

MALWARE DETECTION USING ENSEMBLE METHOD

AMIZAH AIDA BINTI AHMAD



UNIVERSITI TEKNIKAL MALAYSIA MELAKA

MALWARE DETECTION USING ENSEMBLE METHOD

BORANG PENGESAHAN STATUS TESIS

JUDUL: MALWARE DETECTION USING ENSEMBLE METHOD

SESI PENGAJIAN: SEMESTER II 2016/2017

Saya **AMIZAH AIDA BINTI AHMAD (B031510056)** mengaku membenarkan tesis (PSM/Sarjana/Doktor/Falsafah) ini disimpan di Perpustakaan Fakulti Teknologi Maklumat dan Komunikasi dengan syarat-syarat kegunaan seperti berikut:

1. Tesis dan projek adalah hak milik 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

(TANDATANGAN PENULIS)

(TANDATANGAN PENYELIA)

Alamat tetap:

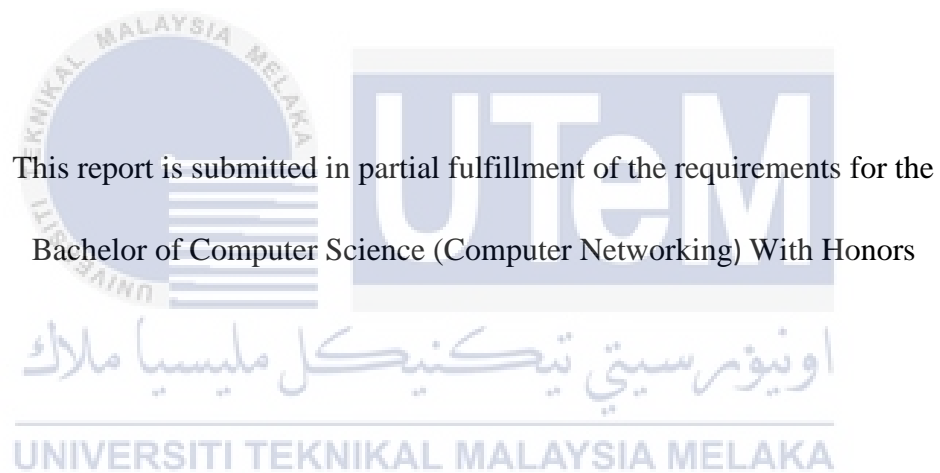
No 34C, Blok F Jalan 17/1A'
Seksyen 17, 46400 Petaling Jaya

Tarikh:

Tarikh: 18/8/2017

MALWARE DETECTION USING ENSEMBLE METHOD

AMIZAH AIDA BINTI AHMAD



FACULTY OF INFORMATION AND COMMUNICATION TECHNOLOGY

UNIVERSITI TEKNIKAL MALAYSIA MELAKA

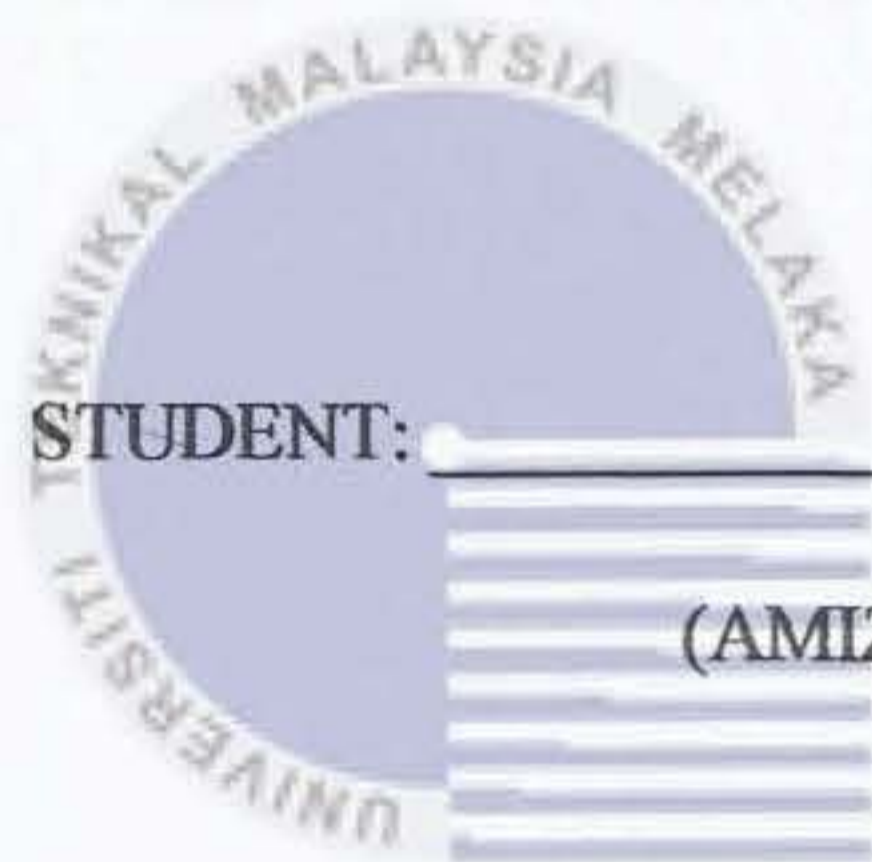
2017

DECLARATION

I hereby declare that this project report entitled

MALWARE DETECTION USING ENSEMBLE METHOD

is written by me and is my own effort and that no part has been plagiarized without citations.



STUDENT: _____

Date: _____

(AMIZAH AIDA BINTI AHMAD)

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
UNIVERSITI TEKNIKAL MALAYSIA MELAKA

Bachelor of Computer Science (Interactive Media) with Honors

SUPERVISOR: _____

Handwritten signature of En. Azman Bin Mat Ariff in black ink.

Date: _____

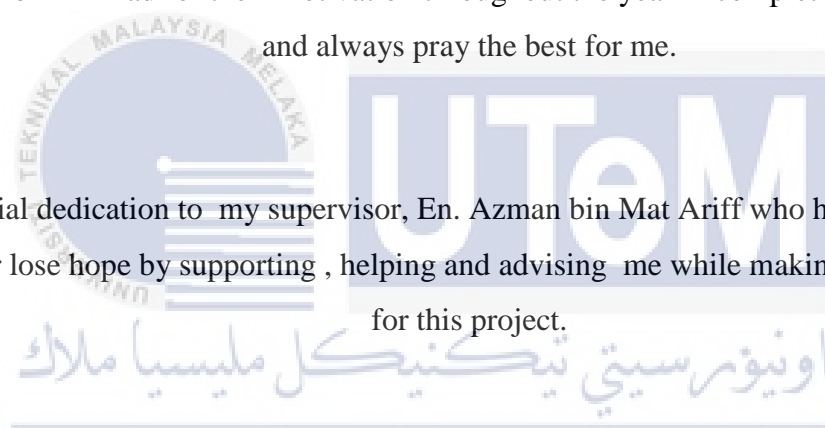
18/8/2017

(EN. AZMAN BIN MAT ARIFF)

DEDICATION

This final project is dedicated to all generous people who have been aided us and supported, encouraged and help me during this project.

Especially to my beloved parents Ahmad bin Jasa and Maimon binti Siron for their endless support and love , both of my siblings Amirah Arina binti Ahmad and Armin Azrai bin Ahmad for their motivation throughout the year in completing my studies and always pray the best for me.



Special dedication to my supervisor, En. Azman bin Mat Ariff who has guided and never lose hope by supporting , helping and advising me while making the progress for this project.

Not forget to my evaluator, P.M. DR. Faizal bin Abdollah for her positive advices, feedback and comment on this project,

Last but not least, to all my beloved friends who always help me from the beginning of this project until the end of it.

ACKNOWLEDGEMENT

Bismillahirrahmanirrahim,

First and foremost, I would like to express my gratitude to Allah S.W.T for giving me strength and for keeping me in the path of righteous while completing the whole process of this project. Without His blessings, I would not have been able to be as I am today and completing this project.

Also a million thanks to my supervisor, En. Azman bin Mat Ariff for his guidance, constant supervision and kindness in completing this project. Without his guide, I am surely would not know the right trail.

I would also like to special thanks to my family members , my dear father and mother for who always support and pray the best for me and constantly give me motivation throughout my project. Besides, thanks to my friends that really play a big role who have been very helpful by giving ideas and suggestion especially Nor Anis Syazana and Sarah Shamimi the during project development.

Thank you.

ABSTRACT

In today's technology driven world, the increasing of malware in the cybercriminals that exploiting the internet and always create and distribute harmful malware has become a serious threat. Malware significantly impact computer's performance and often go unnoticed in our systems and causes several problems to the user. Hence, It's imperative to take the precautions necessary to detect and prevent malware infections. One of the way to detect malware detection is by using machine learning techniques. Malware detection is detected by looking at its behavioural. Behavioural malware detection is a field where malware is detected by its behaviour and the machine learning will look at the pattern of the behavioural. Then it will be analyzed and a report will be generate from the data. Thus, in this project, the behavioral of malware is analyzed and ensemble method is applied in detecting malware. Firstly, the data is collected by a multiple categories of system log and parser chooses from application. Then from the dataset it will classify it to 5 type of n-gram. Secondly, the best features from each of the n-gram are extracted using three feature selection techniques, namely Information Gain, Symmetrical Uncertainty and Chi-Square. SVM classifier is used to train the feature vectors and create a model for each n-gram. Finally, every model from 1-gram to 5-gram is combined using ensemble method. The significant contribution of this project is the effectiveness and efficiently of malware prediction using the state-of-the art techniques named ensemble method.

ABSTRAK

Dalam dunia teknologi terkini, peningkatan malware dalam penjenayah siber yang mengeksploitasi internet dan sentiasa mencipta dan mengedarkan malware berbahaya telah menjadi ancaman yang serius kepada dunia. Malware secara ketara telah memberi kesan kepada prestasi komputer dan sering kali tidak disedari oleh sistem sedia ada yang menyebabkan beberapa masalah kepada pengguna. Oleh itu, sangat penting untuk mengambil langkah berjaga-jaga yang diperlukan untuk mengesan dan mencegah jangkitan malware. Salah satu cara untuk mengesan pengesanan malware adalah dengan menggunakan teknik pembelajaran mesin. Pengesanan perisian hasad dikesan dengan melihat perilakunya. Pengesanan malware tingkah laku adalah bidang di mana malware dikesan oleh tingkah laku dan pembelajaran mesin akan melihat corak tingkah laku. Kemudian ia akan dianalisis dan laporan akan menjana dari data. Oleh itu, dalam projek ini, perilaku malware dianalisis dan kaedah ensemble digunakan untuk mengesan malware. Pertama, data dikumpulkan oleh pelbagai kategori log sistem dan parser yang dipilih dari aplikasi. Kemudian dari dataset ia akan mengklasifikasikannya kepada 5 jenis n-gram. Kedua, ciri-ciri terbaik dari setiap n-gram diekstrak dengan menggunakan tiga teknik pemilihan ciri, iaitu *Information Gain*, Ketidaktentuan Simetris dan Chi-Square. Pengelas SVM digunakan untuk melatih vektor ciri dan membuat model untuk setiap n-gram. Akhir sekali, setiap model dari 1 gram hingga 5 gram digabungkan menggunakan kaedah ensemble. Sumbangan besar projek ini adalah keberkesanan dan kecekapan ramalan malware menggunakan teknik ensiklopedia state-of-the-art yang dinamakan ensemble method.

TABLE OF CONTENTS

CHAPTER	SUBJECT	PAGE
	DECLARATION	ii
	DEDICATION	iii
	ACKNOWLEDGEMENT	iv
	ABSTRACT	v
	ABSTRAK	vi
	TABLE OF CONTENTS	vii
	LIST OF FIGURES	xi
	LIST OF TABLES	Xii
	LIST OF APPENDICES	xiv
CHAPTER I	INTRODUCTION	
	1.1 Introduction	1
	1.2 Problem Statement	2
	1.3 Project Question	3
	1.4 Project Objective	4
	1.5 Project Scope	4
	1.6 Project Contribution	5
	1.7 Thesis Organization	5
	1.8 Conclusion	6
CHAPTER II	LITERATURE REVIEW AND PROJECT METHODOLOGY	
	2.1 Introduction	7
	2.2 Related Work / Previous Work	8
	2.2.1 Malware	9

2.2.1.1 Malware Defination	9
2.2.1.2 Malware Type	9
2.2.1.3 Effects Of Malware	11
2.2.1.4 Malware Detection technique	11
2.2.2 Machine Learning	15
2.2.2.1 Machine Learning Definition or basic concept	15
2.2.2.2 Features	17
2.2.2.3 Features Selection	18
2.2.2.4 Classifier	20
2.2.3 Ensemble method	25
2.2.3.1 Ensemble method Definition	25
2.2.3.2 Ensemble method technique	26
2.2.3.3 Ensemble method advantage	27
2.3 Review of current problem and justification	27
2.3.1 Journal 1	27
2.3.2 Journal II	28
2.3.3 Journal III	29
2.3.4 Journal IV	29
2.4 Conclusion	30
<hr/>	
CHAPTER III PROJECT METHODOLOGY	
3.1 Introduction	31
3.2 Project Methodology	32
3.2.1 Data collection (phase I)	33
3.2.2 Feature Selection (phase II)	35
3.2.2.1 N-gram	35
3.2.3 Machine Learning Classifier (phase III)	37
3.2.4 Ensemble method (phase IV)	39
3.3 Project Schedule and milestone	40
3.3.1 Milestone of project	41
3.3.2 Gantt Chart	42
3.4 Conclusion	42

CHAPTER IV	ANALYSIS & DESIGN	
4.1	Introduction	43
4.2	Problem Analysis	44
4.3	Flow chart	45
4.3.1	Concept of flowchart	45
4.3.2	Flow chart diagram	45
4.4	Use case diagram	46
4.5	Process Design	51
4.5.1	Sequence Diagram concept	51
4.5.2	Sequence Diagram for malware detection	51
4.6	Requirement Analysis	53
4.6.1	Software Requirement	54
4.6.2	Hardware Requirement	55
4.7	Conclusion	56
CHAPTER V	EXPERIMENT ANALYSIS AND RESULT	
5.1	Introduction	57
5.2	Software Development Environment setup	58
5.3	Execution Of Implementation	63
5.3.1	Execution of implementation of N-gram	63
5.3.2	Execution of implementation of Ensemble Method	66
5.3	Result	69
5.3.1	Result 1-gram to 5-gram	69
5.3.2	Result Ensemble method	81
5.4	Conclusion	85
CHAPTER VI	DISCUSSION	
6.1	Introduction	86

6.2 Discussion Of the project	87
6.3 Discussion Of New Purpose Technique	88
6.4 Conclusion	91

CHAPTER VII CONCLUSION

7.1 Introduction	92
7.2 Project Summary	92
7.2.1 Weaknesses and Strength	93
7.3 Project Contribution	93
7.4 Project Limitation	94
7.5 Future Works	94
7.6 Conclusion	94

REFERENCES

95

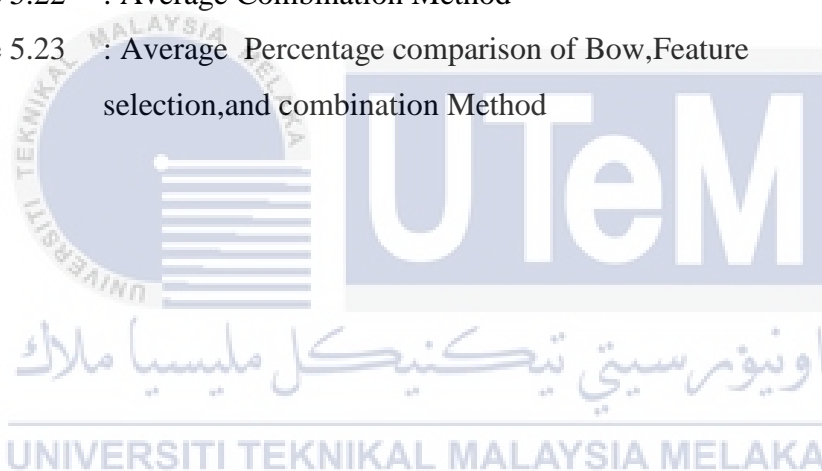
APPENDICES



LIST OF FIGURES

FIGURE	TITLE	PAGE
Figure 2.1	: Hierarchy	8
Figure 2.2	: Type of computer malware	11
Figure 2.3	: Machine Learning algorithm Category	16
Figure 2.4	: Illustration of Support Vector Machine	22
Figure 2.5	: Illustration of Decision Tree	24
Figure 2.6	: A common ensemble architecture	26
Figure 3.1	: Methology of this research.	33
Figure 3.2	: Network design of the experiment.	34
Figure 3.3	: Illustration of Support vector machine separation classifier	38
Figure 3.4	: Diagram process of 1-gram to 5-gram using combining classifier	39
Figure 3.5	Shows diagram of using a specific n-gram with a different classifier in ensemble method	40
Figure 3.6	: Gantt Chart	42
Figure 4.1	: Flow chart of project.	45
Figure 4.2	: Use case diagram for project module.	46
Figure 4.3	: Sequence Diagram phase I	52
Figure 4.4	: Sequence Diagram phase II	53
Figure 5.1	: Software requirement of this project.	58
Figure 5.2	: Implementation of the project	59
Figure 5.3	: Process generating the train and test file	60
Figure 5.4	: Sample Random Number of train and test	61
Figure 5.5	: Step by step verification random number	61
Figure 5.6	: The process create the train .scale.model	62
Figure 5.7	: The process generate the prediction file	63
Figure 5.8	: Step by step to run the .bat file.	65
Figure 5.9	: Process of ensemble method	66
Figure 5.10	: Formula majority voting	67

Figure 5.11	: Process Majority Voting	67
Figure 5.12	: Formula product rules	68
Figure 5.13	: Formula mean	69
Figure 5.14	: Accuracy of 10 run 1-gram	70
Figure 5.15	: Accuracy of 10 run 2-gram	72
Figure 5.16	: Accuracy of 10 run 3-gram	74
Figure 5.17	: Accuracy of 10 run 4-gram	76
Figure 5.18	: Accuracy of 10 run 5-gram	77
Figure 5.19	: Average bow and each feature selection of each n-gram.	79
Figure 5.20	: Total average of bow and other feature selection.	80
Figure 5.21	: Combination Method Of Ensemble Method	82
Figure 5.22	: Average Combination Method	83
Figure 5.23	: Average Percentage comparison of Bow,Feature selection,and combination Method	84



LIST OF TABLES

TABLE	TITLE	PAGE
Table 1.1	: Problem Statements	2
Table 1.2	: Summaries of Project Questions	3
Table 1.3	: Summaries of Project Objectives	4
Table 2.1	: Types of malware detection	14
Table 3.1	: Example 1-gram is function	36
Table 3.2	: Example 2-gram function	36
Table 3.3	: Example 3-gram function	37
Table 3.4	: Milestone of project	40
Table 4.1	: Narrative Use Case for Get Process Dataset	47
Table 4.2	: Narrative Use Case for Pre-processing Dataset	48
Table 4.3	: Narrative for Training Machine Learning	49
Table 4.4	: Predict file apply ensemble method	50
Table 4.5	: Software Requirement	54
Table 4.6	: Personal computer Specification	55
Table 5.1	: Prediction accuracy of total 10 number of for 1-gram)	69
Table 5.2	: Prediction accuracy of total 10 number of for 2-gram)	71
Table 5.3	: Prediction accuracy of total 10 number of for 3-gram)	73
Table 5.4	: Prediction accuracy of total 10 number of for 4-gram)	75
Table 5.5	: Prediction accuracy of total 10 number of for 5-gram)	77
Table 5.6	: Average of bow and each feature selection.	76
Table 5.7	: Average number of each n-gram.	80
Table 5.8	: Percentage of accuracy of all combination method	81
Table 5.9	: Average percentage from the N-gram 2 and average percentage of combination method	83
Table 6.1	: Average and accuracy of each of feature selection with the respective No of Ngram	87
Table 6.2	: Accuracy of 6 No of N-gram using the K-fold cross-validation	88

Table 6.3	: Accuracy of combination method	90
Table 6.4	: Comparison Accuracy of 2-gram with accuracy of Ensemble	90



CHAPTER I

INTRODUCTION



1.1 Introduction

Malware, short for malicious software, is any PC program software used to mess up PC or portable operations, collect sensitive data, have access to privatepc framework , show undesirable promoting. Malware may be stealthy, expected to take data or keep an eye on PC clients for a period time of without their insight and designed to infiltrate and damage computers and often go unnoticed in our systems and causes several problem . Thus, malware is a software programmed to do bad thing for user or user computer, and user don't want it. One of the way to detect malware detection is by using machine learning techniques. Malware detection is detected by looking at its behavioural. Behavioural malware detection is a field where malware is detected by its behaviour and the machine learning will look at the pattern of the behavioural Then it will be analyzed and a report will be generate from the data. Thus, in this project, the behavioral of malware is analyzed and ensemble method is applied in detecting malware. Other than that, The increasing of latest malware and various of malware that exploiting the internet and always create and

distribute harmful malware has become a serious issue. Thus, in this project we will purposed a new malware classification technique by analyzing and look at the pattern of the behavioral of malware and how to detect the the malware detection using an ensemble method . Using this technique we can analyze the malware behavioral and from the results showed that this technique approach can effectively classify and differentiate the malware.

1.2 Problem Statement

The problem that has been identified is summarized in Table 1.1 below

Table 1.1 Problem Statements

PS	Problem Statement
PS1	Malware detection has increasing day by day thus it is needed to detect malware s by using machine learning techniques. Malware detection is detected by looking at its behavioural .

a) PS1: Malware should be detected by looking at its behavioral

Malware be identified and detected so that the proposed to take data or spy on users pc without their knowledge for a period of time will not happen. One of the method to detect the malware is by using Machine Learning techniques. This technique will be used in this project to detect the behavioral attack of malware

1.3 Project Question

In this project, there is three Project Question (PQ) needs to be answered in this project. The summary of project question is shown in Table 1.2.

Table 1.2 Summaries of Project Questions

PQ	Project Questions
PQ1	What are the purposes on this project?
PQ2	What are the approach techniques that will use to detect the malware ?
PQ3	What is the contribution when applying the process and techniques?

PQ1: What are the purposes on this project?

Identify the purposes and goals of this project.

PQ2: What are approach or techniques that will use to malware ?

Distinguish the approach and procedures that will use to recognize malware to have a greater accuracy.

PQ3: What is the contribution when applying the process and techniques?

The enhance on applying the process and techniques on detecting malware to be more accurate .

1.4 Project Objective

Based on this project, the Project Objectives (PO) are developed as follows.

Table 1.3 Summaries of Project Objectives

PQ	Project Objectives
PO1	To study the performance of ensemble method that combine 1-gram to 5-gram feature vector
PO2	To measure the effectiveness of three different combination methods namely majority voting, product rules and mean rules.
PO3	To study the performance of ensemble method that combine specific n-gram feature vectors using different classifier model

1.5 Project Scope


The scope of these projects is to detect the behavioural of malware using ensemble method so that can easily classify it to malware or non-malware and Identify the technique and algorithm that is used to detect the malware. This project is to test the accuracy of the detection of malware using the machine learning technique and to find result of which technique is the most accurate.

1.6 Project Contribution

The contribution of this technique are by doing this project , it will evaluate from five n-gram of feature selection. Each of the n-gram will be evaluate and test based on the support vector machine classifier. Then, from the result each of the model will be combine to predict the accuracy and find out which is the best andmore accurate. The contribution of this finding is it will contribute an idea in finding the vest method to detect the malware and which n-gram has the high accuracy.

1.7 Thesis Organization

Chapter 1: Introduction



In this chapter 1 section it will review the history of project .

Other than that, In this chapter 1 also contain the problem statement, project question, project scope statement , project contribution, thesis organization and conclusion.

Chapter 2: Literature Review

In this chapter 2 it will review about about past work that are knit together to the project and review of previous researcher by study at the literature review. Each of the literature review result will be consider to help in generate this project. In addition . in this chapter also include the conclusion that are clear up to run this project .

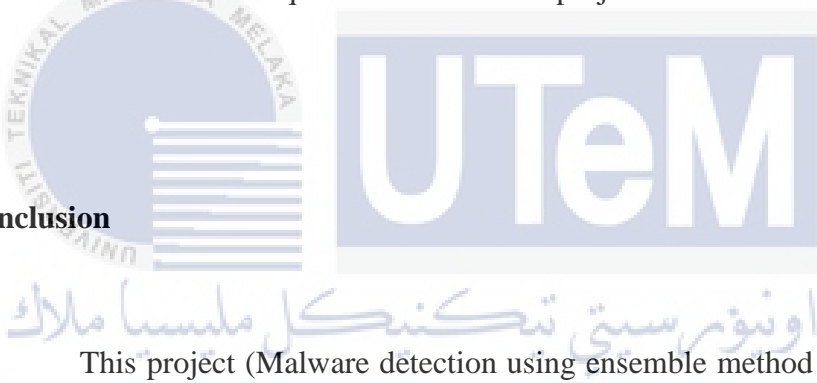
Chapter 3: Project Methodology

In this chapter 3 section it will review about the methodology that are used in these projects. This chapter explain step by step of the methodology in and the project milestone.

Chapter 4: Analysis and Design

In this chapter 4 it will review the discussion about the analysis of the project and the design of the project that had been used in the experiment that consists of the software and hardware requirement to run this project

1.8 Conclusion



This project (Malware detection using ensemble method) is a research project that detect the malware by looking at its behavioural. Malware will significantly impact on your computer performance and often go unnoticed in your system. Thus its imperative to take the precautions necessary to detect and prevent malware infections . In this reasearch it used the machine learning technique to detect the malware using an algorithm. From the algorithm the data will be identified whether it is malware or non-malware. The objective of the project is to help users to prevent from being attack by the user by analysis the behavioural of malware. Thus in next chapter will explains the previous research about malware detection and how the data is conducted.

CHAPTER II



2.1 Introduction

The goals of this Chapter 2 is to summarize, explained and increase knowledge and provide better understanding about malware detection and ensemble method. In this chapter 2 we will explain the definition, approach and elaborate from a technique that were used from previous research from a literature review. From the literature review we can compare the findings and evaluate based on the information. The explanation of literature review are categorize to several topics which consist of malware, detection technique and ensemble method. In malware will focus more on the definition of malware, effect of malware, type of malware and the use of malware. Next for the

detection technique will be focus more on machine learning technique where will explain the definition and benefits of machine learning technique . Lastly in this topic also will explain the detection that are use in this research which is ensemble method. In the ensemble method will focus more on definition of ensemble method and why ensemble method is use in conducting this research. The hierarchy of this topics is based on the figure 2.1 and consists of three phase . Phase one is about malware , phase 2 is detection technique and lastly phase 3 is detection technique.

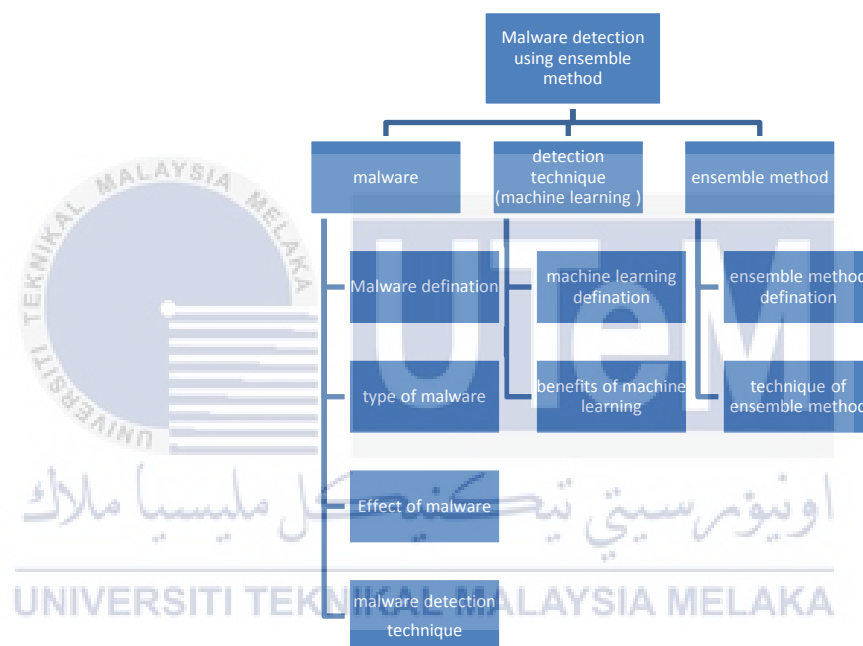


Figure 2.1 Hierarchy

2.2 Related Work / Previous Work

Based on figure 2.1 , this section 2.2.1 is an explanation about malware , section 2.2.2 explanation about detection technique name machine learning and 2.2.3 definition about ensemble method and responsibility of ensemble method. In addition in this chapter also discuss the related work of previous researcher . All of the section will explain based on understanding and previous research..