

Malware Program Classification Using Machine Learning and Principal Component Analysis

Nathan Kueterman
University of Dayton
Dayton, OH
nkueterman1@udayton.edu

Abstract—The field of cyber security involves processing and analyzing massive amounts of data for potential malware. The malicious programs continue to evolve and grow, which means the counter efforts need to be robust and efficient. In this project, the byte files from thousands of different malware programs were converted to images and then classified using machine learning algorithms. To remove unwanted features in the dataset, Principal Component Analysis is utilized leaving the statistically relevant features for classification. The performance of Artificial Neural Network (ANN) algorithms along with K-Nearest Neighbors (k NN) and Support Vector Machine (SVM) are compared in terms of malware classification accuracy. K-fold validation is employed to ensure that each approach is robust and effective. The study uses the Microsoft Malware Classification Challenge (BIG 2015), which is publicly available on the Kaggle database.

Keywords—Malware Detection, Principle Component Analysis, Support Vector Machine, K-Nearest Neighbors, Neural Network

I. INTRODUCTION

The classification of malware data is an essential problem at the forefront of cyber security. Having the capability to find and identify malicious software is increasingly becoming more important as the internet and device interconnectivity rapidly expand in today's society.

The approaches used in this study are scaled back in terms of data usage. The BIG 2015 dataset includes both an assembler source (.asm) file as well as the binary (.bytes) file but only the binary files were used in the classification process. Leveraging the assembly code is much more demanding as it must be decrypted and analyzed via program outputs. Doing this to thousands of files would take immense processing times. If these dynamic code approaches are employed in real-time the malware might have enough time to deploy some or all of the malicious content.

In this study, the approaches are identical to that in [1]. Essentially this is a replication of that study except only a back propagation, single hidden layer neural network was tested instead of the other ANNs that were used in [1].

II. APPROACH

A. Binary File Processing

As mentioned in the introduction, only the byte files were utilized for classification. The general format of the binary files were the same besides the variations of dimensionality. Each row of the file began with an 8 digit, decimal number and then

was followed by 16 pairs of hexadecimal numbers separated by spaces. This only changed on the last line of the file where the number of hexadecimal pairs varied. There was one special case that was common among files where the hexadecimal pair was '??', which was mitigated by replacing it with -1. This was done in order to add priors to the PCA dimensionality reduction which decreases the variance and statistical relevance of this special case.

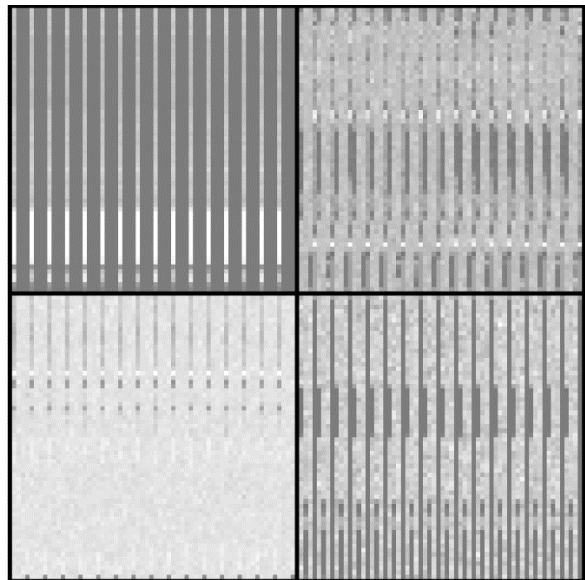


Fig. 1. 4 random samples of binary files converted into 256x16 images.

The procedure for generating the images was to remove the decimal row headers, convert each hexadecimal pair into a decimal number, concatenate them into an array, scale them from 0 to 1, replace the special cases with -1, and reshape the array using bicubic interpolation to create a grayscale image. Four samples of the malware byte images, which were scaled to 256x16, can be seen in Figure 1.

B. PCA Transformation

Once the binary files are converted to images, PCA transformation is applied in order to reduce data dimensionality. PCA transformation maps an N -dimensional vector into an M -dimensional ($N > M$) one without redundancy.

C. k NN Classifier

The k -nearest neighbor, or k NN, classifier utilizes distance in the feature space to determine the class of each sample. This

method has very few parameters but they consist of the weights of the distances, number of neighbors, and distance method. The default parameters used in this experiment for PCA performance analysis were squared inverse, 5 neighbors, and Euclidean distance.

D. SVM

SVM draws decision boundaries in the feature space using optimal hyperplanes. This hyperplane is generated using the training data and is used to classify testing data. There are more parameters and methods to calculate the hyperplane as well as decision metrics compared to the kNN method. The hyperplane was a polynomial of order 2, as well as one k-fold verification with order 1. The other kernel functions, or hyperplane methods, are Gaussian and linear which are better suited for one/two class learning.

E. ANN

The artificial neural network approach focused purely on a single layer network with a logistic, mean-squared error loss function. The number of hidden layers, iterations, and learning rate were varied for performance analysis. The initial weights were initialized randomly in the range $[-0.12, 0.12]$ on each run of the k-fold validation. The class with the highest output from the sigmoid activation function is assigned to the respective class. All pixels were considered during the training function so varying number of PCA did not affect ANN performance.

III. RESULTS

A. PCA Transformation

Experimentation with different amounts of features using PCA reveals a statistical hierarchy of the data. This information can be essential in optimizing classification in terms of speed. In this study the amount of features used varied from 10 to 25 in intervals of 5. After approximately 10 features is where the kNN classifier plateaued in performance. Approximately 25 features is where the SVM classifier started to plateau. The classifier accuracies can be seen in Figure 2.

	10	15	20	25
kNN	96.50%	96.40%	96.40%	96.30%
SVM	94.40%	94.90%	95.10%	95.30%

Fig. 2. Number of PCA features used versus overall 10-fold validation accuracy for each classifier method.

B. kNN Classifier

The kNN classifier had the highest overall 10-fold validation accuracy based on the different parameters and methods that are comparatively included in this study. The parameters that gave the highest performance were using squared inverse, 5 neighbors, Euclidean distance, and 15 PCA features. The confusion matrix output from this model can be seen in Figure 3.

		Confusion Matrix									
Output Class	1	2	3	4	5	6	7	8	9	10	
	1445	3	0	5	6	21	0	34	1	0	95.4%
	13.3%	0.0%	0.0%	0.0%	0.1%	0.2%	0.0%	0.3%	0.0%	0.0%	4.6%
	24	2441	1	1	1	5	1	9	2	0	98.2%
	0.2%	22.5%	0.0%	0.0%	0.0%	0.0%	0.0%	0.1%	0.0%	0.0%	1.8%
	2	2	2937	2	0	2	3	11	0	0	99.3%
	0.0%	0.0%	27.0%	0.0%	0.0%	0.0%	0.0%	0.1%	0.0%	0.0%	0.7%
	10	1	2	460	2	6	0	14	2	0	92.6%
	0.1%	0.0%	0.0%	4.2%	0.0%	0.1%	0.0%	0.1%	0.0%	0.0%	7.4%
	0	2	0	1	19	0	0	0	0	0	86.4%
	0.0%	0.0%	0.0%	0.0%	0.2%	0.0%	0.0%	0.0%	0.0%	0.0%	13.6%
	21	10	0	2	3	705	1	27	2	0	91.4%
	0.2%	0.1%	0.0%	0.0%	0.0%	6.5%	0.0%	0.2%	0.0%	0.0%	8.6%
	2	3	2	2	4	1	392	2	1	0	95.8%
	0.0%	0.0%	0.0%	0.0%	0.0%	0.0%	3.6%	0.0%	0.0%	0.0%	4.2%
	25	4	0	1	7	8	0	1118	3	0	95.9%
	0.2%	0.0%	0.0%	0.0%	0.1%	0.1%	0.0%	10.3%	0.0%	0.0%	4.1%
	12	12	0	1	0	3	1	13	1002	0	96.0%
	0.1%	0.1%	0.0%	0.0%	0.0%	0.0%	0.0%	0.1%	9.2%	0.0%	4.0%
	0	0	0	0	0	0	0	0	0	0	NaN%
	0.0%	0.0%	0.0%	0.0%	0.0%	0.0%	0.0%	0.0%	0.0%	0.0%	NaN%
	93.8%	98.5%	99.8%	96.8%	45.2%	93.9%	98.5%	91.0%	98.9%	NaN%	96.8%
	6.2%	1.5%	0.2%	3.2%	54.8%	6.1%	1.5%	9.0%	1.1%	NaN%	3.2%
Target Class											

Fig. 3. Confusion matrix for highest performing kNN malware classifier.

Some parameter experimentation was performed to determine what methods were optimal for classification in this dataset. Number of neighbors was the largest impact and therefore was closely studied. The highest accuracy resulted from the lowest neighbors, but not below 5. The results of this study can be seen in Figure 4.

C. SVM

The SVM algorithm did not seem too affected from changing the polynomial order. The performance went down approximately 12% when reducing the order from 2 to 1 and did not have any substantial increase when changing the order to 3. However, the number of PCA features had a large impact on performance which is prevalent in Figure 2. The highest performing SVM model, order 2 polynomial with 25 PCA features is presented in Figure 4.

		Confusion Matrix									
Output Class	1	2	3	4	5	6	7	8	9	10	
	1444	9	2	7	21	34	4	69	16	0	89.9%
	13.3%	0.1%	0.0%	0.1%	0.2%	0.3%	0.0%	0.6%	0.1%	0.0%	10.1%
	26	2413	5	2	2	10	12	17	10	0	96.6%
	0.2%	22.2%	0.0%	0.0%	0.0%	0.1%	0.1%	0.2%	0.1%	0.0%	3.4%
	2	0	2932	1	0	0	0	1	0	0	99.9%
	0.0%	0.0%	27.0%	0.0%	0.0%	0.0%	0.0%	0.0%	0.0%	0.0%	0.1%
	4	0	0	454	2	6	0	19	6	0	92.5%
	0.0%	0.0%	0.0%	4.2%	0.0%	0.1%	0.0%	0.2%	0.1%	0.0%	7.5%
	0	1	0	0	6	0	0	0	0	0	85.7%
	0.0%	0.0%	0.0%	0.0%	0.1%	0.0%	0.0%	0.0%	0.0%	0.0%	14.3%
	23	14	1	4	2	667	0	26	5	0	89.9%
	0.2%	0.1%	0.0%	0.0%	0.0%	6.1%	0.0%	0.2%	0.0%	0.0%	10.1%
	9	0	0	1	1	1	382	0	0	0	97.0%
	0.1%	0.0%	0.0%	0.0%	0.0%	0.0%	3.5%	0.0%	0.0%	0.0%	3.0%
	26	8	2	5	8	20	0	1089	4	0	93.7%
	0.2%	0.1%	0.0%	0.0%	0.1%	0.2%	0.0%	10.0%	0.0%	0.0%	6.3%
	7	33	0	1	0	13	0	7	972	0	94.1%
	0.1%	0.3%	0.0%	0.0%	0.0%	0.1%	0.0%	0.1%	8.9%	0.0%	5.9%
	0	0	0	0	0	0	0	0	0	0	NaN%
	0.0%	0.0%	0.0%	0.0%	0.0%	0.0%	0.0%	0.0%	0.0%	0.0%	NaN%
	93.7%	97.4%	99.7%	95.6%	14.3%	88.8%	96.0%	88.7%	96.0%	NaN%	95.3%
	6.3%	2.6%	0.3%	4.4%	85.7%	11.2%	4.0%	11.3%	4.0%	NaN%	4.7%
Target Class											

Fig. 4. Confusion matrix for the highest performing SVM malware classifier.

D. ANN

The ANN approach produced the lowest accuracy across the board but this is most likely because of the shallowness of the network. One aspect that is very important for ensuring proper training for the ANN was to normalize the training data. The accuracy of the same network went from 51.3% to 91.2% just from normalizing the data. Learning rate was varied from 0.05 to 0.2 and 0.1 was consistently performing the highest. Also, increasing iteration up to approximately 10000 increased accuracy but plateaued after that and was not worth the extra computing times. The model parameters that provided the highest accuracy was using normalized features, 25 hidden layers, 10000 iterations, and a learning rate of 0.1 which produced an accuracy of 91.2% which can be seen in Figure 5 below.

Confusion Matrix											
Output Class	1	2	3	4	5	6	7	8	9	10	
	1382 12.7%	35 0.3%	0 0.0%	11 0.1%	5 0.0%	56 0.5%	4 0.0%	75 0.7%	20 0.2%	0 0.0%	87.0%
	35 0.3%	2366 21.8%	4 0.0%	19 0.2%	1 0.0%	17 0.2%	16 0.1%	19 0.2%	62 0.6%	0 0.0%	93.2%
	9 0.1%	1 0.0%	2930 27.0%	1 0.0%	4 0.0%	4 0.0%	0 0.0%	12 0.1%	0 0.0%	0 0.0%	99.0%
	10 0.1%	6 0.1%	1 0.0%	433 4.0%	6 0.1%	29 0.3%	11 0.1%	30 0.3%	7 0.1%	0 0.0%	81.2%
	0 0.0%	0 0.0%	0 0.0%	0 0.0%	0 0.0%	0 0.0%	0 0.0%	0 0.0%	0 0.0%	0 0.0%	NaN%
	48 0.4%	31 0.3%	2 0.0%	4 0.0%	1 0.0%	591 5.4%	3 0.0%	29 0.3%	59 0.5%	0 0.0%	77.0%
	10 0.1%	1 0.0%	0 0.0%	4 0.0%	13 0.1%	2 0.0%	364 3.3%	3 0.0%	0 0.0%	0 0.0%	91.7%
	30 0.3%	7 0.1%	5 0.0%	2 0.0%	12 0.1%	36 0.3%	0 0.0%	1056 9.7%	4 0.0%	0 0.0%	91.7%
	17 0.2%	31 0.3%	0 0.0%	1 0.0%	0 0.0%	16 0.1%	0 0.0%	0 0.0%	861 7.9%	0 0.0%	92.6%
	0 0.0%	0 0.0%	0 0.0%	0 0.0%	0 0.0%	0 0.0%	0 0.0%	0 0.0%	0 0.0%	0 0.0%	NaN%
											91.9%
Target Class											8.1%

Fig. 5. Confusion matrix for the highest performing ANN malware classifier.

IV. CONCLUSION

Overall it seemed that most simple solution, kNN, was the most robust in this situation. There are many other algorithms out there that could outperform it in terms of accuracy but it is among the fastest of classifiers which is an aspect to consider in this realm. The one glaring issue in this study is the very low classification accuracy on class 5. This class only had 42 training samples so it was heavily imbalanced. Some interesting future studies or improvements could include trying different types of interpolation on the images, using different neural networks especially deeper ones, experimenting with class imbalance and possibly augmenting more of the minority class.

REFERENCES

- [1] B. N. Narayanan, O. Djaneye-Boundjou, and T. M. Kebede, "Performance analysis of machine learning and pattern recognition algorithms for Malware classification," *2016 IEEE National Aerospace and Electronics Conference (NAECON) and Ohio Innovation Summit (OIS)*, 2016.