



IJITCE

ISSN 2347- 3657

International Journal of Information Technology & Computer Engineering

www.ijitce.com



Email : ijitce.editor@gmail.com or editor@ijitce.com

MALWARE DETECTION USING MACHINE LEARNING

P.Prashanth, Professor, Department Of CS SICET, Hyderabad

Vanama Tharun Kumar, Kalabandalapati Mounika, Kasoju Akshay, Sura Sai Kiran

UG Student, Department Of CS, SICET, Hyderabad

ABSTRACT:

While aiming to reduce the number of negatives, we propose various functions where different machine learning can be used to get the difference between malware files and clean files. In this paper, we first use cascaded singlesided perceptrons to illustrate the idea behind our framework, and then use cascaded nucleated singlesided perceptrons. After successful testing of the average malware and clean archive, the idea behind the framework was sent to the expansion process, which allowed us to solve large malware and collect clean data.

Introduction

Malware is defined as software that enters or damages a computer without the prior permission of the owner. Malware is a general term for all types of computer threats. Basic classifications of malware include program files and standalone malware. Another way to classify malware is by their specific behavior: worms, backdoors, Trojans, rootkits, spyware, adware, etc. This becomes even more difficult as all malware applications now have multiple polymorphic layers to avoid detection or use a utility to update themselves to new versions on short notice to avoid being caught by antivirus software. For an example of dynamic data analysis for malware detection (tested in a virtual environment), interested readers can refer to [2]. The classical method of detecting metamorphic organisms is described in [3]. Here we provide some information about such procedures. work becomes more efficient. Association rules are also used, but there are known symbols on the honey symbols as they are [7]. No) replacement of previous program files. To achieve similar goals, [9] adopted the Profile Hidden Markov model, which has been previously successful in sequence analysis in bioinformatics. This capability is in [10]. In [11], selfmapUsed to describe the behavior of the virus in Windows executable files. The search software aims to obtain as many parameters as possible by simple and easy multilayer combination (cascade) of different models of the perceptron algorithm [12]. Other automatic classification algorithms [13] can also be used in this framework, but we do not explore this alternative. The main steps taken by this framework are summarized below:

TABLE I
NUMBER OF FILES AND UNIQUE COMBINATIONS OF FEATURE VALUES IN THE TRAINING, TEST, AND SCALE-UP DATASETS.

Database	Files		Unique combinations	
	malware	clean	malware	clean
Training	27475	273133	7822	415
Test	11605	6522	506	130
Scale-up	approx. 3M	approx. 180M	12817	16437

TABLE II
MALWARE DISTRIBUTION IN THE TRAINING AND TEST DATASETS.

	Training Dataset		Test Dataset
Malware type	Files	Unique combinations of feature values	Files
Backdoor	35.52%	40.19%	9.16%
Hacktool	1.53%	1.73%	0.00%
Rootkit	0.09%	0.15%	0.04%
Trojan	48.06%	43.15%	37.17%
Worm	12.61%	12.11%	33.36%
Other malware	2.19%	2.66%	20.26%

Because most features are designed to indicate some aspect of the malware profile Map algorithms on large (large) datasets. Not all malware is actually malware, and not all clean samples are clean. This is because the larger the data, the higher the probability that the samples will not be classified in training. Since our algorithm aims to reduce the number of false positives to 0, the detection value (sensitivity) obtained in large data sets will be less (due to negative problems). In Figure 4, we can see how the detection rate decreases as the data grows. Table X shows that accuracy, specificity, and the number of artifacts generally decrease as data size increases.

Conclusion and future work

Our main goal is to propose a machine learning system that can detect as many types of malware as possible, including the hard limit that zero is not a good value. Although our actual accuracy is still not zero, we are close to our goal. In order for this framework to become part of a highly competitive product, many special exemption procedures need to be added. We believe that malware detection through machine learning will complement, not replace, the detection methods used by antivirus vendors. All antivirus operations are subject to some speed and memory limitations; Therefore, the most reliable algorithms

TABLE X
DETECTION RATE (SE) COMPARISON ON THE SCALE-UP (LARGE) DATASET
WHEN TRAINING THE COS-P ALGORITHM.

Datas et	TP	FP	SE	SP	ACC
S10	170	5	51.76 46.94	97.75%	71.94 69.73
S20	309	5	44.24	98.91%	68.32
S30	438	6	42.13	99.22%	67.18
S40	555	6	39.32	99.36%	65.72
S50	648	5	38.66	99.61%	65.39
S60	764	5	36.55	99.68%	64.29
S70	842	2	36.82	99.89%	64.45
S80	969	2	36.89	99.90%	64.48
S90	1092	3	33.45	99.87%	62.56
S100	1100	3		99.88%	

Presented here are the Cascaded One-Sided Perceptron (COS-P) and its clearly defined model (COS-P-Map). It can be seen that the total detection rate produced by our algorithm increases by 3%~4%, which is very significant. (Please note that training was conducted on malware samples that were not detected by the detection process.)

To date, our framework has proven to be useful research for computer security experts in Bit Defender's anti-malware research department. In the future, we plan to integrate more classification algorithms such as wideedge perceptrons [18] and support vector machines [14], [19], [20].

REFERENCES

1. The authors would like to thank BitDefender management for their support on these issues. Penya, J. Devesa and P. G. Garcia, "N-gram based data signatures for malware detection," 2009. Holz, C. Willems, P. Dussel, and P. Laskov,
2. "Learning and Classification of Malware Behavior," DIMVA '08: Proceedings of the 5th International Conference on Intrusion and Malware Detection, and Vulnerability Assessment. Berlin, Heidelberg: Springer-Verlag, 2008, p. 108–125.

- [3] E. Konstantinou, "Metamorphosis viruses: analysis and detection", 2008, Technical Report RHUL-MA-2008-2, Search Security Award Master's Thesis, 93 Nploo . 6, s. 2669–2672, 2006. Maloof, "Learning to detect and identify malicious behavior in wildlife," Journal of Machine Learning Research, vol. 7, p. 2721–2744, December 2006, Special Issue on Machine Learning in Computer Security. Li and D. Ye, "Imds: An intuitive malware detection system," in KDD, P. Berkhin, R. Caruana, and X. Wu, editors. ACM, 2007, nr 17. 1043-1047 Ib. Upadhyaya, "Spycon: Simulating user activity to detect and bypass spyware," in, IPC CC. IEEE Computer Society, 2007, p. 502-509. Walenstein thiab A. Lakhotia, "Machine-passed variants of malware filtering using Markov chains," Malware and Unwanted Software, 2008. Conference, 2008, p. 77–84. Stamp, S. Attaluri and S. McGhee, "Examining the latent signature model and metamorphic virus detection," Journal of Computational Virology, 2008. Comprehensive detection and classification of polymorphic malware, 2006. and Research Data Proceedings. New York, NY, AB: ACM. 82–89. McGraw-Hill Education (ISE Edition), October 1997. Cambridge University Press, March 2000. J. Smola, Learning with Kernels: Support vector machines, optimization, optimization, and more. MIT Press, 2002. Chauvin, C. A. Andersen, and H. Nielsen, "Accuracy evaluation of distribution prediction algorithms," Bioinformatics, vol. 5, page 5. 412-424, Mart 2000. K. Mishra, "De novo SVM-like i-fixation of precursor microRNAs from genomic pseudohairpins using global and internal folding methods." . Freund and R. E. Schapire, "Classification of large flowers using perceptron algorithm," Machine Learning, vol. 37, 1999, s. 277-296.
- [20] V. N. Vapnik, Information Science and the Nature of Statistics. Springer, November 1999.