

Malware detection using machine learning

Fan Li

CSI-5500 Operating Systems

School of Engineering and Computer Science Oakland University

fanli@oakland.edu

Abstract

In recent years, people have become increasingly inseparable from the Internet in their working and lives, which has led to the existence of underground economy such as malware. Malware can steal our information and cause damage to individuals or businesses. Compared with firewalls, data encryption and other traditional network security protection methods, intrusion detection and malware detection systems have irreplaceable advantages in information mis-information, software leakage, and have become the main approaches of protecting network security. In this paper, I tried to detect malware using machine learning. I extracted static features from binary files which will be used as the dataset for training the model. Comparing the accuracies of different models on different datasets and finally get the best accuracy is 99.79%.

1 Introduction

Malware threats have been with us since the development of computer technology. Malwares have also been iteratively updated with the times, and the number of global cybercrimes and security breaches has increased in recent years. Malwares often steal sensitive information (such as credit card numbers or passwords) from users without their knowledge, or steal their email accounts to send fake emails, thereby damaging users' or business property. Malwares include, but is not limited to, viruses, worms, spywares, adwares, and Trojan horses. In order to better protect the security of information property of individuals and enterprises, I proposed a method of utilizing machine learning techniques to detect the malwares. In this study, I have read the relevant literature and in Section II I present the solutions for malware in recent years. Section III shows the types of malwares and the methods for extracting malware features, which are static feature extraction and dynamic feature extraction. In the experiments I used static analysis to extract features from PE files. In the model se-

lection, I choose several commonly used models, including GaussianNB, MLP, Linear Regression, Decision Tree, and Gradient Boosting. I design three different data preprocessing methods to train the model, the first one retains the extracted feature information, the second one uses feature selection to process the data set and filters out some features before training the model, and the third one uses an AutoEncoder to encode original features to their latent representations. Finally, conclusions are drawn by comparing the performance of different models on different data. I concluded that DecisionTree has the highest accuracy, the accuracy is 99.79%.

2 Related work

A framework called PEMiner is proposed in [1] for malware detection and classification. They extracted several features from PE files: the COFF header, the optional header and some section attributes. In [2], a new approach to monitor malware called Malware Instruction Set (MIST), which uses data mining and machine learning techniques to achieve effective software behavior analysis. A static analysis approach is used in [3] to extract three different features: strings, system resource information, and byte sequences. In [4], n-grams of bytecodes are used as features and machine learning algorithms such as: decision trees [4], naive Bayes, support vector machines and boosting are applied, where boosted decision trees outperform the other methods. In [5], the authors categorize methods according to

- What is the problem they want to solve.
- What are the types of features extracted from the portable files.
- What are the machine learning algorithms they use.

In addition, feature selection can be used with self-supervised learning to improve performance. In this research, I will deal with the malware detect using machine learning from two aspects: obtaining the malware features

through feature extraction, then training the model with the features.

3 Background

3.1 Malware

Malware is a malicious application or code that damages or disrupts the normal use of an endpoint device. When a device is infected with malware, you may experience unauthorized access, data leakage, or the device may be locked until you pay a ransom. There are three main types of malwares on the Web today: viruses, worms, and Trojans, and five widely used types of malwares: rootkits, keyloggers, ransomware, spyware, and adware [6].

Virus: A virus is a type of malware, but not all malware is a virus. Most malware programs are not viruses, and computer viruses try to attach themselves to legitimate programs in order to spread their malicious code to as many other programs and devices as possible. When the execution of the victim file is activated, the virus is also executed and replicates to infect other files on the system.

Typically viruses are spread by exploiting vulnerabilities in software applications and operating systems, they can also infect user devices through

- Infected email attachments
- Malicious scripts and website downloads
- Removable external media

Nowadays, pure computer viruses are not very common, which is considered a good thing for users, but viruses are very difficult to remove because they execute malware from legitimate programs, so in most cases the only way to remove a virus is to quarantine or delete the infected file. **Worms:** Like viruses, worms are a type of malware that can also spread from the host computer to other devices. However, unlike viruses, worms are more insidious because they do not require human intervention to spread and are runnable, standalone malicious programs. Based on their unique characteristics of self-replication, worms therefore spread faster and more widely and are more destructive than viruses. **Trojan horse:** A Trojan horse is a malicious software or code program used by hackers to remotely control a device. It is a piece of malicious code with special functions hidden in a normal program that can perform illegal operations such as monitoring, data modification, and stealing important information from the controlled computer. Trojan horses are usually disguised and hidden in the following legitimate files: Office documents (e.g. Word and Excel, PDF files, executable files. Trojans are usually spread via email or pushed when users visit infected websites. **Rootkit:**

Rootkits are a remote access tool or application provided for use by remote administrators of devices. Hackers can use Rootkits to hide traces of intrusion activity, retain ROOT access, and also hide malicious programs in the operating system. **Spyware:** Spyware is a type of surveillance software that allows someone to observe your activities without your knowledge or authorization, and through this type of software can access your sensitive personal information and data.

3.2 Malware analysis

3.2.1 PE file format

The malware used in this paper is PE files. PE files are the general name of executable files under Windows, such as DLL, EXE, OCX, SYS, etc. Generally speaking, the structure of PE files is shown in the following diagram: DOS header, NT header, section table and specific sections in order from the starting position [5]. The DOS header is used for MS-DOS compatibility to indicate a text when the file runs on MS-DOS, in most cases: This program cannot be run in DOS mode. Another purpose is to indicate the location of the NT header in the file. The NT header contains the main information about the windows PE file, including a 'PE' signature, the PE file header (IMAGE_FILE_HEADER) and the PE optional header (IMAGE_OPTIONAL_HEADER32). section table: is the description of subsequent sections of the PE file, windows loads each section according to the description of the section table. section: each section is actually a container that can contain code, data, etc. Each section can have independent memory permissions, for example, the code section has read/execute permissions by default, and the name and number of sections can be defined by yourself, not necessarily the three in the below figure 1.

3.2.2 PE file execution sequence

When a PE file is executed, the PE loader first checks the offset of the PE header in the DOS header. If it finds it, it jumps directly to the PE header location.

After the PE loader jumps to the PE header, the second step is to check if the PE header is valid. If the PE header is valid, it jumps to the end of the PE header.

After the PE loader executes the second step, it starts to read the section information in the section table and uses file mapping (when executing a PE file, Windows does not read the whole file into memory at the beginning, but uses the memory mapping mechanism, that is, the Windows loader only establishes the virtual address and the PE file when it loads, and only when it actually executes the PE file does the file mapping mechanism work. This means that the Windows loader only establishes the mapping relationship

DOS header	
NT Header	PE Signature
	File Header
	Optional Header
Data Directory	
.idata	
.text	
.data	
...	

Figure 1. PE format

between the virtual address and the PE file at the time of loading, and the page is committed from disk to physical memory only when the instruction in a memory page is actually executed or when the data in a page is accessed; this mechanism makes the speed of file loading not much related to the file size) method of mapping these sections to memory, and attaching the read and write attributes of the specified sections in the section table. After the PE file is mapped into memory, the PE loader will continue to process the logical part of the PE file that resembles an import table.

3.2.3 Static analysis

There are two basic approaches to malware analysis: (1) static analysis and (2) dynamic analysis. On the one hand, static analysis involves checking the malware without running it. On the other hand, dynamic analysis involves running the malware [7]. Static analysis means examining the code or structure of an executable file without executing it. This method can protect the security of the operating system and also check out whether the file is malicious or not. Common static extraction methods are

- Analyzing PE file headers and sections. PE files contain metadata about the file itself, such as the actual section of the file.
- Collecting metadata about the executable's linked libraries and functions, as well as the files contained in the header files. This data provides information about the code base and functions that are common to many programs. The names of these Windows functions can give us an idea of what the executable is.
- Finding sequences of characters or strings. Search-

ing for a program's strings is the easiest way to get hints about its functions. Strings extracted from binary files can contain references to file paths modified or accessed by the executable, URLs accessed by the program, domain names, IP addresses, attack commands, names of loaded Windows dynamic link libraries (dlls), registry entries, etc. Utility Tools String 5 can be used to search for ASCII or Unicode strings while ignoring context and formatting in executable files.

- Search for packaged/encrypted code. Malware writers often use packaging and encryption to make their files more difficult to analyze. Packaged or encrypted software programs usually contain few strings and higher entropy than legitimate programs.
- Translation of machine code into assembly language

3.2.4 Dynamic analysis

Dynamic features are features extracted from the execution of the malware at runtime. Dynamic analysis involves monitoring the malware while it is running (and observing the actual sequence of instructions executed or the sequence of API functions triggered), or examining the system after the malware has been executed. The creation and modification of files and registry values. Generally speaking, dynamic features include the following aspects: operation behavior of files; operation behavior of registry keys; loading behavior of dynamic link libraries; operation behavior of process access; system service behavior; network access requests; and API calls. Compared with static analysis, dynamic analysis is more complex and time-consuming, and uses relatively higher-level features, and is therefore less interpretable.

Example of dynamic feature design, the article was published in 2016 [8] and the PE file was entered directly into the Cuckoo sandbox, which is an open source sandbox, and is more general for extracting dynamic features in academic papers; then a preprocessing operation was performed to convert the text into a vector representation, for example, 200 dynamic features were extracted, which can be represented using a 200-dimensional vector, and the position of each array represents the corresponding API, and then The resulting sequence is fed into the convolutional neural network LSTM for classification, and finally the structure of the family classification is obtained.

4 Approach

The main process of this project is shown in Figure 4, which consists of three steps [9]: input data, feature extraction, and model training. In this project, I am doing static

PEView - C:\Users\vboxuser\Downloads\MobaXterm_Portable_v22.2\MobaXterm_Personal_22.2.exe

File View Go Help

	pFile	Raw Data	Value
MobaXterm_Personal_22.2.exe	00000000	4D 5A 50 00 02 00 00 00 04 00 0F 00 FF FF 00 00	MZP.....
IMAGE_DOS_HEADER	00000010	B8 00 00 00 00 00 00 00 40 00 1A 00 00 00 00 00@.....
MS-DOS Stub Program	00000020	00 00 00 00 00 00 00 00 00 00 00 00 00 00 00 00
IMAGE_NT_HEADERS	00000030	00 00 00 00 00 00 00 00 00 00 00 00 01 00 00 0001.....
IMAGE_SECTION_HEADER CODE	00000040	BA 10 00 0E 1F B4 09 CD 21 B8 01 4C CD 21 90 90!..L..!
IMAGE_SECTION_HEADER DATA	00000050	54 68 69 73 20 70 72 6F 67 72 61 6D 20 6D 75 73	This program mus
IMAGE_SECTION_HEADER BSS	00000060	74 20 62 65 20 72 75 6E 20 75 6E 64 65 72 20 57	t be run under W
IMAGE_SECTION_HEADER .idata	00000070	69 6E 33 32 0D 0A 24 37 00 00 00 00 00 00 00 00	in32..\$7.....
IMAGE_SECTION_HEADER .tls	00000080	00 00 00 00 00 00 00 00 00 00 00 00 00 00 00 00
IMAGE_SECTION_HEADER .rdata	00000090	00 00 00 00 00 00 00 00 00 00 00 00 00 00 00 00
IMAGE_SECTION_HEADER .reloc	000000A0	00 00 00 00 00 00 00 00 00 00 00 00 00 00 00 00
IMAGE_SECTION_HEADER .rsrc	000000B0	00 00 00 00 00 00 00 00 00 00 00 00 00 00 00 00
SECTION CODE	000000C0	00 00 00 00 00 00 00 00 00 00 00 00 00 00 00 00
SECTION DATA	000000D0	00 00 00 00 00 00 00 00 00 00 00 00 00 00 00 00
SECTION BSS	000000E0	00 00 00 00 00 00 00 00 00 00 00 00 00 00 00 00
SECTION .idata	000000F0	00 00 00 00 00 00 00 00 00 00 00 00 00 00 00 00
SECTION .tls	00000100	50 45 00 00 4C 01 08 00 80 32 53 63 00 00 00 00	PE..L.....2Sc
SECTION .rdata	00000110	00 00 00 00 E0 00 AF 81 0B 01 02 19 00 78 6F 00xo.
SECTION .reloc	00000120	00 BE 97 00 00 00 00 00 A8 58 6F 00 00 10 00 00Xo.....
SECTION .rsrc	00000130	00 90 6F 00 00 00 40 00 00 10 00 00 00 02 00 00@.....
CERTIFICATE Table	00000140	04 00 00 00 00 00 00 00 04 00 00 00 00 00 00 005.....
	00000150	00 20 0A 01 00 04 00 00 35 1E 02 01 02 00 00 00@.....
	00000160	00 00 10 00 00 40 00 00 00 00 00 00 10 00 00 00
	00000170	00 00 00 00 10 00 00 00 00 00 00 00 00 00 00 00
	00000180	00 50 7B 00 1A 4B 00 00 00 F0 81 00 00 28 88 00	.P{...K.....{
	00000190	00 00 00 00 00 00 00 00 00 0E 01 01 90 22 00 00
	000001A0	00 00 00 00 00 00 00 00 00 00 00 00 00 00 00 00
	000001B0	00 00 00 00 00 00 00 00 00 00 00 00 00 00 00 00
	000001C0	00 B0 7B 00 18 00 00 00 00 00 00 00 00 00 00 00	..{.....

Figure 2. PE file

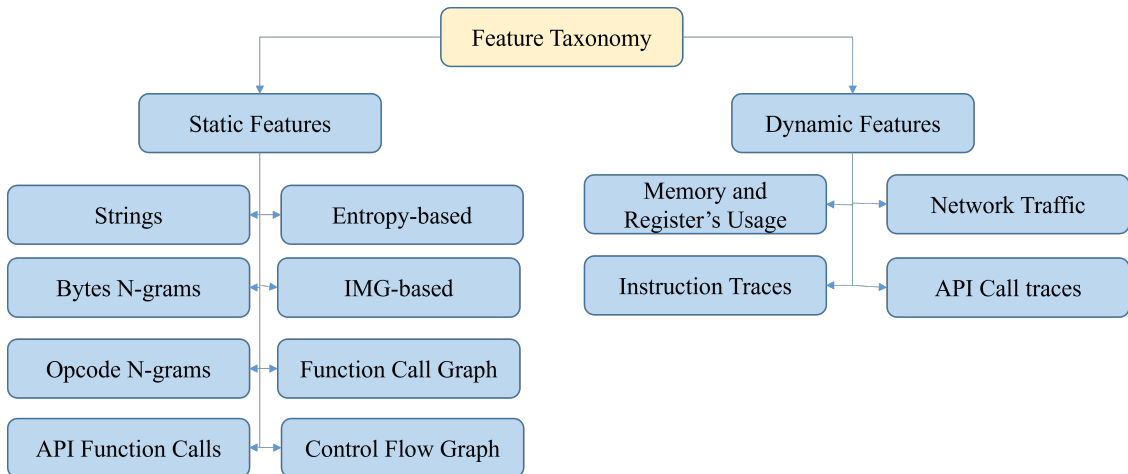


Figure 3. Feature Taxonomy

malware analysis. After downloading the datasets, I am going to extract the features from PE header. The features will be stored csv file and used for training the models.

4.1 Extract features

I extracted the static features from PE files using the pe-file library. The following I will introduce some features that have been extract.

Entropy: entropy measures the randomness of the data in a file and is used to determine whether a file contains hidden data or suspicious scripts. Many malware programs use entropy as a packing technique [10]. Packaging leads to an increase in entropy. So entropy can be used as an indicator of malware programs. Also, some benign programs use packing. In general, a higher percentage of malware programs use packing.

Checksum: The algorithm for computing the checksum is incorporated into IMAGHELP.DLL. The following are checked for validation at load time: all drivers, any DLL loaded at boot time, and any DLL that is loaded into a critical Windows process. According to the studies, it has been observed that the value of this field is zero for many malware programs, while benign programs have significantly larger values for this field.

FileInfo: VS_FIXEDFILEINFO Contains version information for a file. FileVersionLS and ProductVersionLS will contains the version information about the binary file. Each file has a different number of these fields. Since benign programs are valid programs produced by reliable individuals or companies, they usually contain these fields, while many malware programs lack this set or have only a small number of these fields.

4.2 Preprocessing data

There are many redundant features are included after extracting features. It is necessary to select features that are useful for distinguishing software types.

In this step, I choose three ways to preprocessing the data: only preprocessing data simple, using feature selection to filter out some features, and using AutoEncoder to process the data. Dataset One, only simple data processing, only need to remove the information in the file that is not related to the training model, such as file name, file md5 value, etc. Dataset two, the purpose of feature selection is to select the most relevant set of features among these features, while leaving fewer features to improve the classification performance. In this step, I used the Extra Trees Classifier to select the features.

Extremely Randomized Trees Classifier(Extra Trees Classifier) is a type of ensemble learning technique which aggregates the results of multiple de-correlated decision trees

collected in a “forest” to output it’s classification result. In concept, it is very similar to a Random Forest Classifier and only differs from it in the manner of construction of the decision trees in the forest. Each Decision Tree in the Extra Trees Forest is constructed from the original training sample. Then, at each test node, Each tree is provided with a random sample of k features from the feature-set from which each decision tree must select the best feature to split the data based on some mathematical criteria (typically the Gini Index). This random sample of features leads to the creation of multiple de-correlated decision trees.

To perform feature selection using the above forest structure, during the construction of the forest, for each feature, the normalized total reduction in the mathematical criteria used in the decision of feature of split (Gini Index if the Gini Index is used in the construction of the forest) is computed. This value is called the Gini Importance of the feature. To perform feature selection, each feature is ordered in descending order according to the Gini Importance of each feature and the user selects the top k features according to his/her choice. Dataset three, use AutoEncoder to process the data and map from low-dimensional data to high-dimensional data. An autoencoder is a type of artificial neural network used to learn efficient coding of unlabeled data (unsupervised learning) [3]. The encoding is validated and refined by attempting to regenerate the input from the encoding. The autoencoder learns a representation (encoding) for a set of data, typically for dimensionality reduction, by training the network to ignore insignificant data (“noise”). But this project, I used it to ascend dimension. This way may contain hidden information that does not appear in the original features.

4.3 Training models

I choose 5 models in this paper, including GaussianNB, MLP, LinearRegression, DecisionTree, and GradientBoosting. The dataset is generally divided into two parts for the whole dataset, one for training and one for validation. In this paper I choose 10-fold cross-validation to handle the training and test sets. 10-fold cross-validation randomly divides the original data into 10 parts. Then chooses one of them as the testing set and the other nine as the training set each time. Then training the model 10 times. In this way, each subset has a chance to be the training set or the test set. In addition, the five models will be trained on three different datasets, and finally the results on the different datasets will be compared.

Gaussian Naive Bayes (GNB) is a classification technique for machine learning based on probabilistic methods and Gaussian distributions. Plain Bayes assumes that each parameter (also called a feature or predictor variable) has the independent ability to predict the output variable. The com-

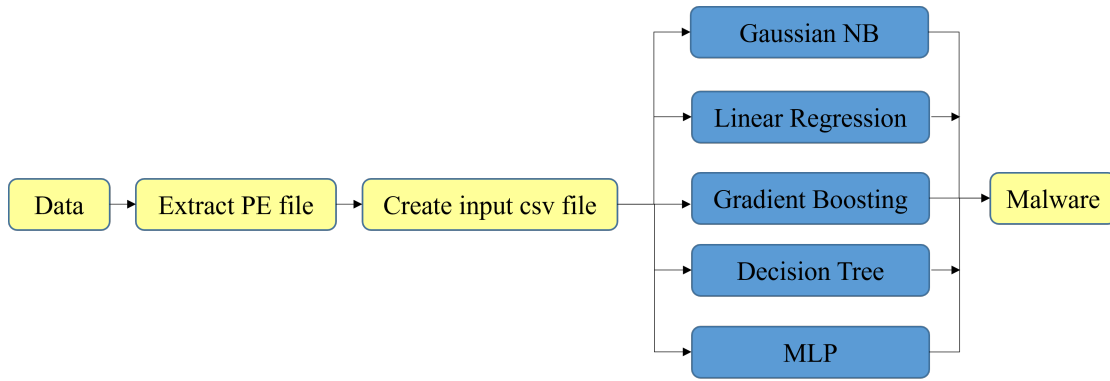


Figure 4. Framework

bination of predictions for all parameters is the final prediction, which returns the probability that the dependent variable will be classified into each group, with the final classification being assigned to the group (class) with the higher probability.

Linear Regression is a linear model for prediction by linear combination of attributes, which aims to find a straight line or a plane or a higher dimensional hyperplane that minimizes the error between the predicted and true values.

Decision tree is a decision analysis method that evaluates the project risk and judges its feasibility by forming a decision tree to find the probability that the expected value of the net present value is greater than or equal to zero on the basis of the known probability of occurrence of various situations, and is a graphical method that intuitively uses probability analysis. Since this decision branch is drawn as a graph resembling the branches of a tree, it is called a decision tree. Decision trees are a basic classification and regression method, and learning usually involves three steps: feature selection, decision tree generation, and decision tree pruning. In machine learning, a decision tree is a predictive model that represents a mapping relationship between object attributes and object values.

The basic idea of Gradient Boosting is to serially generate multiple weak learners, each with the goal of fitting the negative gradient of the loss function of the previously accumulated model so that the cumulative model loss after adding the weak learner is reduced in the direction of the negative gradient. Gradient Boosting can also be understood as gradient descent over the function space. The more familiar gradient descent is usually a gradient descent of values over the parameter space (e.g. training a neural network, where the gradient of the current loss with respect to the parameters is computed in each iteration and the parameters are updated). In Gradient Boosting, on the other hand, a weak learner is generated in each iteration, which fits the gradient of the loss function with respect to the previous cumulative model, and then this weak learner is added to the cumu-

lative model to gradually reduce the loss of the cumulative model. That is, the gradient descent in parameter space uses the gradient information to adjust the parameters to reduce the loss, and the gradient descent in function space uses the gradient to fit a new function to reduce the loss.

MLP consists of three layers: input layer, hidden layer, and output layer, the first layer is called the input layer, the last layer is called the output layer, and the middle layer is called the hidden layer. The number of hidden layers can be chosen according to the respective needs. MLP has good recognition rate and fast classification, but in exchange its training process is time consuming.

5 Experiment

5.1 Project operating environment

- *Language:* Python(3.8)
- *Platform:* Jupyter Notebook
- *Libraries:* Scikit Learn, Pandas, Pickle, Numpy, torch, dfencoder, etc

5.2 Data collection

The malware dataset is obtained from a publicly available website with the address <https://virusshare.com/> and the benign software come from <https://github.com/iosifache/DikeDataset>. The dataset consists of 207781 data entries listing PE information from various sources. 113329 are malicious while 94452 are legitimate. The two figures shown below are the malware downloaded from the virusshare website, the benign software downloaded from github.

Out[37]:

	Name	md5	Machine	SizeOfOptionalHeader	Characteristics	MajorLinkerVer
0	572733582ed2076416e035a1d670d12d984c8d56a1b94d...	e0c8ac0427fe71d9c5913fadb5d19d1b	332	224	258	
1	2999b0709f50ed5026039496142254af172990a50e2a6e...	19611ed77cbac092503767eef94fbc50	332	224	783	
2	ae484f4c0341f050d40d3e148be8d11ce86358349ddf31...	4794245bd7eb172176aa65074895fc0f	332	224	290	
3	dec0efefe7749ab45c541f8c2c0bc4fa946e865025e209...	a84e0f438a5637255096c185fb50bb54	34404	240	34	
4	6b8d1a1f785c103e2d75702c22491d1b3942878bcb9796...	26a13fba483291b5e14329fc1f04837	34404	240	34	
...
207776	VirusShare_b2c592bb21e0a566e512ff697bd73cef	b2c592bb21e0a566e512ff697bd73cef	332	224	271	
207777	VirusShare_d6ec89f00667d2439391c9a37b8ce49e	d6ec89f00667d2439391c9a37b8ce49e	332	224	271	
207778	VirusShare_9dd3ee41cb1b7e8902e38eeb645a772f	9dd3ee41cb1b7e8902e38eeb645a772f	332	224	33166	
207779	VirusShare_ad480c003f5e93a0306ec3387fdb448d	ad480c003f5e93a0306ec3387fdb448d	332	224	271	
207780	VirusShare_2eb062a323b59e0734ae1f40c46a9608	2eb062a323b59e0734ae1f40c46a9608	332	224	271	

207781 rows x 7 columns

Figure 8. Features after extraction

1. feature Subsystem (0.150836)
2. feature SizeOfOptionalHeader (0.140299)
3. feature DllCharacteristics (0.088468)
4. feature Machine (0.085965)
5. feature MajorSubsystemVersion (0.072326)
6. feature SectionsMaxEntropy (0.045958)
7. feature SectionAlignment (0.044875)
8. feature SizeOfHeaders (0.034736)
9. feature LoadConfigurationSize (0.029276)
10. feature MinorSubsystemVersion (0.023359)
11. feature SizeOfStackReserve (0.022674)
12. feature ImageBase (0.020640)
13. feature SectionsNb (0.019799)

Figure 9. Features in dataset two

5.5 Training models & result

There are five models used in my experiments:

- GaussianNB
- MLPClassifier
- LinearRegression
- DecisionTree
- GradientBoosting

5.5.1 Data set one

The highest accuracy model is DecisionTree, its accuracy is 99.79%.

Model	Accuracy
GaussianNB	52.24%
MLPClassifier	97.03%
LinearRegression	78.39%
DecisionTree	99.79%
LinearRegression	98.24%

5.5.2 Data set two

The highest accuracy model is DecisionTree, its accuracy is 99.52%.

Model	Accuracy
GaussianNB	45.51%
MLPClassifier	80.17%
LinearRegression	74.21%
DecisionTree	99.52%
LinearRegression	98.14%

5.5.3 Data set three

The highest accuracy model is MLP, its accuracy is 99.00%.

Model	Accuracy
GaussianNB	71.49%
MLPClassifier	99.00%
LinearRegression	68.22%
DecisionTree	98.39%
LinearRegression	93.28%

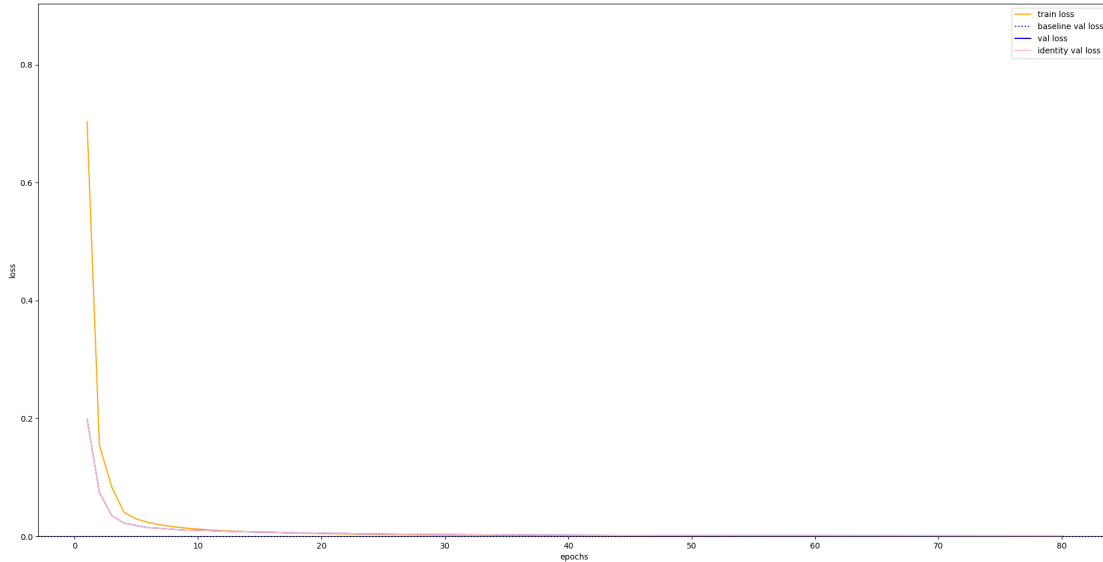


Figure 10. Training loss of dfencoder

6 Conclusion and Future work

My framework consists of two main parts, feature extraction and malware prediction by machine learning technique. Three different ways of data pre-processing were used before training models. According to the final results, performances of all the models trained by dataset one is better than models trained by dataset two. I think this is due to the fact that some relevant information was dropped out during feature selection, therefore impaired model performances. Secondly, by comparing dataset 1 and dataset 3, it is obvious that the accuracies of GaussianNB and MLPClassifier models improved significantly. However, the accuracies of the other three models decreased. It may be due to that GaussianNB and MLPClassifier are more suitable for high-dimensional data, while the other three are more suitable for low-dimensional data. The reason behind this case will be a future work.

Besides the work I have done so far, other static features such as file bytecode features, strings, API call sequences can be used in the feature extraction part. Since there are some limitations of static features extraction, dynamic features extraction is a possible way to improve the performance. Integration of static and dynamic features extraction is promising and will be another future work.

References

- [1] M. Z. Shafiq, S. M. Tabish, F. Mirza, and M. Farooq, "Pe-miner: Mining structural information to detect malicious executables in realtime," in *International workshop on recent advances in intrusion detection*. Springer, 2009, pp. 121–141.
- [2] P. Trinius, C. Willems, T. Holz, and K. Rieck, "A malware instruction set for behavior-based analysis," *None*, 2009.
- [3] M. G. Schultz, E. Eskin, F. Zadok, and S. J. Stolfo, "Data mining methods for detection of new malicious executables," in *Proceedings 2001 IEEE Symposium on Security and Privacy. S&P 2001*. IEEE, 2000, pp. 38–49.
- [4] J. Z. Kolter and M. A. Maloof, "Learning to detect malicious executables in the wild," in *Proceedings of the tenth ACM SIGKDD international conference on Knowledge discovery and data mining*, 2004, pp. 470–478.
- [5] D. Ucci, L. Aniello, and R. Baldoni, "Survey of machine learning techniques for malware analysis," *Computers & Security*, vol. 81, pp. 123–147, 2019.
- [6] M. Belaoued and S. Mazouzi, "Statistical study of imported apis by pe type malware," in *2014 International Conference on Advanced Networking Distributed Systems and Applications*. IEEE, 2014, pp. 82–86.
- [7] D. Gibert, C. Mateu, and J. Planes, "The rise of machine learning for detection and classification of malware: Research developments, trends and challenges," *Journal of Network and Computer Applications*, vol. 153, p. 102526, 2020.

```

dataset_encoder_train = dataset.drop(['Name', 'md5', 'Legitimate'], axis=1).copy()
train = dataset_encoder_train.sample(frac=.9, random_state=36)
test = dataset_encoder_train.loc[~dataset_encoder_train.index.isin(train.index)]

X_train = train.copy()
X_val = test.copy()

dfencoder_model = AutoEncoder(
    encoder_layers = [32, 64, 128], #model architecture
    decoder_layers = [], #decoder optional - you can create bottlenecks if you like
    activation='relu',
    swap_p=0.1, #noise parameter
    lr = 0.01,
    lr_decay=.99,
    batch_size=128,
    logger='ipynb', #special logging for jupyter notebooks
    verbose=False,
    optimizer='sgd',
    scaler='gauss_rank', #gauss rank scaling forces your numeric features into standard normal distributions
    min_cats=1, #Define cutoff for minority categories, default 10
    progress_bar=False
)
dfencoder_model.device = 'cuda'
plt.figure(figsize=(24,10))
plt.rcParams['figure.figsize'] = [24, 12]
dfencoder_model.fit(X_train, epochs=80, val=X_val)

```

Figure 11. Example code of training a dfencoder model

- [8] B. Kolosnjaji, A. Zarras, G. Webster, and C. Eckert, "Deep learning for classification of malware system call sequences," in *Australasian joint conference on artificial intelligence*. Springer, 2016, pp. 137–149.
- [9] J. Saxe and K. Berlin, "Deep neural network based malware detection using two dimensional binary program features," in *2015 10th international conference on malicious and unwanted software (MALWARE)*. IEEE, 2015, pp. 11–20.
- [10] T. Rezaei and A. Hamze, "An efficient approach for malware detection using pe header specifications," in *2020 6th International Conference on Web Research (ICWR)*. IEEE, 2020, pp. 234–239.