

# A Study of Detecting Computer Viruses in Real-Infected Files in the $n$ -gram Representation with Machine Learning Methods

Note, this is an extended version. The LNCS published version is available via [springerlink](#)

Thomas Stibor

Fakultät für Informatik  
Technische Universität München  
*thomas.stibor@in.tum.de*

**Abstract.** Machine learning methods were successfully applied in recent years for detecting new and unseen computer viruses. The viruses were, however, detected in small virus loader files and not in real infected executable files. We created data sets of benign files, virus loader files and real infected executable files and represented the data as collections of  $n$ -grams. Histograms of the relative frequency of the  $n$ -gram collections indicate that detecting viruses in real infected executable files with machine learning methods is nearly impossible in the  $n$ -gram representation. This statement is underpinned by exploring the  $n$ -gram representation from an information theoretic perspective and empirically by performing classification experiments with machine learning methods.

## 1 Introduction

Detecting new and unseen viruses with machine learning methods is a challenging problem in the field of computer security. Signature based detection provides adequate protection against known viruses, when proper virus signatures are available. From a machine learning point of view, signature based detection is based on a prediction model where no generalization exists. In other words, no detection beyond the known viruses can be performed. To obtain a form of generalization, current anti-virus systems use heuristics generated by hand.

In recent years, however, machine learning methods are investigated for detecting unknown viruses [1–5]. Reported results show high and acceptable detection rates. The experiments, however, were performed on *skewed data sets*, that is, not real infected executable files are considered, but small executable virus loader files. In this paper, we investigate and compare machine learning methods on real infected executable files and virus loader files. Results reveal that real infected executable files when represented as  $n$ -grams can *barely* be detected with machine learning methods.

## 2 Data Sets

In this paper, we focus on computer virus detection in DOS executable files. These files are executable on DOS and some Windows platforms only and can be identified by the leading two byte signature “MZ” and the file name suffix “.EXE”. It is important to note that although DOS platforms are obsolete, the principle of virus infection and detection on DOS platforms can be generalized to other OS platforms and executable files.

In total three different data sets are created. The benign data set consists of virus-free executable DOS files which are collected from public FTP servers. The collected files consist of DOS games, compilers, editors, etc. and are preprocessed such that no duplicate files exist. Additionally, files that are compressed with executable packers<sup>1</sup> such as lzexe, pklite, diet, etc. are uncompressed. The preprocessing steps are performed to obtain a clean data set which is supposed to induce high detection rates of the considered machine learning methods. In total, we collected 3514 of such “clean” DOS executable files.

DOS executable virus loader files are collected from the VXHeaven website.<sup>2</sup> The same processing steps are applied, that is, all duplicate files are removed and compressed files are uncompressed. It turned out however, that some files from VXHeaven are real infected executable files and not virus loader files. These files can be identified by inspecting their content with a disassembler/hexeditor and by their large file sizes. Consequently, we removed all files greater than 10 KBytes and collected in total 1555 virus loader files.

For the sake of clarity, the benign data set is denoted as  $\mathcal{B}_{\text{files}}$ , the virus loader data set as  $\mathcal{L}_{\text{files}}$ .

The data set of real infected executable files denoted as  $\mathcal{I}_{\text{files}}$  is created by performing the following steps:

```
1:  $\mathcal{I}_{\text{files}} := \{\}$ 
2: for each file.v in  $\mathcal{L}_{\text{files}}$ 
3:   setup new DOS environment in DOS emulator
4:   randomly determine file.r from  $\mathcal{B}_{\text{files}}$ 
5:   boot DOS and execute file.v and file.r
6:   if (file.r gets infected)
7:      $\mathcal{I}_{\text{files}} := \mathcal{I}_{\text{files}} \cup \{\text{file.r}\}$ 
```

**Fig. 1.** Steps performed to create real infected executable DOS files.

These steps were executed on a GNU/Linux machine by means of a Perl script and a DOS emulator. A DOS emulator is chosen since it is too time consuming to setup a clean DOS environment on a DOS machine whenever a new file gets

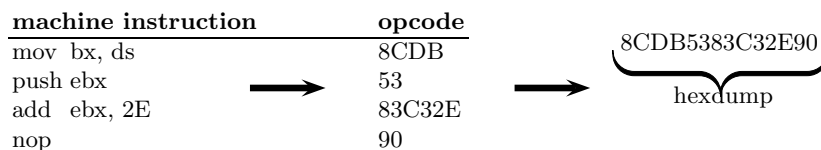
<sup>1</sup> Executable compression is frequently used to confuse and hide from virus scanners.

<sup>2</sup> Available at <http://vx.netlux.org/>

infected. The statement “file.r gets infected” at line 6 in Figure 1 denotes the data labeling process. Since only viruses were considered for which signatures are available, *no* elements in the created data sets are mislabeled. In total 1215 real infected files were collected.<sup>3</sup>

### 3 Transforming Executable Files to n-grams and Frequency Vectors

Raw executable files cannot serve properly as input data for machine learning methods. Consequently, a suitable representation is required. A common and frequently used representation is the hexdump. A hexdump is a hexadecimal representation of an executable file (or of data in general) and is related to machine instructions, termed opcodes. Figure 2 illustrates the connection between machine instructions, opcodes, and a hexdump. It is important to note, that the hexdump representation used and reported in the literature, is computed over a *complete* executable file which consists of a header, relocation tables and the executable code. The machine instructions are located in the code segment only and this implies that a hexdump obtained over a complete executable file may not be interpreted completely as machine instructions.



**Fig. 2.** Relation between machine instructions, opcodes and the resulting hexdump.

In the second preprocessing step, the hexdump is “cut” into substrings of length  $n \in \mathbb{N}$ , denoted as  $n$ -grams. According to the literature an  $n$ -gram can be any set of characters in a string [6]. In terms of detecting viruses in executable files,  $n$ -grams are adjacent substrings created by shifting a window of length  $s \in \mathbb{N}$  over the hexdump. The resulting number of created  $n$ -grams depends on the value of  $n$  and the window shift length  $s$ .

We created for each file an ordered collection<sup>4</sup> of  $n$ -grams by cutting the hexdump from left to right (see Fig. 3). Note that the collection can contain the same  $n$ -grams multiple times. The total number of  $n$ -grams in the collection when given  $n$  and  $s$  is  $\lfloor (l - n) / s + 1 \rfloor$ , where  $l$  denotes the hexdump length. In a final step, the collection is transformed into a vector of dimension<sup>5</sup>  $d = 16^n$ , where

<sup>3</sup> For the sake of verifiability, the three data sets are available at <http://www.sec.in.tum.de/~stibor/ieaaie2010/>

<sup>4</sup> If not otherwise stated we denote an ordered collection as collection.

<sup>5</sup> An unary hexadecimal number can take 16 values.



Let  $X$  be a random variable that represents the outcome of an  $n$ -gram and  $\mathcal{X}$  the set of all  $n$ -grams of length  $n$ . The entropy is defined as

$$H(X) = \sum_{x \in \mathcal{X}} P(X = x) \log_2 \frac{1}{P(X = x)}. \quad (1)$$

Informally speaking, high entropy implies that  $X$  is from a roughly uniform distribution and  $n$ -grams sampled from it are distributed all over the place. Consequently, predicting sampled  $n$ -grams is “hard”. In contrast, low entropy implies that  $X$  is from a varied distribution and  $n$ -grams sampled from it would be more predictable.

#### 4.1 Information Loss in Frequency Vectors

Given a hexdump of a file, the resulting  $n$ -gram collection  $\mathcal{C}$  and the created  $n$ -gram frequency vector  $\mathbf{c}$ . We are interested in minimizing the amount of information loss for reconstructing the  $n$ -gram collection  $\mathcal{C}$  when given  $\mathbf{c}$ . For the sake of simplicity we consider the case  $n = s$  and denote the cardinality of  $\mathcal{C}$  as  $t$ . Assume there are  $r$  unique  $n$ -grams in  $\mathcal{C}$ , where each  $n$ -gram appears  $n_i$  times such that  $\sum_{i=1}^r n_i = t$ . The collection  $\mathcal{C}$  can have

$$\frac{t!}{n_1! n_2! \dots n_r!} = \binom{t}{n_1, n_2, \dots, n_r} \quad (2)$$

many rearrangements of  $n$ -grams. Applying the theorem from information theory on the size of a type class [8], (pp. 350, Theorem 11.1.3) it follows:

$$\frac{1}{(t+1)^{|\Sigma|}} 2^{tH(P)} \leq \binom{t}{n_1, n_2, \dots, n_r} \leq 2^{tH(P)} \quad (3)$$

where  $P$  denotes the empirical probability distribution of the  $n$ -grams,  $H(P)$  the entropy of distribution  $P$  and  $\Sigma$  the used alphabet. Using  $n = s$  it follows from (3) that  $2^{\lfloor l/n \rfloor H(P)}$  is an upper bound of the information loss. To keep the upper bound as small as possible, the value of  $n$  has to be close to the value of  $l$  and the empirical probability distribution  $P$  has to be a varied distribution rather than a uniform distribution.

#### 4.2 Entropy of $n$ -gram Collections

To support the statement in Section 4.1 empirically, the entropy values of the created  $n$ -gram collections are calculated. The values are set in relation to the maximum possible entropy (uniform distribution) denoted as  $H_{\max}$ . An  $n$ -gram collection where the corresponding value of  $H(X)/H_{\max}$  is close to zero, contains more predictive  $n$ -grams than the one where  $H(X)/H_{\max}$  is close to one. The results for different parameter combinations of  $n$  and  $s$  are listed in Table 1. One can see that larger  $n$ -grams lengths are more predictable than short

ones. Additionally, one can observe that window shift length  $s$  has to be a multiple of two, because opcodes most frequently occur as two and four byte values. Moreover, one can observe that the result of  $H(X)/H_{\max}$  for  $n$ -gram collections created from the benign set  $\mathcal{B}_{\text{files}}$  and virus infected set  $\mathcal{I}_{\text{files}}$  have approximately the same values. In other words, discriminating between those two classes is limited realizable. This observation is also identifiable in Section 6 where the classification results are presented and by inspecting Figure 4 which shows histograms of the relative frequencies of the  $n$ -grams. From a practical point of view, however, large values of  $n$  induce a computational complexity which is not manageable. To be more precise, the crucial factor that influences the computational complexity is not just the total number of created  $n$ -grams (see Table 2), but also the number of unique  $n$ -grams, that is, the lower bound of the dimension of the frequency vector. For our data sets we obtained the results listed in Table 3. One can observe that all possible  $n$ -grams occurred in the collections. As a consequence, we focus in the paper on  $n$ -gram lengths 2 and 3. We tried to run the classification experiments for  $n = 4, 5, 6, 7, 8$  in combination with a dimensionality reduction described in the subsequent section, however, due to the resulting space complexity and lack of memory space<sup>6</sup> no results were obtained.

### 4.3 Feature Selection

The purpose of feature selection is to reduce the dimensionality of the data such that the machine learning methods can be applied more efficiently in terms of time and space complexity. Furthermore, feature selection often increases classification accuracy by eliminating noisy features. We apply on our data sets the mutual information method [7]. This method is used in previous work and appears to be practical [5]. A comparison study of other feature selection methods in terms of classifying malicious executables is provided in [9].

**Mutual Information Selection:** Let  $U$  be a random variable that takes values  $e_{ng} = 1$  (the collection contains  $n$ -gram  $ng$ ) and  $e_{ng} = 0$  (the collection does not contain  $ng$ ). Let  $C$  be a random variable that takes values  $e_c = 1$  (the collection belongs to class  $c$ ) and  $e_c = 0$  (the collection does not belong to class  $c$ ).

The mutual information measures how much information the presence/absence of an  $n$ -gram contributes to making the correct classification decision on  $c$ . More formally

$$\sum_{e_{ng} \in \{0,1\}} \sum_{e_c \in \{0,1\}} P(U = e_{ng}, C = e_c) \cdot \log_2 \frac{P(U = e_{ng}, C = e_c)}{P(U = e_{ng})P(C = e_c)}. \quad (4)$$

---

<sup>6</sup> A Quad-Core AMD 2.6 Ghz with 16 GB RAM was used.

$s \backslash n$	2	3	4
	0.923	0.896	0.862
1	0.924	0.895	0.861
	0.738	0.679	0.628
	0.889	0.857	0.818
2	0.888	0.857	0.816
	0.708	0.649	0.593
		0.896	0.862
3		0.895	0.860
		0.679	0.626
			0.816
4			0.813
			0.590

**Table 1.** Results of the expression  $H(X)/H_{\max}$  for different parameter combinations of  $n$  and  $s$ . Smaller values indicate that more predictive  $n$ -grams can be created. The upper value in each row denotes the result of  $n$ -grams created from benign set  $\mathcal{B}_{\text{files}}$ , the middle value from virus infected set  $\mathcal{I}_{\text{files}}$ , the lower value from virus loader set  $\mathcal{L}_{\text{files}}$ .

$s \backslash n$	2	3	4
	378 238 050	378 234 536	378 231 022
1	128 293 929	128 292 714	128 291 499
	8 281 043	8 279 488	8 277 933
	189 120 782	189 117 268	189 117 268
2	64 147 572	64 146 357	64 146 357
	4 141 299	4 139 744	4 139 744
		126 079 338	126 078 183
3		42 764 639	42 764 238
		2 760 368	2 759 815
			94 560 011
4			32 073 535
			2 070 340

**Table 2.** Number of created  $n$ -grams for different  $n$ -gram and window lengths  $s$ . The upper number in each row denotes the number of  $n$ -grams created from benign set  $\mathcal{B}_{\text{files}}$ , the middle number from virus infected set  $\mathcal{I}_{\text{files}}$ , the lower number from virus loader set  $\mathcal{L}_{\text{files}}$ .

$n$	2	3	4	5
$\mathcal{B}_{\text{files}}$	256	4096	65536	1048576
$\mathcal{L}_{\text{files}}$	256	4096	65536	1048576
$\mathcal{I}_{\text{files}}$	256	4096	65536	1048576

**Table 3.** The number of unique  $n$ -grams of the three data sets  $\mathcal{B}_{\text{files}}$ ,  $\mathcal{L}_{\text{files}}$ ,  $\mathcal{I}_{\text{files}}$ . The number determines the dimension of the corresponding frequency vectors.

## 5 Classification Methods

In this section we briefly introduce the classification methods which are used in the experiments. These methods are frequently applied in the field of machine learning.

### 5.1 Cosine Similarity Classification

Given two vectors  $\mathbf{x}, \mathbf{y} \in \mathbb{R}^d$ , the cosine similarity is given as follows

$$\cos \phi = \frac{\langle \mathbf{x}, \mathbf{y} \rangle}{\|\mathbf{x}\| \|\mathbf{y}\|} \quad (5)$$

$$= \frac{x_1 y_1 + \dots + x_d y_d}{\sqrt{x_1^2 + \dots + x_d^2} \sqrt{y_1^2 + \dots + y_d^2}}. \quad (6)$$

The dot product of  $\mathbf{x}$  and  $\mathbf{y}$  denotes as  $\langle \mathbf{x}, \mathbf{y} \rangle$  has a straightforward geometrical interpretation, namely, the length of the projection of  $\mathbf{x}$  onto the unit vector  $\mathbf{y}/\|\mathbf{y}\|$ . Consequently, (5) represents the cosine angle  $\phi$  between  $\mathbf{x}$  and  $\mathbf{y}$ . The result of (5) lies in interval  $[-1, 1]$ , where 1 is the result of equal vectors and  $-1$  of total dissimilar vectors.

For performing classification, the arithmetic mean vector  $\bar{\mathbf{c}}$  of each class is determined. A new vector  $\mathbf{u}$  is assigned to the class with the largest cosine similarity value between  $\bar{\mathbf{c}}$  and  $\mathbf{u}$ .

### 5.2 Naive Bayesian Classifier

A naive Bayesian classifier is well known in the machine learning community and is popularized by applications such as spam filtering or text classification in general.

Given feature variables  $F_1, F_2, \dots, F_n$  and a class variable  $C$ . The Bayes' theorem states

$$P(C | F_1, F_2, \dots, F_n) = \frac{P(C) P(F_1, F_2, \dots, F_n | C)}{P(F_1, F_2, \dots, F_n)}. \quad (7)$$

Assuming that each feature  $F_i$  is conditionally independent of every other feature  $F_j$  for  $i \neq j$  one obtains

$$P(C | F_1, F_2, \dots, F_n) = \frac{P(C) \prod_{i=1}^n P(F_i | C)}{P(F_1, F_2, \dots, F_n)}. \quad (8)$$

The denominator serves as a scaling factor and can be omitted in the final classification rule

$$\operatorname{argmax}_c P(C = c) \prod_{i=1}^n P(F_i = f_i | C = c). \quad (9)$$



### 5.3 Support Vector Machine

The Support Vector Machine (SVM) [10, 11] is a classification technique which has its foundation in computational geometry, optimization and statistics. Given a sample

$$(\mathbf{x}_1, y_1), (\mathbf{x}_2, y_2), \dots, (\mathbf{x}_l, y_l) \in \mathbb{R}^d \times Y, \quad Y = \{-1, +1\}$$

and a (nonlinear) mapping into a potentially much higher dimensional feature space  $\mathcal{F}$

$$\begin{aligned} \Phi : \mathbb{R}^d &\rightarrow \mathcal{F} \\ \mathbf{x} &\mapsto \Phi(\mathbf{x}). \end{aligned}$$

The goal of SVM learning is to find a separating hyperplane with maximal margin in  $\mathcal{F}$  such that

$$y_i(\langle \mathbf{w}, \Phi(\mathbf{x}_i) \rangle + b) \geq 1, \quad i = 1, \dots, l, \quad (10)$$

where  $\mathbf{w} \in \mathcal{F}$  is the normal vector of the hyperplane and  $b \in \mathbb{R}$  the offset. The separation problem between the  $-1$  and  $+1$  class can be formulated in terms of the quadratic convex optimization problem

$$\min_{\mathbf{w}, b, \xi} \frac{1}{2} \|\mathbf{w}\|^2 + C \sum_{i=1}^l \xi_i \quad (11)$$

$$\text{subject to } y_i(\langle \mathbf{w}, \Phi(\mathbf{x}_i) \rangle + b) \geq 1 - \xi_i, \quad (12)$$

where  $C$  is a penalty term and  $\xi_i \geq 0$  slack variables to relax the hard constraints in (10). Solving the dual form of (11) and (12) leads to the final decision function

$$f(\mathbf{x}) = \text{sgn} \left( \sum_{i=1}^l y_i \alpha_i \langle \Phi(\mathbf{x}), \Phi(\mathbf{x}_i) \rangle + b \right) \quad (13)$$

$$= \text{sgn} \left( \sum_{i=1}^l y_i \alpha_i k(\mathbf{x}, \mathbf{x}_i) + b \right) \quad (14)$$

where  $k(\cdot, \cdot)$  is a kernel function satisfies the Mercer's condition [12] and  $\alpha_i$  Lagrangian multipliers obtained by solving the dual form.

In the performed classification experiments we used the Gaussian kernel  $k(\mathbf{x}, \mathbf{y}) = \exp(-\|\mathbf{x} - \mathbf{y}\|^2/\gamma)$  as it gives the highest classification results compared to other kernels. The hyperparameters  $\gamma$  and  $C$  are optimized by means of the inbuilt validation method of the used R package `kernlab` [13].

## 6 ROC Analysis and Results

ROC analysis is performed to measure the goodness of the classification results. More specifically, we are interested in the following quantities:

- true positives (TP): number of virus executable examples classified as virus executable,
- true negatives (TN): number of benign executable examples classified as benign executable,
- false positives (FP): number of benign executable examples classified as virus executable,
- false negatives (FN): number of virus executable examples classified as benign executable.

The *detection rate* (true positive rate), *false alarm rate* (false positive rate), and *overall accuracy* of a classifier are calculated then as follow

$$\begin{aligned} \text{detection rate} &= \frac{TP}{TP + FN} \\ \text{false alarm rate} &= \frac{FP}{FP + TN} \\ \text{overall accuracy} &= \frac{TP + TN}{TP + TN + FP + FN}. \end{aligned}$$

Moreover, to avoid under/overfitting effects of the classifiers, we performed a  $K$ -crossfold validation. That is, the data set is split in  $K$  roughly equal-sized parts. The classification method is trained on  $K - 1$  parts and has to predict the not seen testing part. This is performed for  $k = 1, 2, \dots, K$  and combined as the mean value to estimate the generalization error. Since the data sets are large, a 10-crossfold validation is performed. The classification results are depicted in Tables 4-7, where also the standard deviation result of the crossfold validation is provided.

The classification results evidently show (cf. Tables 4-7) that high detection rates and low false alarm rates can be obtained when discriminating between benign files and virus loader files. This observation supports previously published results [5, 1, 2]. However, for detecting viruses in real infected executables these methods are barely applicable. This observation is not a great surprise, because benign files and real infected files are from a statistical point of view nearly indistinguishable (cf. Figure 4, last page).

It is interesting to note that the SVM gives excellent classification results when discriminating between benign and virus loader files. However, poor results when discriminating between benign and virus infected files. In terms of the overall accuracy, the SVM still outperforms the Bayes and cosine measure method.

Additionally, one can observe that selecting 50% of the features by means of the mutual information criterion does not significantly increase and decrease the classification accuracy of the SVM. The Bayes and cosine measure, however, benefit by the feature selection preprocessing step as a significant increase of the classification accuracy can be observed.

## 7 Conclusion

We investigated the applicability of machine learning methods for detecting viruses in real infected DOS executable files when using the  $n$ -gram representation. For that reason, three data sets were created. The benign data set contains virus free DOS executable files collected from all over the Internet. For the sake of comparison to previous work, we created a data set which consists of virus loader files collected from the VXHeaven website. Furthermore, we created a new data set which consists of real infected executable files. The data sets were transformed to collections of  $n$ -grams and represented as frequency vectors. Due to this transformational step a information loss occurs. We showed empirically and theoretically with arguments from information theory that this loss can be minimized by increasing the lengths of the  $n$ -grams. As a consequence however, this results in a computational complexity which is not manageable.

Moreover, we calculated the entropy of the  $n$ -gram collections and observed that the benign files and real infected files nearly have identical entropy values. As a consequence, discriminating between those two classes is limited realizable. This was also observed by creating histograms of the relative  $n$ -gram frequencies.

Furthermore, we performed classification experiments and confirmed that discriminating between those two classes is limited realizable. More precisely, the SVM gives unacceptable detection rates. The Bayes and the cosine measure gives higher detection rates, however, they also give unacceptable false alarm rates.

In summary, our results doubt the applicability of detecting viruses in *real* infected executable files with machine learning methods when using the (naive)  $n$ -gram representation. However, it has not escaped our mind that learning algorithms for sequential data could be one approach to tackle this problem more effectively. Another promising approach is to learn the *behavior* of a virus, that is, the sequence of system calls a virus is generating. Such a behavior can be collected in a secure sandbox and learned with appropriate machine learning methods.

## References

1. Schultz, M.G., Eskin, E., Zadok, E., Stolfo, S.J.: Data mining methods for detection of new malicious executables. In: Proceedings of the IEEE Symposium on Security and Privacy, IEEE Computer Society (2001) 38–49
2. Abou-Assaleh, T., Cercone, N., Keselj, V., Sweidan, R.: Detection of new malicious code using n-grams signatures. In: Second Annual Conference on Privacy, Security and Trust. (2004) 193–196
3. Reddy, D.K.S., Pujari, A.K.: *N*-gram analysis for computer virus detection. Journal in Computer Virology **2**(3) (2006) 231–239
4. Yoo, I.S., Ultes-Nitsche, U.: Non-signature based virus detection. Journal in Computer Virology **2**(3) (2006) 163–186
5. Kolter, J.Z., Maloof, M.A.: Learning to detect and classify malicious executables in the wild. Journal of Machine Learning Research **7** (2006) 2721–2744
6. Cavnar, W.B., Trenkle, J.M.: N-gram-based text categorization. In: Proceedings of Third Annual Symposium on Document Analysis and Information Retrieval. (1994) 161–175
7. Manning, C.D., Raghavan, P., Schütze, H.: Introduction to Information Retrieval. Cambridge University Press (2008)
8. Cover, T.M., Thomas, J.A.: Elements of Information Theory. 2nd edn. Wiley-Interscience (2006)
9. Cai, D.M., Gokhale, M., James, J.T.: Comparison of feature selection and classification algorithms in identifying malicious executables. Computational Statistics & Data Analysis **51**(6) (2007) 3156–3172
10. Boser, B.E., Guyon, I.M., Vapnik, V.: A training algorithm for optimal margin classifiers. In: Proceedings of the fifth annual workshop on Computational learning theory (COLT), ACM Press (1992) 144–152
11. Cortes, C., Vapnik, V.: Support-vector networks. Machine Learning **20**(3) (1995) 273–297
12. Schölkopf, B., Smola, A.J.: Learning with Kernels: Support Vector Machines, Regularization, Optimization, and Beyond. MIT Press (2001)
13. Karatzoglou, A., Smola, A., Hornik, K., Zeileis, A.: kernlab – an S4 package for kernel methods in R. Journal of Statistical Software **11**(9) (2004) 1–20

benign vs. virus loader				
s	Method	Detection Rate	False Alarm Rate	Overall Accuracy
1	SVM	<b>0.9499</b> ( $\pm 0.0191$ )	<b>0.0272</b> ( $\pm 0.0075$ )	<b>0.9654</b> ( $\pm 0.0059$ )
	Bayes	0.6120 ( $\pm 0.0459$ )	0.0710 ( $\pm 0.0110$ )	0.8315 ( $\pm 0.0194$ )
	Cosine	0.5883 ( $\pm 0.0337$ )	0.1239 ( $\pm 0.0123$ )	0.7873 ( $\pm 0.0151$ )
2	SVM	<b>0.9531</b> ( $\pm 0.0120$ )	<b>0.0261</b> ( $\pm 0.0091$ )	<b>0.9674</b> ( $\pm 0.0083$ )
	Bayes	0.7267 ( $\pm 0.0244$ )	0.0692 ( $\pm 0.0136$ )	0.8684 ( $\pm 0.0122$ )
	Cosine	0.6385 ( $\pm 0.0255$ )	0.1211 ( $\pm 0.0190$ )	0.8050 ( $\pm 0.0194$ )

benign vs. virus infected				
s	Method	Detection Rate	False Alarm Rate	Overall Accuracy
1	SVM	0.0370 ( $\pm 0.0161$ )	<b>0.0731</b> ( $\pm 0.0078$ )	<b>0.6982</b> ( $\pm 0.0155$ )
	Bayes	<b>0.6039</b> ( $\pm 0.0677$ )	0.6113 ( $\pm 0.0679$ )	0.4438 ( $\pm 0.0416$ )
	Cosine	0.5479 ( $\pm 0.1300$ )	0.5674 ( $\pm 0.0906$ )	0.4620 ( $\pm 0.0394$ )
2	SVM	0.0556 ( $\pm 0.0153$ )	<b>0.0888</b> ( $\pm 0.0125$ )	<b>0.6912</b> ( $\pm 0.0206$ )
	Bayes	<b>0.6682</b> ( $\pm 0.0644$ )	0.6533 ( $\pm 0.0736$ )	0.4284 ( $\pm 0.0465$ )
	Cosine	0.5835 ( $\pm 0.1096$ )	0.5938 ( $\pm 0.0897$ )	0.4506 ( $\pm 0.0474$ )

**Table 4.** Classification results for  $n$ -gram length  $n = 2$  and sliding window  $s = 1, 2$  without feature selection.

benign vs. virus loader with feature selection				
s	Method	Detection Rate	False Alarm Rate	Overall Accuracy
1	SVM	<b>0.9443</b> ( $\pm 0.0151$ )	<b>0.0324</b> ( $\pm 0.0057$ )	<b>0.9605</b> ( $\pm 0.0035$ )
	Bayes	0.7897 ( $\pm 0.0297$ )	0.1123 ( $\pm 0.0187$ )	0.8577 ( $\pm 0.0178$ )
	Cosine	0.7688 ( $\pm 0.0292$ )	0.1016 ( $\pm 0.0194$ )	0.8587 ( $\pm 0.0181$ )
2	SVM	<b>0.9354</b> ( $\pm 0.0099$ )	<b>0.0353</b> ( $\pm 0.0129$ )	<b>0.9558</b> ( $\pm 0.0105$ )
	Bayes	0.8844 ( $\pm 0.0301$ )	0.1653 ( $\pm 0.0284$ )	0.8498 ( $\pm 0.0223$ )
	Cosine	0.8643 ( $\pm 0.0391$ )	0.1479 ( $\pm 0.0275$ )	0.8559 ( $\pm 0.0186$ )

benign vs. virus infected with feature selection				
s	Method	Detection Rate	False Alarm Rate	Overall Accuracy
1	SVM	0.0251 ( $\pm 0.0128$ )	<b>0.0702</b> ( $\pm 0.0201$ )	<b>0.6974</b> ( $\pm 0.0213$ )
	Bayes	0.6084 ( $\pm 0.0502$ )	0.5860 ( $\pm 0.0456$ )	0.4637 ( $\pm 0.0287$ )
	Cosine	<b>0.6454</b> ( $\pm 0.0486$ )	0.6408 ( $\pm 0.0210$ )	0.4326 ( $\pm 0.0173$ )
2	SVM	0.0264 ( $\pm 0.0120$ )	<b>0.0782</b> ( $\pm 0.0147$ )	<b>0.6916</b> ( $\pm 0.0184$ )
	Bayes	<b>0.6270</b> ( $\pm 0.0577$ )	0.6120 ( $\pm 0.0456$ )	0.4495 ( $\pm 0.0312$ )
	Cosine	0.6212 ( $\pm 0.0791$ )	0.6216 ( $\pm 0.0695$ )	0.4419 ( $\pm 0.0323$ )

**Table 5.** Classification results for  $n$ -gram length  $n = 2$  and sliding window  $s = 1, 2$  with selecting 50 % of the original features by means of the mutual information method.

benign vs. virus loader				
s	Method	Detection Rate	False Alarm Rate	Overall Accuracy
1	SVM	<b>0.9622</b> ( $\pm 0.0148$ )	<b>0.0267</b> ( $\pm 0.0080$ )	<b>0.9700</b> ( $\pm 0.0080$ )
	Bayes	0.7493 ( $\pm 0.0302$ )	0.0594 ( $\pm 0.0106$ )	0.8820 ( $\pm 0.0076$ )
	Cosine	0.6452 ( $\pm 0.0258$ )	0.1284 ( $\pm 0.0148$ )	0.8021 ( $\pm 0.0150$ )
2	SVM	<b>0.9560</b> ( $\pm 0.0183$ )	<b>0.0266</b> ( $\pm 0.0091$ )	<b>0.9680</b> ( $\pm 0.0095$ )
	Bayes	0.8044 ( $\pm 0.0210$ )	0.0595 ( $\pm 0.0142$ )	0.8985 ( $\pm 0.0108$ )
	Cosine	0.6905 ( $\pm 0.0330$ )	0.1220 ( $\pm 0.0214$ )	0.8202 ( $\pm 0.0137$ )
3	SVM	<b>0.9618</b> ( $\pm 0.0089$ )	<b>0.0289</b> ( $\pm 0.0094$ )	<b>0.9682</b> ( $\pm 0.0071$ )
	Bayes	0.7488 ( $\pm 0.0377$ )	0.0580 ( $\pm 0.0101$ )	0.8828 ( $\pm 0.0146$ )
	Cosine	0.6450 ( $\pm 0.0410$ )	0.1272 ( $\pm 0.0215$ )	0.8025 ( $\pm 0.0224$ )

benign vs. virus infected				
s	Method	Detection Rate	False Alarm Rate	Overall Accuracy
1	SVM	0.0868 ( $\pm 0.0266$ )	<b>0.0997</b> ( $\pm 0.0162$ )	<b>0.6912</b> ( $\pm 0.0239$ )
	Bayes	<b>0.6152</b> ( $\pm 0.0742$ )	0.6110 ( $\pm 0.0661$ )	0.4476 ( $\pm 0.0352$ )
	Cosine	0.5518 ( $\pm 0.1016$ )	0.5624 ( $\pm 0.0804$ )	0.4677 ( $\pm 0.0370$ )
2	SVM	0.1155 ( $\pm 0.0345$ )	<b>0.0979</b> ( $\pm 0.0190$ )	<b>0.7001</b> ( $\pm 0.0186$ )
	Bayes	<b>0.6668</b> ( $\pm 0.0877$ )	0.6386 ( $\pm 0.1169$ )	0.4398 ( $\pm 0.0679$ )
	Cosine	0.5976 ( $\pm 0.0940$ )	0.6032 ( $\pm 0.0798$ )	0.4485 ( $\pm 0.0393$ )
3	SVM	0.0665 ( $\pm 0.0204$ )	<b>0.0944</b> ( $\pm 0.0157$ )	<b>0.6899</b> ( $\pm 0.0195$ )
	Bayes	<b>0.6110</b> ( $\pm 0.0844$ )	0.5995 ( $\pm 0.0980$ )	0.4542 ( $\pm 0.0563$ )
	Cosine	0.5510 ( $\pm 0.0964$ )	0.5609 ( $\pm 0.0594$ )	0.4677 ( $\pm 0.0244$ )

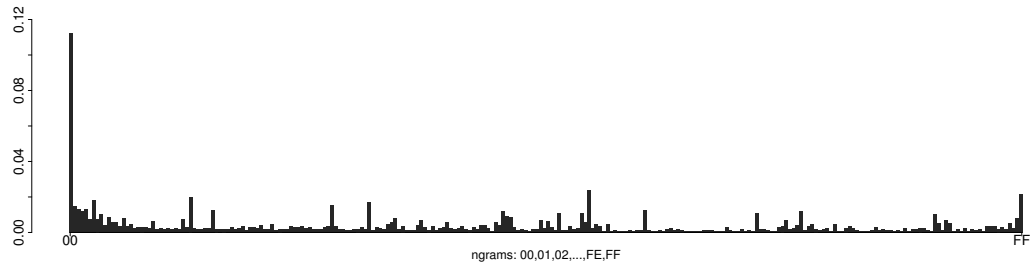
**Table 6.** Classification results for  $n$ -gram length  $n = 3$  and sliding window  $s = 1, 2, 3$  without feature selection.

benign vs. virus loader with feature selection				
s	Method	Detection Rate	False Alarm Rate	Overall Accuracy
1	SVM	<b>0.9619</b> ( $\pm 0.0173$ )	0.0334 ( $\pm 0.0086$ )	<b>0.9648</b> ( $\pm 0.0066$ )
	Bayes	0.9281 ( $\pm 0.0195$ )	0.1702 ( $\pm 0.0242$ )	0.8597 ( $\pm 0.0182$ )
	Cosine	0.1645 ( $\pm 0.0240$ )	<b>0.0065</b> ( $\pm 0.0035$ )	0.7391 ( $\pm 0.0231$ )
2	SVM	<b>0.9620</b> ( $\pm 0.0203$ )	0.0320 ( $\pm 0.0082$ )	<b>0.9660</b> ( $\pm 0.0107$ )
	Bayes	0.8205 ( $\pm 0.0358$ )	0.0701 ( $\pm 0.0151$ )	0.8962 ( $\pm 0.0200$ )
	Cosine	0.2645 ( $\pm 0.0277$ )	<b>0.0196</b> ( $\pm 0.0044$ )	0.7605 ( $\pm 0.0202$ )
3	SVM	<b>0.9549</b> ( $\pm 0.0092$ )	0.0325 ( $\pm 0.0073$ )	<b>0.9636</b> ( $\pm 0.0064$ )
	Bayes	0.7362 ( $\pm 0.0368$ )	0.0429 ( $\pm 0.0094$ )	0.8891 ( $\pm 0.0126$ )
	Cosine	0.2306 ( $\pm 0.02646$ )	<b>0.0096</b> ( $\pm 0.0034$ )	0.7573 ( $\pm 0.0128$ )

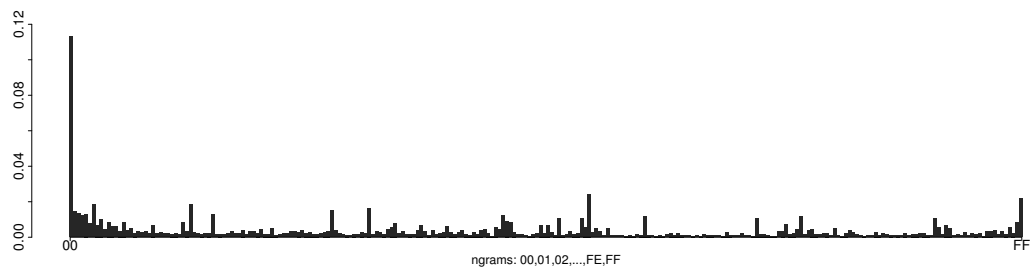
  

benign vs. virus infected with feature selection				
s	Method	Detection Rate	False Alarm Rate	Overall Accuracy
1	SVM	0.1057 ( $\pm 0.0281$ )	<b>0.1024</b> ( $\pm 0.0184$ )	<b>0.6944</b> ( $\pm 0.0196$ )
	Bayes	<b>0.6380</b> ( $\pm 0.0659$ )	0.6159 ( $\pm 0.0829$ )	0.4482 ( $\pm 0.0500$ )
	Cosine	0.4688 ( $\pm 0.1429$ )	0.4473 ( $\pm 0.1491$ )	0.5307 ( $\pm 0.0733$ )
2	SVM	0.1406 ( $\pm 0.0372$ )	<b>0.1010</b> ( $\pm 0.0220$ )	<b>0.7037</b> ( $\pm 0.0276$ )
	Bayes	<b>0.6643</b> ( $\pm 0.0844$ )	0.6287 ( $\pm 0.1157$ )	0.4455 ( $\pm 0.0657$ )
	Cosine	0.5913 ( $\pm 0.1283$ )	0.6035 ( $\pm 0.1265$ )	0.4447 ( $\pm 0.0627$ )
3	SVM	0.0611 ( $\pm 0.0185$ )	<b>0.1001</b> ( $\pm 0.0165$ )	<b>0.6842</b> ( $\pm 0.0253$ )
	Bayes	<b>0.6075</b> ( $\pm 0.0762$ )	0.5668 ( $\pm 0.0823$ )	0.4764 ( $\pm 0.0520$ )
	Cosine	0.5713 ( $\pm 0.0954$ )	0.5557 ( $\pm 0.0711$ )	0.4757 ( $\pm 0.0325$ )

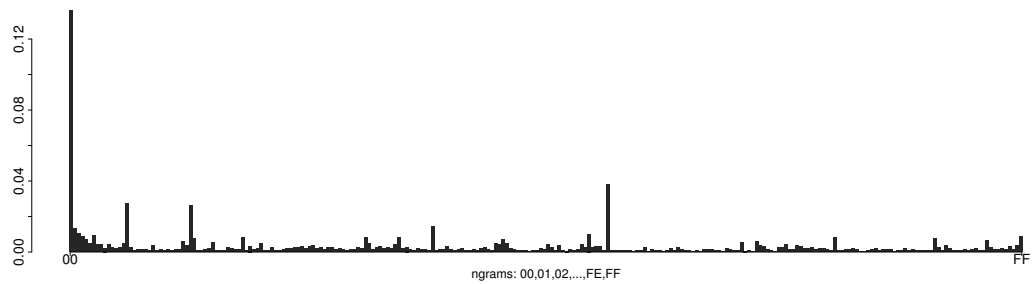
**Table 7.** Classification results for  $n$ -gram length  $n = 3$  and sliding window  $s = 1, 2, 3$  with selecting 50 % of the original features by means of the mutual information method.



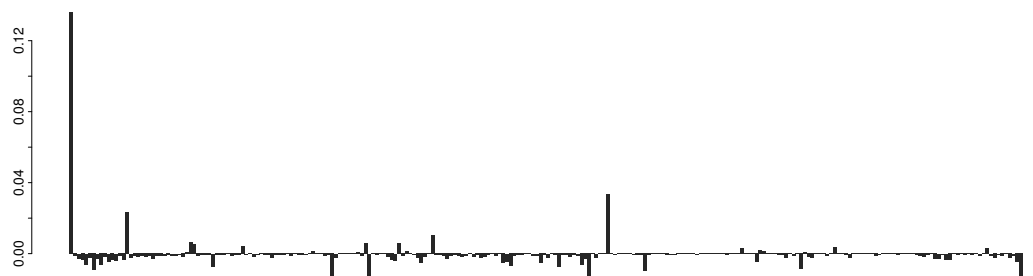
(a) Relative frequencies of  $n$ -grams created from benign files for  $n = 2$  and  $s = 2$ .



(b) Relative frequencies of  $n$ -grams created from virus infected files for  $n = 2$  and  $s = 2$ .



(c) Relative frequencies of  $n$ -grams created from virus loader files for  $n = 2$  and  $s = 2$ . Note that the relative frequency of  $n$ -gram "00" is cropped.



(d) Difference of relative frequency between benign files and virus loader files.



(e) Difference of relative frequency between benign files and virus infected files.

**Fig. 4.** Relative frequencies of  $n$ -grams created from benign files, virus infected files and virus loader files. Observe, that the relative frequencies of the  $n$ -grams created from benign files and files infected are nearly indistinguishable. As a consequence, discriminating  $n$ -grams created from benign files and virus infected files is nearly infeasible.