

Capstone Project

Dileep Tiku

Machine Learning Engineer Nanodegree June 11th, 2016

Definition

Project Overview

Computer Malware poses a security risk for businesses throughout the world. With the advent of Industrial internet of things besides compromising information technology assets, malware can impact operational technology which has potential of causing harm to society at large. Using advanced machine learning algorithms, computer malware can be identified using behavioral analysis of the software. This will allow organizations to use the malware detection model, to assist in advanced threat protection not possible only by traditional methods of malware detection, which are signature based. In this project, I created a model that analyzes the Windows API calls generated by executing a software in a controlled environment and then classifies the software as malware or benign. The model is trained by using the malicious training sets from CSMining¹ website.

Problem Statement

Develop a behavior analysis model that can identify whether a software is malware or not.

- Identify the algorithm to process the Windows API Calls list in the CSMining data.
- Apply the algorithm to extract the features from API Calls for classification
- Train a classifier to use the features extracted from the API Calls.
- Test the classifier trained to use the extracted features from API Calls.

Metrics

Accuracy is a common metric for binary classifiers; it takes into account both true positives and true negatives with equal weight.

accuracy = correct predictions/all predictions

¹ <http://csmining.org/index.php/malicious-software-datasets-.html>

This metric was used when evaluating the classifier because false negatives compromise the security of an entity and false positives are responsible for wasting a lot of time for the scarce security personnel in an organization.

Analysis

Data Exploration

The CSMining dataset used for this project is a collection 388 software traces. The software traces include API Call lists generated by the execution of malware and benign software. Out of a total of 388 traces, 320 are from malware execution and 68 are from benign software execution. The 320 traces generated from different types of malware including worms, trojans and viruses. The API calls are repeated in the dataset for completeness of the calls.

A subset of the Windows APU/System-Calls which are considered informative for differentiating a malware from a benign software are logged by API monitors when a designated program is running in the system¹

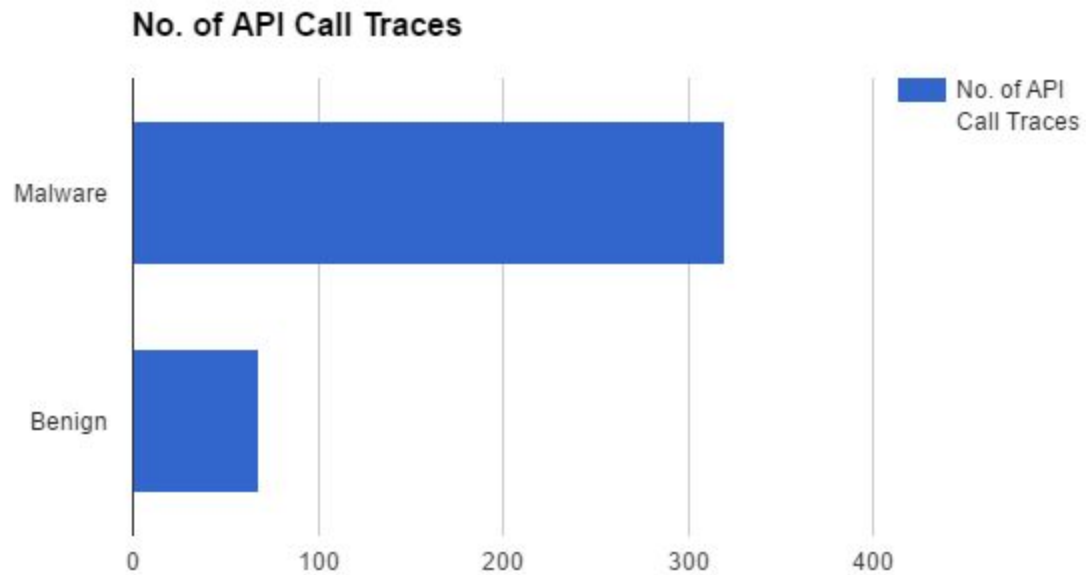
A partial list of API Calls from one malware run taken from one record in the CSMining dataset

LoadLibraryW HeapAlloc HeapAlloc HeapFree HeapAlloc HeapFree HeapFree NtOpenKey
LoadLibraryW GetProcAddress GetProcAddress GetProcAddress GetProcAddress GetProcAddress
GetProcAddress GetProcAddress GetProcAddress GetProcAddress GetCurrentThreadld NtOpenKey
NtQueryValueKey HeapAlloc HeapAlloc HeapFree HeapFree HeapAlloc HeapAlloc HeapFree HeapFree
NtOpenKey GetProcessHeap HeapAlloc NtOpenKey HeapAlloc NtOpenKey NtQueryValueKey
NtQueryValueKey HeapFree HeapAlloc HeapAlloc NtOpenKey NtQueryValueKey HeapAlloc HeapAlloc
RegOpenKeyExW LocalAlloc HeapAlloc HeapAlloc NtOpenKey HeapAlloc GetStartupInfoA GetFileType
GetFileType GetFileType GetEnvironmentStringsW HeapAlloc FreeEnvironmentStringsW
GetCommandLineA GetCommandLineW GetCommandLineW HeapAlloc HeapFree HeapAlloc
VirtualQuery GetModuleFileNameW VirtualQuery NtOpenKey HeapAlloc GetProcessHeap HeapAlloc
HeapAlloc HeapAlloc HeapAlloc HeapAlloc HeapAlloc HeapAlloc HeapAlloc HeapAlloc HeapAlloc
HeapAlloc GetProcessVersion RegOpenKeyExA RegQueryValueExA RegCloseKey RegOpenKeyExA
RegQueryValueExA RegCloseKey GetProcessHeap HeapAlloc HeapAlloc HeapAlloc HeapAlloc
GetProcessVersion HeapAlloc HeapAlloc HeapAlloc HeapAlloc GetProcessVersion HeapAlloc HeapAlloc
HeapAlloc HeapAlloc HeapAlloc HeapAlloc NtOpenKey NtQueryValueKey NtQueryValueKey
NtQueryValueKey NtOpenKey NtQueryValueKey NtQueryValueKey GetModuleFileNameW NtOpenKey
HeapAlloc HeapAlloc HeapAlloc HeapAlloc HeapFree HeapAlloc HeapAlloc HeapAlloc RegOpenKeyExA
HeapAlloc HeapAlloc RegOpenKeyExA HeapAlloc HeapAlloc RegOpenKeyExA NtOpenKey
RegOpenKeyExW RegNotifyChangeKeyValue HeapAlloc HeapAlloc HeapAlloc HeapAlloc HeapAlloc
NtOpenKey HeapAlloc RegOpenKeyExA GetProcAddress GetProcAddress GetProcAddress
GetProcAddress GetProcAddress LocalAlloc LocalAlloc LocalAlloc LocalAlloc HeapAlloc HeapAlloc
HeapAlloc LocalAlloc HeapFree HeapFree HeapFree LocalFree NtOpenKey RegOpenKeyExW

[illegible]

Exploratory Visualization

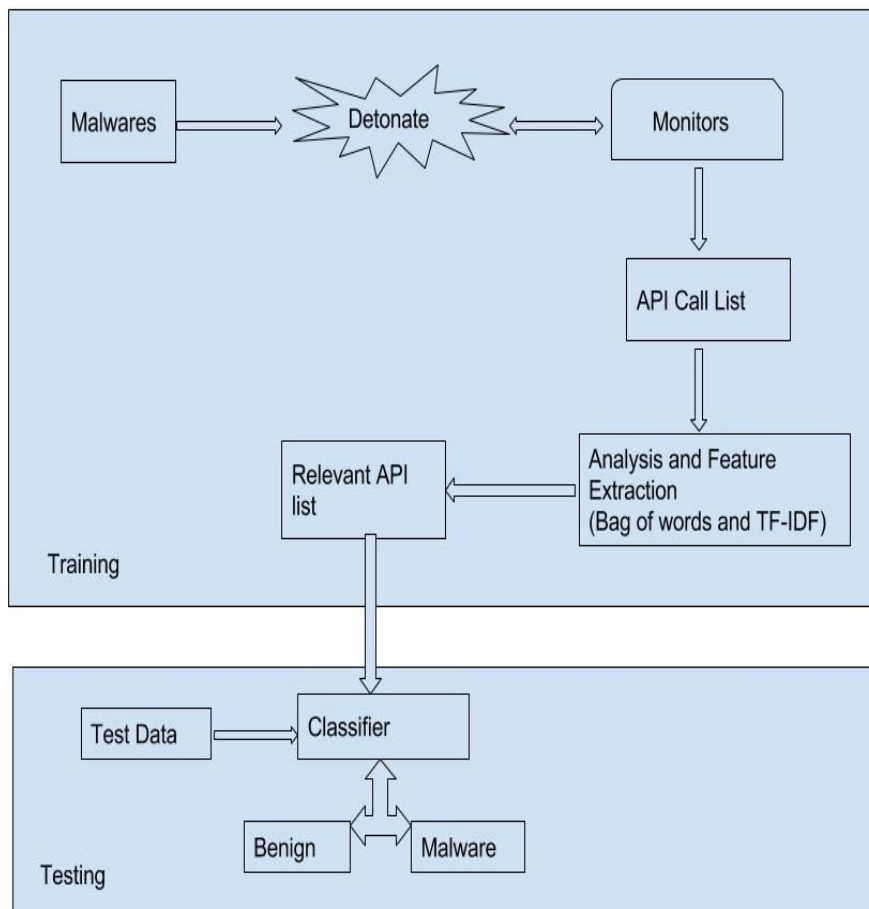
The API Call list for some records is exceptionally long. The number of API Calls for a single software run is approximately equal to 1800.



Algorithms and Techniques

As seen in the figure 1 below, the API Call list had to be analyzed such that relevant features could be extracted. In order to extract the relevant features

Figure 1.



The API Calls are converted to numerical vector format and then analyzed by counting the number of distinct API Calls used by the malware as well as benign software. The count of individual API Calls as a ratio of the total API Calls made in a software run are used for extracting the features from the API Call list.

Tf-idf transformer was used to analyze the words for their importance. TF-IDF (term frequency - inverse document frequency) is a numerical statistic that is intended to reflect how important a word is to a document in a collection or [corpus](#)¹. Inverse document

frequency is a concept which assists in diminishing the importance of words that occur across malware and benign software. TF-IDF transforms the data set with vector representation of the API Calls to a data set of term frequency and inverse document frequency.

Stochastic Gradient Descent method was used to execute linear SVM. Numerous studies (Joachmis 1998, Dumais et al. 1998, Drucker et al. 1999) have shown the superiority of SVM over other machine learning methods for text categorization problems.

Stochastic Gradient Descent is an algorithm which helps in optimizing the parameters for another algorithm in this case a classifier. In case of large data sets the process of choosing parameters for a model can be very slow and sometimes hard to control if the main memory of the computer is insufficient for the whole training data . In such cases gradient descent process helps to improve the speed of computing and control the process. Stochastic Gradient descent computes the gradient of its parameters with only a few training examples unlike the standard gradient descent algorithms. For Stochastic Gradient Descent to work well the training data must be randomly shuffled before each epoch of training. By randomly shuffling the data, poor convergence can be avoided. Since the high number of API Calls can give rise to a large number of features, we need to use Stochastic Gradient Descent algorithm for efficiently training the Support Vector Machine .

Support Vector Machines is a supervised machine learning algorithm, which can be used for both classification and regression problems. For classification problems, the idea behind support vector machine is to find a hyperplane which separates the data into two classes. The data points nearest to the hyperplane are called support vectors. In our case the Support Vector Machines makes perfect sense because the data is linearly separable. Support Vector Machine is widely believed to be the best text classifier.

Benchmark

To create an initial benchmark for the classifier, I used MultiNomial Naive Bayes algorithm. Naive Bayes algorithm was able to achieve approximately 84% accuracy. My goal was to use an algorithm that will get a higher accuracy than Mutinomial Naive Bayes.

Methodology

Data Preprocessing

Since API Call list is textual data, the text analysis algorithms were used to analyze them. Following pre-processing steps were done before the data was sent through the classifier for the final classification.

1) Vectorizer

Scikit learn's feature extraction module offers a CountVectorizer function, which converts the text to numerical vectors. This function combines bag of words generation, tokenization, and filtering of stop words. The countvectorizer also uses sparse matrices to save space in the memory. The bag of words representation indicates that n-features are the distinct words in the corpus. In this case the bag of words is the list of distinct API Calls from all the records. CountVectorizer supports the count of individual API Calls. It produces a vocabulary list, where the index value of the API Call is linked to its frequency.

2) TF-IDF Transformer

Even though occurrence count of an API call would give the indication of type of software, some of the API Calls appear multiple times in a single software execution. To get a better sense of the API Calls occurring during the execution of the software, term frequency is used to get the ratio of the count of a specific API Call and the total number of API Calls from the run. Since some of the API Calls are seen both in malware and benign software runs, inverse document frequency is used to lower the importance of API calls seen in many records. The TF-IDF transformer transforms the data by calculating term frequency and inverse document frequency.

Implementation

I evaluated several approaches before finalizing on the usage of countvectorizer and tfidf transformer for feature extraction. Before finding these built in functions in scikit learn, I was considering writing the algorithms from scratch. Writing the algorithms from scratch would have taken a very long time.

For the classifier, I considered using different sizes of n-grams for the textual analysis of API Calls, but the use of single word provided the best results. This finding was in line with research done in this field by other people like Veeramani R and Nitin Rai¹.

I chose to use SGDClassifier because it is highly efficient and easy to use. Support Vector Machines are a good fit for textual analysis problems because the data generated is sparse and most text classification problems are linearly separable².

Refinement

I ran the grid search to get the best values for following relevant hyper parameters for textual analysis algorithms used in this model.

- 1) Vector_ngram_range (1,1) or (1,2) for countvectorizer
- 2) Use Inverse document frequency or not in tfidf transformer
- 3) Alpha equal to 1e-2 and 1e-3 for SVM

After running the search I found that following were the parameters for the best score.

Vector_ngram_range (1,2)

Tfidf_use_idf = true

Clf_alpha = .001

Results

Model Evaluation

The final model used for classification is shown in Fig 2. This model uses a SVM classifier with SGD training. The loss = 'hinge' as a parameter for SGDClassifier uses a 'hinge' loss function, which give a linear SVM. Since API Calls are textual the data is linearly separable. The penalty also known as regularization term is set to 'l2', which is a standard for SVM.

Fig 2

```
#create a pipeline, which uses a countvectorizer, a tfidftransformer and a SGDclassifier to classify the software  
#SGDClassifier implements linear models with stochastic gradient descent learning. The loss function of Hinge  
#gives a Linear SVM  
#penalty = l2 is a standard regularizer for linear SVM  
#alpha is  
#n_iter is the number of passes over the training set  
  
text_clf = Pipeline([('vect', CountVectorizer()),  
                    ('tfidf', TfidfTransformer()),  
                    ('clf', SGDClassifier(loss='hinge', penalty='l2',  
                                         alpha=1e-3, n_iter=5, random_state=42)),])
```

To evaluate the robustness of model, I used shuffle split to make sure that the training and testing sets are selected randomly. I used ten iterations with the shuffle split to ensure that the training data set in one iteration becomes test data set in another iteration and vice versa. The results from the shuffle split were similar to the results as shown below.

tp: true positive fp: false positive	fn: false negative tn: true negative	Predicted Class	
		Malware	Benign
Actual Class	Malware	tp	fn
	Benign	fp	tn

$$\text{Accuracy} = \frac{tp + tn}{tp + tn + fp + fn}$$

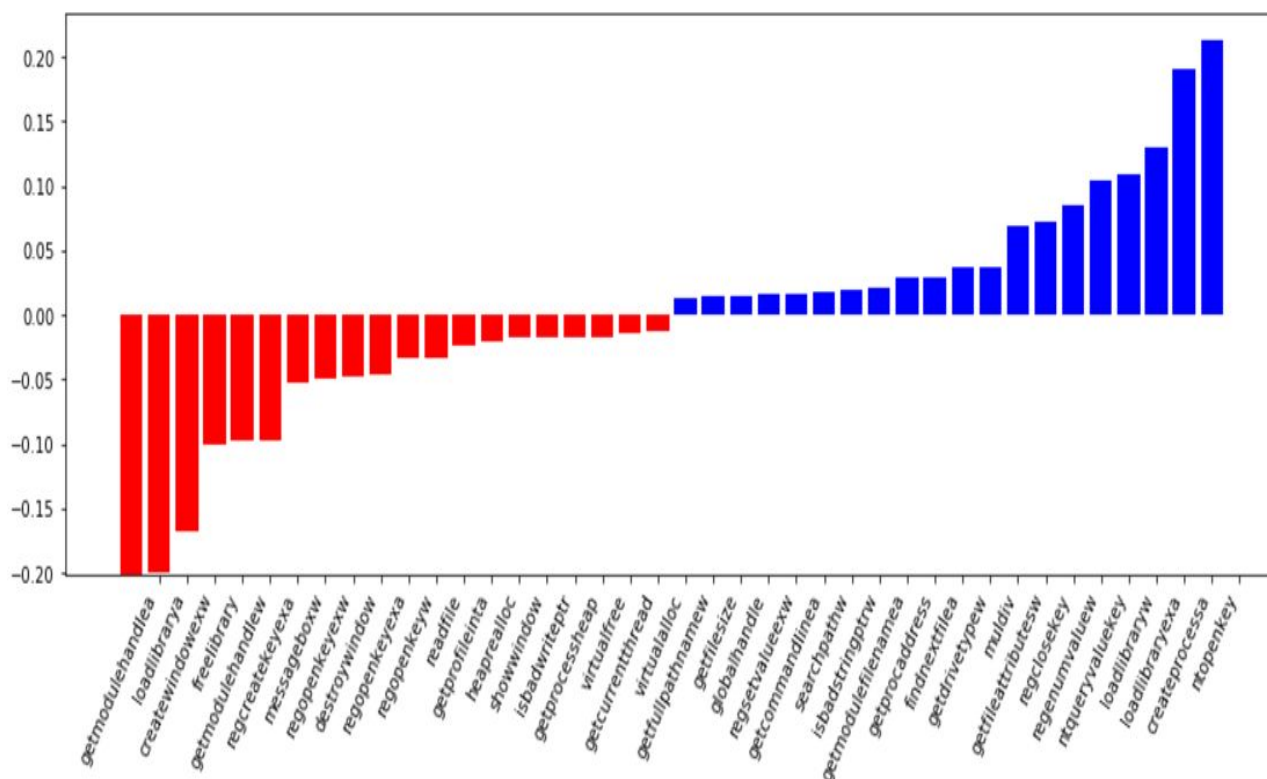
Accuracy with Mutinomial Naive Bayes = 84.6%

Accuracy with SVM classifier = 89.7%

Metrics from SGD Classifier

	precision	recall	f1-score	support
Benign	0.83	0.42	0.56	12
Malware	0.90	0.98	0.94	66
avg / total	0.89	0.90	0.88	78

Conclusion



The plot shows what top twenty features (API Calls) are used by the model to

make the positive and negative decisions from the malware dataset. These features are shown simply based on transformation and calculations executed by count vectorizer. The visualization does not factor in the term frequency and the inverse document frequency transformations performed by tf-idf transformer. The top twenty features will be different when plotted after tf-idf transformation.

I found this project very interesting because when I thought of performing behavior analysis for malware, I was not thinking about text analysis. With this project, I have understood the power of text analysis and the relevant scikit learn modules that offer the textual analysis functionality. Following are the steps, I followed to develop the model.

- 1) Read various white papers to understand how malware is identified using traditional signature based approaches and behavior analysis approach.
- 2) Discussed the use of neural networks for behavior analysis and the complexities involved with it.
- 3) Narrowed down on using textual analysis for Windows API Calls data for malware identification published at CSMining web site.
- 4) Converted the list of API Calls collected from software execution in text form to vectors with countvectorizer.
- 5) Transformed the vectors in step 4 to term frequency and inverse document frequency using tfidftranformer.
- 6) Split the data into test and training data sets.
- 7) Used SVM classifier based on its success with classification from textual data.

In order to augment textual analysis, the malware behavior analysis should also factor in the outbound calls made to URLs, which are black listed. This is based on the fact that a common motive of malware is to steal critical or valuable information. The malware interacts with a command and control host center which is typically managed by malicious people. The analysis has to be further enhanced to classify different malware types like viruses, worms, trojans, bots etc.