# Malware behavioural detection and vaccine development by using a support vector model classifier

## Ping Wang [*], Yu-Shih Wang

*Department of Information Management, Kun Shan University, Tainan, Taiwan*

A B S T R A C T

Most existing approaches for detecting viruses involve signature-based analyses to match the precise patterns of malware threats. However, the problem of classification accuracy regarding unspecified malware detection depends on correct extraction and completeness of training signatures. In practice, malware detection system uses the generalization ability of support vector models (SVMs) to guarantee a small classification error by machine learning. This study developed an automatic malware detection system by training an SVM classifier based on behavioural signatures. A cross-validation scheme was used for solving classification accuracy problems by using SVMs associated with 60 families of real malware. The experimental results reveal that the classification error decreases as the sizing of testing data is increased. For different sizing ($N$) of malware samples, the prediction accuracy of malware detection goes up to 98.7% with $N = 100$. The overall detection accuracy of the SVC is more than 85% for unspecific mobile malware.

© 2014 Elsevier Inc. All rights reserved.

## 1. Introduction

Effective security defence mechanisms involving threat-analysis techniques in open networks are essential for detecting intruder attacks. In implementing network applications, defence mechanisms against network threats must focus on two fundamental security concerns. First, vulnerabilities that are exploited by malware must be identified, and the exploitability must be compared with that of attack scenarios. Second, established methods for detecting malware must be used for classifying malicious executables to respond promptly to cyber attacks [1,2].

The automatic malware detection system (AMDS) [3] is typically used for detecting and evaluating potential attack profiles by incorporating cyber-threat analysis (CTA) [4,5] techniques to assist defenders in determining effective defences against network threats caused by malware infection. CTA of malware attacks typically focuses on examining threats and their exposure by accumulating information on recognised attacks to identify malware signatures associated with system vulnerabilities to estimate detection accuracy and the impact of malware threats, as described in the Common Vulnerabilities and Exposures dictionary.

Current malware detection schemes, such as signature-based and semantic analyses, provide methods for examining the precise patterns of malware threats. Signature-based detection is the most widely employed technique in antivirus software featuring precise comparison. Studies on malware detection have primarily focused on performing static analyses to inspect the code-structure signature of viruses, rather than dynamic behavioural aspects. In other words, automatic analysis of malware behaviour using machine learning (ML) techniques for determining unidentified classes of malware or variants

---

\*  Corresponding author.
   *E-mail address:* pingwang@mail.ksu.edu.tw (P. Wang).

have generally been disregarded. The complete installation of precise patterns for real-time malware detection programs on mobile devices remains a challenging task, restricted by computing power, battery capacity, and limited storage space.

Support vector machines (SVMs) [6,7] are used for clustering data into two categories according to maximum boundary geometry. Solving classification problems by using an SVM classifier (SVC) guarantees few classification errors caused by maximising the generalisation ability of learning incorporating the Lagrange multiplier optimisation algorithm [8]. Generally, the results from using SVM classification algorithms are more accurate than those from using other ML approaches involving non-optimised search methods, such as artificial neural networks, least squares, $k$-nearest neighbour, Bayesian probability, and classification and regression trees [9], particularly when defence systems collect only limited training data.

Previous researches [10–12] have indicated that SVM analysis is useful for discriminating malicious behaviour of malware from the normal behaviour of legitimate applications by training a classifier. SVM is not only used for detecting identified malware, but also for predicting the classification of unidentified malware.

Malware variants generally exhibit similar behaviours to those of malware signatures. Thus, behavioural detection schemes can detect new malware or variants based on existing malware. Furthermore, generalisation is a key benefit of using a behavioural approach instead of payload signatures [10]. This study proposes an improved behaviour-based SVC learning algorithm for use in the AMDS to categorise mobile malware according to collected malware behaviours, enabling the defence system to respond promptly to high-risk security concerns. In particular, the learning signature synthesises both the code features of static analysis and the behavioural patterns of dynamic analysis techniques. A grid search algorithm [13] was used to facilitate defence systems to increase training accuracy by considering the selection of SVC parameters. Furthermore, a digital vaccine (DV) [14] against cyber attacks was developed as a defence solution for preventing malware infection in mobile devices. System validation involved a cross-validation scheme to classify and identify the class of malware by using SVCs associated with 60 families of real mobile malware to test their accuracy.

Three crucial steps were considered in developing the proposed model: 1) perform malware classification by using a behaviour-based SVC based on a heuristic approach incorporating both behavioural [10,15] and code analyses [16,17] of malware to accurately determine the signature of mobile devices, 2) propose a digital vaccine that effectively prevents malware infection and treats infected hosts with a backup and restore approach, and 3) comprehensively determine the classification accuracy of unspecified malware or variants.

The remainder of the paper is organised as follows: Section 2 reviews the basic principles of SVM. Section 3 presents the proposed analytical model used to evaluate the detection accuracy of mobile malware and discusses an effective defence solution for malware infection, and Section 4 describes the proposed approach by presenting two cases of mobile malware attacks on a cloud computing security application. Finally, Section 5 concludes the paper.

## 2. Related work

This section reviews the use of two important issues, namely classification techniques for anomaly-based detection and SVM, in establishing classification rules of malware detection for discriminating the abnormal behaviour from the normal behaviour of the AMDS to solve the threat analysis problem of cyber attacks.

### 2.1. Classification techniques for anomaly-based detection

The classification problem is most frequently discussed regarding data mining or machine learning (ML) techniques. Many classification approaches incorporate ML algorithms for detecting malware [13,18,19]. Machine learning techniques for classification algorithm, including SVMs, least squares (LR), $k$-NN ($k$-nearest neighbours), decision tree (DT), artificial neural networks (ANN), and Bayesian classifiers approaches have been used to facilitate the prediction performance. These classification schemes are summarised in Table 1.

Compared with those of the traditional ML methods such as LR and $k$-NN, SVMs have produced excellent results, and are generally considered as best classifiers by a clear margin as the feature set gets larger provided the sample size is not too small [22]. It implies that SVM can be used to discriminate the abnormal behaviour from the normal behaviour for AMDS by training a binary classifier.

### 2.2. Support vector machine

SVMs, developed by Vapnik in 1995 (AT&T Labs), are supervised learning models associated with learning algorithms and used to analyse data and recognise patterns. In an SVM training algorithm, new examples are assigned to one category or another as either nonlinear or linear binary classifiers obtained from a set of training examples. SVMs have been proven a useful tool for conducting clustering and classification analyses. In particular, SVM theory has been developed gradually from linear SVCs to hyperplane classifiers; that is, SVMs can efficiently perform nonlinear classification by using a kernel function, implicitly mapping their inputs into high-dimensional feature spaces by selecting an appropriate kernel function. Furthermore, a favourable classification result is achieved using a hyperplane that has the largest distance from the nearest training data point of any class [4]. Basic SVM theory is discussed as follows [3,8,23,24].

Given a training dataset $D$, $(x_i, y_i)$, where $x_i$ denotes $n$ observations of malware signatures, $x_i \in R^N$, $i = 1, ..., N$; and $y_i$ is the corresponding class label whose value is either 1 or $-1$ (i.e., malicious or benign), indicating the class to which the