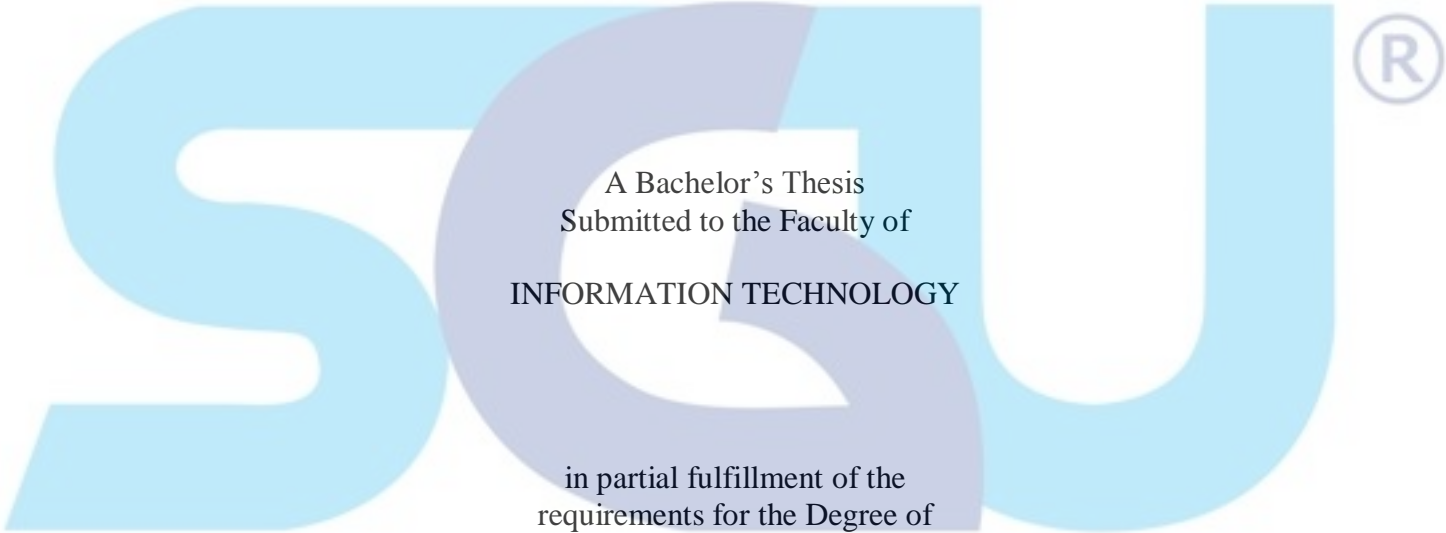


**ANALYSIS OF MACHINE LEARNING TECHNIQUES USED IN
BEHAVIOR-BASED MALWARE DETECTION**

By

Ivan Firdausi



A Bachelor's Thesis
Submitted to the Faculty of
INFORMATION TECHNOLOGY

in partial fulfillment of the
requirements for the Degree of

BACHELOR OF SCIENCES
WITH A MAJOR IN INFORMATION TECHNOLOGY

SWISS GERMAN UNIVERSITY

SWISS GERMAN UNIVERSITY
EduTown BSD City
Tangerang – 15339
Island of Java, Indonesia
www.sgu.ac.id

July 2010

STATEMENT BY THE AUTHOR

I hereby declare that this submission is my own work and to the best of my knowledge, it contains no material previously published or written by another person, not material which to a substantial extent has been accepted for the award of any other degree or diploma at any educational institution, except where due acknowledgement is made in the thesis.

Ivan Firdausi

Date

Approved by:

Charles Lim, Msc., ECSA, ECSP, ECIH, CEH, CEI

Date

Anto Satriyo Nugroho, Dr.Eng

Date

Chairman of the Examination Steering Committee

Date

Ivan Firdausi

ABSTRACT

ANALYSIS OF MACHINE LEARNING TECHNIQUES USED IN BEHAVIOR-BASED MALWARE DETECTION

By

Ivan Firdausi

SWISS GERMAN UNIVERSITY

Bumi Serpong Damai

Charles Lim, Msc., ECSA, ECSP, ECIH, CEH, CEI, Major Lecturer

The increasing of malware that are exploiting the Internet daily has become a serious threat. The manual heuristic inspection of malware analysis is no longer considered effective and efficient compared against the high spreading rate of malware. Conventional signature matching-based antivirus systems fail to detect polymorphic, obfuscated, and new, previously unseen malicious executables. Hence, automated behavior-based malware detection using machine learning techniques is considered a profound solution. The behavior of each malware on an emulated (sandbox) environment will be automatically analyzed and will generate behavior reports. These reports will be preprocessed into sparse vector models for further machine learning (classification). The classifiers used in this research are k -Nearest Neighbors (k NN), Naïve Bayes, Decision Tree, Support Vector Machine (SVM), and Artificial Neural Network (ANN). According to the analysis of the test and experiment results of all the 5 classifiers, the overall best performance goes to J48 with a recall (true positive rate) of 95.9%, a false positive rate of 2.4%, a precision (positive predictive value) of 97.3%, and an accuracy of 96.8%. In summary, it can be concluded that a proof-of-concept based on automatic behavior-based malware analysis and the use of machine learning techniques could detect malware quite effectively and efficiently.

Keywords—malware analysis, dynamic analysis, behavior analysis, data mining, machine learning, classification, malware detection

DEDICATION

I dedicate this thesis to my family, my relatives, and my friends.



ACKNOWLEDGMENTS

The author wishes to give the utmost gratitude to Allah SWT. His blessing and guidance had helped me to give me strength, great health, great ideas, and finally led me to complete this thesis.

The author would also like to give special thanks to Charles Lim, Msc., Anto Satriyo Nugroho, Dr.Eng, and Alva Erwin, Msc., MTI for acting as my thesis advisors. They have had given me support, suggestions, and guidance in completing this thesis.

In addition, the author would also like to give thanks to Attur S. Widjaja (Aat) for giving us (SGU Malware Analysis Team) Indonesian malware samples for our data sets and to Anubis for giving us the permission to use their XML behavior report files of the malware and benign software that we had submitted for our research.

Last but not least, the author would like to give thanks to James Purnama, Msc. and Dipl.-Inf. Kho I Eng for supporting and helping us (SGU Malware Analysis Team) in the establishment of our new malware laboratory in SGU and to IDSIRTII Malware Analysis Team for sharing ideas, discussions, and their cooperation with us.

TABLE OF CONTENTS

STATEMENT BY THE AUTHOR	2
ABSTRACT	3
DEDICATION	4
ACKNOWLEDGMENTS	5
CHAPTER 1 – INTRODUCTION.....	11
1.1 Background	11
1.2 Research Problem.....	11
1.3 Research Objective.....	12
1.4 Scope of Analysis	12
1.5 Research Contributions.....	12
1.6 Methodology	13
1.7 Thesis Organization.....	14
CHAPTER 2 – LITERATURE REVIEW	15
2.1 Malware	15
2.1.1 Virus.....	15
2.1.2 Worm	15
2.1.3 Trojan Horse.....	16
2.1.4 Malicious Rootkit	16
2.1.5 Botnet.....	16
2.1.6 Backdoor	16
2.1.7 Spyware and Adware	17
2.2 Malware Analysis.....	17
2.3 Data Mining and Machine Learning.....	18
2.4 Classification and Supervised Learning	18
2.4.1 Nearest Neighbor	20
2.4.2 Naïve Bayes.....	21
2.4.3 Support Vector Machine (SVM)	23
2.4.3.1 Linear Separable SVM.....	24
2.4.3.2 Linear Non-separable SVM	26
2.4.3.3 Nonlinear SVM	28
2.4.4 Decision Tree.....	31
2.4.5 Artificial Neural Network (ANN).....	35
2.4.5.1 Perceptron	35
2.4.5.2 Multilayer Perceptron	36
2.5 Static Malware Analysis	38
2.5.1 Manual Static Malware Analysis.....	39
2.5.2 Automatic Static Malware Analysis	40
2.6 Dynamic Malware Analysis.....	41
2.6.1 Manual Dynamic Malware Analysis	42

2.6.2	Automatic Dynamic Malware Analysis	43
2.7	Hybrid (Static and Dynamic) Malware Analysis	45
2.8	Previous Related Works	45
CHAPTER 3 – METHODOLOGY		49
3.1	Data Acquisition and Storage	49
3.2	Automatic Behavior Monitoring and Report Generation	54
3.3	Data Preprocessing	55
3.3.1	Create XML File Parser	56
3.3.2	Create Term Dictionary	59
3.3.3	Create Binary-weight and Term Frequency-weight Counters and Vector Models	59
3.3.4	Create Attribute-Relation File Format (ARFF) Files	61
3.3.5	Create Batch Files	62
3.4	Learning and Classification	62
3.5	Results Analysis and Documentation	64
CHAPTER 4 – RESULT AND DISCUSSION		65
4.1	Result	65
4.1.1	Classifier Parameter Tuning	65
4.1.1.1	IBk	66
4.1.1.2	SVM	70
4.1.1.3	J48	71
4.1.1.4	Multilayer Perceptron (MLP)	71
4.1.2	Performance Metrics	72
4.1.3	Without Feature Selection	74
4.1.3.1	Binary-weight Vector Model	75
4.1.3.2	Term Frequency-weight Vector Model	76
4.1.4	With Feature Selection	78
4.1.4.1	Binary-weight Vector Model	80
4.1.4.2	Term Frequency-weight Vector Model	82
4.2	Discussion	83
4.2.1	IBk	83
4.2.2	SVM	85
4.2.3	J48	85
4.2.4	Multilayer Perceptron	86
4.2.5	Performance Metrics	86
4.2.6	Without Feature Selection	87
4.2.7	With Feature Selection	88
4.2.8	Overall Summary	90
CHAPTER 5 – CONCLUSION AND FUTURE WORKS		92
5.1	Conclusion	92
5.2	Future Works	93
GLOSSARY		95
REFERENCES		99

PUBLICATION	102
APPENDICES.....	103
APPENDIX A – SGU MALWARE ANALYSIS TEAM.....	103
APPENDIX B – SAMLPLE ANUBIS XML BEHAVIOR REPORT FILE.....	104
APPENDIX C – SAMLPLE ARFF FILE.....	108
APPENDIX D – BINARY-WEIGHT FEATURE SELECTION LIST	110
APPENDIX E – TERM FREQUENCY-WEIGHT FEATURE SELECTION LIST	114
APPENDIX F – PUBLISHED WORK	115
CURRICULUM VITAE.....	119

