

**MALWARE DETECTION BY USING TWO-LAYER STACKING SVM  
CLASSIFIER**



**UNIVERSITI TEKNIKAL MALAYSIA MELAKA**

## BORANG PENGESAHAN STATUS LAPORAN

JUDUL: MOBILE MALWARE DETECTION USING TWO-LAYER STACKING SVM CLASSIFIER

SESI PENGAJIAN: [2020 / 2021]

Saya RAJA NURUL AINI BINTI RAJA MOHD ANUAR mengaku membenarkan tesis Projek Sarjana Muda ini disimpan di Perpustakaan Universiti Teknikal Malaysia Melaka dengan syarat-syarat kegunaan seperti berikut:

1. Tesis dan projek adalah hakmilik Universiti Teknikal Malaysia Melaka.
2. Perpustakaan Fakulti Teknologi Maklumat dan Komunikasi dibenarkan membuat salinan untuk tujuan pengajian sahaja.
3. Perpustakaan Fakulti Teknologi Maklumat dan Komunikasi dibenarkan membuat salinan tesis ini sebagai bahan pertukaran antara institusi pengajian tinggi.
4. \* Sila tandakan (✓)

\_\_\_\_\_ SULIT (Mengandungi maklumat yang berdarjah keselamatan atau kepentingan Malaysia seperti yang termaktub di dalam AKTA RAHSIA RASMI 1972)

\_\_\_\_\_ TERHAD (Mengandungi maklumat TERHAD yang telah ditentukan oleh organisasi / badan di mana penyelidikan dijalankan)

\_\_\_\_\_ TIDAK TERHAD

*Rajalwani*

(TANDATANGAN PELAJAR)

*Sarim*

(TANDATANGAN PENYELIA)

Alamat tetap: Batu 17, Dusun Tua  
Seberang 43200, Hulu Langat, Selangor

Ts Nor Azman Bin Mat Ariff

Tarikh: 12 SEPTEMBER 2021

Tarikh: 12 SEPTEMBER 2021

CATATAN: \* Jika tesis ini SULIT atau TERHAD, sila lampirkan surat daripada pihak berkuasa.



## DECLARATION

I hereby declare that this project report entitled  
**MOBILE MALWARE DETECTION BY USING TWO-LAYER STACKING SVM  
CLASSIFIER**

is written by me and is my own effort and that no part has been plagiarized

without citations.

STUDENT:  Date : 12 SEPTEMBER 2021

RAJA NURUL AINI BINTI RAJA MOHD ANUAR

UNIVERSITI TEKNIKAL MALAYSIA MELAKA

I hereby declare that I have read this project report and found  
this project report is sufficient in term of the scope and quality for the award of  
Bachelor of [Computer Science (Software Development)] with Honours.

SUPERVISOR :  Date: 12 SEPTEMBER 2021

TS. NOR AZMAN MAT ARIFF

## DEDICATION

Special dedicated to my beloved parents and friend who have encouraged, guided and inspired me throughout this journey of education

To my helpful lecturer, thank you for the guidance from the beginning until the end of this final year project



## ACKNOWLEDGEMENTS

All praise to Allah SWT for giving me the strength to finish my PSM, I manage to complete this Projek Sarjana Muda 1 (PSM1) and Projek Sarjana Muda 2 (PSM2). I would like to thank Ts. Nor Azman Bin Mat Ariff for giving assistant to complete this project successfully. Without his guide, this report and project cannot be completed.

Special thanks to my family for the continuous support throughout this journey from day to day. I would like to thanks to all my friend for their time, concern efforts, and always encourage me when preparing this report. With all the help from involves parties, I manage to complete my report and project which is Mobile Malware Detection Using Machine Learning.

Last but not least, thanks to Universiti Teknikal Malaysia Melaka (UTeM) for the opportunity given.



## ABSTRACT

Mobile device usage increased with new high technology that attracted the attacker to launch the attack, such as mobile malware. Mobile malware is malicious that is loaded into the device system and causes damage. Malware has become more prevalent in recent years. Nowadays, there are many techniques available to detect this attack. In this research, N-gram with opcode sequence dataset was used, focusing on mobile malware detection model. The dataset undergoes feature selection which is Information Gain and Chi-square, to reduce irrelevant data and redundant data. Then, the classifier used in this project is Support Vector Machine (SVM) to develop a model. The developed model will test and verify the accuracy with the test set. The project is giving hope to produce a system that can detect mobile malware attacks.

## ABSTRAK

Penggunaan peranti mudah alih meningkat dengan teknologi tinggi baru yang menarik penyerang untuk melancarkan serangan, seperti perisian hasad mudah alih. Perisian hasad mudah alih berniat jahat yang dimuat ke dalam sistem peranti dan menyebabkan kerosakan. Perisian hasad semakin berleluasa dalam beberapa tahun kebelakangan ini. Pada masa kini, terdapat banyak teknik yang ada untuk mengesan serangan ini. Dalam penyelidikan ini, N-gram dengan dataset urutan opcode digunakan, dengan fokus pada model pengesanan malware mudah alih. Set data menjalani pemilihan ciri iaitu *Information Gain* dan *Chi-Square*, untuk mengurangkan data yang tidak relevan dan data yang berlebihan. Kemudian, pengklasifikasi yang digunakan dalam projek ini adalah Mesin Sokongan Vektor (SVM) untuk membangunkan model. Model yang dibangunkan akan diuji dan mengesahkan ketepatannya dengan set ujian. Projek ini memberi harapan untuk menghasilkan sistem yang dapat mengesan serangan perisian hasad mudah alih.



## TABLE OF CONTENTS

	<b>PAGE</b>
<b>DECLARATION.....</b>	<b>II</b>
<b>DEDICATION.....</b>	<b>III</b>
<b>ACKNOWLEDGEMENTS.....</b>	<b>IV</b>
<b>ABSTRACT.....</b>	<b>V</b>
<b>ABSTRAK.....</b>	<b>VI</b>
<b>TABLE OF CONTENTS.....</b>	<b>VII</b>
<b>LIST OF TABLES.....</b>	<b>XII</b>
<b>LIST OF FIGURES.....</b>	<b>XIII</b>
<b>LIST OF ABBREVIATIONS.....</b>	<b>XV</b>
<b>LIST OF ATTACHMENTS.....</b>	<b>XVI</b>
<b>CHAPTER 1: INTRODUCTION.....</b>	<b>1</b>
1.1 Introduction.....	1
1.2 Problem Statement.....	1
1.3 Project Question (PQ).....	3
1.4 Project Objective (PO).....	4
1.5 Project Scope.....	4
1.6 Project Contribution.....	5
1.7 Thesis Organization.....	5

1.7.1	Chapter 1: Introduction.....	5
1.7.2	Chapter 2: Literature Review.....	5
1.7.3	Chapter 3: Project Methodology.....	6
1.7.4	Chapter 4: Analysis and Design .....	6
1.7.5	Chapter 5: Implementation .....	6
1.7.6	Chapter 6: Discussion .....	6
1.7.7	Chapter 7: Project Conclusion.....	6
1.8	Conclusion .....	7
<b>CHAPTER 2: LITERATURE REVIEW.....</b>		<b>8</b>
2.1	Introduction.....	8
2.2	Mobile Malware.....	9
2.2.1	Mobile Malware Definition .....	9
2.2.2	Mobile Malware Type .....	10
2.2.3	Mobile Malware Analysis Technique.....	12
2.2.4	Mobile Malware Detection Technique .....	12
2.3	Machine Learning .....	14
2.3.1	Machine Learning Definition .....	14
2.3.2	Feature Selection .....	15
2.3.2.1	Type of Feature Selection.....	15
2.4	Classification.....	18
2.4.1	Classification Definition.....	18

2.4.2	Type of Classification.....	18
2.5	Two-Layer Stacking SVM Classifier.....	21
2.6	Critical View.....	22
2.6.1	Previous Research on Mobile Malware.....	22
2.6.2	Previous Research using Machine Learning.....	24
2.6.3	Previous Research using Feature Selection.....	26
2.7	Conclusion.....	28
<b>CHAPTER 3: PROJECT METHODOLOGY.....</b>		<b>30</b>
3.1	Introduction.....	30
3.2	Methodology.....	30
3.2.1	Previous Research.....	31
3.2.2	Information Gathering.....	32
3.2.3	Define Scope.....	32
3.2.4	Design and Implementation.....	32
3.2.5	Testing and Evaluation of Model.....	32
3.2.6	Documentation.....	33
3.3	Project Schedule and Milestones.....	33
3.3.1	Project Flowchart.....	33
3.3.2	Project Milestones.....	34
3.3.3	Project Gantt Chart.....	37
3.4	Conclusion.....	37

<b>CHAPTER 4: ANALYSIS AND DESIGN</b> .....	<b>38</b>
4.1 Introduction.....	38
4.2 Problem Analysis .....	38
4.3 Project Design.....	39
4.3.1 Dataset .....	41
4.3.2 Feature Extraction.....	42
4.3.3 Feature Selection .....	42
4.3.4 Normalization .....	44
4.3.5 Classification .....	44
4.3.6 Two-Layer Stacking SVM Classifier using Output Probabilities .	49
4.4 Requirement Analysis.....	50
4.4.1 Software Requirement .....	50
4.4.2 Hardware Requirement.....	51
4.5 Conclusion .....	53
<b>CHAPTER 5: EXPERIMENT</b> .....	<b>54</b>
5.1 Introduction.....	54
5.2 Software Development Environment Setup.....	54
5.3 Process Module.....	54
5.3.1 Data Collection .....	55
5.3.2 Train and Test Data .....	55
5.3.3 Scale and Model for Train Data.....	57

5.3.4	Generate The Predict File .....	57
5.4	Project Implementation .....	58
5.4.1	Script Execution.....	58
5.5	Result .....	65
5.5.1	Attribute.....	65
5.5.2	10 runs .....	66
5.5.3	Accuracy Table.....	66
<b>CHAPTER 6: DISCUSSION .....</b>		<b>70</b>
6.1	Introduction.....	70
6.2	Discussion of The Project.....	70
6.3	Discussion on The Newly Proposed Method.....	71
6.4	Conclusion .....	72
<b>CHAPTER 7: PROJECT CONCLUSION .....</b>		<b>73</b>
7.1	Introduction.....	73
7.2	Project Summary.....	73
7.3	Project Contribution.....	74
7.4	Project Limitation .....	74
7.5	Future Work.....	74
7.6	Conclusion .....	75
<b>REFERENCES.....</b>		<b>76</b>
<b>APPENDIX .....</b>		<b>79</b>

## LIST OF TABLES

	PAGE
<b>Table 1.1: Problem Statement.....</b>	<b>3</b>
<b>Table 1.2: Summary of Project Question.....</b>	<b>3</b>
<b>Table 1.3: Summary of Project Objective.....</b>	<b>4</b>
<b>Table 2.1: Mobile Malware Detection Techniques.....</b>	<b>12</b>
<b>Table 2.2: Type of Learning.....</b>	<b>14</b>
<b>Table 2.3 Feature Selection Summary .....</b>	<b>16</b>
<b>Table 2.4: Survey on Mobile Malware Detection.....</b>	<b>23</b>
<b>Table 2.5: Previous Research using Machine Learning's Literature .....</b>	<b>25</b>
<b>Table 2.6: Previous Research using Feature Selection's Literature.....</b>	<b>27</b>
<b>Table 3.1: Project Milestones .....</b>	<b>35</b>
<b>Table 4.1: Description of Dataset.....</b>	<b>41</b>
<b>Table 4.2: Type of SVM's Kernal Summary .....</b>	<b>49</b>
<b>Table 4.3: Software Requirement.....</b>	<b>50</b>
<b>Table 4.4: Hardware Requirement.....</b>	<b>51</b>
<b>Table 4.5: Workstation Specification .....</b>	<b>52</b>
<b>Table 5.1: Parameter Description for Batch File .....</b>	<b>58</b>
<b>Table 5.2: Description of each command.....</b>	<b>60</b>
<b>Table 5.3: The Command and Its Description .....</b>	<b>63</b>
<b>Table 5.4: The Description of Comment for Train Data .....</b>	<b>64</b>
<b>Table 5.5: Testing's Command Description .....</b>	<b>64</b>
<b>Table 5.6: The Total Number of Attributes.....</b>	<b>65</b>
<b>Table 6.1: Benchmark.....</b>	<b>71</b>
<b>Table 6.2: Final Result.....</b>	<b>71</b>

## LIST OF FIGURES

	PAGE
<b>Figure 1.1: Total Malware by Year .....</b>	<b>2</b>
<b>Figure 2.1: Literature Review's Structure.....</b>	<b>9</b>
<b>Figure 2.2: Mobile Malware Detection Classification .....</b>	<b>13</b>
<b>Figure 2.3: Classification Techniques in Supervised Learning .....</b>	<b>18</b>
<b>Figure 2.4: Decision Tree.....</b>	<b>19</b>
<b>Figure 2.5: Two-Layer Stacking Spatial Pyramid Classifier .....</b>	<b>22</b>
<b>Figure 3.1: Flowchart for the System Framework.....</b>	<b>31</b>
<b>Figure 3.2: Project Flowchart .....</b>	<b>34</b>
<b>Figure 4.1: Flowchart Design .....</b>	<b>40</b>
<b>Figure 4.2: Part of Dataset 1g .....</b>	<b>42</b>
<b>Figure 4.3: Scatter Plot.....</b>	<b>45</b>
<b>Figure 4.4: Placement of Hyperplane.....</b>	<b>45</b>
<b>Figure 4.5: Hyperplane with Margin .....</b>	<b>46</b>
<b>Figure 4.6: Best Hyperplane.....</b>	<b>46</b>
<b>Figure 4.7: Hyperplane with Outlier .....</b>	<b>47</b>
<b>Figure 4.8: Non-Linear SVM .....</b>	<b>48</b>
<b>Figure 4.9: SVM with Multi-Dimensional Space .....</b>	<b>48</b>
<b>Figure 5.1: Process Module .....</b>	<b>55</b>
<b>Figure 5.2: Sample Train and Test Data.....</b>	<b>56</b>
<b>Figure 5.3: Example Random File Number.....</b>	<b>56</b>
<b>Figure 5.4: Process creating scale and model for train data .....</b>	<b>57</b>
<b>Figure 5.5: Process generating predict file .....</b>	<b>58</b>
<b>Figure 5.6: Convert File Script .....</b>	<b>58</b>

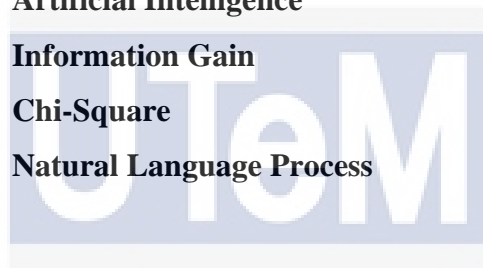
<b>Figure 5.7: Run Batch File Using CMD</b> .....	<b>59</b>
<b>Figure 5.8: Sample run batch program for mobile malware dataset</b> .....	<b>59</b>
<b>Figure 5.9: Execute the batch file program</b> .....	<b>61</b>
<b>Figure 5.10: Sample Predict File</b> .....	<b>62</b>
<b>Figure 5.11: Probability Output</b> .....	<b>62</b>
<b>Figure 5.12: Sample execution for scaling</b> .....	<b>62</b>
<b>Figure 5.13: The Command for Train Data</b> .....	<b>63</b>
<b>Figure 5.14: Testing Command</b> .....	<b>64</b>
<b>Figure 5.15: 10 Run Time</b> .....	<b>66</b>
<b>Figure 5.16: Result Second Experiment</b> .....	<b>68</b>
<b>Figure 5.17: Bar Graph First Experiment</b> .....	<b>68</b>
<b>Figure 5.18: Second Experiment Bar Graph</b> .....	<b>69</b>
<b>Figure 6.1: Final Accuracy</b> .....	<b>72</b>





**LIST OF ABBREVIATIONS**

<b>FYP</b>	-	<b>Final Year Project</b>
<b>IT</b>	-	<b>Information Technology</b>
<b>OS</b>	-	<b>Operating System</b>
<b>ML</b>	-	<b>Machine Learning</b>
<b>AI</b>	-	<b>Artificial Intelligence</b>
<b>IG</b>	-	<b>Information Gain</b>
<b>CS</b>	-	<b>Chi-Square</b>
<b>NLP</b>	-	<b>Natural Language Process</b>



اونيورسيتي تيكنيكل مليسيا ملاك

UNIVERSITI TEKNIKAL MALAYSIA MELAKA

## LIST OF ATTACHMENTS

	<b>PAGE</b>
<b>Appendix A</b>	
<b>Gantt Chart</b>	<b>52</b>



## CHAPTER 1: INTRODUCTION

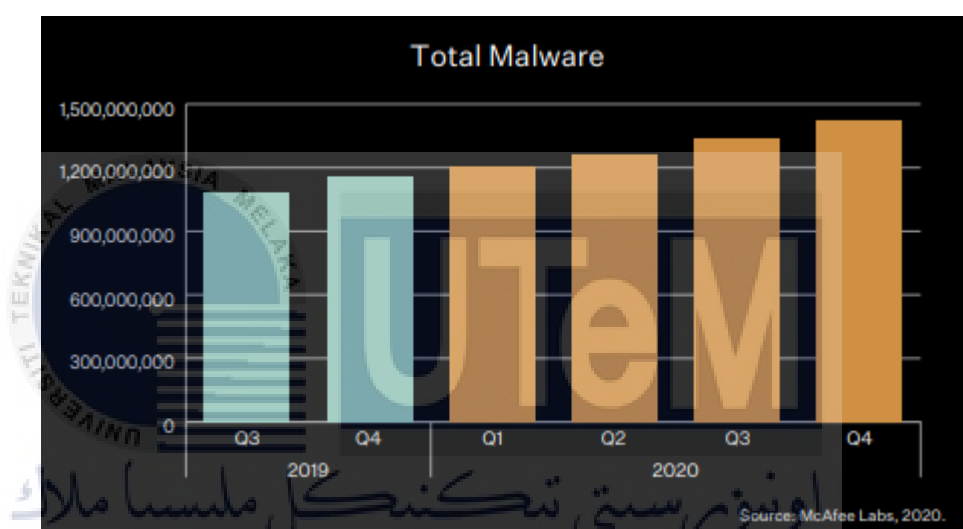
### 1.1 Introduction

Nowadays, mobile device usage is increasing and it is expanded with new advanced technology every year. The continuous development of mobile hardware and technology has created a highly connected world. Mobile devices are used in many fields such as business, healthcare, education, and, Information Technology (IT). This transformation of mobile devices is attractive to the attacker to attempt to steal the information from the devices or system. Malware is malicious software that loaded onto mobile devices and causes damage then destroys it. The more malware attacks are introducing every year and it needs to be overcome before it widespread. This project focuses on mobile malware attack detection using a machine learning algorithm technique. This technique was chosen because it has good performance to track the malware and able to get high accuracy in classification.

### 1.2 Problem Statement

The mobile device becomes a place where intruders try their malicious code to fulfill their intention. There have many researchers proposed to prevent the mobile malware attack by introducing techniques that can protect the devices. Malware developers think it is easy to transfer attackers on mobile devices because mobile applications downloaded from Apple App Stores, Google Play, and websites usually have no time to analyze them (Gyamfi & Owusu, 2019). The current malware can do anything, their threat level becomes bigger and keeps growing. The variation of this malware and the exposed vulnerabilities require new and improved methods for detection.

There are plenty of attack situations in which an attacker can compromise a user's information by exploiting operating system's vulnerabilities. The growth of malware has led to different mobile antiviruses, firewalls, and encoding products, which have been released by various security vendors such as F-Secure, Kaspersky Lab, McAfee, and Symantec. The number of F-Secure mobile malware in 2008 was 401, while McAfee numbered 457. Mobile malware is malicious software that has certain mobile device targets. It first appeared for Symbian Operating System (OS) in 2004 and has now grown exponentially with smartphones in popularity (Gyamfi & Owusu, 2019). Figure 1.1 shows the mobile malware growth from 2010 up to 2018.



**Figure 1.1: Total Malware by Year  
(McAfee 2020)**

COVID-19 pandemic Movement Control Order (MCO) is another alarming issue with mobile malware attacks. According to McAfee (2020), McAfee Lab's average volume of malware threats was 588 per minute, up from 169 per-minute threats in the third quarter of 2020 (40%). The fourth quarter was 648 per minute threats, an increase of 60 per minute threats (10%). Mobile malware grew 118% by SMS Reg driven from Q3 to Q4 of 2020. The best solution to the problem is the detection of mobile malware. The problem was defined as in table 1.1.

**Table 1.1: Problem Statement**

PS	Problem Statement
PS1	An increasing number of users connect to the Internet has been driven by the tremendous growth in mobile and smart device usage, it is hard to discover effective and fast technique in differentiate between benign and malware. Another issue is to identify the suitable machine learning technique to use in this project to process the dataset.

### 1.3 Project Question (PQ)

In this project, there are three Project Question (PQ) needed answers made out from the problem statement. Table 1.2 shows the summary of the project question based on the problem statement.

**Table 1.2: Summary of Project Question**

PS	PQ	Project Question
PS1	PQ1	What is malware in mobile devices?
	PQ2	How mobile malware classified?
	PQ3	Is the machine learning technique give a result model accurately?

#### 1.4 Project Objective (PO)

Project Objective (PO) is a goal for the expected outcome of the research. To fulfill the goals of the project, the problem statement and project question must be assigned. Three objectives in this research to achieve and the relationship between PS, PQ, and PO for this project as shown in table 1.3.

**Table 1.3: Summary of Project Objective**

PS	PQ	PO	Project Objective
	PQ1	PO1	To study the behavior of malware.
PS1	PQ2, PQ3	PO2	To implement machine learning classifier to detect malware attack.
	PQ1, PQ2, PQ3	PO3	To test and verify the accuracy of machine learning to detect the malware.

#### 1.5 Project Scope

The scope of this project is related to the following criteria:

- I. The dataset only focuses on mobile malware.
- II. This project focus on machine learning method.
- III. The effectiveness of the proposed model is measured using accuracy.

## 1.6 Project Contribution

The result of this project is depending on the project implementation to detect the mobile malware attacks and select effective features that will give high accuracy. The following is a brief description of this project's contribution:

- I. Identification of different mobile malware attacks based on their taxonomy.
- II. Propose a model that can detect malware attacks on the mobile device.
- III. Propose to check the accuracy of mobile malware attack detection and test and validate the developed program.

## 1.7 Thesis Organization

This report divided to seven chapter namely Chapter 1: Introduction, Chapter 2: Literature Review, Chapter 3: Project Methodology, Chapter 4: Analysis and Design, Chapter 5: Implementation, Chapter 6: Testing and Validation and Chapter 7: Project Conclusion.

### 1.7.1 Chapter 1: Introduction

This chapter act as the guideline to identifying the important things prior to the project. It will explain what and how this project attempted to accomplish, as well as its relevance.

### 1.7.2 Chapter 2: Literature Review

This chapter consists the related studies and early task from previous researcher about this project and critical analysis of this chapter.

### **1.7.3 Chapter 3: Project Methodology**

Chapter 3 focuses on planning works that must answer the objectives of this project that has been mention previously. This chapter includes the procedures throughout this project.

### **1.7.4 Chapter 4: Analysis and Design**

This section is related to the procedure that correlate with the experiment of this project. The experiment needs to be analyzed the effectiveness by referring to the literature review and design the project clearly to reach the goals.

### **1.7.5 Chapter 5: Implementation**

Chapter 5 is the main part that the procedure from the previous chapter needs to be carried out. The experiment will be processed in this section until has the result and it will be gathered the result to be compared with another method.

### **1.7.6 Chapter 6: Discussion**

This chapter explains the result and analyzed them where the findings are answers the objectives.

### **1.7.7 Chapter 7: Project Conclusion**

This chapter will summarize the project, state the contribution, and highlight the constraints that were encountered throughout the project. This chapter will also describe how to improve the project in the future.