

# Detection of Unknown Computer Worms based on Behavioral Classification of the Host

Robert Moskovitch, Yuval Elovici, Lior Rokach

Deutsche Telekom Laboratories at Ben-Gurion University  
Ben Gurion University, Be'er Sheva, 84105, Israel  
{robertmo, elovici, liorrk}@bgu.ac.il

## Abstract.

Machine learning techniques are widely used in many fields. One of the applications of machine learning in the field of the information security is classification of a computer behavior into malicious and benign. Anti viruses consisting on signature-based methods are helpless against new (unknown) computer worms. This paper focuses on the feasibility of accurately detecting unknown worm activity in individual computers while minimizing the required set of features collected from the monitored computer. A comprehensive experiment for testing the feasibility of detecting unknown computer worms, employing several computer configurations, background applications, and user activity, was performed. During the experiments 323 computer features were monitored by an agent that was developed. Four feature selection methods were used to reduce the amount of features and four learning algorithms were applied on the resulting feature subsets. The evaluation results suggests that using classification algorithms applied on only 20 features the mean detection accuracy exceeded 90%, and for specific unknown worms accuracy reached above 99%, while maintaining a low level of false positive rate.

**Keywords:** Classification Algorithms, Decision Tree, Bayesian Networks, Feature selection; Artificial Neural Networks; Naive Bayes classifier; Malicious Code, Worms detection

## 1 Introduction

Malicious code (malcode) detection, transmitted over computer networks has been researched intensively in recent years (Kabiri and Ghorbani, 2005). One type of abundant malcode is *worms*, which proactively propagate across networks while exploiting vulnerabilities in operating systems and programs. Other types of malcode include computer *viruses*, *Trojan horses*, *spyware*, and *adware*. In this study we focus on worms, though we plan to extend the proposed approach to other types of malcode.

Nowadays, excellent technology (i.e., antivirus software packages) exists for detecting and eliminating *known* malicious code. Typically, *antivirus software packages* inspect each file that enters the system, looking for known signs (signatures) which uniquely identify an instance of known malcode. Nevertheless, antivirus technology is based on prior explicit knowledge of malcode signatures and cannot be used for detecting unknown malcode. Following the appearance of a new *worm*, a

patch is provided by the operating system provider (if needed) and the antivirus vendors update their signatures-base accordingly. This solution is not perfect since worms propagate very rapidly and by the time the local antivirus software tools have been updated, very expensive damage has already been inflicted by the worm (Fosnock, 2005).

Intrusion detection, commonly at the network level, called *network based intrusion detection (NIDS)*, was researched substantially (Kabiri and Ghorbani, 2005). However, *NIDS* are limited in their detection capabilities (like any detection system). In order to detect malcodes which slipped through the *NIDS* at the network level, detection operations are performed locally at the host level. Detection systems at the host level, called *Host-based Intrusion Detection (HIDS)*, are currently very limited in their ability to detect unknown malcode.

Recent studies have proposed methods for detecting unknown malcode using Machine Learning techniques. Given a training set of malicious and benign executables binary code, a classifier is trained to identify and classify unknown malicious executables as being malicious (Schultz et al., 2001; Abou-Assaleh et al., 2004; Kolter and Maloof, 2006; Caia et al., 2007).

Existing methods rely on the analysis of the binary for the detection of unknown malcode. Some less typical worms are left undetectable. Therefore an additional detection layer at runtime is required. The proposed approach assumes that the malicious actions are reflected in the general behavior of the host. Thus, by monitoring the host, one can inexplicitly identify malcodes. This property can be used as an additional protection layer.

In this study, we focus on detecting the presence of a worm based on the computer's (host) behavior. Our suggested approach can be classified under *HIDS*. The main contribution of our approach is that the knowledge is acquired automatically using inductive learning, given a dataset of known worms (avoids the need for *manual* acquisition of knowledge). While the new approach does not prevent infection, it enables a fast detection of an infection which may result in an alert, which can be further reasoned by the system administrator. Further reasoning based on the network-topology can be performed by a network and system administration function, and relevant decisions and policies, such as disconnecting a single computer or a cluster, can be applied.

Generally speaking, malcode within the same category (e.g., *worms*, *Trojans*, *spyware*, *adware*) share similar characteristics and behavior patterns. These patterns are reflected by the infected computer's behavior. Thus, we hypothesize that it is feasible to learn the computer behavior in the presence of a certain type of malcode, which can be measured through the collection of various parameters along time (CPU, Memory, etc.). In the proposed approach, a classifier is trained with computer measurements from infected and not infected computers. Based on the generalization capability of the learning algorithm, we argue that a classifier can further detect previously unknown worm activity. Nevertheless, this approach may be affected by the variance in computer and application configurations as well as user activity (running and using various applications) on each computer. In this study, we investigate whether an unknown worm activity can be detected, at a high level of accuracy, given the variation in hardware and software environmental conditions on individual computers, while minimizing the set of monitored features.

In this paper we introduce three main contributions: We show that current machine learning techniques are capable to detect and classify worms solely by monitoring the host activity. Using feature selection techniques we show that a relatively small set of features are sufficient for solving the problem without sacrifice accuracy. We present empirical results from an extensive study of various machine configurations suggesting that the proposed methods achieve high detection rates on previously unseen worms.

The rest of the paper is structured as follows: in section 2, a survey of the relevant background for this study is presented. The methods used in this study are described in section 3, followed by the description of the experiments design in section 4. In section 5 we present the evaluation results and conclude with summary and conclusions in section 6.

## 2 Background and Related Work

### 2.1 Malicious Code and Worms

The term 'malicious code' (malcode) refers to a piece of code, not necessarily an executable file, intended to harm, whether generally or in particular, a specific owner (host). The approach suggested in this study aims at detecting any malcode activity, whether known or unknown. However, since our preliminary research is on worms, we will focus on them in this section.

Kienzle and Elder (2003) define a *worm* by several aspects through which it can be distinguished from other types of malcode: 1) *Malicious code* – worms are considered malicious in nature; 2) *network propagation or human intervention* – a commonly agreed-upon aspect, that is, worms propagate actively over a network, while other types of malicious codes, such as viruses, commonly require human activity to propagate; 3) *standalone or file infecting* – while viruses infect a file (its host), a worm does not require a host file, and sometimes does not even require an executable file, residing entirely in the memory, as did the *Code Red* (Moore et al., 2002) worm. Different purposes and motivations stand behind worm developers (Weaver et al., 2003) including: *Experimental curiosity* (ILoveYou worm.; CERT, 2000); *pride and power* leading programmers to show off their knowledge and skill through the harm caused by the worm; *commercial advantage, extortion and criminal gain, random and political protest, and terrorism and cyber warfare*. The existence of all these types of motivation indicates that computer worms are here to stay as a network vehicle serving different purposes and implemented in different ways. To address the challenge posed by worms effectively, meaningful experience and knowledge should be extracted by analyzing known worms. Today, given the known worms, we have a great opportunity to learn from these examples in order to generalize. We argue that supervised learning methods can be very useful in learning and generalizing from previously encountered worms, in order to classify unknown worms effectively.

## 2.2 Detecting Malicious Code Using Supervised Learning Techniques

Supervised and unsupervised learning has already been used for detecting and protecting against malicious codes. A recent survey on intrusion detection systems (Kabiri and Ghorbani, 2005, Rokach and Elovici, 2007) summarizes recently proposed applications for recognizing malcodes in single computers and in computer networks. Lee et al. (1999). proposed a framework consisting of set of algorithms for the extraction of anomalies of user normal behavior for use in anomaly detection , in which a normal behavior is learned and any abnormal activity is considered as intrusive. The authors suggest several techniques, such as classification, meta-learning, association rules, and frequent episodes, to extract knowledge for implementation in intrusion detection systems, evaluating their approach on the DARPA98 (Lippmann et al., 1998) benchmark.

A Naïve Bayesian classifier was suggested in (Kabiri and Ghorbani, 2005), referring to its implementation within the ADAM system, developed by Barbara et al. (2001), which had three main parts: (a) a network data monitor listening to TCP/IP protocol; (b) a learning engine which enables acquisition of the association rules from the network data; and (c) a classification module which classifies the nature of the traffic in two possible classes, normal and abnormal, which can later be linked to specific attacks. Other soft computing algorithms were proposed for detecting malicious code: Artificial Neural Networks (ANN) (Zanero and Savaresi, 2004; Kayacik et al., 2003; Lei and Ghorbani; 2004). Self Organizing Maps (SOM) (Hu and Heywood, 2003) and fuzzy logic (Dickerson and Dickerson, 2000; Bridges and Vaughn Rayford, 2000; Botha and von Solms, 2003).

## 3 Methods

The goal of this study was to assess the feasibility of detecting *unknown* malicious code, in particular computer worms, based on the computer's behavior (measurements), using machine learning techniques, and the potential accuracy of such methods. In order to create the datasets we built an isolated local network of computers, simulating a real Internet network which allows worms to propagate. This setup enabled us to inject worms into a controlled environment, while monitoring the computer behavior. The monitoring is performed by an agent, developed specifically for this purpose, that measures various parameters and save their values in log files.

In this study we examine whether a classifier, trained on data collected from a computer having a certain hardware configuration and certain specific background activity, is capable to correctly classify the behavior of a computer having other configurations? In order to answer this question we designed several experiments. We created eight datasets having different configurations, different background applications, and different user activities. Another goal was to select the minimal subset of features which are required to correctly classify new cases. Reducing the number of features used in the model, implies that less monitoring efforts are needed

when the proposed approach is served as the basis for an operational system. Finally, we applied four classification algorithms on the given datasets in a varied series of experiments, starting with detecting known worms in different environments and later detecting completely new, previously unseen worms.

Figure 1 specifies the process that was used in order to perform this study. The upper part refers to the training phase. We collected a set of worms and used them to infect the hosts in the controlled environment. Then an agent, which was installed on each host, recorded the behavior of the host. Based on collected dataset, we trained the classifiers. The bottom part in Figure 1 refers to the test phase. In this phase we examine if the induced classifier can be used to identify the existence of unknown worm.

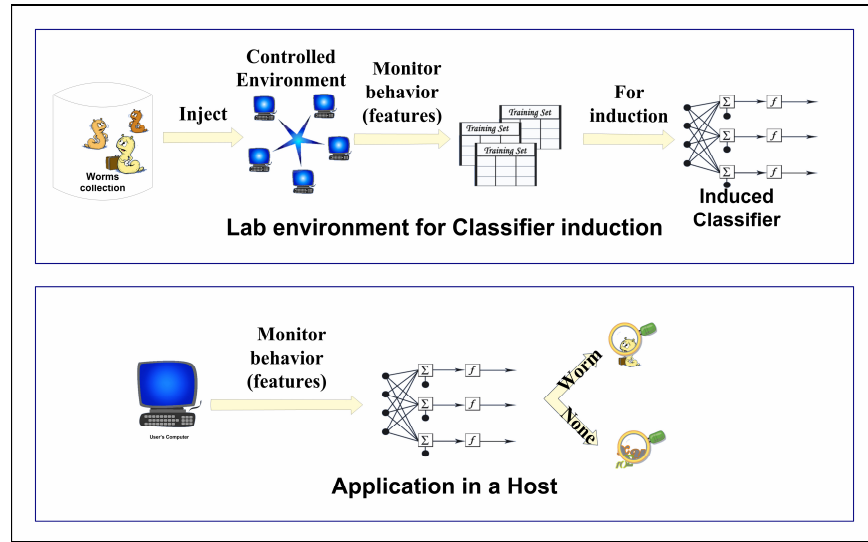


Figure 1: Outline of the Train phase and the Test Phase

### 3.1 Dataset Creation

Since there is no benchmark dataset which could be used for this study, we created our own dataset. A network with various computers (configurations) was deployed, enabling us to inject worms, and monitor the computer behavior and log the measurements.

### Environment Description

The lab network consisted of seven computers, which contained heterogenic hardware, and a server computer simulating the internet. We used the *windows*

*performance counters*<sup>1</sup>, which enable monitoring system features that appear in these main categories (including the number of features in parenthesis): *Internet Control Message Protocol* (27), *Internet Protocol* (17), *Memory* (29), *Network Interface* (17), *Physical Disk* (21), *Process* (27), *Processor* (15), *System* (17), *Transport Control Protocol* (9), *Thread* (12), and *User Datagram Protocol* (5). In addition we used *VTrace* (Lorch and Smith, 2000), a software tool which can be installed on a PC running Windows for monitoring purposes. *VTrace* collects traces of the *file system*, *the network*, *the disk drive*, *processes*, *threads*, *interprocess communication*, *waitable objects*, *cursor changes*, *windows*, and *the keyboard*. The data from the *windows performance counter* were configured to measure the features every second and store them in a log file as vectors. *VTrace* stored time-stamped events, which were aggregated into the same fixed intervals, and merged with the *windows performance* log files. These eventually included a vector of 323 features for every second.

## Injected Worms

While selecting worms from the wild, our goal was to choose worms that differ in their behavior, from among the available worms. Some of the worms have a heavy payload of Trojans to install in parallel to the distribution process upon the network; others focus only on distribution. Another aspect is that they have different strategies for IP scanning which results in varying communication behavior, CPU consumption, and network usage. While all the worms are different, we wanted to find common characteristics by the presence of which it would be possible to detect an unknown worm. We briefly describe here the main characteristics, relevant to this study, of each worm included in this study. The information is based on the virus libraries on the web<sup>2,3,4</sup>. We briefly describe the five worms we used:

(1) **W32.Dabber.A** scans IP addresses randomly. It uses the W32.Sasser.D worm to propagate and opens the FTP server to upload itself to the victim computer. Registering itself enables its execution on the next user login (human based activation). It drops a backdoor, which listens on a predefined port. This worm is distinguished by its use of an external worm in order to propagate.

(2) **W32.Deborm.Y** is a self-carried worm, which prefers local IP addresses. It registers itself as an MS Windows service and is executed upon user login (human based activation). This worm contains three Trojans as a payload: Backdoor.Sdbot, Backdoor.Litmus, and Trojan.KillAV, and executes them all. We chose this worm because of its heavy payload.

(3) **W32.Korgo.X** is a self carrying worm which uses a totally random method for IP addresses scanning. It is self-activated and tries to inject itself as a function to MS Internet Explorer as a new thread. It contains a payload code which enables it to connect to predefined websites in order to receive orders or download newer worm versions.

---

<sup>1</sup>[http://msdn.microsoft.com/library/default.asp?url=/library/en-us/counter/counters2\\_lbfc.asp](http://msdn.microsoft.com/library/default.asp?url=/library/en-us/counter/counters2_lbfc.asp)

<sup>2</sup> Symantec – [www.symantec.com](http://www.symantec.com)

<sup>3</sup> Kasparsky [www.viruslist.com](http://www.viruslist.com)

<sup>4</sup> Macfee <http://vil.nai.com>

(4) *W32.Sasser.D* uses a preference for local addresses optimization while scanning the network. About half the time it scans local addresses, and the other half random addresses. In particular it opens 128 threads for scanning the network, which requires a heavy CPU consumption, as well as significant network traffic. It is a self-carried worm and uses a shell to connect to the infected computer's FTP server and to upload itself.

(5) *W32.Slackor.A*, a self-carried worm, exploits MS Windows sharing vulnerability to propagate. The worm registers itself to be executed upon user login. It contains a Trojan payload and opens an IRC server on the infected computer in order to receive orders.

All the worms perform port scanning and possess different characteristics. Further information about these worms can be found on the web<sup>567</sup>.

### Dataset Description

In order to examine the influence of a computer hardware configuration, background running applications, and user activity, we considered three major aspects: *computer hardware configuration*, *constant background application* consuming extreme computational resources, and *user activity*, being binary variables. (1) *Computer hardware configuration*: Both computers ran on Windows XP, which considered the most widely used operating system, having two hardware configuration types: an "old," having Pentium 3 800Mhz CPU, bus speed 133Mhz and memory 512 Mb, and a "new," having Pentium 4 3Ghz CPU, bus speed 800Mhz and memory 1 Gb. (2) *Background application activity*: We ran an application affecting mainly the following features: *Processor object*, *Processor Time* (usage of 100%); *Page Faults/sec*; *Physical Disk object*, *Avg Disk Bytes/Transfer*, *Avg Disk Bytes/Write*, and *Disk Writes/sec*. (3) *User activity*: several applications, including browsing, downloading and streaming operations through Internet Explorer, Word, Excel, chat through MSN messenger, and Windows Media Player, were executed to imitate user activity in a scheduled order. Appendix A specifies the set of features that was examined in this research.

We created eight datasets (see table I). Each dataset contained monitored samples of each one of the five injected worms separately, and samples of a *normal* computer behavior, without any injected worm. Each worm was monitored for a period of 20 minutes. We collected the values of the features every second. Thus, each record, containing a vector of measurements and a label, presented an activity along a second labeled by a specific *worm*, or a *none* activity label. Each dataset contained a few thousand (labeled) samples of each worm or clean computer. We therefore had three binary aspects, which resulted in eight possible combinations, shown in Table 1, representing a variety of dynamic computer configurations and usage patterns. Each dataset contained monitored samples for each of the five worms injected separately,

---

<sup>5</sup> Symantec – [www.symantec.com](http://www.symantec.com)

<sup>6</sup> Kaspersky [www.viruslist.com](http://www.viruslist.com)

<sup>7</sup> Macfee <http://vil.nai.com>

and samples of a normal computer behavior without any injected worm. Each sample (record) was labeled with the relevant worm (class), or 'none' for clean samples.

**Table 1. The three aspects resulting in eight datasets, representing a variety of situations of a monitored computer.**

Computer	Background Application	User Activity	Dataset Name
Old	No	No	O
Old	No	Yes	Ou
Old	Yes	No	Oa
Old	Yes	Yes	Oau
New	No	No	N
New	No	Yes	Nu
New	Yes	No	Na
New	Yes	Yes	Nau

### 3.2 Feature Selection Methods

In many applications, the large number of features in many domains presents a huge challenge. Typically, some of the features do not contribute to the accuracy of the classification task and may even hamper it. Feature selection is the process of identifying relevant features in the dataset and discarding everything else as irrelevant and redundant. Since feature selection reduces the dimensionality of the data, it enables the classification algorithms to operate more effectively and rapidly. In some cases, classification performance can be improved; in other instances, the obtained classifier is more compact and can be easily interpreted. In host-based detection applications there is an additional motivation. Ideally, we would like to minimize the self-consumption of computer resources required for the monitoring operations (measurements), i.e. minimizing the collection of the features.

In order to compare the performance of the various classification algorithms, we used the *filters* approach, which is applied on the dataset and is independent of any classification algorithm, in which a measure is calculated to quantify the correlation of each feature with the class (the presence or absence of worm activity). Each feature receives a rank which represents its expected contribution in the classification task.

#### 3.2.1 Feature Selection Methods.

We used three feature-selection methods, which resulted in a list of ranked features for each feature-selection method and an ensemble incorporating all three of them. We used Chi-Square (CS), Gain Ratio (GR) and ReliefF implemented in the WEKA environment (Witten and Frank, 2005) and their ensemble.



### Chi-Square

*Chi-Square* measures the lack of independence between a feature  $f$  and a class  $c_i$  (such as W32.Dabber.A) and can be compared to the chi-square distribution with one degree of freedom to judge extremeness. Equation 1 shows how the chi-square measure is defined and computed, where  $N$  is the total number of documents and  $f$  refers to the presence of the feature (and  $\bar{f}$  its absence), and  $c_i$  refers to its membership in  $c_i$ .

$$\chi^2(f, c_i) = \frac{N[P(f, c_i)P(\bar{f}, \bar{c}_i) - P(f, \bar{c}_i)P(\bar{f}, c_i)]^2}{P(f)P(\bar{f})P(c_i)P(\bar{c}_i)} \quad (1)$$

### Gain Ratio

*Gain Ratio* was originally presented by Quinlan in the context of *Decision Trees* (Mitchell, 1997), which was designed to overcome a bias in the *Information Gain* (*IG*) measure, and which measures the expected reduction of entropy caused by partitioning the examples according to a chosen feature. Given entropy  $E(S)$  as a measure of the impurity in a collection of items, it is possible to quantify the effectiveness of a feature in classifying the training data. Equation 3 presents the formula of the entropy of a set of items  $S$ , based on  $C$  subsets of  $S$  (for example, classes of the items), presented by  $S_c$ . *Information Gain* measures the expected reduction of entropy caused by partitioning the examples according to attribute  $A$ , in which  $V$  is the set of possible values of  $A$ , as shown in Equation 2. These equations refer to discrete values; however, it is possible to extend them to continuous values attribute.

$$IG(S, A) = E(S) - \sum_{v \in V(A)} \frac{|S_v|}{|S|} \cdot E(S_v) \quad (2)$$

$$E(S) = \sum_{c \in C} -\frac{|S_c|}{|S|} \cdot \log_2 \frac{|S_c|}{|S|} \quad (3)$$

The *IG* measure favors features having a high variety of values over those with only a few. *GR* overcomes this problem by considering how the feature splits the data (Equations 4 and 5).  $S_i$  are  $d$  subsets of examples resulting from partitioning  $S$  by the  $d$ -valued feature  $A$ .

$$GR(S, A) = \frac{IG(S, A)}{SI(S, A)} \quad (4)$$

$$SI(S, A) = -\sum_{i=1}^d \frac{|S_i|}{|S|} \cdot \log_2 \frac{|S_i|}{|S|} \quad (5)$$

## Relief

*ReslifF* (Pearl, 1986) estimates the quality of the features according to how well their values distinguish between instances that are near each other. Given a randomly selected instance  $x$ , from a dataset  $s$  with  $k$  features, Relief searches the dataset for its two nearest neighbors from the same class, called nearest hit  $H$ , and from a different class, called nearest miss  $M$ . The quality estimation  $W[A_i]$  is stored in a vector of the features  $A_i$ , based on the values of a difference function  $diff()$  given  $x$ ,  $H$  and  $M$  as shown in Equation 6.

$$diff(A_i, x_{1i}, x_{2i}) = \begin{cases} |x_{1i} - x_{2i}| & \text{if } A_i \text{ is numeric,} \\ 0 & \text{if } A_i \text{ is nominal \& } x_{1i} = x_{2i}, \\ 1 & \text{if } A_i \text{ is nominal \& } x_{1i} \neq x_{2i}, \end{cases} \quad (6)$$

## Features Ensembles

Instead of selecting features based on of the feature selection methods, one can use the ensemble strategy (see for instance Rokach et al., 2007) which combines the features subsets that are obtained from several features selection methods. Specifically, we combine several methods by averaging the features ranks as shown in Equation 7:

$$Rank(fi) = \frac{\sum_{j=1}^k rank^j(fi)}{k}. \quad (7)$$

where  $fi$  is a feature,  $filter$  is one of the  $k$  filtering (feature selection) methods. Specifically in our case  $k=3$ .

### 3.2.2 Consolidating Features from Different Environments: Averaged vs. Unified Consolidation.

Often when applying a feature selection method, such as the filters approach, the method is applied on the entire dataset aiming to rank the features based on their measured correlation to the class. However, unlike in the common datasets, our dataset consisted of eight datasets coming from different environments, as explained earlier (see Table 1). Since some features might be more important in specific environments and less in others it is not clear how this has to be considered. In this study we propose two approaches to considering and integrating the aspects of the datasets. In the first approach, termed *unified* dataset, we unified all the eight datasets into a single dataset and applied the filter on the unified dataset.

Alternatively, we examined the approach termed *averaged* in which we applied the filter on each one of the eight datasets and computed the average rank for each

feature. Note that for averaging the features ranks obtained from the different environments, we used Equation 7 again.

Figure 2 illustrates both approaches, in which the top refers to the unified approach, where the feature selection (FS) is applied on the unified dataset 'all' and the averaged approach, at the bottom, in which the feature selection is applied on each dataset and averaged into a rank list.

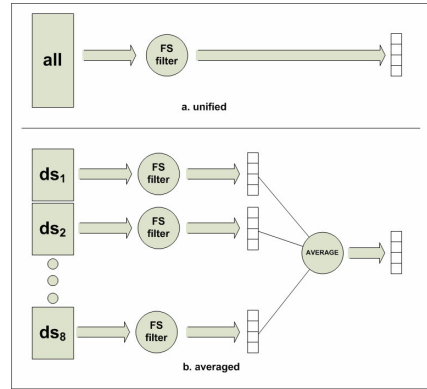


Figure 2 – Unified versus the averaged approach for environment features consolidation.

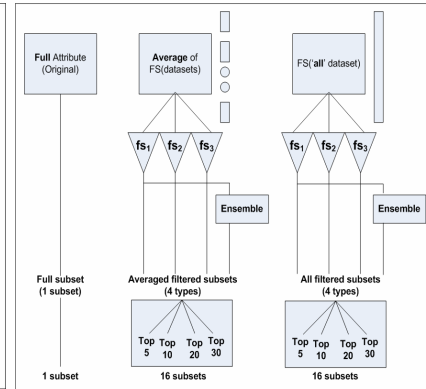


Figure 3 - The creation of 33 features sets.

After applying both approaches we extracted the top ranked features. We took the highest ranked (top) features 5, 10, 20 and 30 from the output of each feature selection method. Finally, we had four features sets (Top 5, 10, 20, 30) for each of the four filters (fs1, fs2, fs3, ensemble), for each feature consolidation (unified, averaged). On top of that we also examined the *full* features set (with no feature selection). This totally results with 33 features sets ( $4 \times 4 \times 2 + 1$ ) as shown in Figure 3.

### 3.3 Classification algorithms

One of the goals of this study was to pinpoint the classification algorithm that provides the highest level of detection accuracy. We employed four commonly used Machine Learning algorithms: *Decision Trees*, *Naïve Bayes*, *Bayesian Networks* and *Artificial Neural Networks*, in a *supervised learning* approach, in which the classification algorithm learns from a provided training set, containing labeled examples.

While the focus of this paper is not on *classification algorithm* techniques, but on their application in the task of detecting worm activity, we briefly describe the classification algorithms we used in this study.

### **Decision Trees**

Decision tree learners (Quinlan, 1993) are a well-established family of learning algorithms. Classifiers are represented as trees whose internal nodes are tests on individual features, and leaves are classification decisions. Typically, a greedy top-down search method is used to find a small decision tree that correctly classifies the training data. The decision tree is induced from the dataset by splitting the variables based on the *expected information gain*. Modern implementations include pruning, which avoids over-fitting. In this study we evaluated J48, the WEKA version of the commonly used C4.5 algorithm (Quinlan, 1993). An important characteristic of Decision Trees is the explicit form of their knowledge which can be represented as a set of if-then rules. This set of rules can be then easily embedded in any existing IDS.

### **Naïve Bayes**

The Naïve Bayes classifier is based on the *Bayes theorem*, which in the context of classification states that the posterior probability of a class is proportional to its prior probability as well as to the conditional likelihood of the features, given this class. If no independent assumptions are made, a Bayesian algorithm must estimate conditional probabilities for an exponential number of feature combinations. Naïve Bayes simplifies this process by making the assumption that features are *conditionally independent* given the class, and requires that only a linear number of parameters be estimated. The prior probability of each class and the probability of each feature, given each class, is easily estimated from the training data and used to determine the posterior probability of each class, given a set of features. Naïve Bayes has been shown empirically to produce good classification accuracy across a variety of problem domains. In this study, we evaluated Naïve Bayes, the standard version that comes with WEKA.

### **Bayesian Networks**

Bayesian networks are a form of the probabilistic graphical model (Pearl, 1986). Specifically, a Bayesian network is a directed acyclic graph of nodes with variables and arcs representing dependence among the variables. Like Naïve Bayes, Bayesian networks are based on the Bayes Theorem; however, unlike Naïve Bayes they do not assume that the variables are independent. Actually Bayesian Networks are known for their ability to represent conditional probabilities, which are the relations between variables. A Bayesian network can thus be considered a mechanism for automatically constructing extensions of Bayes Theorem to more complex problems. Bayesian

networks were used for modeling knowledge and implemented successfully in different domains. We evaluated the Bayesian Network standard version which comes with WEKA.

### Artificial Neural Networks

An Artificial Neural Network (ANN) (Bishop, 1995) is an information processing paradigm that is inspired by the way biological nervous systems (i.e., the brain) are modeled with regard to information processing. The key element of this paradigm is the structure of the information processing system. It is a network composed of a large number of highly interconnected processing elements, called neurons, working together in order to approximate a specific function. An ANN is configured for a specific application, such as pattern recognition or data classification, through a *learning process* during which the weights of the inputs in each neuron are updated. The weights are updated by a *training algorithm*, such as back-propagation, according to the examples the network receives, in order to reduce the value of *error function*. The power and usefulness of ANN have been demonstrated in numerous applications including speech synthesis, medicine, finance, and many other pattern recognition problems. For some application domains, neural models show more promise in achieving human-like performance than do more traditional artificial intelligence techniques. All ANN manipulations in this study have been performed within a MATLAB(r) environment using Neural Network Toolbox (Demuth and Beale, 1998).

## 4 Experimental Design

Our main goal in this study was to investigate whether the approach presented here, in which *unknown* malicious code is detected, based on the computer behavior (measurements), is feasible and enables a high level of accuracy when applied to a variety of computers. We defined four research questions accordingly:

**Q<sub>1</sub>:** In the detection of *known* malicious code, based on a computer's measurements, using machine learning techniques, what is the achievable level of accuracy?

**Q<sub>2</sub>:** Is it possible to reduce the amount of features to below 30, while maintaining a high level of accuracy (compared to the full set of features). Which feature consolidation approach (unified versus averaged) and feature selection method is superior?

**Q<sub>3</sub>:** Will the *computer configuration* and the *computer background activity*, from which the training sets were taken, have a significant influence on the detection accuracy?

**Q<sub>4</sub>:** Is the detection of *unknown* worms possible, based on a training set of known worms?

In addition to these research questions, we wanted to identify the best classification algorithms and the best combination of top ranked features and classification

algorithm. We start with the definition of the evaluation measures and continue with the experiments we designed for this study.

#### 4.1 Evaluation Measures

For evaluation purposes, we measured the *True Positive Rate (TPR)* measure, which is the number of *positive* instances classified correctly, as shown in Equation 8, *False Positive Rate (FPR)*, which is the number of *negative* instances misclassified (Equation 8), and the *Total Accuracy*, which measures the number of absolutely correctly classified instances, either positive or negative, divided by the entire number of instances shown in Equation 9. Additionally, we calculated the ROC curves, but we don't present them because of lack of room.

$$TPR = \frac{|TP|}{|TP| + |FN|}; \quad FPR = \frac{|FP|}{|FP| + |TN|} \quad (8)$$

$$Total\ Accuracy = \frac{|TP| + |TN|}{|TP| + |FP| + |TN| + |FN|}. \quad (9)$$

We also measured a *confusion matrix*, which depicts the number of instances from each class which were classified in each one of the classes (ideally all the instances would be in their actual class).

In the first part of the study, we wanted to identify the best feature consolidation approach (unified or averaged) and feature selection method, the best classification algorithm and the minimal features required to maintain a high level of accuracy. In the second part we wanted to measure the capability of classifying unknown worms based on a training set of known worms. In order to answer these questions we designed two experimental plans, based on 33 datasets, as will be described later. After evaluating all the classification algorithms on the 33 datasets, we selected the best feature selection and the top features to evaluate the unknown worms detection.

#### 4.2 Experiment I

To determine the best combination of feature selection method, number of features, and classification algorithm, we performed a wide set of experiments, in which we evaluated all the combinations of feature selection method, classification algorithm, and number of top features.

In this experiment, called *e1*, we trained each classifier on a single dataset *i* and tested on each one (*j*) of the *eight* datasets. Thus, we had a set of eight iterations in which a dataset was used for training, and eight corresponding evaluations which were done on each one of the datasets, resulting in 64 evaluation runs. When  $i = j$ , we used *10 fold cross validation*, in which the dataset is randomly partitioned into ten partitions and repeatedly the classifier is trained on nine partitions and tested on the tenth. Each evaluation run (out of the 64) was repeated for each one of the combinations of feature selection method, classification algorithm, and number of top

features. Thus, each evaluation run was repeated for the 33 features set described earlier in Figure 3 (in each repetition different features are extracted from the datasets). Note that the task was to classify specifically the exact worm out of the five or a none (worm) activity, and not to generate a general binary classification of “worm” or a “none” activity, which was our final goal in the context of an unknown worm detection. Such conditions, while being more challenging, were expected to bring more insights.

### 4.3 Experiment II

To estimate the potential of the suggested approach in classifying an *unknown worm* activity, which was the main objective of this study, we designed an additional experiment, called *e2*, in which we trained classifiers based on *part* of the (five) *worms* and the *none* activity, and tested on the *excluded worms* (from the training set) and the *none* activity, in order to measure the capability to detect an *unknown worm* and the *none* activity accurately.

In this experiment the training set consisted of  $5-k$  worms and the testing set contained the  $k$  excluded worms, while the *none* activity appeared in both datasets. This process repeated for all the possible combinations of the  $k$  worms ( $k = 1$  to  $4$ ). In each combination a classifier was trained on the training set and tested on all the remaining seven datasets. The test set included *only* the excluded worms and not the worms presented in the training set since we wanted to measure specifically the detection rate of the unknown. Note that in these experiments, unlike in *e1*, there were two classes: (generally) *worm*, for any type of worm, and *none* activity. This experiment was evaluated on each classification algorithm, using the outperforming top selected features found in *e1*.

## 5 Results

### Experiment I

Our objective in *e1* was to determine the best: feature selection approach, feature selection method, number of top features, and classification algorithms. We ran 132 (four classification algorithms applied to 33 features sets) evaluations (each comprises 64 runs), summing up to 8448 evaluation runs.

Figure 4 shows the mean accuracy (of all the classification algorithms) achieved for each environments features consolidation (unified or averaged), each feature selection method, top 5, 10, 20, 30 features and for the full set of features as a baseline. While the feature selection method and number of top features aren’t relevant for the FULL features set curve (blue lines in Figure 4) we presented the curve for comparison purposes. In general, all the feature subsets having less than 30 features achieved a mean performance quite similar to the full set of features (including 323 features). The *unified* consolidation approach outperforms the

*averaged* consolidation approach for most of the cases, especially when the Gain Ratio feature selection method is used. Additionally, unlike the averaged consolidation approach, which in most of the cases is below the full set performance, the unified consolidation approach for most of the cases is above it. Additionally, the *Top20* features delivered the best in most of the cases.

Based on the mean accuracy of the four classification algorithms GainRatio feature selection method outperformed the other feature selection methods for most of the top number of features, while the ensemble feature selection method outperformed for *Top5*. Unlike the independent measures, in which there was a monotonic growth when features were added, in the ensemble a monotonic slight decrease was observed as more features were used. The *Top20* features outperformed in general (by averaging) and when using GainRatio feature selection method in particular.

Figure 5 shows the same results, but presents the mean accuracy of the classification algorithms for several numbers of top features. *Bayesian Networks* outperforms for any number of top features, and on average the 20 top features *Top20* outperformed the other number of top features.

For example, *Top5* features for *GainRatio* feature selection method included: in the category of *ICMP*: (1) *Sent\_Echo\_sec* – the rate of ICMP Echo messages sent; (2) *Messages Sent/sec* – the rate, in incidents per second, at which the server attempted to send. The rate includes those messages sent in error; (3) *Messages/sec* – the total rate, in incidents per second, at which ICMP messages were sent and received by the target entity. The rate includes messages received or sent in error. In the category of *TCP*: (4) *Connections Passive* – the number of times TCP connections made a direct transition to the SYN-RCVD state from the LISTEN state; (5) *Connection Failures* – the number of times TCP connections made a direct transition to the CLOSED state from the SYN-SENT state or the SYN-RCVD state, plus the number of times TCP connections made a direct transition to the LISTEN state from the SYN-RCVD state. The list of the top twenty ranked features is presented in the Appendix B.

Based on the results achieved in *e1*, in which the unified consolidation approach, *Top20* features, *GainRatio* feature selection method, outperformed on average, we used only this features subset in the second part of the evaluation (*e2*).



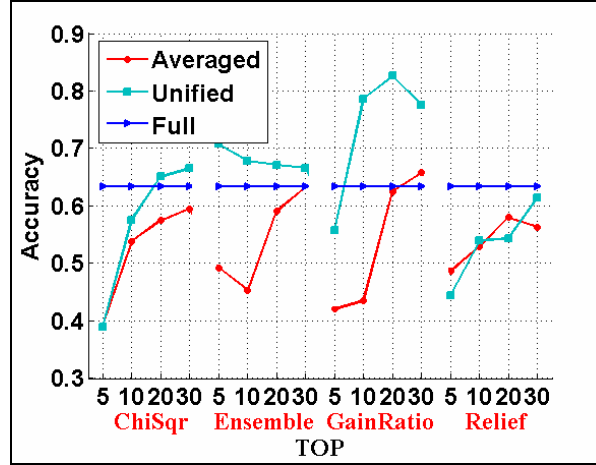


Figure 4 – The mean accuracy for various feature selection methods and the number of top features

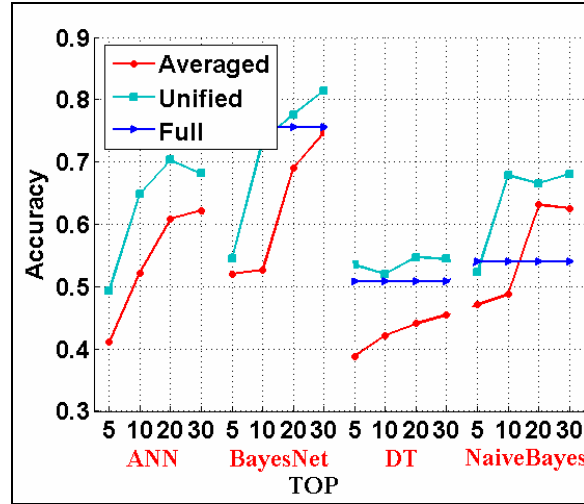


Figure 5 – The mean accuracy for various classification algorithms and number of top features

In  $Q_3$  we wanted to estimate the performance sensitivity of the suggested approach given several training sets sampled from a variety of computers, represented by the eight datasets. Thus, we tested whether the accuracy obtained by training a classifier on a training set sampled from a given computer will vary significantly when evaluated on a variety of test sets. To perform this test we designed two experiments ( $e1_1, e1_2$ ).

Table 2 presents the mean accuracy and the resulting standard deviation in each experiment, including 128 evaluation runs. In the first main column the results of  $e1_1$

are presented in two columns, when the dataset  $di$  was used as a training set (left column), and when used as test sets (right column). The results of  $eI_1$  when  $di$  was the training set were not homogenous. The results of  $eI_1$  when  $di$  was the testset were not homogenous too, though we found that the training on datasets created in the 'old' computer was significantly better ( $\alpha = 0.01$ ). In  $eI_2$ , in which all datasets except  $di$  were training set and  $di$  was the testset, the results were statistically significant ( $\alpha = 0.05$ ) homogenous. Note that  $eI_1$  and  $eI_2$  experiments are based on all the learning algorithms, and not specifically on each algorithm. The classification accuracy in  $eI_2$  outperformed the accuracy in  $eI_1$ , since the training sets in  $eI_2$  included several datasets.

**Table 2. The results achieved in  $eI_1$  and  $eI_2$ .**

Experiment	$eI_1$		$eI_2$
Dataset $di$	<i>Training: <math>di</math></i>	<i>Testset: <math>di</math></i>	<i>Training: all datasets except <math>di</math></i> <i>Testset: <math>di</math></i>
o	$0.68 \pm 0.23$	$0.73 \pm 0.23$	$0.78 \pm 0.22$
ou	$0.76 \pm 0.22$	$0.73 \pm 0.22$	$0.82 \pm 0.18$
oa	$0.73 \pm 0.21$	$0.73 \pm 0.23$	$0.81 \pm 0.18$
oau	$0.77 \pm 0.21$	$0.72 \pm 0.21$	$0.81 \pm 0.19$
n	$0.61 \pm 0.24$	$0.64 \pm 0.22$	$0.71 \pm 0.20$
nu	$0.76 \pm 0.21$	$0.72 \pm 0.22$	$0.82 \pm 0.22$
na	$0.70 \pm 0.21$	$0.73 \pm 0.24$	$0.86 \pm 0.17$
nau	$0.73 \pm 0.22$	$0.71 \pm 0.23$	$0.79 \pm 0.20$
Average	$0.72 \pm 0.22$	$0.72 \pm 0.22$	$0.80 \pm 0.19$

## Experiment II

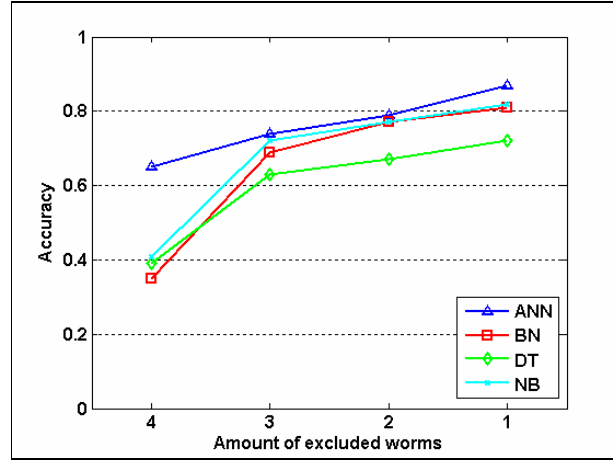
In  $Q_4$  we wanted to estimate the possibility of classifying an unknown worm when training on data collected from a single computer. In this set of experiments we used only the *Top20* features, which outperformed in  $eI$ . The training set included four worms out of the five and the none activity samples, and the test set included the excluded worm and the none activity samples. This process was done for each worm repeating in five iterations. Note that in these experiments, unlike in  $eI$ , in which each worm class was defined separately, there were two classes: (generally) *worm* and *none* activity.

Table 3 presents the results of  $e2$ . On average the *Decision Trees* and *Bayesian Networks* outperformed the others classification algorithms. The table shows also the *true positive* (TP) and *false positive* (FP). *Decision Trees* and *Bayesian Networks* achieved high level of accuracy and maintained a low *false positive* rate.

**Table 3 - The results of  $e2$ . There is a difference in the detection accuracy of each classifier for each type of worm. On average, Decision Trees outperformed the other classifiers, while maintaining a low false positive rate.**

	<i>ANN 20</i>			<i>BN 20</i>			<i>DT 20</i>			<i>NB 20</i>		
Worm	Acc	TP	FP	Acc	TP	FP	Acc	TP	FP	Acc	TP	FP
1	0.985	0.985	0.014	0.554	0.557	0.443	0.936	0.937	0.063	0.497	0.500	0.499
2	0.494	0.499	0.500	0.992	0.992	0.007	0.999	0.999	0.0005	0.997	0.997	0.002
3	0.952	0.952	0.047	0.992	0.993	0.007	0.680	0.678	0.3215	0.844	0.843	0.156
4	0.994	0.994	0.005	0.998	0.998	0.002	0.968	0.968	0.032	0.998	0.998	0.002
5	0.636	0.637	0.362	0.990	0.991	0.008	0.999	0.999	0.0005	0.975	0.974	0.026
Average	0.81	0.81	0.19	0.91	0.91	0.09	0.92	0.92	0.08	0.86	0.86	0.14
StdDev	0.05	0.05	0.05	0.04	0.04	0.04	0.02	0.02	0.02	0.05	0.05	0.05

Figure 6 presents the results of  $e2$ , in which a monotonic increase in the accuracy is shown, as more worms are included in the training set. Note that the number of worms in the  $x$  axis refers to the number of worms excluded from the training set, and were included in the test set. In general the *ANN* outperformed all the other algorithms, while the *BN* kept on showing very good results. Note that testing on the seven datasets separately decreased the mean accuracy slightly. In addition, when only one worm was excluded, in specific worms we achieved 99% accuracy and a very low false positive rate of 0.005.



**Figure 5. The performance monotonically increases as fewer worms are excluded from the training set**

## 6. Discussion and Conclusions

In this paper we explored the feasibility of detecting unknown worm activity in individual computers, at a high level of accuracy, given the variation in hardware and software environmental conditions, while minimizing the set of features collected from the monitored

computer. Four research questions were investigated, referring to the feasibility of the approach, the best settings, and the level of achieved accuracy, for which a dataset was created and several corresponding experiments were designed. In the first experiment we showed that the detection of *known* worms is feasible at a very high level of accuracy. To reduce the computational resources in the classification task we wanted to reduce the number of features. Two consolidation approaches to integrating the eight datasets for the task of features selection were proposed: unified, in which all the datasets were unified into a single dataset, and averaged, in which we first applied the feature selection method on each dataset and averaged the ranked features into a single rank. Our results showed that the mean performance of the unified approach outperformed the averaged approach. Based on the evaluation results, in general *Bayesian Networks* outperformed the other algorithms; and using the *Top 20* ranked features from the *GainRatio* was the best. The reduction in the amount of features and the improvement in accuracy compared well to the baseline of above 300 features in the full set, since it reduces the computer's resources consumption needed for monitoring its behavior. We investigated the influence of the variance in the training phase and detection phase on the configuration of a computer and its programs, to determine whether this method can be generalized. We found that training on seven unified datasets was significantly homogenous ( $\alpha=0.01$ ) based on a homogeneity test, unlike training on a single dataset (as in  $e1_1$ ). This is a very encouraging result, since we assume that, when applying such an approach in the real world, a training set that consists of samples from several types of computer activity in several environments is a reasonable requirement.

To examine the possibility of classifying *unknown* worms, unlike in previous experiments, two classes were defined in the dataset, a worm type consisting of the worms' samples and 'none' type. The training sets had four worms and the 'none' activity and the test set consisted only of the excluded *worm* and the *none-activity*. We found that the level of detection accuracy for each worm varies from algorithm to algorithm. Finally, in  $e2$  above 85% accuracy was achieved in general; Decision Trees achieved 92%, while specific algorithms exceeded the 95% level of accuracy for specific worms. We noticed that the detection of each worm varied within each algorithm, while being different among algorithms, and thus we suggest using an ensemble of classifiers to achieve a higher level of accuracy for instances of all potential worm classes. In general *Bayesian Networks* resulted constantly in very good results, which might be explained by the consideration of the dependency within features, unlike other classifiers. Later we reduced the amount of worms in the training set and increased the amount of unknown worms in the test set. We found an increase in accuracy as more worms were presented in the training set.

The limitations of this study are the number of worms and the variety of computer configurations. Note that the worms were selected to provide a reasonable variety and the computers which were used were dramatically different. However, this was enough to achieve statistically significant results.

To conclude, we have shown that it is possible to detect previously un-encountered computer worms using our novel approach, which is based on monitoring the computer "behavior" (features). In order to attain a high level of accuracy in different types of computers, which is an essential requirement in real life, the training set

should include samples taken from several computers types (i.e., different configurations).

**Acknowledgments.** This work was supported by Deutsche Telekom Co. We would like to thank the undergraduate students Shai and Ido and Clint Feher who contributed in the preparations of the dataset.

## References

1. Abou-Assaleh, T., Cercone, N., Keselj, V., and Sweidan, R. (2004). N-gram Based Detection of New Malicious Code, Proceedings of the 28th Annual International Computer Software and Applications Conference (COMPSAC'04)
2. Barbara, D., Wu, N., Jajodia, S. (2001). Detecting Novel Network Intrusions using Bayes Estimators, in Proceedings of the First SIAM International Conference on Data Mining (SDM 2001), Chicago, USA
3. Bishop, C. (1995). Neural Networks for Pattern Recognition. Clarendon Press, Oxford.
4. Botha, M. and von Solms, R. (2003). Utilising Fuzzy Logic and Trend Analysis for Effective Intrusion Detection,” Computers & Security, vol. 22, no. 5, pp. 423–434, 2003.
5. Bridges, S.M. and Vaughn Rayford, M. (2000). Fuzzy Data Mining and Genetic Algorithms Applied to Intrusion Detection, in Proceedings of the Twenty-third National Information Systems Security Conference. National Institute of Standards and Technology, Oct. 2000.
6. Caia D. M., Gokhaleb M., Theiler J., Comparison of feature selection and classification algorithms in identifying malicious executables, Computational Statistics & Data Analysis 51 (2007) 3156 – 3172.
7. CERT. CERT Advisory CA-2000-04, Love Letter Worm, <http://www.cert.org/advisories/ca-2000-04.html>
8. Demuth, H. and Beale, M. (1998). Neural Network Toolbox for Use with Matlab. The Mathworks Inc., Natick, MA.
9. Dickerson, J.E. and Dickerson, J.A. (2000). Fuzzy Network Profiling for Intrusion Detection,” in Proceedings of NAFIPS 19th International Conference of the North American Fuzzy Information Processing Society, pp. 301–306, Atlanta, USA, July 2000.
10. Fosnock, C. (2005). Computer Worms: Past, Present and Future. East Carolina University
11. Hu, P. Z. and Heywood, M.I. (2003). Predicting Intrusions with Local Linear model, in Proceedings of the International Joint Conference on Neural Networks, vol. 3, pp. 1780–1785. IEEE, IEEE, July 2003.
12. Kabiri, P., Ghorbani, A.A. (2005). Research on Intrusion Detection and Response: A Survey. International Journal of Network Security, vol. 1(2), pp. 84-102.

13. Kayacik, H.G., Zincir-Heywood, A.N., and Heywood, M.I. (2003) On the Capability of a Som Based Intrusion Detection System, in Proceedings of the International Joint Conference on Neural Networks, vol. 3, pp. 1808–1813. IEEE, IEEE, July 2003.
14. Kienzle, D.M. and Elder, M.C. (2003). Recent Worms: a Survey and Trends. In Proceedings of the 2003 ACM Workshop on Rapid Malcode, pages 1--10. ACM Press, October 27, 2003.
15. Kolter, J.Z. and Maloof, M.A. (2006). Learning to Detect and Classify Malicious Executables in the Wild, *Journal of Machine Learning Research*, 7, 2721-2744.
16. Lee, W., Stolfo, S.J. and Mok, K.W. (1999). A Data Mining Framework for Building Intrusion Detection Models. In Proceedings of the 1999 IEEE Symposium on Security and Privacy, May 1999
17. Lei, J. Z. and Ghorbani, A. (2004). Network Intrusion Detection Using an Improved Competitive Learning Neural Network,” in Proceedings of the Second Annual Conference on Communication Networks and Services Research (CNSR04), pp. 190–197. IEEE-Computer Society, IEEE, May 2004.
18. Lippmann, R.P., Graf, I., Wyszogrod, D., Webster, S.E., Weber, D.J., and Gorton, S. (1998). The 1998 DARPA/AFRL Off-Line Intrusion Detection Evaluation, First International Workshop on Recent Advances in Intrusion Detection (RAID), Louvain-la-Neuve, Belgium, 1998.
19. Lorch, J. and Smith, A. J. (2000) The VTrace Tool: Building a System Tracer for Windows NT and Windows 2000. *MSDN Magazine*, 15(10):86–102, October 2000.
20. Mitchell T. (1997). *Machine Learning*, McGraw-Hill.
21. Moore, D., Shannon, C., and Brown, J. (2002). Code Red: a Case Study on the Spread and Victims of an Internet Worm, Proceedings of the Internet Measurement Workshop 2002, Marseille, France, November 2002.
22. Not USED: Domingos, P., and Pazzani, M. (1997). On the Optimality of Simple Bayesian Classifier under Zero-one Loss, *Machine Learning*, 29:103-130.
23. Not Used: Kohavi, R., (1995). A Study of Cross-Validation and Bootstrap for Accuracy Estimation and Model Selection, *International Joint Conference in Artificial Intelligence*, 1137-1145.
24. Pearl J., (1986). Fusion, Propagation, and Structuring in Belief Networks. *Artificial Intelligence* 29(3):241–288.
25. Quinlan, J.R. (1993). *C4.5: Programs for Machine Learning*. Morgan Kaufmann Publishers Inc., San Francisco, CA, USA.
26. Rokach L., Chizi B., Maimon O., A Methodology for Improving the Performance of Non-ranker Feature Selection Filters, *International Journal of Pattern Recognition and Artificial Intelligence* , 21, 5, 1-22, (2007).
27. Rokach L., Elovici Y., Data Mining for Intrusion Detection, *Encyclopaedia of Cyber terrorism and Cyber Warfare*, (Editors: Lech J. Janczewski and Andrew M. Colarik), Idea Group Reference, Winter 2007.

28. Schultz, M., Eskin, E., Zadok, E., and Stolfo, S. (2001). Data Mining Methods for Detection of New Malicious Executables, Proceedings of the IEEE Symposium on Security and Privacy, 2001, pp. 178--184.
29. Weaver, N. Paxson, V. Staniford, and S. Cunningham, R. (2003). A Taxonomy of Computer Worms, Proceedings of the 2003 ACM workshop on Rapid Malcode, Washington, DC, October 2003, pages 11-18
30. Witten, I.H. and Frank E., (2005). Data Mining: Practical Machine Learning Tools and Techniques, 2nd Edition, Morgan Kaufmann, San Francisco, 2005.
31. Zanero. S. and Savaresi, S.M. (2004) Unsupervised Learning Techniques for an Intrusion Detection System," in Proceedings of the 2004 ACM Symposium on Applied Computing, pp. 412-419, Nicosia, Cyprus, Mar. 2004. ACM Press.

## Appendix A - Operating System Measurements

The following table includes all features mapping in the data set. For further information about the objects and their meaning in the windows counters tools, please refer to [http://msdn.microsoft.com/library/default.asp?url=/library/en-us/counter/counters2\\_lbf.asp](http://msdn.microsoft.com/library/default.asp?url=/library/en-us/counter/counters2_lbf.asp), for more information about VTrace please refer to <http://msdn.microsoft.com/msdnmag/issues/1000/VTrace/>.

ID	Feature Name
1	<u>A_1ICMPMessages_sec</u>
2	<u>A_1ICMPMessages_Received_sec</u>
3	<u>A_1ICMPMessages_Received_Errors</u>
4	<u>A_1ICMPReceived_Dest_Unreachable</u>
5	<u>A_1ICMPReceived_Time_Exceeded</u>
6	<u>A_1ICMPReceived_Parameter_Problem</u>
7	<u>A_1ICMPReceived_Source_Quench</u>
8	<u>A_1ICMPReceived_Redirect_sec</u>
9	<u>A_1ICMPReceived_Echo_sec</u>
10	<u>A_1ICMPReceived_Echo_Reply_sec</u>
11	<u>A_1ICMPReceived_Timestamp_sec</u>
12	<u>A_1ICMPReceived_Timestamp_Reply_sec</u>
13	<u>A_1ICMPReceived_Address_Mask</u>
14	<u>A_1ICMPReceived_Address_Mask_Reply</u>
15	<u>A_1ICMPMessages_Sent_sec</u>
16	<u>A_1ICMPMessages_Outbound_Errors</u>
17	<u>A_1ICMPSent_Destination_Unreachable</u>
18	<u>A_1ICMPSent_Time_Exceeded</u>
19	<u>A_1ICMPSent_Parameter_Problem</u>
20	<u>A_1ICMPSent_Source_Quench</u>
21	<u>A_1ICMPSent_Redirect_sec</u>
22	<u>A_1ICMPSent_Echo_sec</u>
23	<u>A_1ICMPSent_Echo_Reply_sec</u>
24	<u>A_1ICMPSent_Timestamp_sec</u>
25	<u>A_1ICMPSent_Timestamp_Reply_sec</u>
26	<u>A_1ICMPSent_Address_Mask</u>
27	<u>A_1ICMPSent_Address_Mask_Reply</u>
28	<u>A_1IPDatagrams_sec</u>
29	<u>A_1IPDatagrams_Received_sec</u>
30	<u>A_1IPDatagrams_Received_Header_Errors</u>



31	A 1IPDatagrams_Received_Address_Errors_
32	A 1IPDatagrams_Forwarded_sec_
33	A 1IPDatagrams_Received_Unknown_Protocol_
34	A 1IPDatagrams_Received_Discarded_
35	A 1IPDatagrams_Received_Delivered_sec_
36	A 1IPDatagrams_Sent_sec_
37	A 1IPDatagrams_Outbound_Discarded_
38	A 1IPDatagrams_Outbound_No_Route_
39	A 1IPFragments_Received_sec_
40	A 1IPFragments_Re_assembled_sec_
41	A 1IPFragment_Re_assembly_Failures_
42	A 1IPFragmented_Datagrams_sec_
43	A 1IPFragmentation_Failures_
44	A 1IPFragments_Created_sec_
45	A 1MemoryPage_Faults_sec_
46	A 1MemoryAvailable_Bytes_
47	A 1MemoryCommitted_Bytes_
48	A 1MemoryCommit_Limit_
49	A 1MemoryWrite_Copies_sec_
50	A 1MemoryTransition_Faults_sec_
51	A 1MemoryCache_Faults_sec_
52	A 1MemoryDemand_Zero_Faults_sec_
53	A 1MemoryPages_sec_
54	A 1MemoryPages_Input_sec_
55	A 1MemoryPage_Reads_sec_
56	A 1MemoryPages_Output_sec_
57	A 1MemoryPool_Paged_Bytes_
58	A 1MemoryPool_Nonpaged_Bytes_
59	A 1MemoryPage_Writes_sec_
60	A 1MemoryPool_Paged_Allocs_
61	A 1MemoryPool_Nonpaged_Allocs_
62	A 1MemoryFree_System_Page_Table_Entries_
63	A 1MemoryCache_Bytes_
64	A 1MemoryCache_Bytes_Peak_
65	A 1MemoryPool_Paged_Resident_Bytes_
66	A 1MemorySystem_Code_Total_Bytes_
67	A 1MemorySystem_Code_Resident_Bytes_
68	A 1MemorySystem_Driver_Total_Bytes_
69	A 1MemorySystem_Driver_Resident_Bytes_
70	A 1MemorySystem_Cache_Resident_Bytes_
71	A 1Memory__Committed_Bytes_In_Use_

72	A_1MemoryAvailable_KBytes
73	A_1MemoryAvailable_MBytes
74	A_1Network_Interface_Packet_Scheduler_Miniport_Bytes_Total_sec
75	A_1Network_InterfaceTX_Packet_Scheduler_Miniport_Packets_sec
76	A_1Network_InterfacePacket_Scheduler_Miniport_Packets_Received_sec
77	A_1Network_Interface_Packet_Scheduler_Miniport_Packets_Sent_sec
78	A_1Network_InterfacePacket_Scheduler_Miniport_Current_Bandwidth
79	A_1Network_Interfacepacket_Scheduler_Miniport_Bytes_Received_sec
80	A_1Network_Interfaceduler_Miniport_Packets_Received_Unicast_sec
81	A_1Network_Interfacer_Miniport_Packets_Received_Non_Unicast_sec
82	A_1Network_Interfaceduler_Miniport_Packets_Received_Discarded
83	A_1Network_Interface_Scheduler_Miniport_Packets_Received_Errors
84	A_1Network_InterfaceScheduler_Miniport_Packets_Received_Unknown
85	A_1Network_Interface_Packet_Scheduler_Miniport_Bytes_Sent_sec
86	A_1Network_InterfaceScheduler_Miniport_Packets_Sent_Unicast_sec
87	A_1Network_Interfaceduler_Miniport_Packets_Sent_Non_Unicast_sec
88	A_1Network_Interfaceduler_Miniport_Packets_Outbound_Discarded
89	A_1Network_Interface_Scheduler_Miniport_Packets_Outbound_Errors
90	A_1Network_Interfacepacket_Scheduler_Miniport_Output_Queue_Length
91	A_1Network_Interface_MS_TCP_Loopback_interface_Bytes_Total_sec
92	A_1Network_Interface_MS_TCP_Loopback_interface_Packets_sec
93	A_1Network_InterfaceTCP_Loopback_interface_Packets_Received_sec
94	A_1Network_Interface_MS_TCP_Loopback_interface_Packets_Sent_sec
95	A_1Network_InterfaceMS_TCP_Loopback_interface_Current_Bandwidth
96	A_1Network_InterfaceS_TCP_Loopback_interface_Bytes_Received_sec
97	A_1Network_Interfaceback_interface_Packets_Received_Unicast_sec
98	A_1Network_Interface_interface_Packets_Received_Non_Unicast_sec
99	A_1Network_Interfaceopback_interface_Packets_Received_Discarded
100	A_1Network_Interface_Loopback_interface_Packets_Received_Errors
101	A_1Network_InterfaceLoopback_interface_Packets_Received_Unknown
102	A_1Network_Interface_MS_TCP_Loopback_interface_Bytes_Sent_sec
103	A_1Network_InterfaceLoopback_interface_Packets_Sent_Unicast_sec
104	A_1Network_Interfaceback_interface_Packets_Sent_Non_Unicast_sec
105	A_1Network_Interfaceopback_interface_Packets_Outbound_Discarded
106	A_1Network_Interface_Loopback_interface_Packets_Outbound_Errors
107	A_1Network_Interface_TCP_Loopback_interface_Output_Queue_Length
108	A_1PhysicalDisk_Total_Disk_Read_Time
109	A_1PhysicalDisk_Total_Disk_Time
110	A_1PhysicalDisk_Total_Disk_Write_Time
111	A_1PhysicalDisk_Total_Idle_Time
112	A_1PhysicalDisk_Total_Avg_Disk_Bytes_Read

113	A_1PhysicalDisk_Total_Avg_Disk_Bytes_Transfer_
114	A_1PhysicalDisk_Total_Avg_Disk_Bytes_Write_
115	A_1PhysicalDisk_Total_Avg_Disk_Queue_Length_
116	A_1PhysicalDisk_Total_Avg_Disk_Read_Queue_Length_
117	A_1PhysicalDisk_Total_Avg_Disk_sec_Read_
118	A_1PhysicalDisk_Total_Avg_Disk_sec_Transfer_
119	A_1PhysicalDisk_Total_Avg_Disk_sec_Write_
120	A_1PhysicalDisk_Total_Avg_Disk_Write_Queue_Length_
121	A_1PhysicalDisk_Total_Current_Disk_Queue_Length_
122	A_1PhysicalDisk_Total_Disk_Bytes_sec_
123	A_1PhysicalDisk_Total_Disk_Read_Bytes_sec_
124	A_1PhysicalDisk_Total_Disk_Reads_sec_
125	A_1PhysicalDisk_Total_Disk_Transfers_sec_
126	A_1PhysicalDisk_Total_Disk_Write_Bytes_sec_
127	A_1PhysicalDisk_Total_Disk_Writes_sec_
128	A_1PhysicalDisk_Total_Split_IO_Sec_
129	A_1Process_Total_Privileged_Time_
130	A_1Process_Total_Processor_Time_
131	A_1Process_Total_User_Time_
132	A_1Process_Total_Creating_Process_ID_
133	A_1Process_Total_Elapsed_Time_
134	A_1Process_Total_Handle_Count_
135	A_1Process_Total_ID_Process_
136	A_1Process_Total_IO_Data_Bytes_sec_
137	A_1Process_Total_IO_Data_Operations_sec_
138	A_1Process_Total_IO_Other_Bytes_sec_
139	A_1Process_Total_IO_Other_Operations_sec_
140	A_1Process_Total_IO_Read_Bytes_sec_
141	A_1Process_Total_IO_Read_Operations_sec_
142	A_1Process_Total_IO_Write_Bytes_sec_
143	A_1Process_Total_IO_Write_Operations_sec_
144	A_1Process_Total_Page_Faults_sec_
145	A_1Process_Total_Page_File_Bytes_
146	A_1Process_Total_Page_File_Bytes_Peak_
147	A_1Process_Total_Pool_Nonpaged_Bytes_
148	A_1Process_Total_Pool_Paged_Bytes_
149	A_1Process_Total_Priority_Base_
150	A_1Process_Total_Private_Bytes_
151	A_1Process_Total_Thread_Count_
152	A_1Process_Total_Virtual_Bytes_
153	A_1Process_Total_Virtual_Bytes_Peak_

154	A 1Process Total Working Set
155	A 1Process Total Working Set Peak
156	A 1Processor Total C1 Time
157	A 1Processor Total C2 Time
158	A 1Processor Total C3 Time
159	A 1Processor Total DPC Time
160	A 1Processor Total Idle Time
161	A 1Processor Total Interrupt Time
162	A 1Processor Total Privileged Time
163	A 1Processor Total Processor Time
164	A 1Processor Total User Time
165	A 1Processor Total C1 Transitions sec
166	A 1Processor Total C2 Transitions sec
167	A 1Processor Total C3 Transitions sec
168	A 1Processor Total DPC Rate
169	A 1Processor Total DPCs Queued sec
170	A 1Processor Total Interrupts sec
171	A 1SystemFile Read Operations sec
172	A 1SystemFile Write Operations sec
173	A 1SystemFile Control Operations sec
174	A 1SystemFile Read Bytes sec
175	A 1SystemFile Write Bytes sec
176	A 1SystemFile Control Bytes sec
177	A 1SystemContext Switches sec
178	A 1SystemSystem Calls sec
179	A 1SystemFile Data Operations sec
180	A 1SystemSystem Up Time
181	A 1SystemProcessor Queue Length
182	A 1SystemProcesses
183	A 1SystemThreads
184	A 1SystemAlignment Fixups sec
185	A 1SystemException Dispatches sec
186	A 1SystemFloating Emulations sec
187	A 1System Registry Quota In Use
188	A 1TCPSegments sec
189	A 1TCPConnections Established
190	A 1TCPConnections Active
191	A 1TCPConnections Passive
192	A 1TCPConnection Failures
193	A 1TCPConnections Reset
194	A 1TCPSegments Received sec

195	A_1TCPSegments_Sent_sec
196	A_1TCPSegments_Retransmitted_sec
197	A_1Thread_Total_Total_Privileged_Time
198	A_1Thread_Total_Total_Processor_Time
199	A_1Thread_Total_Total_User_Time
200	A_1Thread_Total_Total_Context_Switches_sec
201	A_1Thread_Total_Total_Elapsed_Time
202	A_1Thread_Total_Total_ID_Process
203	A_1Thread_Total_Total_ID_Thread
204	A_1Thread_Total_Total_Priority_Base
205	A_1Thread_Total_Total_Priority_Current
206	A_1Thread_Total_Total_Start_Address
207	A_1Thread_Total_Total_Thread_State
208	A_1Thread_Total_Total_Thread_Wait_Reason
209	A_1UDPDatagrams_sec
210	A_1UDPDatagrams_Received_sec
211	A_1UDPDatagrams_No_Port_sec
212	A_1UDPDatagrams_Received_Errors
213	A_1UDPDatagrams_Sent_sec

#### VTRACE features

214	Process_create
215	Process_destroy
216	Thread_create
217	Thread_destroy
218	Thread_switch
219	Process_set_priority
220	Thread_set_priority
221	Message_get_nf_call
222	Message_get_f_call
223	Message_get_return
224	Message_peek_r_nf_call
225	Message_peek_r_f_call
226	Message_peek_nr_nf_call
227	Message_peek_nr_f_call
228	Message_peek_return
229	Message_dispatch_call
230	Message_dispatch_return
231	Message_trans_accel_call
232	Message_trans_accel_ret
233	Message_translate_call

234	Message_translate_return
235	Message_trans_mdi_call
236	Message_trans_mdi_return
237	Message_set_timer
238	Message_cancel_timer
239	Message_wait_call
240	Message_wait_return
241	Message_get_input_state
242	Message_get_queue_status
243	Key_press
244	Cursor_load
245	Cursor_set
246	Window_create
247	Dialog_create
248	File_complete_operation
249	File_open
250	File_read
251	File_write
252	File_close
253	File_query_info
254	File_set_info
255	File_directory_control
256	File_name
257	File_information
258	File_open_query_close
259	File_rename
260	File_delete
261	File_flush
262	File_lock
263	File_unlock
264	File_set_position
265	File_link
266	File_complete_mdl_op
267	File_fcb_information
268	Message_post_window
269	Message_post_thread
270	Message_send
271	Netobj_complete_op
272	Netobj_open
273	Netobj_close
274	Netobj_connect

275	Netobj_disconnect
276	Netobj_send
277	Netobj_send_datagram
278	Netobj_receive
279	Netobj_receive_datagram
280	Netobj_listen
281	Netobj_accept
282	Netobj_associate_address
283	Netobj_disassociate_addr
284	Netobj_notify_connect
285	Netobj_notify_disconnect
286	Netobj_notify_receive
287	Netobj_notify_rcv_dgram
288	Netobj_notify_rcv_exped
289	VTrace_action
290	Device_hook
291	VTrace_version
292	Local_time
293	VTrace_set_mask
294	Beat
295	Device_unhook
296	High_timestamp
297	Ignore_activity
298	Set_cpu_speed
299	User_changed
300	Waitobj_signal
301	Waitobj_set_timer
302	Waitobj_cancel_timer
303	Waitobj_wait_call
304	Waitobj_wait_return
305	Waitobj_wait_gen_call
306	Waitobj_wait_gen_return
307	Msg__waitobj_wait_call
308	Msg__waitobj_wait_ret
309	Rawdisk_read
310	Rawdisk_write
311	Rawdisk_complete_op
312	Section_create
313	Section_open
314	Section_map_view
315	Section_unmap_view

316	Section_get_mdl
317	Testing
318	Flush_icache
319	Flush_write_buffer
320	Terminate_process
321	Terminate_thread
322	Write_request_data
323	Write_VM



## Appendix B – Top Twenty Features Selected by Each Feature Selection Method

All (V All)	
ChiSqr	ReliefF
A_1MemoryCache_Bytes_Peak	A_1ICMPReceived_Dest_Unreachable
A_1Process_Total_Virtual_Bytes_Peak	A_1ICMPSent_Destination_Unreachable
A_1MemoryFree_System_Page_Table_Entries	A_1SystemFile_Control_Bytes_sec
A_1Process_Total_Virtual_Bytes	A_1Process_Total_IO_Other_Bytes_sec
A_1Process_Total_Pool_Nonpaged_Bytes	A_1ICMPMessages_Outbound_Errors
A_1MemoryPool_Nonpaged_Bytes	A_1MemorySystem_Code_Total_Bytes
A_1Process_Total_Thread_Count	Netobj_disconnect
A_1SystemThreads	A_1ICMPSent_Echo_sec
A_1Process_Total_Pool_Paged_Bytes	A_1ICMPMessages_Sent_sec
A_1TCPConnections_Active	A_1Process_Total_Handle_Count
A_1Network_Interface_Packet_Scheduler_Miniport_Bytes_Sent_sec	A_1ICMPMessages_sec
A_1TCPConnection_Failures	A_1Processor_Total_Processor_Time
A_1MemoryPool_Nonpaged_Allocs	A_1SystemException_Dispatches_sec
A_1Process_Total_Handle_Count	A_1TCPConnections_Reset
A_1Network_InterfaceTX_Packet_Scheduler_Miniport_Packets_sec	A_1Processor_Total_Idle_Time
A_1Network_Interface_Packet_Scheduler_Miniport_Bytes_Total_sec	A_1Processor_Total_User_Time
A_1Process_Total_Page_File_Bytes_Peak	A_1Process_Total_User_Time
A_1IPDatagrams_sec	A_1Thread_Total_Total_User_Time
A_1SystemFile_Control_Bytes_sec	A_1Processor_Total_Interrupts_sec
A_1Process_Total_IO_Other_Bytes_sec	A_1Memory_Committed_Bytes_In_Use
GainRatio	Ensemble
A_1ICMPSent_Echo_sec	A_1ICMPSent_Echo_sec
A_1ICMPMessages_Sent_sec	A_1ICMPMessages_Sent_sec
A_1ICMPMessages_sec	A_1Process_Total_IO_Other_Bytes_sec
A_1TCPConnections_Passive	A_1SystemFile_Control_Bytes_sec
Netobj_disconnect	A_1ICMPMessages_sec
A_1TCPConnection_Failures	A_1MemoryCache_Bytes_Peak
A_1TCPConnections_Active	A_1TCPConnection_Failures
A_1TCPSegments_Retransmitted_sec	A_1TCPConnections_Active
Write_request_data	A_1MemoryFree_System_Page_Table_Entries
A_1MemoryFree_System_Page_Table_Entries	A_1Process_Total_Virtual_Bytes_Peak
A_1ICMPReceived_Echo_Reply_sec	A_1TCPConnections_Passive
A_1ICMPMessages_Received_sec	A_1Process_Total_Handle_Count
A_1ICMPReceived_Echo_sec	A_1MemoryPool_Nonpaged_Bytes
A_1ICMPSent_Echo_Reply_sec	A_1ICMPReceived_Dest_Unreachable
A_1UDPDatagrams_No_Port_sec	A_1Process_Total_Thread_Count
A_1Process_Total_Thread_Count	A_1SystemThreads

A_1SystemThreads_	A_1MemoryPool_Nonpaged_Allocs_
A_1Network_Interface_Packet_Scheduler_Miniport_Bytes_Sent_sec_	A_1Process_Total_Pool_Nonpaged_Bytes_
A_1Network_Interface_Scheduler_Miniport_Packets_Outbound_Errors_	Netobj_disconnect
A_1Network_InterfaceTX_Packet_Scheduler_Miniport_Packets_sec_	A_1Process_Total_Virtual_Bytes_