# Feature Analysis of Encrypted Malicious Traffic

Anish Singh Shekhawat[*][†]  Fabio Di Troia[*][‡]  Mark Stamp[*][§]

**Abstract**

In recent years there has been a dramatic increase in the number of malware attacks that use encrypted HTTP traffic for self-propagation or communication. Antivirus software and firewalls typically will not have access to encryption keys, and therefore direct detection of malicious encrypted data is unlikely to succeed. However, previous work has shown that traffic analysis can provide indications of malicious intent, even in cases where the underlying data remains encrypted. In this paper, we apply three machine learning techniques to the problem of distinguishing malicious encrypted HTTP traffic from benign encrypted traffic and obtain results comparable to previous work. We then consider the problem of feature analysis in some detail. Previous work has often relied on human expertise to determine the most useful and informative features in this problem domain. We demonstrate that such feature-related information can be obtained directly from machine learning models themselves. We argue that such a machine learning based approach to feature analysis is preferable, as it is more reliable, and we can, for example, uncover relatively unintuitive interactions between features.

# 1  Introduction

Malicious software, or malware, can be defined as a program that is designed to damage a computer system (Aycock, 2006). Malware is, arguably, the greatest threat to information security today.

It is estimated that more than 90% of small-to-medium sized businesses have recently experienced an increase in the number of malware detected, with some experiencing an increase of 500% in March 2017 alone (Malwarebytes, 2017). Real-time malware detection based on network traffic has the potential to greatly reduce malware propagation on the network.

[*]Department of Computer Science, San Jose State University, San Jose, California, USA
[†]anish.s.shekhawat@gmail.com
[‡]fabio.ditroia@sjsu.edu
[§]mark.stamp@sjsu.edu

One approach to detect network-based malware is to use deep packet inspection (DPI). In DPI, packets are aggregated and the content analyzed to check for signatures or other characteristics that can be used to classify the data as malicious or benign (Sen, Spatscheck, & Wang, 2004). Unfortunately, due to the widespread use of the HyperText Transfer Protocol Secure (HTTPS), or HTTP over Secure Socket Layer (SSL), straightforward deep packet inspection methods can be inadequate to classify network traffic. It is estimated that HTTPS is used in more than 70% of Internet traffic today (Google, 2017).

Since HTTPS traffic is encrypted, it cannot be analyzed to the same degree as plaintext traffic. This is a benefit if the analyzer is a potential eavesdropper or attacker, but it is harmful when firewalls are unable to analyze the traffic, since malware can leverage encryption to evade detection. According to a recent report (Anderson, Paul, & McGrew, 2016), there is a steady 10% to 12% annual increase in encrypted malicious network traffic over HTTPS. The 2017 Global Application & Network Security Report (Radware, 2018) states that 35% of the organizations surveyed faced TLS or SSL based attacks, which represents an increase of 50% over the previous year.

The purpose of this research is to analyze various features that are commonly used to distinguish encrypted malicious network traffic from encrypted benign traffic, in cases where the decryption keys are unavailable. Specifically, we employ machine learning to analyze encrypted network traffic features. Again, the emphasis here is on feature analysis, based on trained machine learning models. We find that feature analysis based on machine learning models can provide at least as much useful information—with respect to individual features—as human experts can provide. We also show that automated feature analysis can easily uncover less intuitive aspects of features.

The remainder of this paper is organized as follows. Section 2 gives an overview of selected previous work related to the problem of detecting malicious traffic, with the emphasis on encrypted traffic. In Section 3, we introduce the datasets used in our experiments, while Section 4 covers our proposed methodology. Section 5 gives our experimental results and analysis, where the focus is on feature analysis. Finally, in Section 6, we conclude the paper and outline some areas for future work.

## 2   Related Work

Malicious network communication detection typically relies on either port-based classification or deep packet inspection and signature matching. Port-based methods inspect Transmission Control Protocol (TCP) and User Datagram Protocol (UDP) port numbers under the assumption that applications use well-known port numbers (Yoon, Park, Park, Oh, & Kim, 2009), which are

assigned by the Internet Assigned Numbers Authority (IANA) (Touch et al., 2018). Not surprisingly, malicious applications frequently use non-standard ports in an effort to evade such network intrusion detection systems (NIDS) and firewalls that rely on similar information (Dreger & Feldmann, 2006). Even legitimate applications such as Skype use dynamic port numbers to overcome restrictive firewall policies (Baset & Schulzrinne, 2006). In (Madhukar & Williamson, 2006) it is shown that port-based classification misclassifies network flow traffic at an astonishingly high rate, estimated to be at least 30%, and possibly as high as 70%.

The authors of (Etienne, 2009) use deep packet inspection to detect malicious traffic by considering payload content and using traditional pattern matching or signature based techniques. This research relies on Snort (Cisco, 2018), an IDS, to detect malicious traffic, and uses signature or string matching on the packet contents. Snort also hosts a popular Intrusion Protection System (IPS) rule set. However, only a minuscule percentage of the rules in Snort are TLS-specific, which indicates that such pattern matching techniques are not likely to be effective for TLS based malware (Bakhdlaghi, 2017).

The paper (Sen et al., 2004) demonstrates the use of deep packet inspection to reduce false positive and false negative rates when classifying peer-to-peer (P2P) traffic. In (Moore & Papagiannaki, 2005), the authors achieve 100% accuracy when identifying network applications, based on an analysis of the entire packet payload. A major limitation of such methods is the overhead of decrypting and analyzing each packet. Of course, these techniques are of not applicable if the decryption key is not available to the IDS or IPS.

BotFinder (Tegeler, Fu, Vigna, & Kruegel, 2012) is a network flow information analyzer that is used to detect bot infections. The system relies on chronologically-ordered flows (or traces) to find irregularities in the network behavior between two endpoints. This information, along with other network metadata, is used as features in a clustering based algorithm (Wang, Qiu, & Zamar, 2007). In (Prasse, Machlica, Pevný, Havelka, & Scheffer, 2017) a neural network based malware detection approach is developed, which relies on various network flow features. The authors of (Lokoc, Kohout, Cech, Skopal, & Pevný, 2016) present a $k$-nearest neighbor ($k$-NN) based classification strategy that is claimed to accurately identify malware utilizing HTTPS traffic, at least in some specific cases.

The problem of detecting encrypted Skype traffic is considered in the paper (Di Mauro & Di Sarno, 2018). In this work, detection is based on a majority vote of three classifiers, namely, a decision tree-based classifier, a logistic classifier, and a Bayesian network classifier. The system is practical and it appears that relatively strong detection results are obtained. Although Skype traffic is not malware, some of the issues that arise when trying to distinguish encrypted Skype traffic are relevant when attempting to detect encrypted malicious traffic.

In (Anderson & McGrew, 2017), a technique is proposed to identify encrypted malware traffic based on network flow metadata, using supervised machine learning. These authors rely on a complex demilitarized zone (DMZ) architecture to collect the necessary data for training their machine learning algorithms, The traffic collected may not be representative of general network traffic, as this data represents only enterprise users. Interestingly, (Anderson & McGrew, 2017) rely heavily on human expertise to determine the most important features, which we view as a significant limitation. Human expertise is not always available, humans are fallible, and it is possible that non-intuitive features (or non-obvious combinations of features) can be as significant to machine learning models—if not more so—than a collection of seemingly intuitive features.

This paper further explores the use of network flow information as a basis for machine learning models that are designed to detect attacks that rely on encrypted traffic. The emphasis here is squarely on machine learning—to not only detect and distinguish encrypted attack traffic, but to analyze the information that is available from the trained models. We show that we can obtain detection results that are comparable to previous work. We then consider feature analysis in some depth, where our feature analysis is based on trained machine learning models. This is in contrast to previous work, which tends to rely primarily on domain-specific expertise provided by humans for feature analysis. We believe that there is significant value in a completely domain-free approach, such as we consider here, where essentially no human intervention or domain-specific expertise is required. A fully automated system that can rapidly adapt to new and changing circumstances would almost certainly require such a domain-free approach. In fact, we find that the most informative features differ depending on the machine learning model used, and some of these features are not among the most intuitive. Feature interactions are also difficult to account for when relying on human expertise.

# 3   Dataset

In this research, we have used the following two previously published network capture dataset, each of which contains both malicious and benign traffic.

- The CTU-13 dataset (García, Grill, Stiborek, & Zunino, 2014) was captured as part of research project at the Czech Technical University. This dataset includes 13 malware traffic captures, consisting of both benign and malware traffic. The malware traffic was captured by executing selected malware in a Windows virtual machine and recording the network traffic generated on the host. The benign traffic was captured on benign hosts, i.e., hosts that were not infected with malware. These network captures are available as `pcap` files.

- The Malware Capture Facility Project dataset is from another research project carried out by the Czech Technical University ATG Group to capture, analyze, and publish real malware network traffic (Erquiaga & Garcia, 2013). In this case, the malware was executed with two restrictions, namely, a bandwidth limit and spam interception. An interesting characteristic of this dataset is that the malware was allowed to execute over a long period of time—up to several months in some cases. Again, the traffic is stored in `pcap` files, and is labeled for ease of use.

Our combined dataset includes a total of 72 captures, of which 59 represent malware while the remaining 13 are benign. Tables 1 and 2 give basic statistics for the connections and flows, respectively, in this combined dataset.

**Table 1:** Dataset connections

| Feature | Count |
| --- | --- |
| Benign connection 4-tuples | 8828 |
| Malicious connection 4-tuples | 52898 |
| Total connection 4-tuples | 61726 |

**Table 2:** Dataset flows

| Feature | Count |
| --- | --- |
| Benign flows | 69358 |
| Malicious flows | 1067273 |
| Total flows | 1136631 |

Each capture is contained in a `pcap` file, and includes a list of infected and benign hosts, as well as Bro IDS logs that were generated from the `pcap` files. Bro (Bro Project Team, 2014) is a powerful open-source network analysis tool that supports various features for traffic inspection, log recording, and attack detection. We use Bro to generate network traffic logs that include network flow information and other metadata. This information is then used to extract various features related to traffic flows. We use the resulting features to train and test our machine learning models.

Bro generates several types of log files. In this research, we only need to make use of the following Bro log files.

**conn.log** — This log file contains information about TCP, UDP, and ICMP connections.

**ssl.log** — The ssl.log file contains information related to SSL/TLS certificates and sessions.

(a) conn.log



(b) ssl.log



(c) x509.log

**Figure 1:** Interconnection of log records using unique keys in Bro

**x509.log** — As the name suggests, this log file contains relevant information about X.509 certificates.

## 3.1 Feature Extraction

We extract several features from Bro logs generated from the network captures under consideration. Note that features related to a single connection are generally spread over different log files. For example, if there is an SSL connection to a specific server, the connection features (e.g., source and destination IP, ports, protocols, duration) are stored in the connection log (conn.log). On the other hand, many SSL-specific features (e.g., cipher used, server name) are stored in the SSL log (ssl.log), while certificate features (e.g., key lengths, common names) are stored in the certificate log (x509.log). Bro enables us to deal effectively with the interconnections between the various logs, as illustrated in Figure 1. Thus, it is relatively easy to extract the features that correspond to a given flow.

Bro tracks every incoming and outgoing connection in conn.log, and hence each record in this log file provides information about a specific connection. Since we are interested only in encrypted network connections, we only consider connection records that are related to HTTPS. And, since HTTPS connections use SSL/TLS to establish an encrypted link between the client and the server, we only extract connection records that have corresponding entries in the ssl.log file.

Every SSL/TLS connection requires a server certificate (Freier, Karlton, & Kocher, 2011), and hence every record in ssl.log contains at least one unique certificate ID that the server uses to validate its signing chain. Unique certificate IDs are used to represent the certificate record in the x509.log file. We extract only the first certificate ID from the ssl.log file, since it corresponds to the end-user certificate, while the remaining IDs correspond to intermediate

and root certificates.

Every connection record can be identified by a 4-tuple of the form

$$(\text{source IP}, \text{destination IP}, \text{destination port}, \text{protocol}).$$

We use these 4-tuples as keys to group network features by connection.

## 3.2   Features

We use Bro to extract a wide variety of features from the connection, SSL, and certificate log files. At this point, we want as many features as we can reasonable collect. Subsequently, we perform feature analysis to determine the relative importance of these features, as well as considering interactions between features.

Table 3 lists the specific features we extract from the conn.log, the ssl.log, and x509.log files. Note that in this table, we use $\mu$ to represent the mean and $\sigma$ for the standard deviation. Again, all features are determined over a single connection, and aggregated based on a specific 4-tuple, as discussed above.

## 3.3   Labels

The datasets from (García et al., 2014) and (Erquiaga & Garcia, 2013) contains IP addresses of infected and benign hosts. We use these IP addresses to label our dataset. That is, if a connection record has an infected source IP address, then the record is labeled as malware; otherwise it is labeled as benign. Of course, these labels are used in both the training phase and in the testing phase.

# 4   Methodolgy

We experiment with three machine learning algorithms. Specifically, we consider experiments involving support vector machines (SVM), random forests (RF), and extreme gradient boosting (XGBoost). Next, we briefly discuss each of these machine learning techniques, followed by an introduction to the evaluation metrics that we use to quantify our experiments results.

## 4.1   Support Vector Machines

According to (Bennett & Campbell, 2000), support vector machines "are a rare example of a methodology where geometric intuition, elegant mathematics, theoretical guarantees, and practical algorithms meet." In particular, the

**Table 3:** Extracted features from conn.log, ssl.log, and x509.log

| Num | Feature Name | Log | Description |
|---|---|---|---|
| 1 | no_of_flows | conn | Number of records in 4-tuples |
| 2 | avg_of_duration | conn | $\mu$ duration of connections |
| 3 | standard_deviation_duration | conn | $\sigma$ of connections |
| 4 | percent_sd_of_duration | conn | Percent exceeding $\sigma$ |
| 5 | size_of_orig_flows | conn | Bytes sent by the originator |
| 6 | size_of_resp_flows | conn | Bytes sent by the responder |
| 7 | ratio_of_sizes | conn | Ratio of responder bytes |
| 8 | percent_of_established_states | conn | Percent established connections |
| 9 | inbound_pckts | conn | Number of incoming packets |
| 10 | outbound_pckts | conn | Number of outgoing packets |
| 11 | periodicity_average | conn | $\mu$ periodicity of connection |
| 12 | periodicity_standard_deviation | conn | $\sigma$ of connection periodicity |
| 13 | ssl_ratio | ssl | Ratio SSL to non-SSL records |
| 14 | avg_key_len | ssl | Average key length |
| 15 | tls_version_ratio | ssl | Ratio of records with TLS |
| 16 | avg_of_certificate_len | x509 | Average length |
| 17 | standart_deviation_cert_len | x509 | $\sigma$ certificate length |
| 18 | is_valid_certificate | x509 | 1 if certificate is valid |
| 19 | amount_diff_certificates | x509 | Number of certificates |
| 20 | no_of_domains_in_cert | x509 | Number of domains |
| 21 | no_of_cert_path | x509 | Number of signed paths |
| 22 | x509_ssl_ratio | x509 | Ratio of SSL logs with x509 |
| 23 | SNI_ssl_ratio | ssl | Ratio with SNI in SSL record |
| 24 | self_signed_ratio | ssl | Ration with self-signed certificate |
| 25 | is_SNIs_in_SAN_dns | x509 | Check if SNI is SAN DNS |
| 26 | is_CNs_in_SAN_dns | x509 | 1 if all common names in SAN |
| 27 | differ_SNI_in_ssl_log | ssl | Ratio SSL with different SNI |
| 28 | differ_subject_in_ssl_log | ssl | Ratio SSL with different subject |
| 29 | differ_issuer_in_ssl_log | ssl | Ratio SSL with different issuer |
| 30 | differ_subject_in_cert | x509 | Ratio subjects |
| 31 | differ_issuer_in_cert | x509 | Ratio issuers |
| 32 | differ_sandns_in_cert | x509 | Ratio SAN DNS |
| 33 | ratio_of_same_subjects | ssl | Ratio SSL with same subject |
| 34 | ratio_of_same_issuer | ssl | Ratio SSL with same issuer |
| 35 | is_same_CN_and_SNI | x509 | Checks if CN is same as SNI |
| 36 | average_certificate_expo | x509 | Average exponent |
| 37 | is_SNI_top_level_domain | ssl | 1 if SNI is a top level domain |
| 38 | ratio_certificate_path_error | x509 | Check if path is valid |

geometric intuition that underlies the topic is especially helpful for understanding the basic concepts behind SVMs.

For binary classification, the key ideas behind SVMs are the following.

- Separating hyperplane — The goal is to separate the labeled training data into two classes based on a hyperplane.
- Maximize the margin — When constructing the separating hyperplane, we maximize the "margin," that is, we maximize the minimum separation between the two classes.
- Work in a higher dimensional space — We often try to reduce the dimensionality of data, due to the so-called curse of dimensionality. However, when training an SVM, it is usually beneficial to work in a higher dimensional space. By moving the problem to a higher dimension, we

have more space available, and hence there is a better chance of finding a separating hyperplane.

- Kernel trick — In an SVM, we use a kernel function to transform the data, typically, to a higher dimensional space, with the goal of obtaining better separation. The "trick" lies in the fact that we pay essentially no performance penalty for working in this higher dimensional space.

By choosing a separating hyperplane that maximizes the margin, we give ourselves the largest possible margin for error, with respect to the training data. This assumes that errors are equally likely (and equal in magnitude) in either direction.

SVM training relies on two techniques to deal with training data that is not linearly separable. A so-called "soft" margin allows for some classification errors when determining a separating hyperplane. The more classification errors we can tolerate, in general, the larger the margin. When training an SVMs, there is a user-defined parameter that specifies the allowable "softness" of the margin.

Another technique employed in SVM training is to map the input data to a feature space, where the problem of constructing a separating hyperplane is more tractable. As previously mentioned, this generally involves transforming the input data to a feature space of higher dimension. In a higher dimensional space, it is more likely that a separating hyperplane can be found.

For more information on SVMs, see, for example, (Cristianini & Shawe-Taylor, 2000; Stamp, 2017b). In particular, (Bennett & Campbell, 2000) is highly recommended.

## 4.2   Random Forest

A random forest can be viewed as a generalization of a decision tree. To illustrate a decision tree, suppose that we have a labeled training set consisting of malware samples and benign samples. From this training set, we observe that malware samples tend to be smaller in size and have higher entropy, as compared to benign samples. We could use this information to construct the decision tree in Figure 2, where the thresholds for "large" versus "small" (size) and "high" versus "low" (entropy) would be based on the training data. This decision tree could then be used to classify any sample as either malware or benign, based on its size and entropy.

We might want to consider the features in a different order. For example, the two features of file size and entropy, as illustrated in Figure 2, could instead be considered in the opposite order, as illustrated in Figure 3.

In general, splits made closer to the root of the tree will tend to have more impact on the final classification. Therefore, we want to make decisions that are based on the most distinguishing features closer to the root of the decision tree. In this way, the decisions for which the training data is less useful are
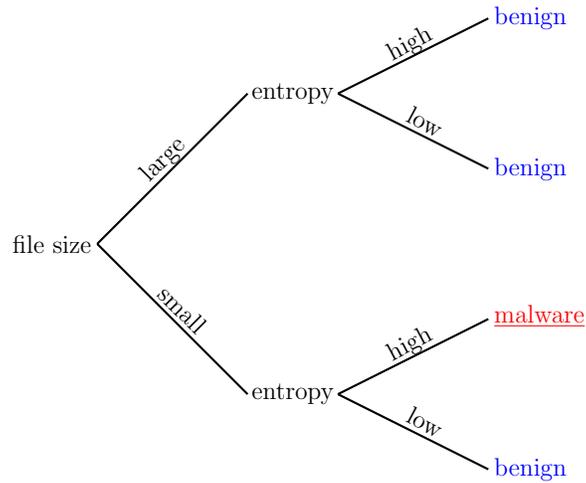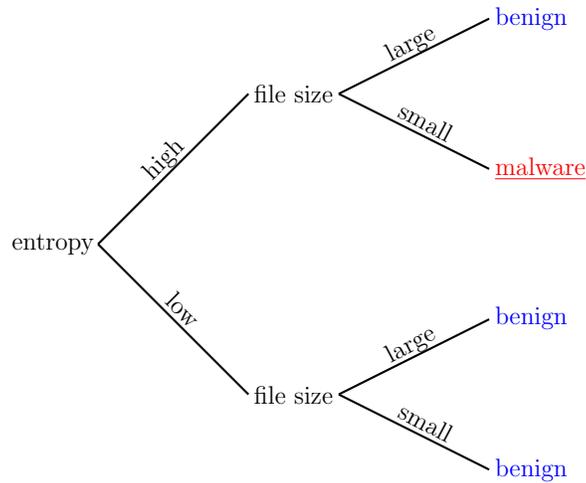
**Figure 2:** Decision tree example



**Figure 3:** Features in different order

made later in the process, where such decisions will have less influence on the final classification.

Information gain is defined as the expected reduction in entropy when we branch on a given feature. In the context of a decision tree, information gain can be computed as the entropy of the parent node minus the average weighted entropy of its child nodes. We can measure the information gain for each feature, and select features in a greedy manner. In this way, features with the highest gain will be closest to the root. This is desirable, since the resulting tree will reduce the entropy as rapidly as possible, which enables us to simplify the tree by trimming features that provide little or no gain.

Among the advantages of decision trees are simplicity and clarity, where "clarity" means that the tree itself is informative—such is not always the case for machine learning models. Of course, there are also some disadvantages to decisions trees, chief of which is a tendency to overfit the training data. That is, the simple and intuitive nature of the tree structure captures the information in the training data too well, in the sense that it does not generalize. In machine learning, this is undesirable, as the training data is only a representative sample, and we want our models to capture the significant properties of this data.

One way to improve on a decision tree is to train multiple trees. We can select different (overlapping) subsets of the training data and construct a tree for each subset, then use a majority vote of the resulting trees to determine the classification. This process is known as bagging the observations, and it is less likely to overfit, since it tends to better generalize the training data.

In a random forest, the idea of bagging is taken one step further—in addition to bagging the observations, we also bag the features. That is, we construct multiple decision trees from selections (with replacement) of the training data, and also for various selections of the features (i.e., various subsets and orderings of the features). Then for classification, we combine the output from all of the resulting decision trees to obtain the final classification. For example, we could generate $t$ decision trees, where $t$ is odd, with each based on a different subset of features and data. Then we could use a simple majority vote of the resulting trees for classification.

As previously mentioned, a significant advantage of a random forest is that it is not prone to overfitting, as compared to a simple decision tree. However, random forests do lose some of the inherent simplicity and intuitiveness of decision trees.

We note in passing that there is a deep connection between $k$-nearest neighbors ($k$-NN) and random forests. Both $k$-NN and decision trees are neighborhood-based classification algorithms, but with different neighborhood structures. Since a random forest is a collection of decision trees, a random forest is also a neighborhood-based classification technique, but with a fairly complex neighborhood structure (Stamp, 2017b).

For the basics on random forests, a good source is (Breiman & Cutler, 2002), while (Liaw & Wiener, 2002) also provides a gentle and practical introduction to the topic. For an in-depth discussion of the connection between random forests and $k$-NN (and related algorithms), see (Lin & Jeon, 2002) or (Breiman & Cutler, 2002).

## 4.3  XGBoost

Boosting is a process whereby multiple (weak) classifiers are combined into one, much stronger classifier (Stamp, 2017b). Of course, many machine learn-
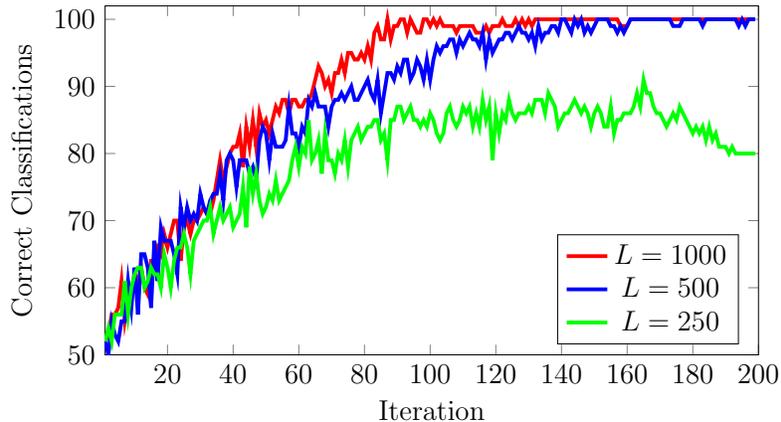
**Figure 4:** Correct classifications vs iteration ($n = 100$)

ing techniques can be applied in a somewhat analogous fashion. For example, SVMs are frequently used to construct a classifier based on a collection of other machine learning based scores (Singh, Di Troia, Visaggio, Austin, & Stamp, 2016). The advantage of boosting is that the individual classifiers can be extremely weak—anything that is better than a coin flip can be used. And, provided that we have a sufficient number of usable classifiers, boosting enables us to construct an arbitrarily strong classifier.

The best-known boosting algorithm is AdaBoost which is shorthand for "adaptive boosting." At each iteration of AdaBoost, we use a greedy strategy, in the sense that we select the individual classifier, and its associated weight, that improves our overall classifier the most. AdaBoost is an adaptive approach, since we build each intermediate classifier based on the classifier that was determined at the previous step of the algorithm. It is worth noting that AdaBoost is not a hill climb—we select the best available classifier at each step in a greedy manner, but there is no guarantee that this selection will improve our overall classifier.

The results of a simple boosting experiments are given in (Stamp, 2017a). In this particular experiment, 1000 classifiers are available, each of which is extremely weak, with each classifier only providing the correct classification 51% to 52% of the time. The results of the boosting experiments in (Stamp, 2017a) are reproduced here in Figure 4, where $L$ is the number of classifiers used. The red line illustrates the case where $L = 1000$ classifiers are used, that is, all of the available classifiers are used. The blue line represents the case where $L = 500$ of the 1000 available classifiers are used, and the green line is the case where only $L = 250$ of the 1000 classifiers are used. In each case, the classification accuracy of the intermediate (boosted) classifiers are graphed for 200 iterations of AdaBoost.

From Figure 4, we see that with $L = 1000$ classifiers, we can obtain ideal

accuracy using only about 100 of the available classifiers. On the other hand, with $L = 500$ classifiers available, we need about 140 iterations before we achieve ideal classification, and with only $L = 250$ weak classifiers, we never achieve more than about 90% accuracy. This illustrates that may need a large number of (weak) classifiers to achieve impressive results from a boosting algorithm.

Gradient boosting relies on a gradient descent approach to speed convergence and thereby improve performance. Extreme gradient boosting, or XGBoost, is a highly optimized version of gradient boosting (Chen & Guestrin, 2016). While AdaBoost is well-known and relatively simple, XGBoost has some potential advantage, in terms of efficiency, and additional flexibility in selecting a cost function (Chen & Guestrin, 2016).

## 4.4    Cross Validation

Cross validation consists of partitioning the available data into $n$ equal sized subsets and training $n$ models, where a different subset is reserved for testing in each of these $n$ "folds." Cross validation serves to minimize the effect of bias in the training data, while also maximizing the number of independent tests on the available data. In this paper, we use 10-fold cross validation for all experiments.

## 4.5    Evaluation Metrics

We use accuracy as an evaluation metric for our classification experiments. For a given experiment on a labeled data set,

$$\text{accuracy} = \frac{\text{TP} + \text{TN}}{\text{TP} + \text{TN} + \text{FP} + \text{FN}}$$

where

$$\text{TP} = \text{true positives}$$
$$\text{TN} = \text{true negatives}$$
$$\text{FP} = \text{false positives}$$
$$\text{FN} = \text{false negatives}$$

That is, the accuracy is simply the ratio of correct classifications to the total number of classifications.

We also employ receiver operating characteristic (ROC) curves in our data analysis. ROC curves originated with radar engineers during World War II, but they are now widely used in many fields. Given a binary classifier, we construct an ROC curve by plotting the TPR versus the FPR as the threshold

varies through the range of data values. Equivalently, we plot $1 -$ specificity versus sensitivity as the threshold varies.

For any ROC curve, the area under the curve (AUC) varies between 0.0 and 1.0. An AUC of 1.0 indicates ideal separation, i.e., a threshold exists such that no false positives or false negatives occur. On the other hand, an AUC of 0.5 indicates that the corresponding binary classifier is no better than flipping a coin. The AUC can be interpreted as the probability that a randomly selected match case scores higher than a randomly selected nomatch case (Bradley, 1997). Whenever we have an AUC of $x < 0.5$, we can simply reverse the match and nomatch criteria to obtain a classifier with an AUC of $1.0 - x > 0.5$. That is, no binary classifier can do worse than flipping a fair coin.

We note in passing that it is sometimes claimed that the AUC tends to overstate the success of experiments. However, as previously mentioned, the AUC measures the probability that a randomly selected positive instance scores higher than a randomly selected negative instance—nothing more nor less. An advantage of the AUC is that it allows us to directly compare different experimental results, without reference to any specific threshold value.

# 5 Experiments

In this section, we present our experimental results. Our primary focus is on feature analysis, but as a verification of the soundness of our approach, we also show that we can obtain detection and classification results that are comparable to previous work. First, we train an SVM classifier and determine the resulting detection accuracy, and we perform feature analysis based on linear SVMs. Then we consider the accuracy of a random forest classifier and again perform feature analysis. As a third test case, we apply XGBoost to our data and determine the accuracy for this technique, and perform feature analysis yet again. As a final experiment, we consider a multiclass problem, where we attempt to classify malicious samples into their respective families. We also consider pairs of features and show that there are some interesting interactions among the features.

## 5.1 SVM Experiments

We conducted a variety of experiment based on support vector machines (SVM). For our first set of experiments, we use a linear kernel and achieve a detection accuracy slightly above 92%. This experiment shows that by using a linear SVM, we are able to successfully separate encrypted malicious network traffic from encrypted benign traffic with reasonably high accuracy.

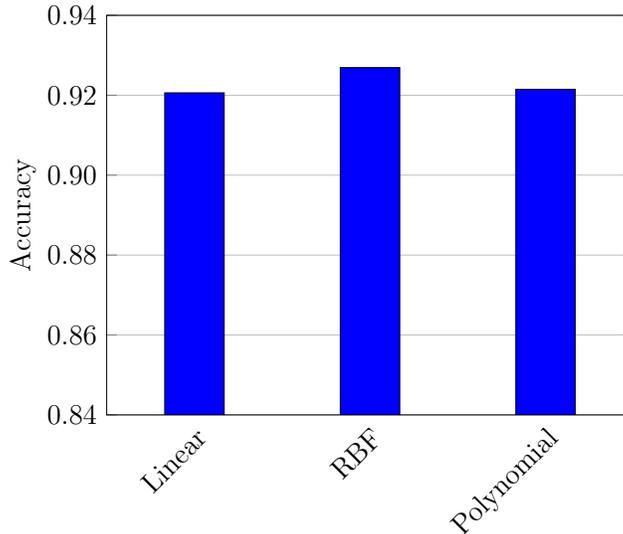We also experimented with other kernel functions, specifically, radial basis

**Figure 5:** SVM with different kernels

function (RBF) and polynomial kernels. These results—along with the linear kernel result—are presented in the form of a bar graph in Figure 5. We observe that the results are not substantially different for any of the kernels. In the remainder of this section, we only consider a linear kernel, as this greatly simplifies feature analysis.

A linear SVM assigns weights to each input feature, where the weight indicate the significance that the SVM places on that feature (Stamp, 2017b). In Figure 6, we see that the linear SVM assigned the highest weight to feature 12, where the features are numbered as in Table 3. Table 4 provides a listing of the top 15 of the 38 features, based on the linear SVM weights. From this ranking, we see that average certificate length, periodicity, and the average public key length are the most significant features, at least from the perspective of a linear SVM. From (Anderson & McGrew, 2016), we expect that malware uses weaker encryption techniques, such as shorter key lengths and that the traffic tends to be more periodic than other applications. Our SVM feature analysis confirms that the SVM does indeed find these same types of features to be most significant. In addition, the validity of certificate is the sixth highest rank feature, which tell us that malware often does not use a valid certificate. This is also consistent with manual analysis of the data.

Next, we consider recursive feature elimination (RFE). In RFE, we iteratively eliminate the lowest ranking feature, then train a new model on the reduced feature set. In this way, interactions between the features are accounted for, which is not the case with the feature ranking in Table 4.

From Figure 7, we see that by using just the top 6 features obtained from an RFE (based on linear SVMs), we can obtain results that are within 2% of
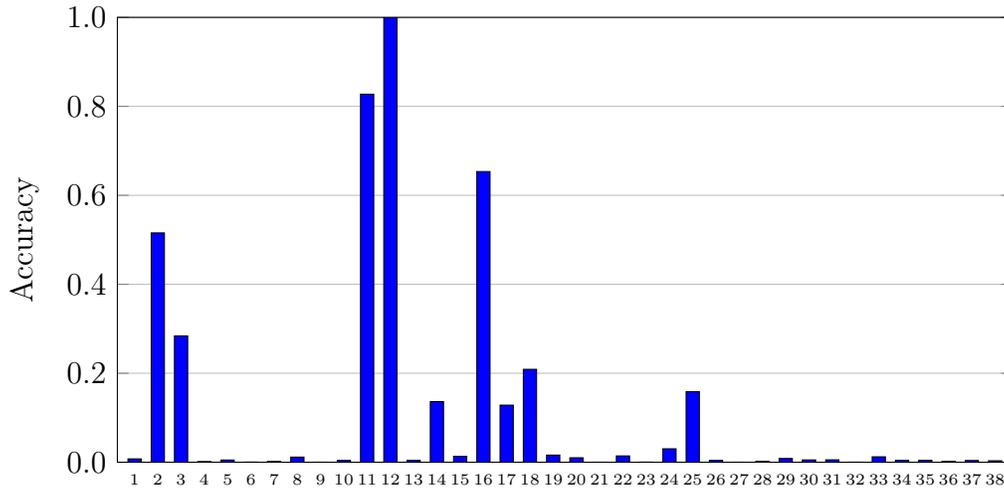
**Figure 6:** Weights assigned to features by SVM

**Table 4:** Top 15 features as ranked by linear SVM

| Rank | Feature |
| --- | --- |
| 1 | periodicity_standard_deviation |
| 2 | periodicity_average |
| 3 | avg_of_certificate_length |
| 4 | avg_of_duration |
| 5 | standard_deviation_duration |
| 6 | is_valid_certificate |
| 7 | is_SNIs_in_SAN_dns |
| 8 | avg_key_len |
| 9 | standard_deviation_cert_length |
| 10 | self_signed_ratio |
| 11 | amount_diff_certificates |
| 12 | x509_ssl_ratio |
| 13 | tls_version_ratio |
| 14 | ratio_of_same_subjects |
| 15 | percent_of_established_states |

those obtained using the full set of 38 features. Furthermore, if we use the top 10 features, our accuracy is within 1% of the the results obtained when using all 38 features. These results demonstrate that for this problem, a small number of features contain most of the discriminatory strength.

In Table 5, we have listed the top 15 of the 38 features, as determined by RFE, based on a linear SVM. Comparing these results to the ranked features
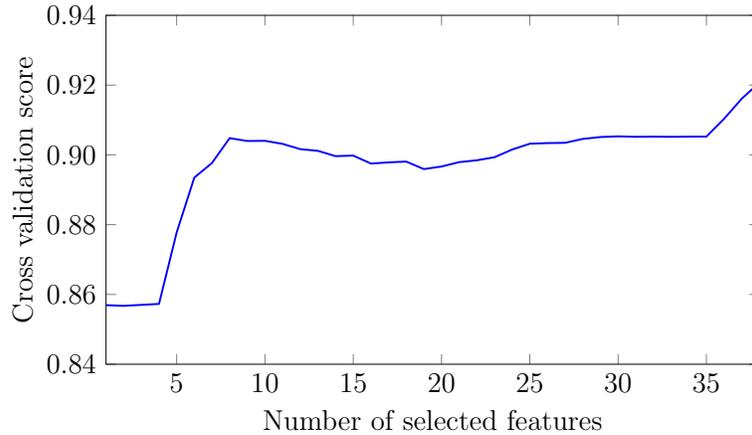
**Figure 7:** Recursive feature elimination with SVM

in Table 4, we see significant differences—for example, the top ranked feature by the linear weight ranking is only the fifth highest ranked feature according to the RFE. In addition, the 8th ranked feature according to the RFE does not appear among the to 15 features, based on a simple feature ranking. In general, we expect the RFE ranking to be more meaningful, as the RFE better accounts for interactions between features.

**Table 5:** Top 15 features by SVM RFE ranking

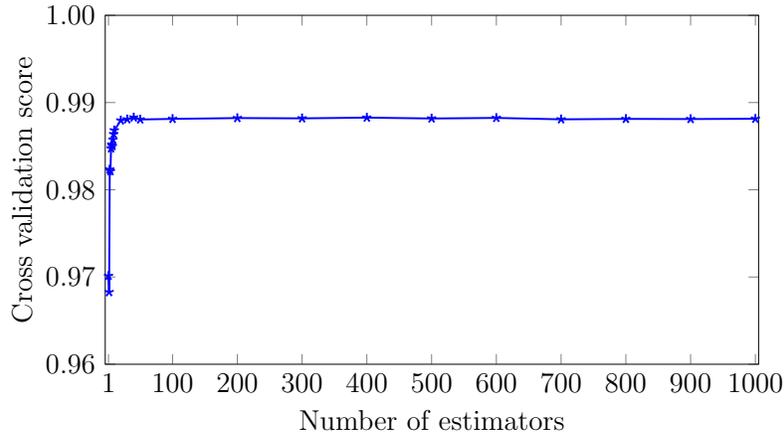| Rank | Feature |
|------|---------|
| 1 | `periodicity_standard_deviation` |
| 2 | `periodicity_average` |
| 3 | `avg_of_duration` |
| 4 | `is_SNIs_in_SAN_dns` |
| 5 | `avg_of_certificate_length` |
| 6 | `avg_key_len` |
| 7 | `self_signed_ratio` |
| 8 | `is_same_CN_and_SNI` |
| 9 | `percent_of_established_states` |
| 10 | `differ_issuer_in_ssl_log` |
| 11 | `differ_subject_in_cert` |
| 12 | `ratio_certificate_path_error` |
| 13 | `ratio_of_same_issuer` |
| 14 | `size_of_resp_flows` |
| 15 | `ratio_of_sizes` |

**Figure 8:** Random forest accuracy as a function of the number of estimators
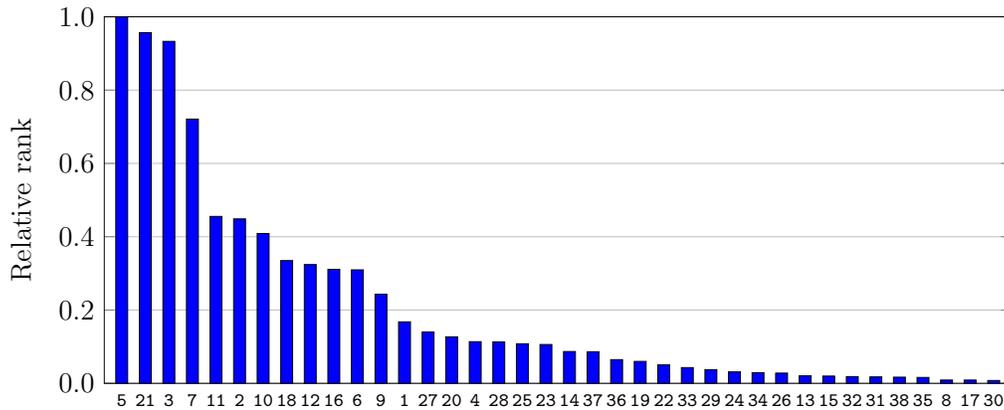


**Figure 9:** Random forest RFE feature rank

## 5.2 Random Forest Experiments

Next, we consider a random forest classifier. The key parameter for a random forest is the number of estimators, that is, the number of bagged decision trees (see Section 4.2 for a discussion of bagging). In Figure 8, we have graphed the accuracy of the random forest as a function of the number of estimators. Interestingly, only a relatively small number of estimators are needed to achieve good accuracy. However, to ensure optimal results, and since there is little additional training cost associated with a larger number of estimators, we use 500 estimators for all random forest experiments presented in this section.

We again performed recursive feature elimination, but based on a random forest classifier rather than a linear SVM. The features rankings obtained from this random forest RFE are given in Figure 9.
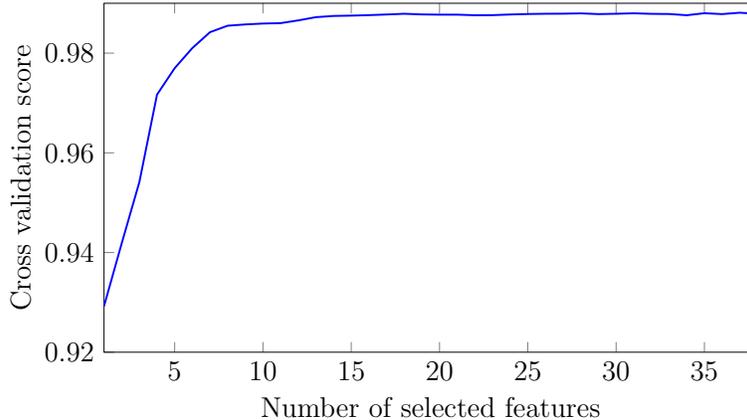
18

**Figure 10:** RFE results for random forest

We observe that the feature ranking in the case of random forest differ significantly from the rankings generated using a linear SVM, as discussed in Section 5.1. From Figure 9, we see that the top four features from the random forest classifier stand out from the remaining features—these four features are listed in Table 6. Interestingly, of these four dominant random forest features, only the fourth appears in the list of top ranked SVM features, as given in Table 5, and this feature only ranks as the $15^{\text{th}}$ best according to our SVM-based analysis.

**Table 6:** Top 4 features from random forest RFE

| Rank | Feature |
|------|---------|
| 1 | size_of_orig_flows |
| 2 | no_of_cert_path |
| 3 | standard_deviation_duration |
| 4 | ratio_of_sizes |

Using the top ranked features from the RFE, as listed in Figure 9, we obtain the accuracy results in Figure 10. In this case, we obtain a maximum accuracy of nearly 99%, which is significantly better than a linear SVM, and we observe that we have an accuracy in excess of 98.5% with just 6 features.

We note that the highly ranked random forest features are less intuitive than the SVM features, yet we can obtain stronger results with the random forest based on these features, as compared to an SVM. This nicely illustrates the importance of letting the model dictate features selection, rather than relying on the most intuitively appealing features.
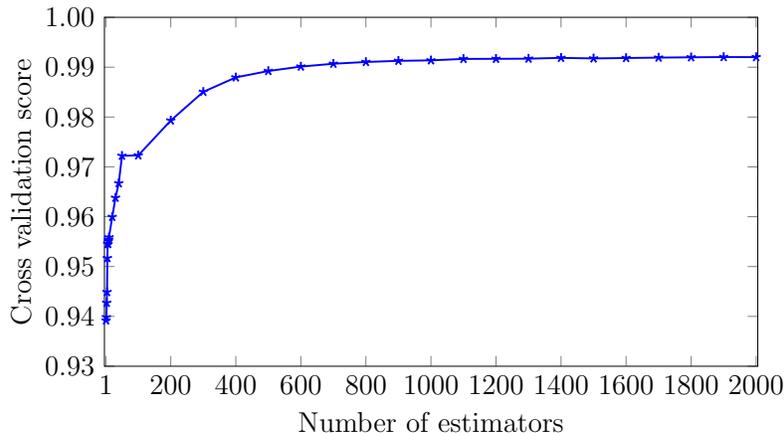
**Figure 11:** XGBoost accuracy as a function of the number of estimators

## 5.3   XGBoost Experiments

Finally, we construct classifiers using XGBoost, based on decision trees. In Figure 11, we give our XGBoost results as a function of the number of estimators (i.e., decision tress). In light of the discussion of boosting in Section 4.3, it should not be surprising that a large number of estimators are required before we obtain near-optimal accuracy. In all of the XGBoost experiments discussed in the remainder of this section, we use 1000 estimators.

We performed recursive feature elimination using XGBoost—these results are summarized in Figure 12. Note that the top ranked XGBoost features differ somewhat from those obtained with the random forest, as give in Figure 10. However, the feature rankings provided by our random forest and SGBoost experiments are more similar to each other than either is to the SVM rankings. That is, the SVM feature rankings can be viewed as the outlier among our three sets of experiments.

In Figure 13, we give the accuracy of XGBoost as a function of the top ranked features. As with a random forest, we see that a small number of features is sufficient to obtain near-optimal accuracy.

The results in Figure 13 serve to reinforce the point that the highly intuitive features determined by the SVM classifier are not the optimal set for constructing a classifier, particularly in case we want to minimize the number of features. In practice, we would want to reduce the number of features that need to be collected, as this will enable more efficient scoring. Efficiency is particularly important for the problem at hand, since we are dealing with network traffic that would need to be analyzed in real time.

Before briefly turning our attention to the multiclass problem, we summarize our SVM, random forest, and boosting experiments. A comparison of the classification accuracy for these machine learning algorithms is given in the
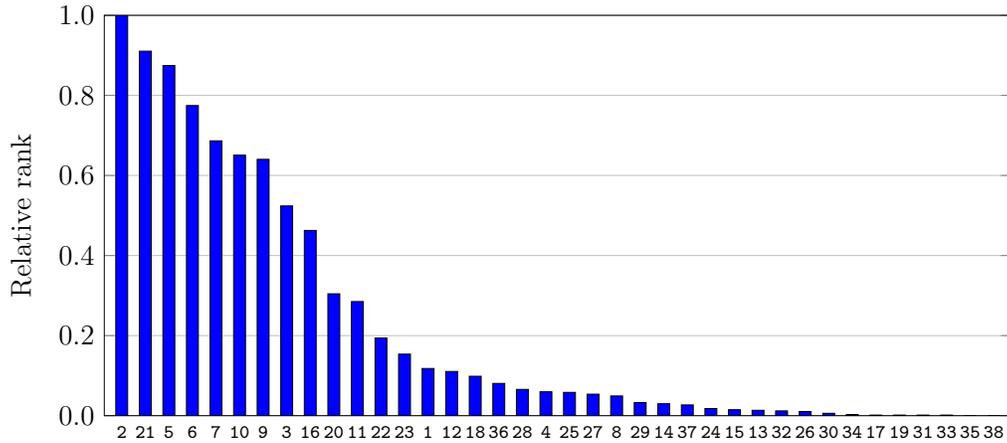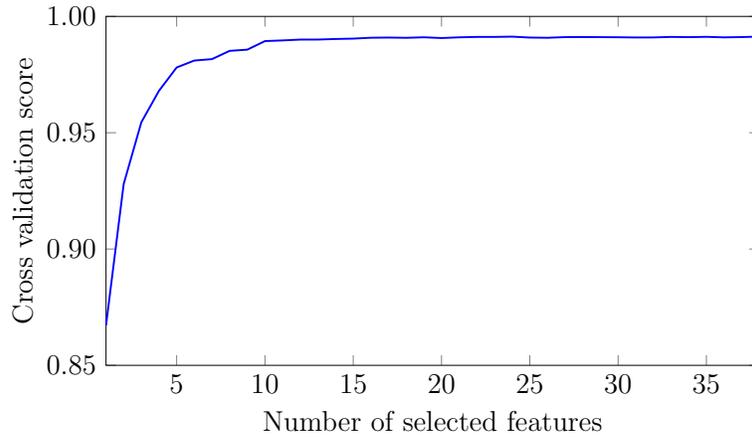
**Figure 12:** XGBoost feature importance



**Figure 13:** Recursive feature elimination with XGBoost

form of a bar graph in Figure 14. We see that XGBoost gives the highest accuracy, at 99.15%, while random forest is virtually indistinguishable at 98.78% accuracy, with a linear SVM performing significantly worse.

## 5.4 Multiclass Family Classification

For our multiclass experiments, we attempt to classify samples from four malware families, namely, Dridex, Trickbot, WannaCry, and Zbot. The data for these experiments is summarized in Table 7.

We trained models for all six pairs of the four malware families listed in Table 7, and we tested each model independently. These results are summarized in Figure 15, where we see that a random forest performs better than either an SVM or XGboost in most cases, with XGBoost giving us higher
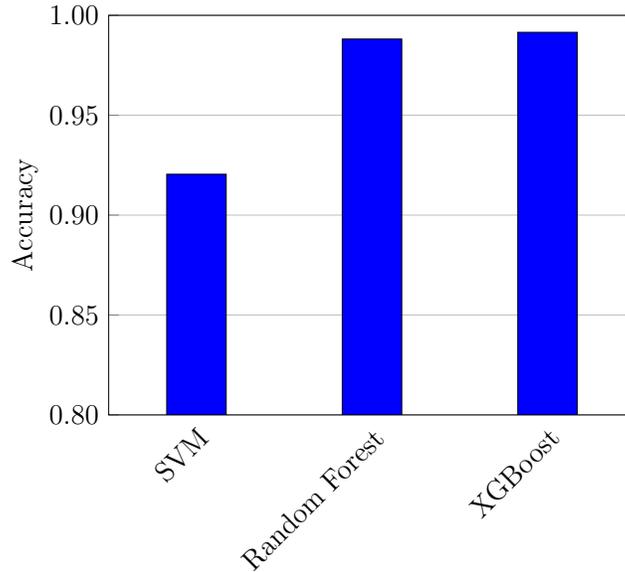
**Figure 14:** Accuracy comparison of machine learning algorithms

**Table 7:** Malware families

| Family | Connection 4-tuples | Flows |
|--------|:---:|---:|
| Dridex | 24 | 65465 |
| Trickbot | 217 | 465289 |
| WannaCry | 34 | 785 |
| Zbot | 45 | 1788 |

accuracy only for the Dridex versus Trickbot case. However, the differences are relatively small in every case.

We also trained true multiclass models, that is, we comsidered data from all four families together for training. The results of these experiments are given in Figure 16, where we see that a random forest performs best, but only marginally better than XGBoost. Perhaps the most interesting aspect of this experiment is the overall high accuracy, given that we are classifying samples from four diverse malware families. These results indicate that although the families are significantly different from each other, they are generally more different from the benign traffic than they are from each other. This is a significant observation, since it implies that we are likely to be able to filter malicious encrypted traffic from benign traffic with one model, instead of needing to rely on a model for each individual family. This would make detection based on machine learning models far more practical than a situation where a multitude of different models are required.
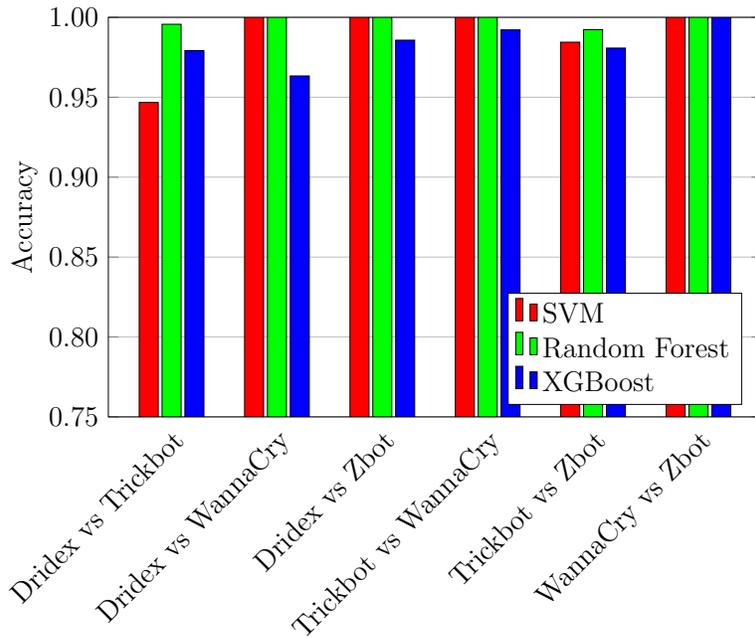
**Figure 15:** Accuracy for pairwise classification
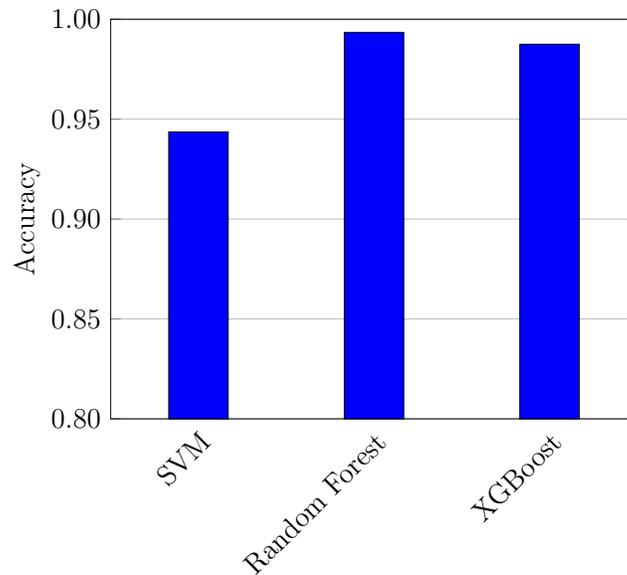


**Figure 16:** Accuracy for the multiclass problem

## 5.5  Discussion

The accuracy results in the previous section show that XGBoost performs better than the other two algorithms considered, but only marginally better
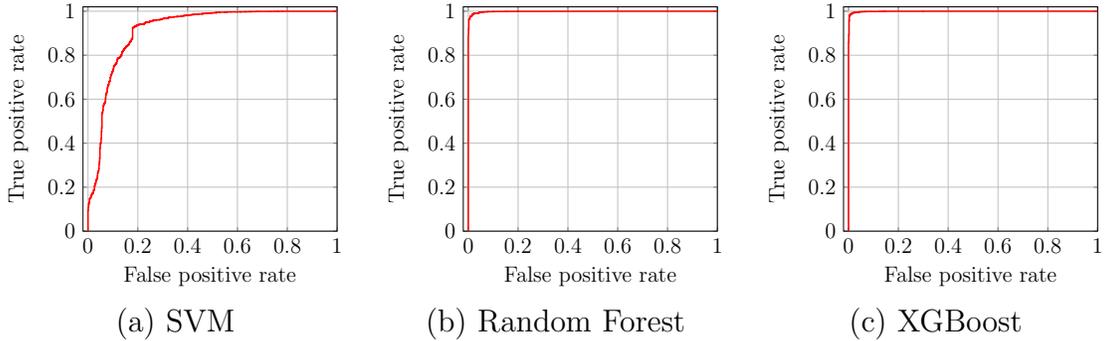
**Figure 17:** ROC curves

than a random forest. In addition to the accuracy numbers cited above, we have computed the area under the ROC curve (AUC). For the AUC statistic, we obtain a value of 0.9122 for the linear SVM, we have 0.9980 for a random forest, while XGBoost achieves an AUC of 0.9988. The corresponding ROC curve for these SVM, random forest, and XGBoost experiments are given in Figure 17 (a), (b), and (c), respectively. These AUC values serve to further reinforce the point that the random forest and XGBoost offer virtually the same performance on our dataset, while also providing a more detailed view of the relative weakness of the SVM.

In Figure 18, we give a summary of several measures of success for our experiments. These results confirm that our results are comparable to—if not better than—those achieved in any of the relevant previous work cited in Section 2.

Feature analysis is the main point of our work, and the experiments above show that we can determine strong sets of features directly from the machine learning models. It is interesting to note that the strongest features for our best models are not the most intuitive features. This illustrates the strength of a domain-free approach to feature selection, as opposed to relying on the most intuitively appealing features.

In an attempt to gain further insight into the various features, we briefly consider feature interactions. Figure 19 gives the Pearson correlation coefficient (Pearson, 1895) for all pairs of features in the form of a heatmap. Note that the feature numbers in this heatmap correspond to the feature numbers in Table 3.

From the heatmap in Figure 19 we see that, for example, features 14 through 18 are all highly correlated. These features include the length of the encryption key, the certificate length, and the validity of the certificate. The fact that these features are correlated provides evidence that substandard encryption is generally employed by encrypted malware. We also note that periodicity properties, which appear as features 11 and 12, are uncorrelated

24

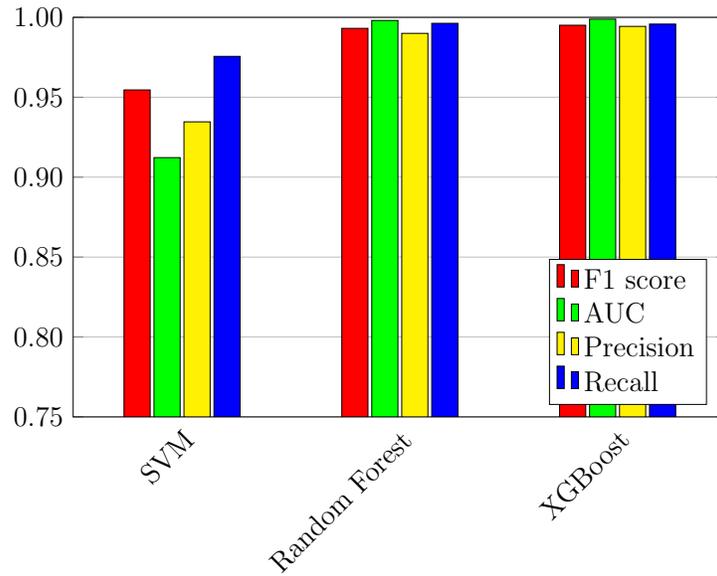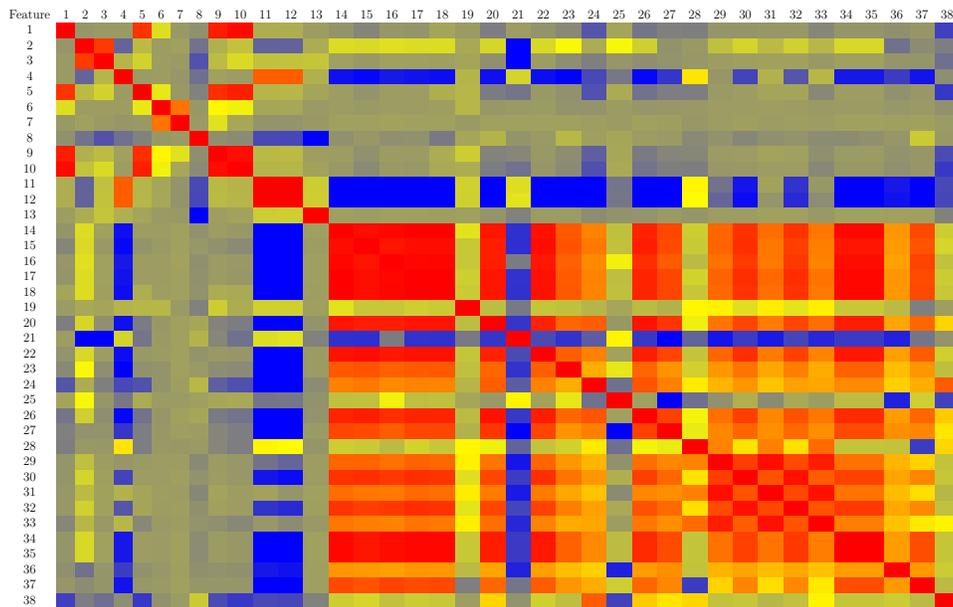**Figure 18:** Comparison based on various measures of success



**Figure 19:** Feature correlation heat map

with the encryption features, as given in features 14 through 18. In fact, these periodicity features are uncorrelated with almost all other features.

It is also interesting to consider, for example, the four strongest features of the random forest, as given in Table 6. These four features are, in order,

numbers 5, 21, 3, and 7. From the heatmap in Figure 19, we see that these features are only weakly correlated with each other. This provides further evidence that the random forest RFE is able to determine a strong set of features that are relatively independent of each other, and hence these features are providing independent information to the random forest.

On the other hand, if we look closely at the more intuitive features determined by the SVM RFE, we find some strong correlations. For example, the periodicity features (numbers 11 and 12) are the top two features in the SVM, yet they are highly correlated, and hence including both of these features provides very little new information, as compared to only including either one. This is significant as, again, fewer features will enable faster and more efficient scoring. This also points out a risk of relying on domain-specific knowledge for feature extraction—correlations are difficult to account for in such an approach, and hence some level of redundancy is highly likely.

# 6    Conclusion and Future Work

With the widespread use of HTTPS and advancement in malware detection techniques, there has been a rapid increase in the number of malware samples using HTTPS encryption to evade detection—a trend that is sure to continue into the future. This is worrying because encryption disrupts the most popular and effective malware detection techniques available today.

In this research, we considered the challenging problem of classifying encrypted network traffic as malicious or benign, without any decryption or deep packet inspection. Our focus here has been squarely on the important problem of feature analysis—a topic that we believe has been neglected in previous research in this field. We considered three machine learning algorithms, namely, SVM, XGBoost, and random forest. These algorithms were used to train and test models. We then performed feature analysis using RFE in each case. Our results show that XGBoost performed slightly better than a random forest, with both being about 99% accurate, while an SVM performed relatively poorly. Our main contribution is that we provided a thorough analysis of the various features considered. We found that a small number of features suffice, and that the optimal set of features are not necessarily the most intuitive. This is significant, as a minimal set of features is necessary for a practical and efficient implementation of any such detection system. We also argued that feature selection based on the models themselves—which we refer to as a domain-free approach—is preferable to a domain-specific approach that relies on human expertise to select features.

Security research is often a double-edged sword, and the feature analysis considered here is no exception. For efficient detection, it is necessary to know which features contribute the most, but this same analysis also points

the way toward "new-and-improved" malware. That is, malware writers can use the feature analysis presented here to make their creations more difficult to detect. This observation points to the need for more research in this field, so that detection techniques can stay ahead of the inevitable improvements in malware development.

For further research, it would be useful to have a larger dataset that we could mine for more subtle features. Another direction for related research would be to field a system based on the machine learning techniques considered here, so that the effectiveness and robustness of such an approach could be analyzed under real-world and changing conditions. In particular, it would be interesting to consider the evolution of such a system over time, as human adversaries improve their attacks and the detection system then adapts to these changes in attack strategies.

# References

Anderson, B., & McGrew, D. A. (2016). Identifying encrypted malware traffic with contextual flow data. In D. M. Freeman, A. Mitrokotsa, & A. Sinha (Eds.), *Proceedings of the 2016 ACM workshop on artificial intelligence and security* (pp. 35–46). ACM.

Anderson, B., & McGrew, D. A. (2017). Machine learning for encrypted malware traffic classification: Accounting for noisy labels and non-stationarity. In *Proceedings of the 23rd ACM international conference on knowledge discovery and data mining* (pp. 1723–1732). ACM.

Anderson, B., Paul, S., & McGrew, D. (2016). *Deciphering malware's use of TLS (without decryption).* https://arxiv.org/abs/1607.01639.

Aycock, J. (2006). *Computer Viruses and Malware.* Springer.

Bakhdlaghi, Y. (2017). *Snort and SSL/TLS inspection.* SANS Institute InfoSec Reading Room. https://www.sans.org/reading-room/whitepapers/vpns/snort-ssl-tls-inspection-37735.

Baset, S., & Schulzrinne, H. (2006). An analysis of the Skype peer-to-peer internet telephony protocol. In *25th IEEE international conference on computer communications.* IEEE. http://ieeexplore.ieee.org/xpl/mostRecentIssue.jsp?punumber=4146652.

Bennett, K. P., & Campbell, C. (2000). Support vector machines: Hype or hallelujah? *SIGKDD Explorations, 2,* 1–13.

Bradley, A. P. (1997). The use of the area under the roc curve in the evaluation of machine learning algorithms. *Pattern Recognition,*

*30*, 1145–1159.

Breiman, L., & Cutler, A. (2002). *Random forests^{TM}*. https://www
    .stat.berkeley.edu/~breiman/RandomForests/cc_home.htm.

Bro Project Team. (2014). *Bro network security monitor.* https://
    www.bro.org.

Chen, T., & Guestrin, C. (2016). *XGBoost: A scalable tree boosting
    system.* http://arxiv.org/abs/1603.02754.

Cisco. (2018). *Snort — Network intrusion detection and prevention
    system.* https://www.snort.org.

Cristianini, N., & Shawe-Taylor, J. (2000). *An introduction to support
    vector machines and other kernel-based learning methods.* Cam-
    bridge University Press.

Di Mauro, M., & Di Sarno, C. (2018). Improving siem capabilities
    through an enhanced probe for encrypted skype traffic detection.
    *Journal of Information Security and Applications*, *38*, 85–95.

Dreger, H., & Feldmann, A. (2006). Dynamic application-layer protocol
    analysis for network intrusion detection. In A. D. Keromytis
    (Ed.), *Proceedings of the 15th USENIX security symposium.*
    USENIX Association. https://www.usenix.org/conference/
    15th-usenix-security-symposium/dynamic-application
    -layer-protocol-analysis-network.

Erquiaga, M. J., & Garcia, S. (2013). *Malware capture facility project.*
    https://mcfp.weebly.com.

Etienne, L. (2009). *Malicious traffic detection in local networks with
    snort.* https://infoscience.epfl.ch/record/141022/files/
    pdm.pdf.

Freier, A. O., Karlton, P., & Kocher, P. C. (2011). *RFC 6101: The
    secure sockets layer (SSL) protocol version 3.0.*

García, S., Grill, M., Stiborek, J., & Zunino, A. (2014). An empirical
    comparison of botnet detection methods. *Computers & Security*,
    *45*, 100–123.

Google. (2017). *HTTPS encryption on the web.* https://
    transparencyreport.google.com/https/.

Liaw, A., & Wiener, M. (2002). Classification and regression by random-
    Forest. *R News*, *2*, 18–22. http://www.bios.unc.edu/~dzeng/
    BIOS740/randomforest.pdf.

Lin, Y., & Jeon, Y. (2002). *Random forests and adaptive nearest
    neighbors* (Tech. Rep. No. 1055). Madison, Wisconsin: Univer-
    sity of Wisconsin, Department of Statistics. https://www.stat
    .wisc.edu/sites/default/files/tr1055.pdf.

Lokoc, J., Kohout, J., Cech, P., Skopal, T., & Pevný, T. (2016). *k*-NN classification of malware in HTTPS traffic using the metric space approach. In M. Chau, G. A. Wang, & H. Chen (Eds.), *Proceedings of intelligence and security informatics* (pp. 131–145). Springer.

Madhukar, A., & Williamson, C. L. (2006). A longitudinal study of P2P traffic classification. In *14th international symposium on modeling, analysis, and simulation of computer and telecommunication systems* (pp. 179–188). IEEE Computer Society. http://ieeexplore.ieee.org/xpl/mostRecentIssue.jsp?punumber=11154.

Malwarebytes. (2017). *Analysis of malware trends for small and medium businesses Q1 2017.* https://www.malwarebytes.com/pdf/white-papers/MalwareTrendsForSMBQ12017.pdf.

Moore, A. W., & Papagiannaki, K. (2005). Toward the accurate identification of network applications. In C. Dovrolis (Ed.), *Proceedings of 6th international workshop on passive and active network measurement* (pp. 41–54). Springer.

Pearson, K. (1895). Notes on regression and inheritance in the case of two parents. *Proceedings of the Royal Society of London*, *58*, 240–242.

Prasse, P., Machlica, L., Pevný, T., Havelka, J., & Scheffer, T. (2017). Malware detection by analysing network traffic with neural networks. In *2017 ieee security and privacy workshops* (pp. 205–210).

Radware. (2018). *Global application and network security report 2017–18.* https://www.datacom.cz/userfiles/radware_ert_report_2017_2018_final.pdf.

Sen, S., Spatscheck, O., & Wang, D. (2004). Accurate, scalable in-network identification of p2p traffic using application signatures. In S. I. Feldman, M. Uretsky, M. Najork, & C. E. Wills (Eds.), *Proceedings of the 13th international conference on world wide web* (pp. 512–521). ACM.

Singh, T., Di Troia, F., Visaggio, C. A., Austin, T. H., & Stamp, M. (2016). Support vector machines and malware detection. *Journal of Computer Virology and Hacking Techniques*, *12*(4), 203–212. Retrieved from https://doi.org/10.1007/s11416-015-0252-0 doi: 10.1007/s11416-015-0252-0

Stamp, M. (2017a). *Boost your knowledge of AdaBoost.* https://www.cs.sjsu.edu/~stamp/ML/files/ada.pdf.

Stamp, M. (2017b). *Introduction to machine learning with applications in information security.* Chapman & Hall/CRC Press.

Tegeler, F., Fu, X., Vigna, G., & Kruegel, C. (2012). Botfinder: Find-

ing bots in network traffic without deep packet inspection. In C. Barakat, R. Teixeira, K. K. Ramakrishnan, & P. Thiran (Eds.), *Conference on emerging networking experiments and technologies* (pp. 349–360). ACM. http://dl.acm.org/citation.cfm?id=2413176.

Touch, J., Lear, E., Mankin, A., Kojo, M., Ono, K., Stiemerling, M., . . . Nishida, Y. (2018). *IANA: Service name and transport protocol port number registry.* http://www.iana.org/assignments/port-numbers.

Wang, X., Qiu, W., & Zamar, R. H. (2007). CLUES: A non-parametric clustering method based on local shrinking. *Computational Statistics & Data Analysis*, *52*(1), 286–298.

Yoon, S., Park, J., Park, J., Oh, Y., & Kim, M. (2009). Internet application traffic classification using fixed IP-port. In C. S. Hong, T. Tonouchi, Y. Ma, & C. Chao (Eds.), *Proceedings of 12th Asia-Pacific network operations and management symposium* (pp. 21–30). Springer.