



Republic of Iraq
Ministry of Higher Education
and Scientific Research
University of Diyala
College of Science



Detection of Malware under Android Mobile Application

A Thesis

Submitted to the Department of Computer Science\ College
of Science \ University of Diyala as a Partial Fulfillment of
the Requirements for the Degree of Master in Computer
Science

By

Saja Ibraheem Hani Ismail

Supervised By

Prof. Naji M. Sahib

2020 A.D.

1441 A.H.

بِسْمِ اللَّهِ الرَّحْمَنِ الرَّحِيمِ

يَا أَيُّهَا الَّذِينَ آمَنُوا إِذَا قِيلَ لَكُمْ تَفَسَّحُوا فِي الْمَجَالِسِ
فَانْفَسِحُوا يَفْسَحِ اللَّهُ لَكُمْ وَإِذَا قِيلَ انشُرُوا فَانشُرُوا يَرْفَعِ
اللَّهُ الَّذِينَ آمَنُوا مِنْكُمْ وَالَّذِينَ أُوتُوا الْعِلْمَ دَرَجَاتٍ وَاللَّهُ بِمَا
تَعْمَلُونَ خَبِيرٌ ﴿١١﴾

صَدَقَ اللَّهُ الْعَظِيمُ

سورة المجادلة آية ﴿١١﴾

Dedication

*To those who taught me how to stand firmly on
the ground*

Dear father

To the source of love, altruism and generosity

My Mother

To the closest people to myself

My faithful husband

To my soul, my eyesight, and my heartbeat

My children

*I present to you a summary of my scientific
efforts*



Saja Ibraheem

Acknowledgment

At the end of my thesis, I am pleased to thank and show my gratitude to those who have helped me to achieve my study whether through assisting to provide resources or with guidance and presenting me advice.

First, I would like to thank my honorable prof. Naji Muttar Sahib who supervised my thesis and offered me many scientific advices and guidelines; to him I give all of my appreciation.

In addition, I want to thank my respected professors in Department of Computer Science with whom I completed a very important stage of education within the academic journey through their scientific experiences and helping us to complete it.

Also, I would love to thank prof. Azhar Hasan Nussief who provided me with solid resources, so I give my gratitude to her.

In addition, I want to thank my family who supported me and stood by me during my study journey.



Saja Ibraheem

Supervisors' Certification

I certify that this thesis entitled "**Detection of Malware under Android Mobile Application**" was prepared by "**Saja Ibraheem Hani**" under my supervisions at the University of Diyala collage of Science Department of Computer Science, as a partial fulfillment of the requirements needed to award the degree of Master of Science in Computer Science.

(Supervisor)

Name: Prof. Najj M. Sahib

Signature: 

Date: / / 2020

Approved by University of Diyala\ College of Science\
Department of Computer Science.

Signature: 

Name: Asst. Prof. Dr. Taha Mohammed Hasan

Date:

Head of Computer Science Department

(Linguistic Certification)

I certify that this research entitled “**Detection of Malware under Android Mobile Application**” was prepared by “**Saja Ibraheem Hani**” and was reviewed linguistically. Its language was amended meet the style of English language.

Signature:

Name:

Date:

Scientific Amendment

I certify that the thesis entitled “**Detection of Malware under Android Mobile Application**” was prepared by "**Saja Ibraheem Hani**" has been evaluated scientifically; therefore, it is suitable for debate by examining committee.

Signature :

Name :

Date : / / 2020

Abstract

Smartphones have become essential in our daily life. It also can do a lot of work and can browse the Internet, and download many applications for each device, through the available store. As a result, the number of malware applications downloaded also increases.

This malware carries out various activities behind the scenes; Such as confidentiality, breach of privacy, loss of confidentiality, system breakdown, theft of sensitive information, etc.

Many types of research and studies that proposed different techniques to detect malicious programs, but they contained weak points, which are illustrated by efficiency, speed, and lack of comprehensiveness.

In this thesis, a proposed system is developing implementing to detect malware in smartphones, and contains two parts:

In the first part, access control is initially started upon the system launch. The user authentication algorithm adopts the user's permission to detect a threat factor after applying a user's permission policy by improving the method with which the user's activities are extracted. While, **in the second part**, anomaly detection technology begins to extract the important features that play an effective role in detecting malicious code.

The proposed system has been tested by using a hybrid genetic algorithm, and the SVM data has been registered an accuracy of (0.9282).

The experimental results indicate that the proposed system has a high average accuracy rate compared to other existing methods where it (0.8848) average accuracy using PNN, while the average accuracy is (0.8835) and (0.8715) with SVM and K-NN respectively.

Table of Contents

<i>subject</i>	<i>Page NO.</i>
Chapter 1: Introduction	1-14
1.1 Overview	1- 5
1.2 Literature review	5-8
1.3 Problem Statement	8
1.4 Aim of Thesis	9
1.5 Contribution	9
1.6 Thesis Outline	9-10
Chapter 2: Theoretical Background	11-34
2.1 introduction	11
2.2 Android and Application Definition	11-13
2.3 Application Programming Interface call	13
2.4 Android Features Extraction	14-16
2.5 Android Permission and Security Model	16-18
2.6 Malware Detection Methods	19
2.6.1 Types of Malware and How they Affect the System	20-21
2.6.2 Malware Detection Techniques	22-23
2.6.3 Classification Detection Techniques	23-26
2.7 Classification Algorithms	26-27
2.7.1 k-Nearest Neighbor (k-NN) Algorithm	27
2.7.2 Random forest algorithm	27-28
2.7.3 Support Vector Machine (SVM) Algorithm	28-29
2.7.4 Genetic Algorithms (GAs)	29-31
2.7.5 Probabilistic Neural Networks Algorithm	32
2.8 Parameters Used to Evaluate Classification	32-34
Chapter 3: The Propose Detection Malware System	35-53
3.1 Introduction	35
3.2 Proposed System Architecture	35-36
3.3 Access Control Detection	36-37
3.3.1 User Permission	37-39
3.3.2 Logical Rule Based	39-42

3.4 Anomaly Detection	42-45
3.5 Machine Learning Technique for Classifying Dataset	45-46
3.5.1 Data Cleaning Dataset	47-48
3.5.2 Filling in the Missing Value in the Dataset	48-49
3.5.3 Classification with SVM	49-50
3.5.3 Classification with Genetic Algorithm SVM	50-51
3.5.4 Classification with PNN	51-52
3.5.5 Classification with KNN	52-53
Chapter 4: Implementation and Experimental Results	54-68
4.1 introduction	54
4.2 Result Presentation	54-59
4.3 Results about Datasets and Configuration for Algorithms	60
4.3.1 The probabilistic Neural Network Detection Malware in Android	60-61
4.3.2 The KNN Detection Malware in Android	61-63
4.3.3 The SVM Detection Malware in Android	63-64
4.3.4 The Hybrid Genetic Algorithm and SVM Detection Malware in Android	65
4.4 Evaluation of Official Market Metadata	66-68
Chapter 5: Conclusions and Recommendations	69-70
5.1 Conclusion	69
5.2 Future Work	69-70
References	71-77

List of Figures

<i>Figure No.</i>	<i>Caption</i>	<i>Page No.</i>
Figure 1.1	Rate of Mobile Application Downloaded	2
Figure 2.1	System Architecture of Malware Detection Model	14
Figure 2.2	Malware Detection Techniques	25
Figure 2.3	Examples of Crossover	31
Figure 2.4	Two-point crossover	31
Figure 2.5	Mutation	31
Figure 2.6	PNN general architecture	32
Figure 3.1	general flowchart of the proposed system	36
Figure 3.2	Access Control Detection	37
Figure 3.3	Authentication Flowchart	38
Figure 3.4	General Block Diagram (anomaly detector)	46
Figure 3.5	Dataset	47
Figure 4.1	User Login	55
Figure 4.2	User Doesn't have an account	56
Figure 4.3	Login Process	57
Figure 4.4	Shows all Users that Match with Registration and Login process.	57
Figure 4.5	URL page download	58
Figure 4.6	Enter URL	59
Figure 4.7	Download the Page	59
Figure 4.8	Testing Accuracy using PNN	61
Figure 4.9	Testing Accuracy using K-NN	62
Figure 4.10	Representation of classification accuracy rate for each Value of Parameter k Based on Random Rule.	63
Figure 4.11	Testing accuracy using SVM	64

Figure 4.12	Testing accuracy Using Hybrid Genetic Algorithm and SVM	65
Figure 4.13	The Compare of Algorithms	68

List of Table

Table page	Table title page	page
2.1	Advantage and Disadvantage	29-30
2.2	Sample Confusion Matrix	37
4.1	Classification of Training PNN	61
4.2	Classification of Test KNN	62
4.3	Classification of Test SVM	64
4.4	Classification of Test hybrid Genetic algorithm and SVM	65
4.5	The Result of Performance Evaluation of machine learning technique	67
4.6	Comparison with Some Related work	68

List of Abbreviations

Abbreviations	Meaning
API	Application Programming Interface
APK	Android Application Package
AV	Anti-Virus
GA	Genetic Algorithm
IOS	IPhone Operating System
JVM	Java Virtual Machine
K-NN	K-Nearest Neighbors
PNN	Probabilistic Neural Networks
RF	Random Forests
SVM	Support Vector Machine
PE	Portable Executable
DEX	Dalvik Executable
XML	Extensible markup Language
HTML	Hypertext Markup Language
SQL	Structured Query Language

Chapter One

Introduction

Chapter one

Introduction

1.1 Overview

In the past few years, it was clear that smartphone users have increased exponentially. Besides, the operating systems for smartphones are Symbian, iPhone Operating System (IOS), Android, and Blackberry. The smartphone is viewed as a portable Personal Computer system, PCs, as they have all the functionalities of a desktop PC integrated into them. Just as there are hackers/attackers releasing malware for PCs, there are attackers who are now targeting smartphones. The main reason for this is that mobile security is still in its initial stages and the lack of user awareness regarding how their devices can be undermined by using if they are not careful enough. Google's is open-source operating systems. Android is among the most popular smartphone operating systems. Android is a Linux-based operating system that also includes key applications and middleware. In order to fully benefit from and explore the functionality of Android, Google allows third-party developers to create applications and release them to the Android [1].

A recent work indicates the number of mobile applications is increased extremely which also increases malware application as shown in Figure (1.1) [2].

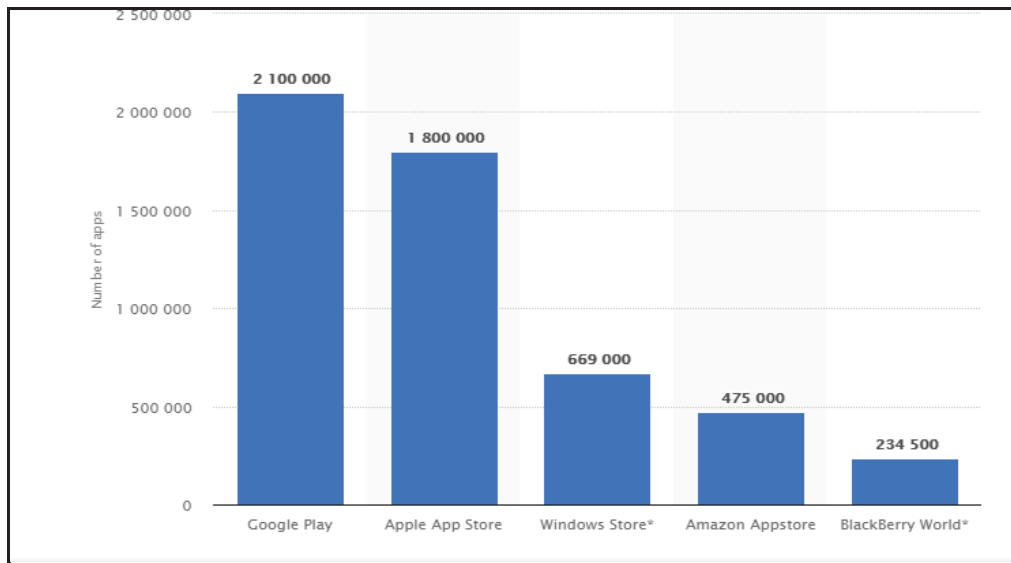


Figure (1.1): Rate of mobile Application downloaded [2]

Malware produced as malicious content that damages the stand-alone computer or networked computers damaged by its harmful effects. It can come in many forms; Spyware, Worm, Trojan or Virus, etc. But whatsoever the form is; its function is the same that is to harm the computers. This intrusive virus can be executable code or non-executable code. Many malware types have the ability to multiply so strict precautions are needed [3].

There are various types of malicious software with different structures, properties, and effects. They also vary in the intensity of the threats they pose [4].

There are different types of malware, the most common one is the virus (an infected code), which after execution multiplies itself and infects other files. It also adds malicious code to other files in order to attack more vigorously. Although viruses are hosted and controlled by a third party, worms create all the damage without being controlled by anyone. They can propagate themselves by infecting the other files. Another side

we have that worms are simply standalone malicious software, which creates the same damage as viruses [5].

Another most common type of malware is adware. It is supposed to be used solely for advertisement and generating revenue. However, nowadays adware has been combined with spyware, it keeps track of user's activities. Addition, Spyware attacks to steal user's sensitive information. On the other hand, adware is simply the popping up of ads on websites or applications. Normally this comes up with free software. Intruders are making use of it and transform adware into spyware stealing user's information and misusing it [6].

Spyware is one of the dangerous malwares, which keeps an evil eye on the user's activities by recording those keystrokes and personal information. Furthermore, personal information can vary from login credentials to sensitive bank account details. Not only stealing users' information, but also spyware can intrude into the user's computer in order to change some software's security or privacy settings or browser's setting making the networks public. The other type of malware is the bot which is an automatic malicious code that intrudes exclusively network of computers. The bots are casually used for positive purposes, but now malicious harmful attacks have been imposed on them. Denial-of-service-attack is an attack on the host computer that transfers its virus to all the networked computers, which is an output of bots. Also, this is a malicious code, which works automatically. Spambots are one of its types, which spam the internet with malicious websites. It is difficult to get rid of such malware but not impossible [7]. Ransomware is also a type of malware, which takes over the hard drive of the user, and the user has to pay some ransom to regain access. It is a crucial kind of malware, which restricts user access to its own computer. Furthermore, it spreads through the downloaded file or any vulnerability in the network service. Trojan horse

is a trick to the users; it presents itself as an authentic file which the user can download. Afterward, the attacker gets access to the infected computer remotely. Now the Trojan can add more malware to the infected computer, control the security configurations, monitor user key logs, steal sensitive information, etc. Such infected computers can be used in botnets [8].

The user authentication algorithm of user permission is dependent in order to detect a threat actor after applying the user permission to approach via the improvement of the user activities extraction method.

In addition, to adapt the unknown sessions and use the rule-based on integration with the attacks, the heuristic sequence will be used. A very important aspect is the identification of the classification algorithm of the most reliable detection accuracy. Therefore, the strongest classifiers are identified through the evaluation of the activeness and detection accuracy of all machine-learning algorithms. Modifying the default input values can enhance the efficiency and accuracy of a classifier. However, enabling an equivalent comparison between the classifiers dictated the implementation of the classifiers with their default input values. Many studies have been emerged to discover and treat malicious programs based on the artificial intelligence algorithm by many researches, such as K-Nearest Neighbors (k-NN), Naive Bayes, Random Forests (RF), Support Vector Machines (SVM), and Genetic algorithm [9].

The k-NN algorithm belongs to the family of methods known as instance-based methods. These methods are based on the principle that observations (instances) within a dataset are usually placed close to other observations that have similar attributes; this method selects the closest observations from the dataset in such a way to minimize the distance.

In machine learning, random forest (RF) is already widely used in bioinformatics the best available methods and superior to most methods

in common use. As the name suggests, RF combines many classification trees to produce more accurate classifications. By-products of the RF calculations include measures of variable importance and measures of similarity of data points that may be used for clustering, multidimensional scaling, graphical representation, and missing value imputation.

The machine learning techniques are Support Vector Machine (SVM) which is used for binary classification. It is a very general technique that can be applied in a wide variety of situations. Also, it has special characteristics that are used to implement efficient parallel algorithms in terms of time and memory. One characteristic is that the solution to the classification problem is obtained by only a few samples called Support Vectors (SV) that determine the maximum margin separating hyperplane. Another characteristic of SVM is to perform the nonlinear mapping without knowing the mapping function using predefined functions called kernels for calculating the inner product of mapping functions [10].

Genetic algorithm (GA) is a class of stochastic global search techniques based on biological evolution principles; several have applied the genetic algorithm to geophysical optimization problems such as seismic attributes. The genetic algorithm represents parameters as an encoded binary string and works with the binary strings to minimize the cost, while the other works with the continuous parameters themselves to minimize cost [11].

1.2 Related work

Several researchers have shown their interest in the detection of malware in a smartphone, the following are some of the published works that are relevant to the current work:

- **Lu, et al., 2013 [12]:** compared Bayesian method alone and Bayesian method combined with Chi Square feature selection method results are compared to evaluate the performance of the two ML algorithms. The study concluded that Bayesian method with Chi Squared yielded an accuracy of 89% while Bayesian method alone yielded 80%.
- **Nuray Baltaci, et al., 2014 [13]:** The main purpose of the study is to investigate the contribution of other application market metadata to the detection of malicious applications in addition to requested permissions. Hence, the information of applications presented on the official market when a user wants to download them was used as the feature set for training supervised classification algorithms.
- **Kurniawan, et al., 2015 [14]:** used Logger, a default application which is inbuilt in Android was used to extract the sum of Internet traffic, percentage of battery used and battery temperature for every minute. This information collected as set of features and is fed into weka, an open source learning library for testing and training with Naive Bayes, J48 decision tree and Random Forest algorithms. The author concluded that Random Forest has high accuracy of 85.6% with these features and proposes other features that can be combined with existing system to improve the accuracy.
- **Weng, et al. [15]:** in their work, published in 2017, propounded a model that classifies using machine learning techniques such as SVM, the nearest neighbors, by extracting 11 different static features. This model can be used in the management of large application markets. A correct classification rate of 99% for the malicious set and 82% for the benign set was obtained working with 100,000 benign and 8,000 malicious application set.

- **Mohsen Kakavand, et al.,2018 [16]:** this work involves in static analysis of applications, which checks for the presence and frequency of keywords in Android application' manifest file and drives the static feature sets from a 400- application dataset to produce better malware detection results. The classification performance of the ML algorithms is measured in terms of accuracy and true positive rate and interpreted to determine which algorithm is more applicable for the Android malware detection. The experimental results for a dataset of real malware and benign applications indicate the average accuracy rate of 79.80% and 80.50% with average true positive rate of over 67% and 80% using SVM and KNN, respectively.

- **Matthew Leeds, et al., 2017 [17]:** Malware is a current threat facing Android users. As users have come to depend on these devices for communication and information, it is essential to make sure they are secure. Therefore, developing and testing new sophisticated malware detection techniques must be a priority. This paper compared two prominent features used to detect Android malware, permissions and system calls, and applied machine learning to both. The results showed that permissions data was better at detecting malware than system call data. An average classification accuracy rate of 80% was achieved when using permissions data to determine malicious activity on Android devices. Therefore, it is a reliable way to detect malware.

- **Michal Kedziora, et al., 2018[18]:** In this study, an overview of Android malware analysis was presented, and a unique set of features was chosen that was later used in the study of malware classification. Five classification algorithms (Random Forest, SVM, K-NN, Nave Bayes, and Logistic Regression) and three attribute

selection algorithms were examined in order to choose those that would provide the most effective malware detection.

- **H al-kaaf1, et al., 2019 [19]:** In this study, proposed feature selection methods to identify clean and malicious applications based on selecting a set combination of permission patterns using different classification algorithms such as sequential minimal optimization (SMO), decision Tree (J48), and Naive Bayes. The experimental results show that sequential minimal optimization (SMO) combining with the SymmetricalUncertAttributeEval method achieved the highest accuracy rate of 0.88, with the lowest false positive rate of 0.085 and the highest precision of 0.910. And the findings prove that feature selection methods enhanced the result of classification.

1.3 Problem Statement

Develop a new method for more classifications performance serve to protect Mobil Application malware.

Applications on smartphones must take Care when Downloading, Due to that, many malicious attacks target them. The majority of operating systems in the smartphone business are operating using the Android OS. However, around 97% of mobile malware targets Android phones. In the Therefore, these incidents motivated us to study mobile application security, especially in Android because the viruses pose risk to the applications as well as the operator.

1.4 Aim of Thesis

The current study aims to find whether or not Google Play, Google's app, and Android's official application market, metadata of Android

applications assist in explaining the malicious behaviors when joined with user's permissions analysis.

1.5 Contribution

Using a hybrid system to detect malware based on static and dynamic approaches. Hence, this contribution will provide protection for the Android mobile based-on access control and anomaly detection.

1.6 Thesis Outline

Beside this chapter, the remaining parts of this thesis include the following chapters:

Chapter Two: Theoretical background

The start with an overview of the Android system architecture and describe the implementation design of Android. Also, discuss an overview of the core components which are found in android applications and it's included in the concepts. That concept is used in this thesis, where the methods used in the malware app such as classifying android apps, SVM, and Genetic Algorithm.

Chapter Three: The propose Detection Malware System

In this chapter, the discussed the proposed system for the Authentication process and the check authentication method with Algorithms SVM as well as the Genetic algorithm.

Chapter Four: Implementation and Experimental Results

This chapter involves studies and results, which are obtained from the system running as well as the performance measures of the results of the test, and comparisons.

Chapter Five: Conclusion and Future Work

In this chapter, the present a list of conclusions from the results of the presented work and some suggestions for future works.

Chapter Two

Theoretical Background