

Malware Prediction and Classification Using Advanced Modeling Techniques

Bharath Ratna Karumudi, Sriganesh Chandrasekaran, Barry Armour, Izzat Alsmadi

College of Engineering and Computer Science
Syracuse University
Syracuse, NY USA - 13244
(bhkarumu,srchandr,blarmour,imalsmad)@syr.edu

Abstract — Antivirus and malware detection software's primarily uses a signature based approach. These software's use a set of rules from a library, based on previously known malware and published vulnerabilities from the software manufacturer. The rules are very specific to the vulnerability and malware, thus making them brittle and very context and malware specific. So a new malware release still has the potential to take advantage of the same vulnerability that has already been addressed for a different malware. The Antivirus software provider releases updates to the library of rules on a periodic basis, or as critical vulnerabilities are discovered. One risk with the traditional Antivirus or malware detection software is zero-day vulnerabilities, where the hackers discover it before either the software manufacturer or the antivirus provider.

In this implementation, we will be attempting to build a malware prediction model using Machine Learning with which we can overcome the above mentioned issues with traditional malware detection software, using the publicly available malware dataset.

Keywords—*Malware, n-grams, Neural Networks, MalConv, Stacking, Random Forest, Machine Learning.*

I. GOAL/HYPOTHESIS

The goal of our project is that we are given a data set [1] that contain malware data in bytes and asm files. We would like to predict the best possible match for the malware to a set of target based on the program characteristics (like Header, bytes file info, IDAT disassembled code features). The overall research hypothesis for this dataset is to parse the data, extract the features, model and then classify the malware into one of the nine categories: *Ramnit, Lollipop, Kelihos_ver3, Vundo, Simda, Tracur, Kelihos_ver1, Obfuscator.ACY, Gatak.*

We will initially experiment with numerous algorithms and formally move further with the one that

accomplishes the outcome with best accuracy. The main challenge with the dataset is the domain knowledge and we need to do domain research as well to do some feature analysis which will help us determine potential features to use.

There are two types [4] of Malware analysis:

- *Static analysis* analyses the malware attributes by executing a static parser and understanding underlying patterns on the suspect malware fingerprints.
- *Dynamic Analysis* is basically executing the suspect malware in a controlled environment and understanding its behaviour.

In our approach, we will be extending the static analysis with Machine Learning and model predictions through the end goal of mapping the dataset with high degree of accuracy to the provided fingerprints. We also intend to understand and learn on the static analysis of malwares (pros and cons) as an additional item. If computing power is not sufficient and scalable with our limited resources to perform the analysis, we may end up limiting the data set.

The entire Machine Learning processing involved:

- Data collection and Preprocessing.
- Feature Selection/reduction.
- Classification analysis.
- Modeling
- Ensembling - techniques to improve the overall performance metrics including stacking and voting based techniques for multivariate classes.

II. DATASET

The dataset is from Kaggle shared by Microsoft Corporation [1] and is about 500 Gigabytes, which includes:

- 10,868 malware files (Bytes file+ASM disassembled version) in training dataset.
- 10,873 malwares in test dataset.
- trainLabels.csv - which has the labels for training dataset.

Every malware has two related files:

- Hexadecimal byte contents.
- ASM files produced by a commercial Disassembler IDA pro.

Figure 1 and Figure 2 shows a sample of byte and asm files respectively.

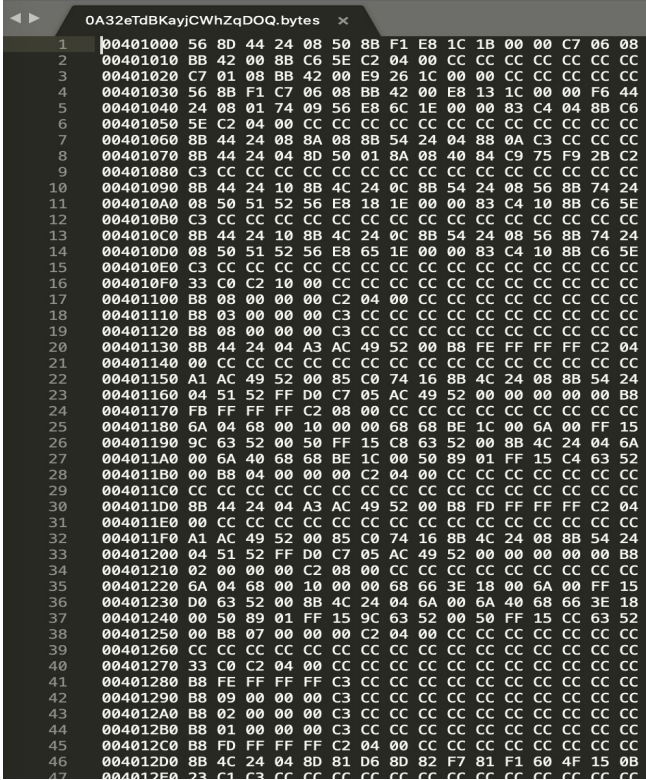


Fig.1. Sample bytes file

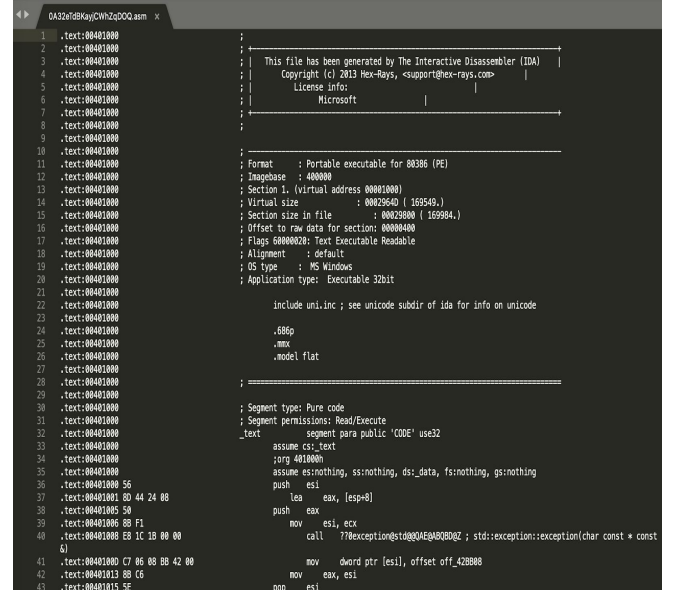


Fig.2. Sample asm file

III. RELATED CONTRIBUTION

We have identified several implementations on this Kaggle competition dataset [1]. This section will summarize their implementation and our understanding of their approaches. We referred to some of the solution approaches from Kaggle and how they had proceeded from thought process to give us some ideas. Almost all major approaches used some sort of XGBoost or gradient boosting techniques to improve performance as the dataset is humongous to train and validate on.

Out of them, one is by the first place winner [2] whose approach is on intensive feature engineering and used the Gradient boosting (XGBoost), Ensembling, semi-supervised learning. The main model implemented is Gradient Boosting and also cross validation helped in avoiding the overfitting. They extracted three kinds of features: opcode N-gram count, segment line count and asm file pixel intensity features and extracted the features by writing the code in Python. For supervised modeling, XGBoost and ensemble are used. The figure 4 can depict the high level overview of their approach.

Also an interesting implementation was ASM file pixel intensity feature, where malwares can be visualized as gray-scale images. So they extracted the gray-scale image from asm files. The model was built on XGBoost to perform a multiclass classification. The implementer also tried Random Forest, Naive Bayes, but XGBoost gave better results.

Code Execution:

The implementer's code[2] is available in GitHub repository, but it is almost five year old and a lot of libraries are obsolete now and running on Python 2.7. We tried our best to refactor the code to rerun and figure 3 is the code execution to create a single model.

```
(base) bharath@melodic-park $ ./single_model.sh
Gathering 4 grams, Class 1 out of 9...
Gathering 4 grams, Class 2 out of 9...
Gathering 4 grams, Class 3 out of 9...
Gathering 4 grams, Class 4 out of 9...
Gathering 4 grams, Class 5 out of 9...
Gathering 4 grams, Class 6 out of 9...
Gathering 4 grams, Class 7 out of 9...
Gathering 4 grams, Class 8 out of 9...
Gathering 4 grams, Class 9 out of 9...

DONE 4 gram features!
DONE bytes count!
DONE instruction count!
DONE DLL functions!
DONE asm image features!
mkdir: cannot create directory 'op_train': File exists
mkdir: cannot create directory 'jump_train': File exists
mkdir: cannot create directory 'jump_map_train': File exists
progress 0.0 457
0.0
0.111111111111
0.222222222222
0.333333333333
0.444444444444
0.555555555555
0.666666666667
0.777777777778
0.888888888889
01kcfWA9K2B0xQeS5Rju .text:10001591 FF 25 70 20 00 10
p ds:_onexit
0.0
01azqd4InC7m9JpocGv5 .text:004113B2 FF 25 CC 25 41 00
p ds:_putenv
01azqd4InC7m9JpocGv5 .text:004113BE FF 25 D4 25 41 00
p ds:_mbscat
01azqd4InC7m9JpocGv5 .text:004113C4 FF 25 C4 25 41 00
p ds:_scalb
01azqd4InC7m9JpocGv5 .text:004113D6 FF 25 B8 25 41 00
p ds:_mbsinc
01azqd4InC7m9JpocGv5 .text:004113DC FF 25 B4 25 41 00
p ds:_wcsdup
01azqd4InC7m9JpocGv5 .text:004113E8 FF 25 AC 25 41 00
p ds:_heapset
01azqd4InC7m9JpocGv5 .text:004113F4 FF 25 A4 25 41 00
p ds:_cabs
01azqd4InC7m9JpocGv5 .text:00411410 FF 25 94 25 41 00
p ds:_assert
01azqd4InC7m9JpocGv5 .text:00411440 FF 25 70 25 41 00

{'push': 16225, 'call': 12706, 'cmp': 24268, 'mov': 40183, 'lea': 10748, 'jmp': 3433, '
2556, 'add': 14530, 'jz': 12161, 'pop': 1389, 'xor': 381, 'jge': 336, 'jle': 296, '
7, 'db': 7772, 'ja': 9956, 'dec': 9955}
{'push': 16225, 'call': 12706, 'cmp': 24268, 'mov': 40183, 'lea': 10748, 'jmp': 3433, '
2556, 'add': 14530, 'jz': 12161, 'pop': 1389, 'xor': 381, 'jge': 336, 'jle': 296, '
7, 'db': 7772, 'ja': 9956, 'dec': 9955}
0.0 142
0.111111111111 2042207
0.222222222222 227812
0.333333333333 3278
0.444444444444 144751
0.555555555555 8717
0.666666666667 3175
0.777777777778 424908
0.888888888889 94351
```

Fig.3. Code execution after refactoring

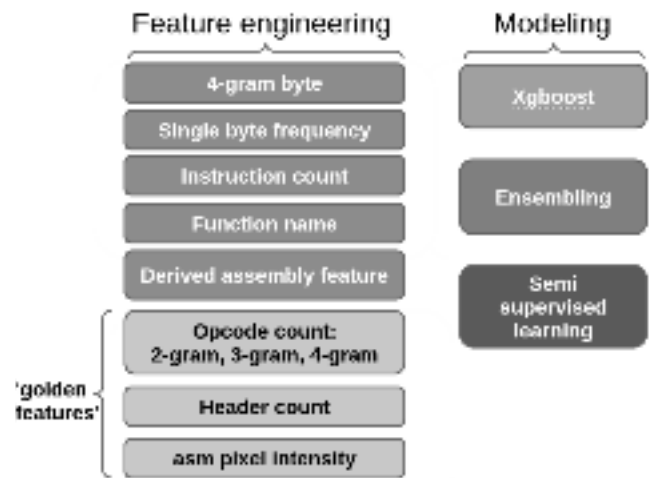


Fig.4. Overview of the approach

The second approach, who was the second place winner [5], used iterative method of feature extraction and model training. In this approach, the researchers had vetted both bytes file and the ASM file for certain unique characteristics that can be considered as features.

Bytes file: One interesting thing to note was that the approach actually took the physical file size properties (in terms of the bytes file and the ratio with corresponding ASM files) which is pretty interesting. This was also done for the ASM file to extract as much features as possible.

There is also the usage of finding out the numeric values of 1, 2 and 4 byte n-grams and calculate ratios of these. The way N-grams were calculated is as below:

Let's consider the *OA32eTdBKayCWhzqDOQ.bytes* file as an example and in figure 5:

```

nder File Edit View Go Window Help
OA32eTdBKayCWhzqDOQ.bytes - wzQL3hi6zy Add License
1 00401000 56 80 44 24 08 50 88 F1 E8 1C 1B 00 00 C7 06 08
2 00401010 88 42 00 88 C6 5E C2 04 00 CC CC CC CC CC CC
3 00401020 C7 01 08 B8 42 00 E9 26 1C 00 CC CC CC CC CC
4 00401030 56 88 F1 C7 06 08 B8 42 00 E8 13 1C 00 00 F6 44
5 00401040 24 08 01 74 09 56 E8 6C 1E 00 00 83 C4 04 8B C6
6 00401050 5E C2 04 00 CC CC CC CC CC CC CC CC CC CC
7 00401060 88 44 24 08 8A 08 88 54 24 04 88 0A C3 CC CC CC
8 00401070 88 44 24 04 8D 50 01 8A 08 40 84 C9 75 F9 28 C2
9 00401080 C3 CC CC CC CC CC CC CC CC CC CC CC CC CC

```

Fig.5. OA32eTdBKayCWhzqDOQ.bytes file

The way 1 gram, 2 gram and 4 grams are calculated are as below -

```

00401000 56 - 1 gram
00401000 44 24 - 2 gram
00401000 F1 E8 1C 1B - 4 gram

```

The research then goes into collating the weightage of each of these n-grams and consolidating based on popular or repeatable n grams that can be used as feature for prediction.

ASM file - In addition to the file size and compression size characteristics, within the ASM file interpunction characters like ‘,’ ‘?’ ‘.’ ‘+’ ‘-’ etc. were also looked at to understand the density of these characters. They also looked at the proportion of lines or characters in each section namely the HEADER, .data, .rdata. This approach though does not need a lot of domain knowledge behind the assembler coding itself. The model had used XGBoost and gradient boosting trees.

The whole analysis is done in two parts -

1. First, evaluate the set of new features using untuned modelling
2. Predictions from untuned model is then used as features in the meta bagging model for the final accuracy evaluation.

The other approach [6] combined an examination of the software fingerprint of each ASM file to look at the type of code used in the ASM file as well as the bytes file ngrams/characteristics. This helped to come up with a pool of features and categorize them into relevant feature sets. The total number of features pulled in were 28,000 and which were grouped into 9 feature sets based on unigrams, bigrams in bytes files and opcodes frequency on n-grams in ASM files at a high level. The method to classify used Gradient Boosting with a depth of 3 and a K-fold validation where K=6.

One interesting thing that we found out is that there will be a huge challenge in parsing the ASM files and properly classifying them into sections like .data, .rdata etc. We are analysing multiple ways of coding this parser tool and are trying out few examples to ensure we can select the right features, correlation between the features and then doing the predictive modelling tests.

One common theme used across multiple solutions is ensemble boosting and stacking methods which improves the performances of the models and help in scalability.

IV. FEATURE SELECTION/ENGINEERING

Each malware file has an identifier, a 20 character hash value uniquely identifying the file, and a class label, which is an integer representing one of the nine family names to which the malware belongs to. There are multiple ways that we are planning to approach the feature selection for both bytes file and the ASM files.

We used *RapidMiner* predominantly for feature selection and engineering techniques. It is critical to gather the tokenized list of n-grams and opcodes for rich and varied subset of features, which we describe in the below section. In addition, since the volume of feature set was high, it was critical to do feature reduction using known techniques like PCA to get out of the ‘Curse of Dimensionality’ and avoid noise in the datasets.

Data collection and Preprocessing: The data preprocessing included writing elaborate processes to get useful information out of the large training dataset provided, reading them as text, tokenizing and generating n-grams/filtering tokens as shown in the process snapshot below.

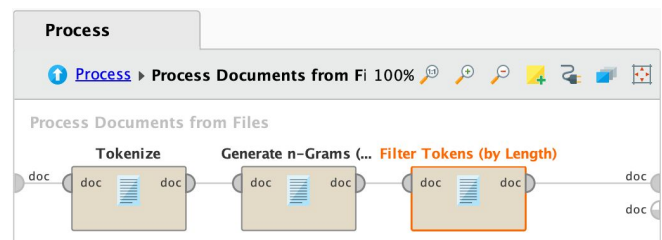


Fig.6. Collection of n-grams from Bytes files

Bytes file - Extract the features of 1 gram, 2 gram and n-gram hash codes. For our experiment, due to computational limitations we restricted our modelling to 1 gram and 2 gram on the Bytes file which still yielded the necessary accuracy. During the selection of attributes we used the below steps to process documents from files in loop, generate and select attributes as shown in figure 7.

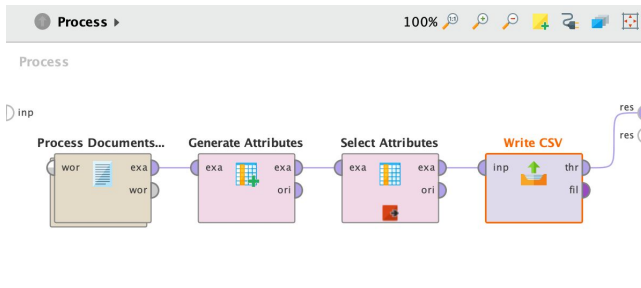


Fig.7. Final extraction/selection of bigram attributes to n-gram files

The mode used for bi-gram collection from Hex bytes file was to process the entire set of bytes file as a Process Document loop in RapidMiner. Series of extraction steps involve tokenizing the entire hex file for the 2 byte grams (modifying the filter criteria to extract 2-byte grams), their TF-IDF frequencies, generating the attributes including the metadata file attributes and consolidating into a single CSV for further processing steps. The base n-grams were 2-byte opcode as explained by Zak. et.al [7]. This yielded close to 65,000+ attributes which were further reduced to the most relevant attributes using the variance based selection/correlation techniques to 1000+ n-gram features on the collection dataset.

ASM file - We looked at the base set of opcode operators with particular focus on opcodes that impact memory processing like arithmetic/branching instructions (sub, mov, jmp, jnz) for tokenization and filtering. We processed the ASM files as text files, tokenized based on the opcodes n-gram and processed as a CSV file for the initial step.

The process step for ASM files was constructed similar to that of the bytes file and shown in figure 8.

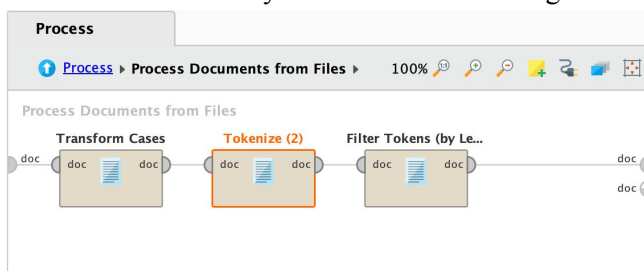


Fig.8. ASM files Tokenization of Opcode N-grams

The parsing was done through RegEx based tokenization and extraction which yielded good results in the data collection step. This is similar to using a custom dictionary /wordlist to pull the relevant tokens from the ASM files as the opcodes may not substitute for english words directly in the dictionary. This ended

up with 150+ core opcode based attributes to be processed on the ASM files.

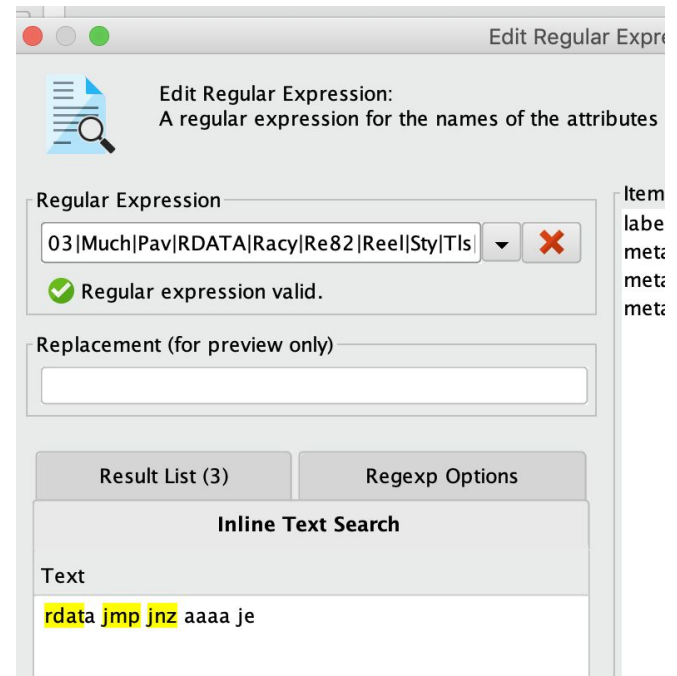


Fig.9. List of Opcodes in RapidMiner RegEx operator

A quick visualization on the bytes n-gram as shown in figure 10 clearly indicates there are outliers that can cause additional noise for modelling. We may have to remove them further using the normalization/feature reduction step before running through the modelling techniques.

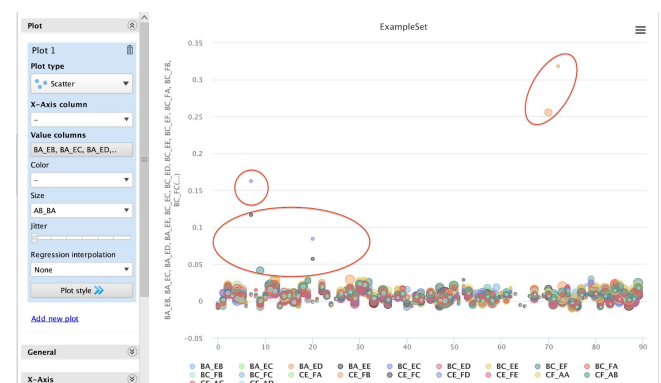


Fig.10.: Outliers in bigrams extracted from bytes files

On a similar approach, ASM file features were also extracted and figure 11 shows a visualization of the extracted dataset for further processing. On close analysis of the extracted features, we found a dense scatter on some of the memory handling/branching opcodes like mov, jnz etc as can be seen below. These seem to add more information to the feature dataset

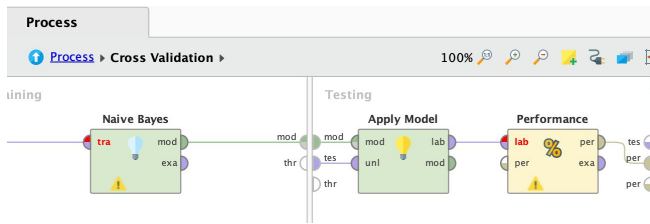


Fig.15. Naive Bayes based classification on the dataset

While AdaBoost on Random forest did not yield improved results, we tried expanding the modelling to Gradient boosting trees.

As a next step to improve the overall accuracy of the predictive modelling, we did stacking ensembling technique with decision tree, random forest and generalized linear model with random forest as the stacking learner. With ensembling technique the accuracy of the model improved to 85%.

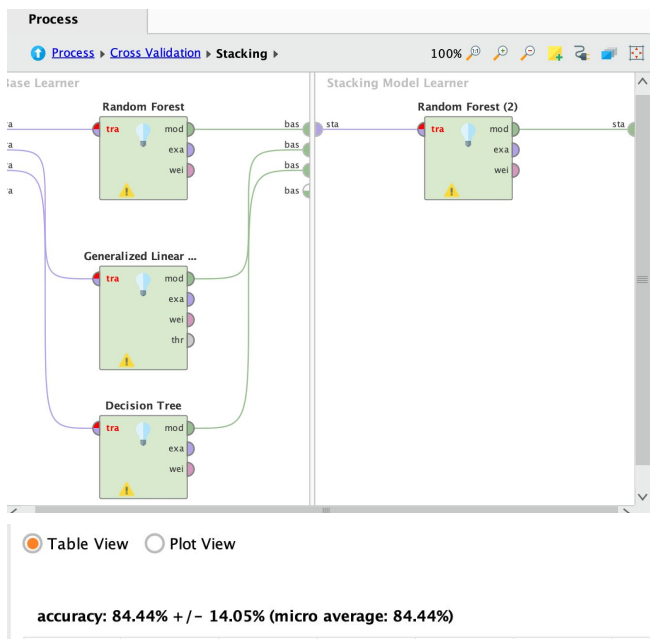


Fig.16. Ensembling using Stacking method on ASM Opcode n-gram and its accuracy

The features selected included key system specific opcode n-grams which were obtained by feature selection (through removing correlated attributes). We kept the remove correlation at 0.85, which yielded in 65 significant attributes for feature dataset. Due to the significant bias in the class 5, we feel that the results when extrapolated for large values can be slightly skewed against this class.

As we did additional modelling/ensemble techniques, with the varied feature set and high dimensions, we needed a better mechanism for dimensionality reduction and getting a higher prediction score.

Malconv - A Neural Network based Architecture

There were many implementations on this dataset during competition, but none with Malconv, so we have selected this architecture and built the model using Malconv. The Malconv model was developed by Nvidia by Raff et al.[3] as shown in figure 17, specially for Malware detection and the model aimed for three features:

1. The ability to scale efficiently in computation and memory usage with sequence length.
2. The ability to consider both local and global context while examining the entire file.
3. An explanatory ability to aid analysis of flagged malware.

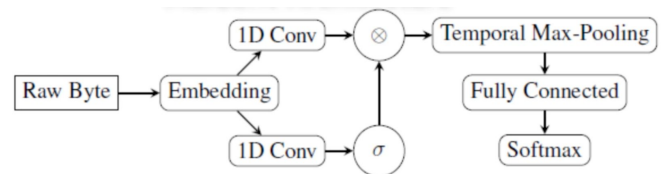


Fig.17. Malconv Architecture

The model does not need to perform any feature engineering ahead.

Malconv Model Execution:

We intend to split the Train dataset into subsets - Train and Test for modelling. For this we have a shell script that will segregate the files into two different directories from where our program will read.

```
(base) karumudi7@melodic-park $ ./train_test_split.sh
Script started at: Fri Nov 29 20:20:33 EST 2019
cp: cannot stat '/mnt/disks/MLProject/data/train_dataset/Id.bytes': No such file or directory
cp: cannot stat '/mnt/disks/MLProject/data/train_dataset/Id.bytes': No such file or directory

Script ended at: Fri Nov 29 20:46:59 EST 2019
(base) karumudi7@melodic-park $
(base) karumudi7@melodic-park $
```

Fig.18. Running the Train-Test split process

The program was written in Python 3 by using Keras and TensorFlow libraries. The approach starts by splitting the given Train dataset into two Train and Test by randomizing the given TrainLabels.csv - for this the above mentioned shell script will do the operations on the file system and prepares the dataset.

We then clean the byte sequences for any unknown characters and do the embedding by converting the input to a numerical format - each byte into an eight dimensional vector. Using Keras functional API, we then captured the Conv1D as conv1 and conv2 and then activate one of the convolutions, in our case conv2, sigmoid. We multiply the Activated convolution and non activated convolution and the result was activated - relu. With these, we build the model and the built model summary as shown in figure 19.

```
model.summary()
```

Layer (type)	Output Shape	Param #	Connected to
input_1 (InputLayer)	(None, 8, 30000)	0	
conv1d_2 (Conv1D)	(None, 1, 32)	122880032	input_1[0][0]
conv1d_1 (Conv1D)	(None, 1, 32)	122880032	input_1[0][0]
sigmoid (Activation)	(None, 1, 32)	0	conv1d_2[0][0]
multiply_1 (Multiply)	(None, 1, 32)	0	conv1d_1[0][0] sigmoid[0][0]
relu (Activation)	(None, 1, 32)	0	multiply_1[0][0]
global_max_pooling1d_1 (GlobalM	(None, 32)	0	relu[0][0]
dense_1 (Dense)	(None, 16)	528	global_max_pooling1d_1[0][0]
dense_2 (Dense)	(None, 9)	153	dense_1[0][0]
Total params: 245,760,745			
Trainable params: 245,760,745			
Non-trainable params: 0			

Fig.19. Malconv Model Summary

The built model is trained with our training dataset and predicted with an accuracy of 97.6%.

```
#Train Model Accuracy
str(sum(pred_label == Train_Y_np)/len(Train_Y_np)*100)+"%"
'97.62269938650306%'

print(datetime.datetime.now())
2019-12-01 21:19:11.688610
```

Fig.20. Training Accuracy

The model is also tested with our test dataset and predicted accuracy of 89.4%.

```
#Test Model Accuracy
str(sum(Test_Y_pred_label == Test_Y_np)/len(Test_Y_np)*100)+"%"
'89.44342226310947%'

print(datetime.datetime.now())
2019-12-01 22:16:15.829232
```

Fig.21. Test Accuracy

The given test dataset do not have any labels, so we performed a prediction on the whole test dataset and model classified the malwares as shown in figure 22.

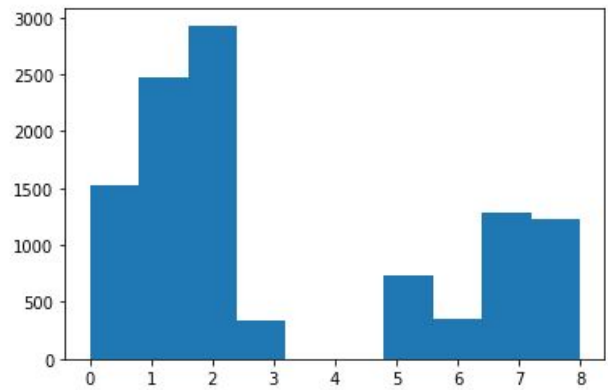


Fig.22. Histogram showing the predicted label distribution on 10,873 test malwares dataset

VI. SUMMARY

We have built our machine learning model using multiple models and below is a summary of our work and observations:

The initial classification and training on the RapidMiner was done through RandomForest. This provided reasonable results on the opcode n-gram and byte n-gram relevant features dataset (reduced to 100+ attributes). Overall accuracy was 81%. Other classifiers like Linear Regression yielded similar results as well. This took an extrapolated training time of close to 6 hours due to volume of dataset and tokenization timeframe. The precision was at 74% and Recall at 81%. We then created an ensemble model for training with Stacking (Random Forest + Generalized Linear Model + Decision Tree) with an accuracy of 88%, which was a 4% increase in accuracy. We also did a AdaBoost based boosting model in RapidMiner which yielded similar results. With Stacking Ensemble, the Precision improved to 79% and Recall improved to 84%.

The Malconv model that was built took 5 hours to process end-to-end and with all the 21K+ byte files in the train and test datasets. Out of 5 hours, it took 2 hours to build and train the model with an accuracy of 97.6% and testing with cross validation took 1 hr yielding an accuracy of ~90%. The model was also used to predict the test dataset of 10K+ byte files which took approximately 2 hours to complete. Over all the model performed better in both running time and accuracy than our initial expectations. It looks like none of the previous users of this dataset used Malconv model to compare our results with them.

Some of the winning contributions had close to 95%+ accuracy. Our results are in-line with the earlier results however done through different techniques. Malconv is more contemporary framework/solution compared to some of the original submissions and different approach to the same problem statement.

VII. ARTIFACTS AND TECHNOLOGY USED

- Code base: [GitHub](#)
- Environment:
 - Malconv Model: 8 vCPU, 64GB RAM, Ubuntu 16.04.6 LTS, 600 GB, Python 3.7 - Anaconda 2019.10, Jupyter IDE on GCP.
Libraries: Keras and Tensorflow.
 - RapidMiner 9.4: 2 CPU, 8 GB RAM, Mac OSX.
Extensions: Text Processing, Statistics, Community extension.

VIII. CONCLUSION

The traditional antivirus and malware detection software's approach is insufficient in growing malwares every day. Because every environment is unique and has specific binaries and these detection softwares might never seen those before and with millions of new malware samples every day it will be difficult. So there is a need to develop a detection system that can adapt to the rapidly changing malware ecosystem is seemingly a perfect fit for Machine Learning. As implemented in this project with Malconv[3] architecture, we could build an anti-virus system without feature engineering and that would save efforts and allows to detect the malware across a variety of operating systems and hardware. Also, as the cybersecurity threat landscape keeps changing, more popular and contemporary modelling techniques like Malconv keep evolving in the Cybersecurity Data Science research space.

IX. REFERENCES

- [1] Microsoft Corporation. (2015). Microsoft Malware Classification Challenge (BIG 2015). Retrieved from <https://www.kaggle.com/c/malware-classification/data>.
- [2] Xiaozhou Wang, Jiwei Liu, Xueer Chen, "Say No to Overfitting" Winner of Microsoft Malware Classification Challenge (BIG 2015).
- [3] Raff, E., Barker, J., Sylvester, J., Brandon, R., Catanzaro, B., & Nicholas, C.K. (2017). Malware Detection by Eating a Whole EXE. ArXiv, abs/1710.09435.
- [4] M. Sharif, V. Yegneswaran, H. Saidi, P. Porras, W. Lee, Eureka: A Framework for Enabling Static Malware Analysis, in: Proceedings of the 2008 European Symposium on Research in Computer Security (ESO-RICS).
- [5] Michailidis, M., & Jacobusse, G. (2015). Microsoft Malware Classification Challenge (BIG 2015). Retrieved October 2019, from <https://www.kaggle.com/c/malware-classification/forums/t/13863/2nd-place-code-and-documentation>
- [6] Microsoft Malware Classification Challenge 6th place solution 21.04.2015, Oleksandr Lysenko
- [7] Zak, R., Raff, E., & Nicholas, C. (2017). What can N-grams learn for malware detection? 2017 12th International Conference on Malicious and Unwanted Software (MALWARE). doi: 10.1109/malware.2017.8323963
- [8] Metamorphic Virus Detection in Portable Executables Using Opcodes Statistical Feature, Babak Bashari Rad, Maslin Masrom.
- [9] Detecting unknown malicious code by applying classification techniques on OpCode patterns. Asaf Shabtai, Robert Moskovitch, Clint Feher, Shlomi Dolev and Yuval Elovic.
- [10] Malware Images: Visualization and Automatic Classification. L. Nataraj, S. Karthikeyan, G. Jacob, B. S. Manjunath.
- [11] Mitchell Mays, Noah Drabinsky, and Stephan Brandle. 2017. Feature Selection for Malware Classification.. In MAICS. 165–170.