

# Homework I

Francesco Terenzi, 1702217

Fall 2018

## Contents

<b>1</b>	<b>Introduction</b>	<b>2</b>
<b>2</b>	<b>Dataset</b>	<b>3</b>
2.1	Prefixes . . . . .	3
2.2	Encoded Features . . . . .	4
2.3	Load Drebin . . . . .	4
<b>3</b>	<b>Classifier</b>	<b>6</b>
3.1	fit and predict . . . . .	6
3.2	scores . . . . .	6
3.3	print_results . . . . .	7
<b>4</b>	<b>Main</b>	<b>8</b>
<b>5</b>	<b>Results</b>	<b>10</b>

# 1 Introduction

In this report it's described the solution implemented by the student (ID: 1702217) to solve the homework request. The programming language used for this project is Python, for its pseudocode-*like* syntax and for the number of useful libraries for Machine Learning: scikit-learn, numpy, pickle, tqdm and other are the main. Their functionality require an installation (via `pip` command or `conda` command if you use the Anaconda framework).

Each section describes a different file present in the working directory: in `dataset.py` are given tools for modelling the *drebin* dataset, in the `malware_classifier.py` is described the classifier and with `main.py` is illustrated a simply test. All the code described in this paper is available on Github platform: <https://github.com/fratere/Malware-Analysis>.

## 2 Dataset

The drebin dataset contains about 130.000 applications. 5560 of them are malicious samples (as reported in a useful csv file) and so, testing the whole dataset requires a huge memory space and a long execution time. Then, the drebin dataset is modified and "filtered" in accordance to the type of test you want to accomplish. In the `main.py` file, reported in Section 4, this idea will be clearer.

In this chapter are described the most important functions in `dataset.py` file.

### 2.1 Prefixes

```
29 def get_prefixes():
30     all_prefixes = ['real_permission', 'feature',
                     'api_call', 'call', 'permission', 'provider',
                     'activity', 'intent', 'service_receiver',
                     'activity', 'url']
31     return all_prefixes[:6]
```

Every feature in a file is composed by a `[prefix]::[rest of the feature]` syntax. The `get_prefixes()` function returns all the prefixes categories, ordered following the non-decreasing number of frequency in the drebin applications, as shown in the figure below. How you can see, the last 5 categories are very frequently, so, probably they are insignificant in a malware analysis because they can represent features common both malicious and not malicious files: they are omitted from the output list. There is, obviously, a gain, computationally speaking.

Listing 1: Number of features ordered by prefixes

```
13 sizes: {'real_permission': 70,
14         'feature': 72,
15         'api_call': 315,
16         'call': 733,
17         'permission': 3812,
18         'provider': 4513,
19         'intent': 6379,
20         'service_receiver': 33222
21         'activity': 185729,
```

```

22         'url': 310488,
23     }

```

## 2.2 Encoded Features

```

72 def get_encoded_features(prefixes):
73     files_features = pref_clean(get_files_features(),
74                                 prefixes)
75     features = get_features(files_features)
76     print("number of features: %d " % len(features))
77     res = {}
78     for file_name in tqdm(files_features, desc='creating
79                             encoded features file'):
80         file_features = files_features[file_name]
81         l = [1 if f in file_features else 0 for f in
              features ]
82         res[file_name] = l
83     return res

```

Every file contains a set of features, that are strings. Our classifier, however, requires a numerical representation so, a transformation is needed. In the line 73 all files with the respective features are taken and filtered from those features which prefix is not in the **prefixes** input list. With the cycle of line 77 a one-hot encoding operation is done: for every file is created an array with '0' values in those positions that represent features not present in the file, '1' values otherwise. Obviously, the length of the array depends on the number of the filtered features and represents the function cost input in terms of space complexity.

## 2.3 Load Drebin

```

83 def load_drebin(prefixes):
84     if os.path.exists('drebin.pickle'):
85         with open('drebin.pickle', 'rb') as all_data:
86             data = pickle.load(all_data)
87             return data[0].todense(), data[1]
88     else:
89         files = get_encoded_features(prefixes)

```

```

90     X = [v for v in files.values()]
91     hashes = malware_hashes()
92     y = []
93     for f in files:
94         if f in hashes:
95             y.append(1)
96         else:
97             y.append(0)
98     print('creating sparse matrix...')
99     X = csr_matrix(X)
100    print('done')
101    """
102    with open('drebin.pickle', 'wb') as all_data:
103        print('saving data in drebin.pickle...')
104        pickle.dump([X, y], all_data)
105        print('done')
106    """
107    return X, y

```

The `load_drebin` function returns the encoded drebin dataset, represented by 'X', a matrix and 'y', a vector. The 'X' represents the encoded files given by the `get_encoded_features()`, explained in the Section 2.2, and the 'y' labels a malicious file with 1 and a non-malicious file with 0. An important observation: the X matrix is a sparse matrix, because it's requires less space in memory than an *array-like* matrix. Anyway, the classifier accepts both types of matrix.

In the lines 102-105 thanks to the `pickle` library methods it's possible to save the dataset on disk (in a pickle file) to avoid repeating the encoding operation when the `load_drebin` is called subsequently. However, in the test reported in the `main.py` file of Section 4, different versions of the dataset are needed and so, in this context, the pickle operation is insignificant (and the 102-105 lines are commented).

### 3 Classifier

The `malware_classifier.py` is a Python Class which methods, following the Machine Learning libraries syntax, offer useful kits for a correct binary classification of a file (i.e. malicious or non-malicious).

```
7 classifiers = {  
8     'svm'      : SVC(kernel='linear'),  
9     'random'   : RandomForestClassifier(n_estimators=100),  
10    'multinomial' : MultinomialNB()  
11 }
```

The three classifiers are taken from the scikit-learn library and are reported in a dictionary. As we'll see in the next section, the final user can select one of these simply instantiating a new `MalwareClassifier` object using 'svm', 'random' or 'multinomial' as parameter.

```
15 def __init__(self, classifier):  
16     self.classifier = classifiers[classifier]
```

#### 3.1 fit and predict

The `fit` function simply trains the selected classifier according to the [X,y] input, the `predict` function return a predicted vector of labels given a test matrix X.

```
18 def fit(self, X, y):  
19     self.classifier.fit(X, y)
```

```
21 def predict(self, X):  
22     return self.classifier.predict(X)
```

#### 3.2 scores

In the `score` method is returned a dictionary with three different evaluation metrics: the detection score (i.e. the accuracy score), the precision score and the false positive rate (given by the confusion matrix here not reported) calculated in the private auxiliary file `_false_positive_rate()`.

```

29 def scores(self, y_pred, y_true):
30     d = {}
31     d['detection_score'] = accuracy_score(y_true, y_pred)
32     d['precision_score'] = precision_score(y_true, y_pred)
33     d['false_positive_rate'] =
        self.__false_positive_rate(y_true, y_pred)
34     return d

```

### 3.3 print\_results

This method simply define provide a useful print function to standard output of the scores calculate in the previous function. For a better accuracy, the numbers are represented with the first three digits before the comma.

```

36 def print_results(self, all_scores):
37     print("\n"+30*"-"+"\nResults:")
38     for s in all_scores:
39         print(s + ': %0.3f' % all_scores[s])
40     print(30*"-"+"\n\n")

```

## 4 Main

Listing 2: main.py file

```
1 import numpy as np
2 import random
3 from malware_classifier import MalwareClassifier
4 from sklearn.model_selection import train_test_split
5 from dataset import load_drebin, get_prefixes
6
7 prefixes = get_prefixes()
8 random.shuffle(prefixes)
9 chunks = [prefixes[x:x+3] for x in range(0,
10         len(prefixes), 3)]
11
12 for p in chunks:
13     print("training with %s" % p)
14
15     X, y = load_drebin(p)
16
17     X_train, X_test, y_train, y_test = train_test_split(X,
18         y, test_size=0.33)
19     md = MalwareClassifier('random')
20
21     print('training...')
22     md.fit(X_train, y_train)
23
24     print('evaluating...')
25     y_pred = md.predict(X_test)
26
27     res = md.scores(y_pred, y_test)
28     md.print_results(res)
```

That's the complete code of the main file of our test. It can be executed from terminal with the `python3 main.py` command.

In line 7 the prefixes are taken and shuffled with the `random.shuffle` method to ensure a better impartiality of the final result. In fact, as already explained in 2.1, the output list of the `get_prefixes()` function is sorted according to non decreasing number of the prefixes. In line 9 the prefixes list is separated in chunks (sub-lists with three elements each).



The core of the `main.py` file is represented by the cycle starting at line 11: for each sub-list the drebin-adapted dataset is given and used as the input for the Random-Forest classifier; every predicted label is compared with the real label using the `md.scores` method. The final results are reported in Section 5.

## 5 Results

Listing 3: results.log file

```
1 training with ['feature', 'api_call', 'provider']
2 4900 features
3 training...
4 evaluating...
5
6 -----
7 Results:
8 detection_score: 0.989
9 precision_score: 0.930
10 false_positive_rate: 0.003
11 -----
12
13
14 training with ['call', 'real_permission', 'permission']
15 4615 features
16 training...
17 evaluating...
18
19 -----
20 Results:
21 detection_score: 0.993
22 precision_score: 0.974
23 false_positive_rate: 0.001
24 -----
```

In this section are reported the results of our test. With the command `python3 main.py > results.log`, a new output file is created in the working directory. Using the Random Forest classifier and the selected prefixes as input, good result are obtained in an acceptable time.

It's very important to underline the test is executed on a PC with i7 processor and 8 GB memory and that in the worst case the program need a memory capacity higher than 6 GB. The random forest emerges as the fastest among the available methods of the `malware_classifier.py` class.