

Data Mining Methods for Detection of New Malicious Executables

Matthew G. Schultz and Eleazar Eskin

*Department of Computer Science
Columbia University*

{mgs,eeskin}@cs.columbia.edu

Erez Zadok

*Department of Computer Science
State University of New York at Stony Brook*

ezk@cs.sunysb.edu

Salvatore J. Stolfo

*Department of Computer Science
Columbia University*

sal@cs.columbia.edu

Abstract

A serious security threat today is malicious executables, especially new, unseen malicious executables often arriving as email attachments. These new malicious executables are created at the rate of thousands every year and pose a serious security threat. Current anti-virus systems attempt to detect these new malicious programs with heuristics generated by hand. This approach is costly and oftentimes ineffective. In this paper, we present a data-mining framework that detects new, previously unseen malicious executables accurately and automatically. The data-mining framework automatically found patterns in our data set and used these patterns to detect a set of new malicious binaries. Comparing our detection methods with a traditional signature-based method, our method more than doubles the current detection rates for new malicious executables.

1 Introduction

A malicious executable is defined to be a program that performs a malicious function, such as compromising a system's security, damaging a system or obtaining sensitive information without the user's permission. Using data mining methods, our goal is to automatically design and build a scanner that accurately detects malicious executables before they have been given a chance to run.

Data mining methods detect patterns in large amounts of data, such as byte code, and use these patterns to detect future instances in similar data. Our framework uses *classifiers* to detect new malicious executables. A classifier is a rule set, or detection model, generated by the data mining algorithm that was trained over a given set of training data.

One of the primary problems faced by the virus community is to devise methods for detecting new malicious programs that have not yet been analyzed [26]. Eight to ten malicious programs are created every day and most cannot be accurately detected until signatures have been generated for them [27]. During this time period, systems protected by signature-based algorithms are vulnerable to attacks.

Malicious executables are also used as attacks for many types of intrusions. In the DARPA 1999 intrusion detection evaluation, many of the attacks on the Windows platform were caused by malicious programs [19]. Recently, a malicious piece of code created a hole in a Microsoft's internal network [23]. That attack was initiated by a malicious executable that opened a back-door into Microsoft's internal network resulting in the theft of Microsoft's source code.

Current virus scanner technology has two parts: a signature-based detector and a heuristic classifier that detects new viruses [8]. The classic signature-based detection algorithm relies on *signatures* (unique telltale strings) of known malicious executables to generate detection models. Signature-based methods create a unique tag for each malicious program so that future examples of it can be correctly classified with a small error rate. These methods do not generalize well to detect new malicious binaries because they are created to give a false positive rate as close to zero as possible. Whenever a detection method generalizes to new instances, the tradeoff is for a higher false positive rate. *Heuristic* classifiers are generated by a group of virus experts to detect new malicious programs. This kind of analysis can be time-consuming and oftentimes still fail to detect new malicious executables.

We designed a framework that used data mining algorithms to train multiple classifiers on a set of malicious and benign executables to detect new examples. The binaries were first statically analyzed to extract properties of the binary, and then the classifiers trained over a subset of the data.

Our goal in the evaluation of this method was to simulate the task of detecting new malicious executables. To do this we separated our data into two sets: a *training set* and a *test set* with standard cross-validation methodology. The training set was used by the data mining algorithms to generate classifiers to classify previously unseen binaries as malicious or benign. A test set is a subset of dataset that had no examples in it that were seen during the training of an algorithm. This subset was used to test an algorithms' performance over similar, unseen data and its performance

over new malicious executables. Both the test and training data were malicious executables gathered from public sources.

We implemented a traditional signature-based algorithm to compare with the the data mining algorithms over new examples. Using standard statistical cross-validation techniques, our data mining-based method had a detection rate of 97.76%—more than double the detection rate of a signature-based scanner over a set of new malicious executables. Comparing our method to industry heuristics cannot be done at this time because the methods for generating these heuristics are not published and there is no equivalent or statistically comparable dataset to which both techniques are applied. However, the framework we provide is fully automatic and could assist experts in generating the heuristics.

2 Background

Detecting malicious executables is not a new problem in security. Early methods used signatures to detect malicious programs. These signatures were composed of many different properties: filename, text strings, or byte code. Research also centered on protecting the system from the security holes that these malicious programs created.

Experts were typically employed to analyze suspicious programs by hand. Using their expertise, signatures were found that made a malicious executable example different from other malicious executables or benign programs. One example of this type of analysis was performed by Spafford [24] who analyzed the Internet Worm and provided detailed notes on its spread over the Internet, the unique signatures in the worm’s code, the method of the worm’s attack, and a comprehensive description of system failure points.

Although accurate, this method of analysis is expensive, and slow. If only a small set of malicious executables will ever circulate then this method will work very well, but the *Wildlist* [22] is always changing and expanding. The *Wildlist* is a list of malicious programs that are currently estimated to be circulating at any given time.

Current approaches to detecting malicious programs match them to a set of known malicious programs. The anti-virus community relies heavily on known byte-code signatures to detect malicious programs. More recently, these byte sequences were determined by automatically examining known malicious binaries with probabilistic methods.

At IBM, Kephart and Arnold [9] developed a statistical method for automatically extracting malicious executable signatures. Their research was based on speech recognition algorithms and was shown to perform almost as good as a human expert at detecting known malicious executables. Their algorithm was eventually packaged with IBM’s anti-virus software.

Lo et al. [15] presented a method for filtering mali-

cious code based on “tell-tale signs” for detecting malicious code. These were manually engineered based on observing the characteristics of malicious code. Similarly, filters for detecting properties of malicious executables have been proposed for UNIX systems [10] as well as semi-automatic methods for detecting malicious code [4].

Unfortunately, a new malicious program may not contain any known signatures so traditional signature-based methods may not detect a new malicious executable. In an attempt to solve this problem, the anti-virus industry generates heuristic classifiers by hand [8]. This process can be even more costly than generating signatures, so finding an automatic method to generate classifiers has been the subject of research in the anti-virus community. To solve this problem, different IBM researchers applied *Artificial Neural Networks* (ANNs) to the problem of detecting boot sector malicious binaries [25]. An ANN is a classifier that models neural networks explored in human cognition. Because of the limitations of the implementation of their classifier, they were unable to analyze anything other than small boot sector viruses which comprise about 5% of all malicious binaries.

Using an ANN classifier with all bytes from the boot sector malicious executables as input, IBM researchers were able to identify 80–85% of unknown boot sector malicious executables successfully with a low false positive rate (< 1%). They were unable to find a way to apply ANNs to the other 95% of computer malicious binaries.

In similar work, Arnold and Tesauro [1] applied the same techniques to Win32 binaries, but because of limitations of the ANN classifier they were unable to have the comparable accuracy over new Win32 binaries.

Our method is different because we analyzed the entire set of malicious executables instead of only boot-sector viruses, or only Win32 binaries.

Our technique is similar to data mining techniques that have already been applied to Intrusion Detection Systems by Lee et al. [13, 14]. Their methods were applied to system calls and network data to learn how to detect new intrusions. They reported good detection rates as a result of applying data mining to the problem of IDS. We applied a similar framework to the problem of detecting new malicious executables.

3 Methodology

The goal of this work was to explore a number of standard data mining techniques to compute accurate detectors for new (unseen) binaries. We gathered a large set of programs from public sources and separated the problem into two classes: *malicious* and *benign* executables. Every example in our data set is a Windows or MS-DOS format executable, although the framework we present is applicable to other formats. To standardize our data-set, we used an updated MacAfee’s [16] virus scanner and labeled our pro-

grams as either malicious or benign executables. Since the virus scanner was updated and the viruses were obtained from public sources, we assume that the virus scanner has a signature for each malicious virus.

We split the dataset into two subsets: the *training set* and the *test set*. The data mining algorithms used the training set while generating the rule sets. We used a test set to check the accuracy of the classifiers over unseen examples.

Next, we automatically extracted a binary profile from each example in our dataset, and from the binary profiles we extracted *features* to use with classifiers. In a data mining framework, features are properties extracted from each example in the data set—such as byte sequences—that a classifier can use to generate detection models. Using different features, we trained a set of data mining classifiers to distinguish between benign and malicious programs. It should be noted that the features extracted were static properties of the binary and did not require executing the binary.

The framework supports different methods for feature extraction and different data mining classifiers. We used system resource information, strings and byte sequences that were extracted from the malicious executables in the data set as different types of features. We also used three learning algorithms:

- an inductive rule-based learner that generates boolean rules based on feature attributes.
- a probabilistic method that generates probabilities that an example was in a class given a set of features.
- a multi-classifier system that combines the outputs from several classifiers to generate a prediction.

To compare the data mining methods with a traditional signature-based method, we designed an automatic signature generator. Since the virus scanner that we used to label the data set had signatures for every malicious example in our data set, it was necessary to implement a similar signature-based method to compare with the data mining algorithms. There was no way to use an off-the-shelf virus scanner and simulate the detection of new malicious executables because these commercial scanners contained signatures for all the malicious executables in our data set. Like the data mining algorithms, the signature-based algorithm was only allowed to generate signatures over the set of training data. This allowed our data mining framework to be fairly compared to traditional scanners over new data.

To quantitatively express the performance of our method we show tables with the counts for *true positives* (TP), *true negatives* (TN), *false positives* (FP), and *false negatives* (FN). A true positive is a malicious example that is correctly tagged as malicious, and a true negative is a benign example that is correctly classified. A false positive is a benign program that has been mislabeled by an algorithm as a malicious program, while a false negative is a malicious executable that has been misclassified as a benign program.

To evaluate the performance, we compute the false positive rate and the detection rate. The false positive rate is the number of benign examples that are mislabeled as malicious divided by the total number of benign examples. The detection rate is the number of malicious examples that are caught divided by the total number of malicious examples.

3.1 Dataset Description

Our data set consisted of a total of 4,266 programs split into 3,265 malicious binaries and 1,001 clean programs. There were no duplicate programs in the data set and every example in the set is labeled either malicious or benign by the commercial virus scanner.

The malicious executables were downloaded from various FTP sites and were labeled by a commercial virus scanner with the correct class label (malicious or benign) for our method. 5% of the data set was composed of Trojans and the other 95% consisted of viruses. Most of the clean programs were gathered from a freshly installed Windows 98 machine running MSOffice 97 while others are small executables downloaded from the Internet. The entire data set is available from our Web site <http://www.cs.columbia.edu/ids/mef/software/>.

We also examined a subset of the data that was in *Portable Executable* (PE) [17] format. The data set consisting of PE format executables was composed of 206 benign programs and 38 malicious executables.

After verification of the data set the next step of our method was to extract features from the programs.

4 Feature Extraction

In this section we detail all of our choices of features. We statically extracted different features that represented different information contained within each binary. These features were then used by the algorithms to generate detection models.

We first examine only the subset of PE executables using LibBFD. Then we used more general methods to extract features from all types of binaries.

4.1 LibBFD

Our first intuition into the problem was to extract information from the binary that would dictate its behavior. The problem of predicting a program's behavior can be reduced to the halting problem and hence is undecidable [2]. Perfectly predicting a program's behavior is unattainable but estimating what a program can or cannot do is possible. For instance if a Windows executable does not call the User Interfaces Dynamically Linked Library (USER32.DLL), then we could assume that the program does not have the standard Windows user interface. This is of course an oversimplification of the problem because the author of that ex-

ample could have written or linked to another user interface library, but it did provide us with some insight to an appropriate feature set.

To extract resource information from Windows executables we used GNU’s Bin–Utils [5]. GNU’s Bin–Utils suite of tools can analyze PE binaries within Windows. In PE, or *Common Object File Format* (COFF), program headers are composed of a COFF header, an Optional header, an MS-DOS stub, and a file signature. From the PE header we used libBFD, a library within Bin–Utils, to extract information in *object format*. Object format for a PE binary gives the file size, the names of DLLs, and the names of function calls within those DLLs and Relocation Tables. From the object format, we extracted a set of features to compose a feature vector for each binary.

To understand how resources affected a binary’s behavior we performed our experiments using three types of features:

1. The list of DLLs used by the binary
2. The list of DLL function calls made by the binary
3. The number of different function calls within each DLL

The first approach to binary profiling (shown in Figure 1) used the DLLs loaded by the binary as features. The feature vector comprised of 30 boolean values representing whether or not a binary used a DLL. Typically, not every DLL was used in all of the binaries, but a majority of the binaries called the same resource. For example, almost every binary called GDI32.DLL, which is the Windows NT Graphics Device Interface and is a core component of WinNT.

$$\neg advapi32 \wedge avicap32 \wedge \dots \wedge winmm \wedge \neg wsock32$$

Figure 1: First Feature Vector: A conjunction of DLL names

The example vector given in Figure 1 is composed of at least two unused resources: ADVAPI32.DLL, the Advanced Windows API, and WSOCK32.DLL, the Windows Sockets API. It also uses at least two resources: AVICAP32.DLL, the AVI capture API, and WINMM.DLL, the Windows Multimedia API.

The second approach to binary profiling (seen in Figure 2) used DLLs and their function calls as features. This approach was similar to the first, but with added function call information. The feature vector was composed of 2,229 boolean values. Because some of the DLL’s had the same function names it was important to record which DLL the function came from.

The example vector given in Figure 2 is composed of at least four resources. Two functions were called in ADVAPI32.DLL: AdjustTokenPrivileges() and GetFileSecu-

$$\begin{aligned} & advapi32.AdjustTokenPrivileges() \\ \wedge & advapi32.GetFileSecurityA() \wedge \dots \\ \wedge & wsock32.recv() \wedge wsock32.send() \end{aligned}$$

Figure 2: Second Feature Vector: A conjunction of DLL’s and the functions called inside each DLL

urityA(), and two functions in WSOCK32.DLL: recv() and send().

The third approach to binary profiling (seen in Figure 3) counted the number of different function calls used within each DLL. The feature vector included 30 integer values. This profile gives a rough measure of how heavily a DLL is used within a specific binary. Intuitively, in the resource models we have been exploring, this is a macro-resource usage model because the number of calls to each resource is counted instead of detailing referenced functions. For example, if a program only called the recv() and send() functions of WSOCK32.DLL, then the count would be 2. It should be noted that we do not count the number of times those functions might have been called.

$$\begin{aligned} & advapi32 = 2 \quad \wedge avicap32 = 10 \wedge \dots \\ \wedge & winmm = 8 \quad \wedge wsock32 = 2 \end{aligned}$$

Figure 3: Third Feature Vector: A conjunction of DLL’s and a count of the number of functions called inside each DLL

The example vector given in Figure 3 describes an example that calls two functions in ADVAPI32.DLL, ten functions in AVICAP32.DLL, eight functions in WINMM.DLL and two functions from WSOCK32.DLL.

All of the information about the binary was obtained from the program header. In addition, the information was obtained without executing the unknown program but by examining the static properties of the binary, using libBFD.

Since we could not analyze the entire dataset with libBFD we found another method for extracting features that works over the entire dataset. We describe that method next.

4.2 GNU Strings

During the analysis of our libBFD method we noticed that headers in PE-format were in plain text. This meant that we could extract the same information from the PE-executables by just extracting the plain text headers. We also noticed that non-PE executables also have strings encoded in them. We theorized that we could use this information to classify the full 4,266 item data set instead of the

small libBFD data set.

To extract features from the first data set of 4,266 programs we used the GNU *strings* program. The strings program extracts consecutive printable characters from any file. Typically there are many printable strings in binary files. Some common strings found in our dataset are illustrated in Table 1.

kernel	microsoft	windows	getmodulehandlea
getversion	getstartupinfo	win	getmodulefilenamea
messageboxa	closehandle	null	dispatchmessagea
library	getprocaddress	advapi	getlasterror
loadlibrarya	exitprocess	heap	getcommandlinea
reloc	createfilea	writefile	setfilepointer
application	showwindow	time	regclosekey

Table 1: Common strings extracted from binaries using GNU strings

Through testing we found that there were similar strings in malicious executables that distinguished them from clean programs, and similar strings in benign programs that distinguished them from malicious executables. Each string in the binary was used as a feature. In the data mining step, we discuss how a frequency analysis was performed over all the byte sequences found in our data set.

The strings contained in a binary may consist of reused code fragments, author signatures, file names, system resource information, etc. This method of detecting malicious executables is already used by the anti-malicious executable community to create signatures for malicious executables.

Extracted strings from an executable are not very robust as features because they can be changed easily, so we analyzed another feature, byte sequences.

4.3 Byte Sequences Using Hexdump

Byte sequences are the last set of features that we used over the entire 4,266 member data set. We used *hexdump* [18], a tool that transforms binary files into hexadecimal files. The byte sequence feature is the most informative because it represents the machine code in an executable instead of resource information like libBFD features. Secondly, analyzing the entire binary gives more information for non-PE format executables than the strings method.

After we generated the hexdumps we had features as displayed in Figure 4 where each line represents a short sequence of machine code instructions.

We again assumed that there were similar instructions in malicious executables that differentiated them from benign programs, and the class of benign programs had similar byte code that differentiated them from the malicious executables. Also like the string features, each byte sequence in a binary is used as a feature.

```
1f0e 0eba b400 cd09 b821 4c01 21cd 6854
7369 7020 6f72 7267 6d61 7220 7165 6975
6572 2073 694d 7263 736f 666f 2074 6957
646e 776f 2e73 0a0d 0024 0000 0000 0000
454e 3c05 026c 0009 0000 0000 0302 0004
0400 2800 3924 0001 0000 0004 0004 0006
000c 0040 0060 021e 0238 0244 02f5 0000
0001 0004 0000 0802 0032 1304 0000 030a
```

Figure 4: Example Hexdump

5 Algorithms

In this section we describe all the algorithms presented in this paper as well as the signature-based method used for comparison. We used three different data mining algorithms to generate classifiers with different features: RIPPER, Naive Bayes, and a Multi-Classifer system.

We describe the signature-based method first.

5.1 Signature Methods

We examine signature-based methods to compare our results to traditional anti-virus methods. Signature-based detection methods are the most commonly used algorithms in industry [27]. These signatures are picked to differentiate one malicious executable from another, and from benign programs. These signatures are generated by an expert in the field or an automatic method. Typically, a signature is picked to illustrate the distinct properties of a specific malicious executable.

We implemented a signature-based scanner with this method that follows a simple algorithm for signature generation. First, we calculated the byte-sequences that were only found in the malicious executable class. These byte-sequences were then concatenated together to make a unique signature for each malicious executable example. Thus, each malicious executable signature contained only byte-sequences found in the malicious executable class. To make the signature unique, the byte-sequences found in each example were concatenated together to form one signature. This was done because a byte-sequence that is only found in one class during training could possibly be found in the other class during testing [9], and lead to false positives in testing.

The method described above for the commercial scanner was never intended to detect unknown malicious binaries, but the data mining algorithms that follow were built to detect new malicious executables.

5.2 RIPPER

The next algorithm we used, RIPPER [3], is an inductive rule learner. This algorithm generated a detection model composed of resource rules that was built to detect future

examples of malicious executables. This algorithm used libBFD information as features.

RIPPER is a rule-based learner that builds a set of rules that identify the classes while minimizing the amount of error. The error is defined by the number of training examples misclassified by the rules.

An inductive algorithm learns what a malicious executable is given a set of training examples. The four features seen in Table 2 are:

1. “Does it have a GUI?”
2. “Does it perform a malicious function?”
3. “Does it compromise system security?”
4. “Does it delete files?”

and finally the class question “Is it malicious?”

Has a GUI?	Malicious Function?	Compromise Security?	Deletes Files?	Is it malicious?
yes	yes	yes	no	yes
no	yes	yes	yes	yes
yes	no	no	yes	no
yes	yes	yes	yes	yes

Table 2: Example Inductive Training Set. Intuitively all malicious executables share the second and third feature, “yes” and “yes” respectively.

The defining property of any inductive learner is that no a priori assumptions have been made regarding the final concept. The inductive learning algorithm makes as its primary assumption that the data trained over is similar in some way to the unseen data.

A hypothesis generated by an inductive learning algorithm for this learning problem has four attributes. Each attribute will have one of these values:

1. \top , truth, indicating any value is acceptable in this position,
2. a value, either yes, or no, is needed in this position, or
3. a \perp , falsity, indicating that no value is acceptable for this position

For example, the hypothesis $\langle \top, \top, \top, \top \rangle$ and the hypothesis $\langle yes, yes, yes, no \rangle$ would make the first example true. $\langle \top, \top, \top, \top \rangle$ would make any feature set true and $\langle yes, yes, yes, no \rangle$ is the set of features for example one.

The algorithm we describe is *Find-S* [20]. *Find-S* finds the most specific hypothesis that is consistent with the training examples. For a positive training example the algorithm replaces any attribute in the hypothesis that is inconsistent with the training example with a more general attribute. Of all the hypotheses values 1 is more general

than 2 and 2 is more general than 3. For a negative example the algorithm does nothing. Positive examples in this problem are defined to be the malicious executables and negative examples are the benign programs.

The initial hypothesis that *Find-S* starts with is $\langle \perp, \perp, \perp, \perp \rangle$. This hypothesis is the most specific because it is true over the fewest possible examples, none. Examining the first positive example in Table 2, $\langle yes, yes, yes, no \rangle$, the algorithm chooses the next most specific hypothesis $\langle yes, yes, yes, no \rangle$. The next positive example, $\langle no, no, no, yes \rangle$, is inconsistent with the hypothesis in its first and fourth attribute (“Does it have a GUI?” and “Does it delete files?”) and those attributes in the hypothesis get replaced with the next most general attribute, \top .

The resulting hypothesis after two positive examples is $\langle \top, yes, yes, \top \rangle$. The algorithm skips the third example, a negative example, and finds that this hypothesis is consistent with the final example in Table 2. The final rule for the training data listed in Table 2 is $\langle \top, yes, yes, \top \rangle$. The rule states that the attributes of a malicious executable, based on training data, are that it has a malicious function and compromises system security. This is consistent with the definition of a malicious executable we gave in the introduction. It does not matter in this example if a malicious executable deletes files, or if it has a GUI or not.

Find-S is a relatively simple algorithm while RIPPER is more complex. RIPPER looks at both positive and negative examples to generate a set of hypotheses that more closely approximate the target concept while *Find-S* generates one hypothesis that approximates the target concept.

5.3 Naive Bayes

The next classifier we describe is a Naive Bayes classifier [6]. The naive Bayes classifier computes the likelihood that a program is malicious given the features that are contained in the program. This method used both strings and byte-sequence data to compute a probability of a binary’s maliciousness given its features.

Nigam et al. [21] performed a similar experiment when they classified text documents according to which newsgroup they originated from. In this method we treated each executable’s features as a text document and classified based on that. The main assumption in this approach is that the binaries contain similar features such as signatures, machine instructions, etc.

Specifically, we wanted to compute the class of a program given that the program contains a set of features F . We define C to be a random variable over the set of classes: *benign*, and *malicious* executables. That is, we want to compute $P(C|F)$, the probability that a program is in a certain class given the program contains the set of features F . We apply Bayes rule and express the probability as:

$$P(C|F) = \frac{P(F|C) * P(C)}{P(F)} \quad (1)$$

To use the naive Bayes rule we assume that the features occur independently from one another. If the features of a program F include the features $F_1, F_2, F_3, \dots, F_n$, then equation (1) becomes:

$$P(C|F) = \frac{\prod_{i=1}^n P(F_i|C) * P(C)}{\prod_{j=1}^n P(F_j)} \quad (2)$$

Each $P(F_i|C)$ is the frequency that string F_i occurs in a program of class C . $P(C)$ is the proportion of the class C in the entire set of programs.

The output of the classifier is the highest probability class for a given set of strings. Since the denominator of (1) is the same for all classes we take the maximum class over all classes C of the probability of each class computed in (2) to get:

$$\text{Most Likely Class} = \max_C \left(P(C) \prod_{i=1}^n P(F_i|C) \right) \quad (3)$$

In (3), we use \max_C to denote the function that returns the class with the highest probability. Most Likely Class is the class in C with the highest probability and hence the most likely classification of the example with features F .

To train the classifier, we recorded how many programs in each class contained each unique feature. We used this information to classify a new program into an appropriate class. We first used feature extraction to determine the features contained in the program. Then we applied equation (3) to compute the most likely class for the program.

We used the Naive Bayes algorithm and computed the most likely class for byte sequences and strings.

5.4 Multi-Naive Bayes

The next data mining algorithm we describe is Multi-Naive Bayes. This algorithm was essentially a collection of Naive Bayes algorithms that voted on an overall classification for an example. Each Naive Bayes algorithm classified the examples in the test set as malicious or benign and this counted as a vote. The votes were combined by the Multi-Naive Bayes algorithm to output a final classification for all the Naive Bayes.

This method was required because even using a machine with one gigabyte of RAM, the size of the binary data was too large to fit into memory. The Naive Bayes algorithm required a table of all strings or bytes to compute its probabilities. To correct this problem we divided the problem into smaller pieces that would fit in memory and trained a Naive Bayes algorithm over each of the subproblems.

We split the data evenly into several sets by putting each i th line in the binary into the $(i \bmod n)$ th set where n is

the number of sets. For each set we trained a Naive Bayes classifier. Our prediction for a binary is the product of the predictions of the n classifiers. In our experiments we use 6 classifiers ($n = 6$).

More formally, the Multi-Naive Bayes promotes a vote of confidence between all of the underlying Naive Bayes classifiers. Each classifier gives a probability of a class C given a set of bytes F which the Multi-Naive Bayes uses to generate a probability for class C given F over all the classifiers.

We want to compute the likelihood of a class C given bytes F and the probabilities learned by each classifier $NaiveBayes_i$. In equation (4) we computed the likelihood, $L_{NB}(C|F)$, of class C given a set of bytes F .

$$L_{NB}(C|F) = \prod_{i=1}^{|NB|} P_{NB_i}(C|F)/P_{NB_i}(C) \quad (4)$$

where NB_i is a Naive Bayes classifier and NB is the set of all combined Naive Bayes classifiers (in our case 6). $P_{NB_i}(C|F)$ (generated from equation (2)) is the probability for class C computed by the classifier $NaiveBayes_i$ given F divided by the probability of class C computed by $NaiveBayes_i$. Each $P_{NB_i}(C|F)$ was divided by $P_{NB_i}(C)$ to remove the redundant probabilities. All the terms were multiplied together to compute $L_{NB}(C|F)$, the final likelihood of C given F . $|NB|$ is the size of the set NB such that $\forall NB_i \in NB$.

The output of the multi-classifier given a set of bytes F is the class of highest probability over the classes given $L_{NB}(C|F)$ and $P_{NB}(C)$ the prior probability of a given class.

$$\text{Most Likely Class} = \max_C (P_{NB}(C) * L_{NB}(C|F)) \quad (5)$$

Most Likely Class is the class in C with the highest probability hence the most likely classification of the example with features F , and \max_C returns the class with the highest likelihood.

6 Rules

Each data mining algorithm generated its own rule set to evaluate new examples. These detection models were the final result of our experiments. Each algorithm's rule set could be incorporated into a scanner to detect malicious programs. The generation of the rules only needed to be done periodically and the rule set distributed in order to detect new malicious executables.

6.1 RIPPER

RIPPER's rules were built to generalize over unseen examples so the rule set was more compact than the signature

based methods. For the data set that contained 3,301 malicious executables the RIPPER rule set contained the five rules in Figure 5.

```

malicious := ¬user32.EndDialog() ∧
             kernel32.EnumCalendarInfoA()
malicious := ¬user32.LoadIconA() ∧
             ¬kernel32.GetTempPathA() ∧ ¬advapi32.
malicious := shell32.ExtractAssociatedIconA()
malicious := msvbvm.
benign    := otherwise

```

Figure 5: Sample Classification Rules using features found in Figure 2

Here, a malicious executable was consistent with one of four hypotheses:

1. it did not call `user32.EndDialog()` but it did call `kernel32.EnumCalendarInfoA()`
2. it did not call `user32.LoadIconA()`, `kernel32.GetTempPathA()`, or any function in `advapi32.dll`
3. it called `shell32.ExtractAssociatedIconA()`,
4. it called any function in `msvbvm.dll`, the Microsoft Visual Basic Library

A binary is labeled benign if it is inconsistent with all of the malicious binary hypotheses in Figure 5.

6.2 Naive Bayes

The Naive Bayes rules were more complicated than the RIPPER and signature based hypotheses. These rules took the form of $P(F|C)$, the probability of an example F given a class C . The probability for a string occurring in a class is the total number of times it occurred in that class’s training set divided by the total number of times that the string occurred over the entire training set. These hypotheses are illustrated in Figure 6.

$$\begin{aligned}
 P(\text{"windows"}|\text{benign}) &= 45/47 \\
 P(\text{"windows"}|\text{malicious}) &= 2/47 \\
 P(\text{"*.COM"}|\text{benign}) &= 1/12 \\
 P(\text{"*.COM"}|\text{malicious}) &= 11/12
 \end{aligned}$$

Figure 6: Sample classification rules found by Naive Bayes.

Here, the string “windows” was predicted to more likely to occur in a benign program and string “*.COM” was more than likely in a malicious executable program.

However this leads to a problem when a string (e.g. “CH2OH-CHOH-CH2OH”) only occurred in one set, for example only in the malicious executables. The probability of “CH2OH-CHOH-CH2OH” occurring in any future benign example is predicted to be 0, but this is an incorrect assumption. If a Chemistry TA’s program was written to print out “CH2OH-CHOH-CH2OH” (glycerol) it will always be tagged a malicious executable even if it has other strings in it that would have labeled it benign.

In Figure 6 the string “*.COM” does not occur in any benign programs so the probability of “*.COM” occurring in class benign is approximated to be 1/12 instead of 0/11. This approximates real world probability that any string could occur in both classes even if during training it was only seen in one class [9].

6.3 Multi-Naive Bayes

The rule sets generated by our Multi-Naive Bayes algorithm are the collection of the rules generated by each of the component Naive Bayes classifiers. For each classifier, there is a rule set such as the one in Figure 6. The probabilities in the rules for the different classifiers may be different because the underlying data that each classifier is trained on is different. The prediction of the Multi-Naive Bayes algorithm is the product of the predictions of the underlying Naive Bayes classifiers.

7 Results and Analysis

We estimate our results over new data by using 5-fold cross validation [12]. Cross validation is the standard method to estimate likely predictions over unseen data in Data Mining. For each set of binary profiles we partitioned the data into 5 equal size groups. We used 4 of the partitions for training and then evaluated the rule set over the remaining partition. Then we repeated the process 5 times leaving out a different partition for testing each time. This gave us a very reliable measure of our method’s accuracy over unseen data. We averaged the results of these five tests to obtain a good measure of how the algorithm performs over the entire set.

To evaluate our system we were interested in several quantities:

1. True Positives (TP), the number of malicious executable examples classified as malicious executables
2. True Negatives (TN), the number of benign programs classified as benign.
3. False Positives (FP), the number of benign programs classified as malicious executables

- False Negatives (FN), the number of malicious executables classified as benign binaries.

We were interested in the detection rate of the classifier. In our case this was the percentage of the total malicious programs labeled malicious. We were also interested in the false positive rate. This was the percentage of benign programs which were labeled as malicious, also called false alarms.

The Detection Rate is defined as $\frac{TP}{TP+FN}$, False Positive Rate as $\frac{FP}{TN+FP}$, and Overall Accuracy as $\frac{TP+TN}{TP+TN+FP+FN}$. The results of all experiments are presented in Table 3.

For all the algorithms we plotted the detection rate vs. false positive rate using *Receiver Operating Characteristic* (ROC) curves [11]. ROC curves are a way of visualizing the trade-offs between detection and false positive rates.

7.1 Signature Method

As is shown in Table 3, the signature method had the lowest false positive rate, 0%. This algorithm also had the lowest detection rate, 33.75%, and accuracy rate, 49.28%.

Since we use this method to compare with the learning algorithms we plot its ROC curves against the RIPPER algorithm in Figure 7 and against the Naive Bayes and Multi-Naive Bayes algorithms in Figure 8.

The detection rate of the signature-based method is inherently low over new executables because the signatures generated were never designed to detect new malicious executables. Also it should be noted that although the signature based method only detected 33.75% of new malicious programs, the method did detect 100% of the malicious binaries that it had seen before with a 0% false positive rate.

7.2 RIPPER

The RIPPER results shown in Table 3 are roughly equivalent to each other in detection rates and overall accuracy, but the method using features from Figure 2, a list of DLL function calls, has a higher detection rate.

The ROC curves for all RIPPER variations are shown in Figure 7. The lowest line represents RIPPER using DLLs only as features, and it was roughly linear in its growth. This means that as we increase detection rate by 5% the false positive would also increase by roughly 5%.

The other lines are concave down so there was an optimal trade-off between detection and false alarms. For DLL's with Counted Function Calls this optimal point was when the false positive rate was 10% and the detection rate was equal to 75%. For DLLs with Function Calls the optimal point was when the false positive rate was 12% and the detection rate was less than 90%.

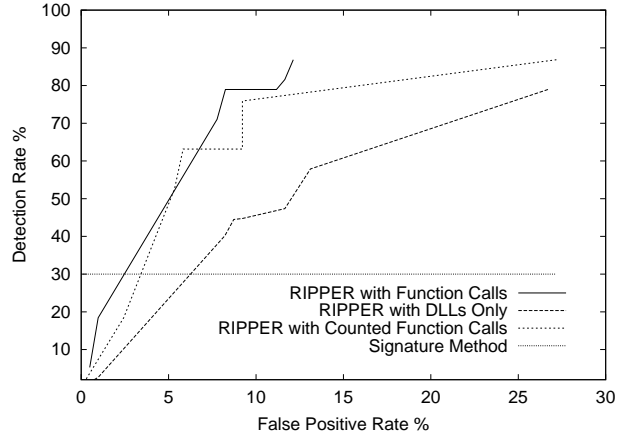


Figure 7: RIPPER ROC. Notice that the RIPPER curves have a higher detection rate than the comparison method with false-positive rates greater than 7%.

7.3 Naive Bayes

The Naive Bayes algorithm using strings as features performed the best out of the learning algorithms and better than the signature method in terms of false positive rate and overall accuracy (see Table 3). It is the most accurate algorithm with 97.11% and within 1% of the highest detection rate, Multi-Naive Bayes with 97.76%. It performed better than the RIPPER methods in every category.

In Figure 8, the slope of the Naive Bayes curve is initially much steeper than the Multi-Naive Bayes. The Naive Bayes with strings algorithm has better detection rates for small false positive rates. Its results were greater than 90% accuracy with a false positive rate less than 2%.

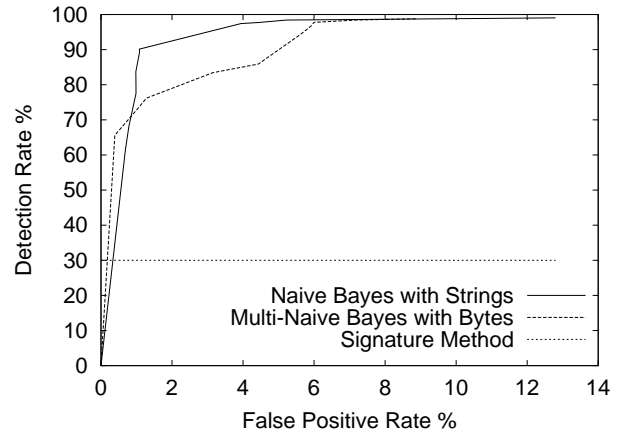


Figure 8: Naive Bayes and Multi-Naive Bayes ROC. Note that the Naive Bayes and Multi-Naive Bayes methods have higher detection rate than the signature method with a greater than 0.5% false positive rate.

Profile Type	True Positives (TP)	True Negatives (TN)	False Positives (FP)	False Negatives (FN)	Detection Rate	False Positive Rate	Overall Accuracy
Signature Method — Bytes	1102	1000	0	2163	33.75%	0%	49.28%
RIPPER							
— DLLs used	22	187	19	16	57.89%	9.22%	83.62%
— DLL function calls	27	190	16	11	71.05%	7.77%	89.36%
— DLLs with counted function calls	20	195	11	18	52.63%	5.34%	89.07%
Naive Bayes							
— Strings	3176	960	41	89	97.43%	3.80%	97.11%
Multi-Naive Bayes							
— Bytes	3191	940	61	74	97.76%	6.01%	96.88%

Table 3: These are the results of classifying new malicious programs organized by algorithm and feature. Multi-Naive Bayes using Bytes had the highest Detection Rate, and Signature Method with strings had the lowest False Positive Rate. Highest overall accuracy was the Naive Bayes algorithm with strings. Note that the detection rate for the signature-based methods are lower than the data mining methods.

7.4 Multi-Naive Bayes

The Multi-Naive Bayes algorithm using bytes as features had the highest detection rate out of any method we tested, 97.76%. The false positive rate at 6.01% was higher than the Naive Bayes methods (3.80%) and the signature methods (< 1%).

The ROC curves in Figure 8 show a slower growth than the Naive Bayes with strings method until the false positive rate climbed above 4%. Then the two algorithms converged for false positive rates greater than 6% with a detection rate greater than 95%.

7.5 Same Model, Different Applications

The ROC curves in Figures 7 and 8 also let security experts understand how to tailor this framework to their specific needs. For example, in a secure computing setting, it may be more important to have a high detection rate of 98.79%, in which case the false positive rate would increase to 8.87%. Or if the framework were to be integrated into a mail server, it may be more important to have a low false positive rate below 1% (0.39% FP rate in our case) and a detection rate of 65.62% for new malicious programs.

8 Defeating Detection Models

Although these methods can detect new malicious executables, a malicious executable author could bypass detection if the detection model were to be compromised.

First, to defeat the signature-based method requires removing all malicious signatures from the binary. Since these are typically a subset of a malicious executable’s total data, changing the signature of a binary would be possible although difficult.

Defeating the models generated by RIPPER would require generating functions that would change the resource

usage. These functions do not have to be called by the binary but would change the resource signature of an executable.

To defeat our implementation of the Naive Bayes classifier it would be necessary to change a significant number of features in the example. One way this can be done is through encryption, but encryption will add overhead to small malicious executables.

We corrected the problem of authors evading a strings-based rule set by initially classifying each example as malicious. If no strings that were contained in the binary had ever been used for training then the final class was malicious. If there were strings contained in the program that the algorithm had seen before then the probabilities were computed normally according to the Naive Bayes rule from Section 4.3. This took care of the instance where a binary had encrypted strings, or had changed all of its strings.

The Multi-Naive Bayes method improved on these results because changing every line of byte code in the Naive Bayes detection model would be an even more difficult proposition than changing all the strings. Changing this many of the lines in a program would change the binary’s behavior significantly. Removing all lines of code that appear in our model would be difficult and time consuming, and even then if none of the byte sequences in the example had been used for training then the example would be initially classified as malicious.

The Multi-Naive Bayes is a more secure model of detection than any of the other methods discussed in this paper because we evaluate a binary’s entire instruction set whereas signature methods looks for segments of byte sequences. It is much easier for malicious program authors to modify the lines of code that a signature represents than to change all the lines contained in the program to evade a Naive Bayes or Multi-Naive Bayes model. The byte sequence model is the most secure model we devised in our test.

A further security concern is what happens when malicious software writers obtain copies of the malicious binaries that we could not detect, and use these false negatives to generate new malicious software. Presumably this would allow them to circumvent our detection models, but in fact having a larger set of similar false negatives would make our model more accurate. In other words, if malicious binary authors clone the undetectable binaries, they are in effect making it easier for this framework to detect their programs. The more data that the method analyzes, and the more false positives and false negatives that it learns from, the more accurate the method becomes at distinguishing between benign and malicious programs.

9 Conclusions

The first contribution that we presented in this paper was a method for detecting previously undetectable malicious executables. We showed this by comparing our results with traditional signature-based methods and with other learning algorithms. The Multi-Naive Bayes method had the highest accuracy and detection rate of any algorithm over unknown programs, 97.76%, over double the detection rates of signature-based methods. Its rule set was also more difficult to defeat than other methods because all lines of machine instructions would have to be changed to avoid detection.

The first problem with traditional anti-malicious executable detection methods is that in order to detect a new malicious executable, the program needs to be examined and a signature extracted from it and included in the anti-malicious executable software database. The difficulty with this method is that during the time required for a malicious program to be identified, analyzed and signatures to be distributed, systems are at risk from that program. Our methods may provide a defense during that time. With a low false positive rate, the inconvenience to the end user would be minimal while providing ample defense during the time before an update of models is available.

Virus Scanners are updated about every month. 240–300 new malicious executables are created in that time (8–10 a day [27]). Our method would catch roughly 216–270 of those new malicious executables without the need for an update whereas traditional methods would catch only 87–109. Our method more than doubles the detection rate of signature based methods for new malicious binaries.

The methods discussed in this paper are being implemented as a network mail filter. We are implementing a network-level email filter that uses our algorithms to catch malicious executables before users receive them through their mail. We can either wrap the potential malicious executable or we can block it. This has the potential to stop some malicious executables in the network and prevent Denial of Service (DoS) attacks by malicious executables. If a malicious binary accesses a user's address book and mails

copies of itself out over the network, eventually most users of the LAN will clog the network by sending each other copies of the same malicious executable. This is very similar to the old Internet Worm attack. Stopping the malicious executables from replicating on a network level would be very advantageous.

Since both the Naive Bayes and Multi-Naive Bayes methods are probabilistic we can also tell if a binary is *borderline*. A borderline binary is a program that has similar probabilities for both classes (i.e., could be a malicious executable or a benign program). If it is a borderline case we have an option in the network filter to send a copy of the malicious executable to a central repository such as CERT. There, it can be examined by human experts.

9.1 Future Work

Future work involves extending our learning algorithms to better utilize byte-sequences. Currently, the Multi-Naive Bayes method learns over sequences of a fixed length, but we theorize that rules with higher accuracy and detection rates could be learned over variable length sequences. There are some algorithms such as Sparse Markov Transducers [7] that can determine how long a sequence of bytes should be for optimal classification.

We are planning to implement the system on a network of computers to evaluate its performance in terms of time and accuracy in real world environments. We also would like to make the learning algorithms more efficient in time and space. Currently, the Naive Bayes methods have to be run on a computer with one gigabyte of RAM.

Finally, we are planning on testing this method over a larger set of malicious and benign executables. Only when testing over a significantly larger set of malicious executables can we fully evaluate the method. In that light, our current results are preliminary. In addition with a larger data set, we plan to evaluate this method on different types of malicious executables such as macros and Visual Basic scripts.

References

- [1] William Arnold and Gerald Tesauro. Automatically Generated Win32 Heuristic Virus Detection. *Proceedings of the 2000 International Virus Bulletin Conference*, 2000.
- [2] Fred Cohen. *A Short Course on Computer Viruses*. ASP Press, 1990.
- [3] William Cohen. Learning Trees and Rules with Set-Valued Features. *American Association for Artificial Intelligence (AAAI)*, 1996.
- [4] R. Crawford, P. Kerchen, K. Levitt, R. Olsson, M. Archer, and M. Casillas. Automated Assistance

- for Detecting Malicious Code. *Proceedings of the 6th International Computer Virus and Security Conference*, 1993.
- [5] Cygnus. GNU Binutils Cygwin. *Online publication*, 1999. <http://sourceware.cygwin.com/cygwin>.
- [6] D.Michie, D.J.Spiegelhalter, and C.C.TaylorD. *Machine learning of rules and trees. In Machine Learning, Neural and Statistical Classification*. Ellis Horwood, 1994.
- [7] Eleazar Eskin, William Noble Grundy, and Yoram Singer. Protein Family Classification using Sparse Markov Transducers. *Proceedings of the Eighth International Conference on Intelligent Systems for Molecular Biology*, 2000.
- [8] Dmitry Gryaznov. Scanners of the Year 2000: Heuristics. *Proceedings of the 5th International Virus Bulletin*, 1999.
- [9] Jeffrey O. Kephart and William C. Arnold. Automatic Extraction of Computer Virus Signatures. *4th Virus Bulletin International Conference*, pages 178–184, 1994.
- [10] P. Kerchen, R. Lo, J. Crossley, G. Elkinbard, and R. Olsson. Static Analysis Virus Detection Tools for UNIX Systems. *Proceedings of the 13th National Computer Security Conference*, pages 350–365, 1990.
- [11] Zou KH, Hall WJ, and Shapiro D. Smooth non-parametric ROC curves for continuous diagnostic tests. *Statistics in Medicine*, 1997.
- [12] R Kohavi. A study of cross-validation and bootstrap for accuracy estimation and model selection. *IJCAI*, 1995.
- [13] W. Lee, S. J. Stolfo, and P. K. Chan. Learning patterns from UNIX processes execution traces for intrusion detection. *AAAI Workshop on AI Approaches to Fraud Detection and Risk Management*, pages 50–56. AAAI Press, 1997.
- [14] Wenke Lee, Sal Stolfo, and Kui Mok. A Data Mining Framework for Building Intrusion Detection Models. *IEEE Symposium on Security and Privacy*, 1999.
- [15] R.W. Lo, K.N. Levitt, and R.A. Olsson. MCF: a Malicious Code Filter. *Computers & Security*, **14**(6):541–566., 1995.
- [16] MacAfee. Homepage - MacAfee.com. *Online publication*, 2000. <http://www.mcafee.com>.
- [17] Microsoft. Portable Executable Format. *Online publication*, 1999. <http://support.microsoft.com/support/kb/articles/Q121/4/60.asp>.
- [18] Peter Miller. Hexdump. *Online publication*, 2000. <http://www.pcug.org.au/millerp/hexdump.html>.
- [19] MIT Lincoln Labs. 1999 DARPA intrusion detection evaluation.
- [20] Tom Mitchell. *Machine Learning*. McGraw Hill, 1997.
- [21] Kamal Nigam, Andrew McCallum, Sebastian Thrun, and Tom Mitchell. Learning to Classify Text from Labeled and Unlabeled Documents. *AAAI-98*, 1998.
- [22] Wildlist Organization. Virus descriptions of viruses in the wild. *Online publication*, 2000. <http://www.f-secure.com/virus-info/wild.html>.
- [23] REUTERS. Microsoft Hack Shows Companies Are Vulnerable. *New York Times*, October 29, 2000.
- [24] Eugene H. Spafford. The Internet worm program: an analysis. *Tech. Report CSD-TR-823*, 1988. Department of Computer Science, Purdue University.
- [25] Gerald Tesauro, Jeffrey O. Kephart, and Gregory B. Sorkin. Neural Networks for Computer Virus Recognition. *IEEE Expert*, **11**(4):5–6. IEEE Computer Society, August, 1996.
- [26] Steve R. White. Open Problems in Computer Virus Research. *Virus Bulletin Conference*, 1998.
- [27] Steve R. White, Morton Swimmer, Edward J. Pring, William C. Arnold, David M. Chess, and John F. Morar. Anatomy of a Commercial-Grade Immune System. *IBM Research White Paper*, 1999. <http://www.av.ibm.com/ScientificPapers/White/Anatomy/anatomy.html>.