection against first day launched harmful programs or apps

PROBLEM FORMULATION

Understanding malware behaviors is useful information for computer system. But it is difficult to identify its behavior and classification. The aim of this thesis is to design a automatic malware classification system based on the reports generated by the executable files. The malwares will be classified into three classifications. First task is to separate the clean files from the malwares and then a text mining and pattern matching algorithm needs to be applied for classifying these malwares into normal malwares and poly/metamorphic malwares. An improved Hybrid DBSCAN algorithm needs to be designed which will take into account the various parameters like icmp counts, number of host calls, affected system part, binary registry key etc.

| S.no. | Author name | Year of publication | Conclusion |
|-------|--|---------------------|---|
| 1 | Mohamad Fadli AmanJantan | 2011 | 1) Presented An Approach for Malware Behavior Identification and Classification. |
| 2 | Xufang Li, Peter K.K. Loh <i>et. al.</i> | 2011 | 1) Mechanisms of Polymorphic and Metamorphic Viruses have been proposed in this paper. 2)This technique is used distinctive bytes string forever fails to take unseen malware |
| 3 | M. Bauer, M.J.G. van Eeten | 2012 | 1)Proposed Obfuscation: The Hidden Malware 2) It applied machine learning rule to their malware detection method |
| 4 | BabakBashari Rad, Maslin Masrom <i>et. al.</i> | 2012 | 1) In this paper Camouflage in Malware: from Encryption to Metamorphism has been presented. 2) This method works on unstructured input, that is, raw executables, with Associate in underlying math compression model. |
| 5 | Vinod P., V.Laxmi <i>et. al.</i> | 2012 | 1)presents Metamorphic Malware Exploration Techniques Using MSA signatures |
| 6 | M. ZubairRafique <i>et al.</i> | 2014 | 1)Evolutionary Algorithms for Classification of Malware Families through Different Network Behaviors |
| 7 | Hira agrawal <i>et al.</i> | 2013 | 1)Proposed Detection of Global, Metamorphic Malware Variants Using Control and Data Flow Analysis 2)This approach is used to boost the efficiency of dynamic malware analysis systems |
| 8 | Seyed Emad Armoun <i>et al.</i> | 2012 | 1) proposed harmful programs or apps detection that applied data mining which is based on the analysis of byte level file content |

III. PROPOSED METHODOLOGY

The thesis proposes a novel algorithm to classify malwares as clean, normal malwares and polymorphic/metamorphic malwares. The approach is to generate pydasm report of each malware and then collect the other parameters like binary count,icmp,opcode and machine cycles. Internet Control Message Protocol (**ICMP**) is one of the main protocols of the Internet Protocol Suite. It is used by network devices, like routers, to send error messages indicating.It is a computer file that is not a text file. Many binary file formats contain parts that can be interpreted as text; for example some computer document files containing formatted text, such as older Microsoft Word document files, contain the text of the document but also contain formatting information in binary form

The instruction sets will be extracted from the report via text mining and text preprocessing will be done for various processes like comment removal, function extraction etc. Once the features are extracted then a database corresponding to various files and function list will be created. The similarity score will be calculated between any two files using signature based technique and then a Hybrid DBSCAN algorithm will be applied on the selected feature for classification.

Pattern Matching

The text report of the remaining files is generated and instruction sets are fetched from it for pattern extraction. A pattern matching algorithm is designed on the instruction sets by calculating the scores of matching between two documents. The matching criteria is set by creating a moving window of the instruction set and sliding the instruction set of one document against the other. All the instructions are given a numerical value and the moving window is subtracted from each other. The Scores are calculated on the basis of number of zeros in the result of subtraction.

Simulation Results and Discussion

All the simulation has been done in a computer with 4 GB of RAM and 1.7 GHz processor. The language utilized for scripting is Python. Various Natural Language processing tool kits have been utilized.

Total number of files taken for testing purpose is 25 and by using DBSCAN algorithm for malware classification the 23 files are correct means malware free out of 25 files and 2 files are wrong as shown in below figure. The value of K is taken as 3.

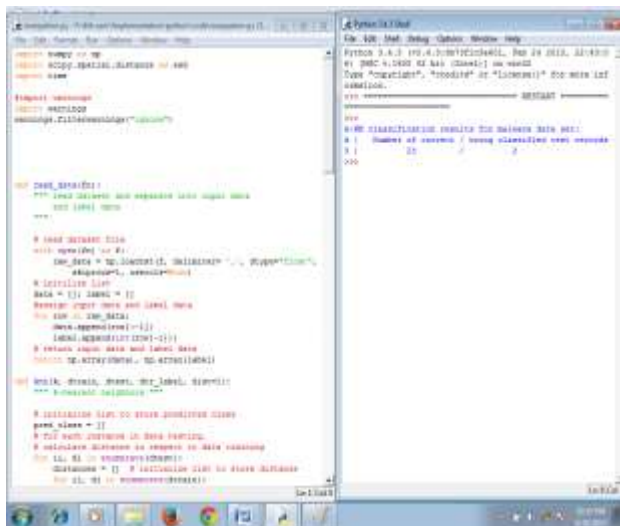


Figure 5: Representation of malware classification by using DBSCAN algorithm

Figure below represents the correctly classified documents vs the incorrectly classified documents. As observed the number of correctly classified documents forms a high percentage of the incorrectly classified ones.

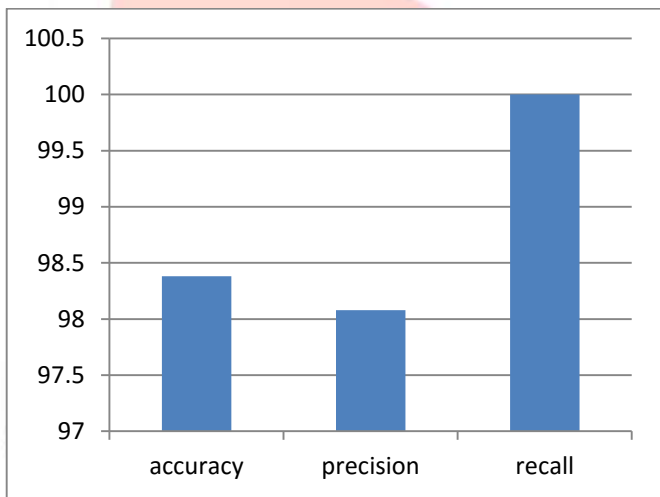


Figure 6: Graphical representation of malware classification by using DBSCAN algorithm

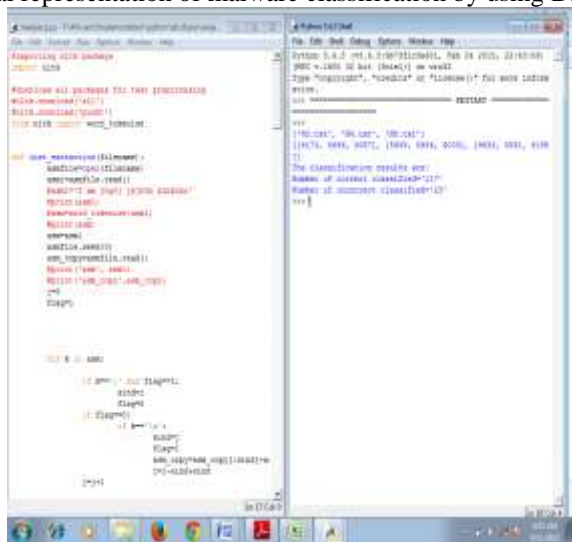


Figure 7: Pre-processing of documents

Figure shows the results of pre-processing of the documents and as it is observed the pre-processing steps enable us to correctly classify the documents.

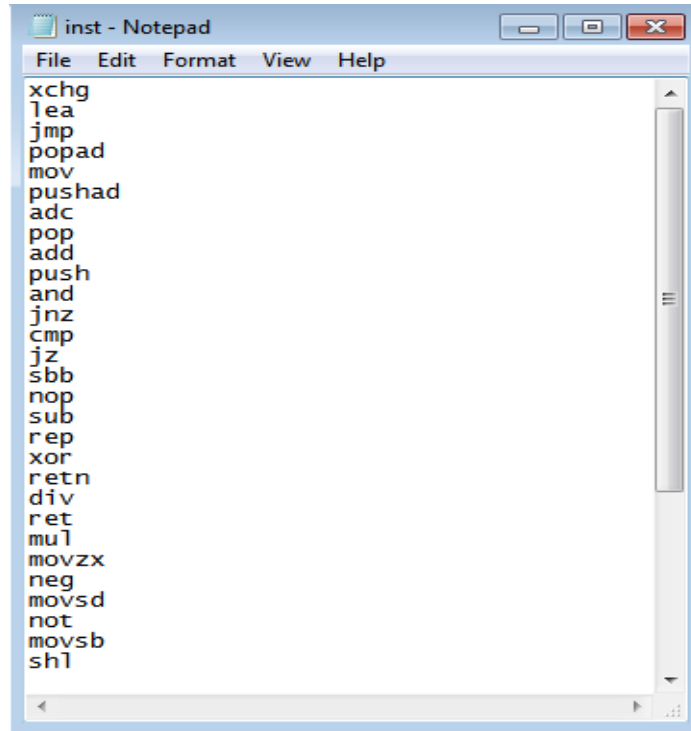


Figure.8: Instruction list

Figure shows the list of instructions which are utilized for pattern matching. These instructions are stored in a text file and read in the code for pattern matching.

Table 1:

| | Clean | Malware | Polydest |
|----------|-------|---------|----------|
| Clean | 26 | 0 | 0 |
| Malware | 3 | 49 | 0 |
| Polydest | 0 | 0 | 105 |

Table 2: Confusion Matrix

| | True | False |
|----------|------|-------|
| Positive | 154 | 3 |
| Negative | 26 | 0 |

$$Accuracy = \frac{154 + 29}{154 + 29 + 3}$$

$$= \frac{183}{186}$$

$$= 98.38\%$$

$$Precision = \frac{Tp}{Tp + Fp}$$

$$Precision = \frac{154}{154 + 3}$$

$$= 98.08\%$$

$$Recall = \frac{Tp}{Tp + Fn}$$

$$Recall = \frac{154}{154}$$

$$= 100 \%$$

REFERENCES

- [1] MohamadFadliAmanJantan “An Approach for Malware Behavior Identification and Classification” ,IEEE 2011.
- [2] Xufang Li, Peter K.K. Loh, Freddy Tan “Mechanisms of Polymorphic and Metamorphic Viruses” European Intelligence and Security Informatics Conference 2011.
- [3] Nitesh Kumar Dixit¹ ,Lokesh mishra² , Mahendra Singh Charan³ and Bhabesh Kumar Dey⁴, “The new age of computer virus and their detection” International Journal of Network Security & Its Applications (IJNSA), May 2012.
- [4] AbdellatifBerkat “Metamorphic computer virus detection by Case Based Reasoning (CBR) methods” International Journal of Software Engineering and Applications (IJSEA) octber 2011
- [5] .M. Bauer, M.J.G. van Eeten “Obfuscation: The Hidden Malware” Security & Privacy IEE 2011.
- [6] BabakBashari Rad, Maslin Masrom,Suhaimi Ibrahim “Camouflage in Malware: from Encryption to Metamorphism” IJCSNS International Journal of Computer Science and Network security August 2012.
- [7] Qinghua Zhang, Douglas S. Reeves “MetaAware: Identifying Metamorphic Malware” National Science Foundation (NSF).
- [8] Vinod P.1, V.Laxmi², M.S.Gaur³, GrijeshChauhan “Metamorphic Malware Exploration Techniques Using MSA signatures” International Conference on Innovations in Information Technology (IIT) 2012.
- [9] Ming Yao “Research on Learning Evidence Improvement for DBSCANBased Classification Algorithm”, IJDTA, 2014.
- [10] M. ZubairRafique, Ping Chen, Christophe Huygens, WouterJoosen Minds-DistriNet, KU Leuven, “Evolutionary Algorithms for Classification of Malware Families through Different Network Behaviors” July 12-16, 2014.
- [11] HiraAgrawal, Lisa Bahler, Josephine Micallef, Shane Snyder, and AlexandrVirodov “Detection of Global, Metamorphic Malware Variants Using Control and Data Flow Analysis”, IEEE 2013.
- [12] SeyedEmadArmoun ,SattarHashemi “A General Paradigm for Normalizing Metamorphic Malwares” IEEE 2012

