

Configuration Manual

MSc Academic Internship
MSc Cyber Security

Viraj Kudtarkar
Student ID: 18178499

School of Computing
National College of Ireland

Supervisor: Ross Spelman

National College of Ireland
MSc Project Submission Sheet
School of Computing



Student Name: Viraj Kudtarkar
Student ID: 18178499
Programme: MSc Cyber Security **Year:** 2020
Module: MSc Internship
Supervisor: Ross Spelman
Submission Due Date: 17/08/2020
Project Title: Android botnet detection using signature data and Ensemble Machine Learning.
Word Count: 1173 **Page Count:** 12

I hereby certify that the information contained in this (my submission) is information pertaining to research I conducted for this project. All information other than my own contribution will be fully referenced and listed in the relevant bibliography section at the rear of the project.

ALL internet material must be referenced in the bibliography section. Students are required to use the Referencing Standard specified in the report template. To use other author's written or electronic work is illegal (plagiarism) and may result in disciplinary action.

I agree to an electronic copy of my thesis being made publicly available on NORMA the National College of Ireland's Institutional Repository for consultation.

Signature: Viraj Kudtarkar

Date: 17 August 2020

PLEASE READ THE FOLLOWING INSTRUCTIONS AND CHECKLIST

Attach a completed copy of this sheet to each project (including multiple copies)	<input type="checkbox"/>
Attach a Moodle submission receipt of the online project submission, to each project (including multiple copies).	<input type="checkbox"/>
You must ensure that you retain a HARD COPY of the project, both for your own reference and in case a project is lost or mislaid. It is not sufficient to keep a copy on computer.	<input type="checkbox"/>

Assignments that are submitted to the Programme Coordinator Office must be placed into the assignment box located outside the office.

Office Use Only	
Signature:	
Date:	
Penalty Applied (if applicable):	

Configuration Manual

Viraj Kudtarkar
18178499

1 System Configurations:

OS-MAC Catalina
Processor – Intel Core i7 2.8GHz
Ram 16 Gb
Python – 3.8
IDE – Pycharm 2020.1
Pip 20.2.1

2 Project Libraries:

Database = Sqlite3

Joblib v0.16.0

CSV Reader = pandas v1.1.0

ML = scikit-learn v0.23.2

3 Walkthrough:

3.1 Virtual Environment:

First of all, we setup a virtual environment “Pycharm” IDE to installing all the dependences/Libraries used in our project. We used “Pip3 v20.3.1” package manager for installing different libraries. The most common libraries are pandas, scikit-learn, Joblib and sqlite3.

pandas	1.1.0	1.1.0
pip	20.2.1	▲ 20.2.2
plac	1.1.3	▲ 1.2.0
prshed	2.0.1	▲ 3.0.2
protobuf	3.12.4	▲ 3.13.0
pyasn1	0.4.8	0.4.8
pyasn1-modules	0.2.8	0.2.8
pymongo	3.11.0	3.11.0
pyarsing	2.4.7	2.4.7
python-dateutil	2.7.5	▲ 2.8.1
pytz	2020.1	2020.1
regex	2020.7.14	2020.7.14
requests	2.22.0	▲ 2.24.0
rfc3986	1.4.0	1.4.0
rsa	4.6	4.6
scikit-learn	0.23.2	0.23.2
scipy	1.5.2	1.5.2
selenium	2.53.6	▲ 3.141.0
setuptools	49.2.1	▲ 49.6.0
six	1.12.0	▲ 1.15.0
sniffio	1.1.0	1.1.0
soupsieve	1.9.2	▲ 2.0.1
spacy	2.3.2	2.3.2

Figure 1: Installed Packages

3.2 Dataset:

We used the dataset, for our ML work, which includes the data of botnet and normal applications. Dataset in csv file contains feature set of 15 permissions, 3 intents, family, botnet-category, App Package-Name/MD5 hash keys and labeled column.

Permissions	Intents
android.permission.INTERNET	android.intent.action.BOOT_COMPLETED
android.permission.READ_PHONE_STATE	android.intent.action.POWER_CONNECTED
android.permission.ACCESS_NETWORK_STATE	android.intent.action.BATTERY_LOW
android.permission.WRITE_EXTERNAL_STORAGE	
android.permission.RECEIVE_BOOT_COMPLETED	
android.permission.READ_SMS	
android.permission.SEND_SMS	
android.permission.WRITE_SMS	
android.permission.RECEIVE_SMS	
android.permission.ACCESS_COARSE_LOCATION	
android.permission.ACCESS_FINE_LOCATION	
android.permission.ACCESS_WIFI_STATE	
android.permission.CALL_PHONE	
android.permission.WAKE_LOCK	
android.permission.READ_CONTACTS	

Table 1: Feature set of Permissions and Intents

	A	B	C	D	E	F	G	H	I	J	K	L	M	N	O	P	Q	R	S	T	Results							
1	android.pern	android.pern	android.pern	android.pern	android.pern	android.pern	android.pern	android.pern	android.pern	android.pern	android.pern	android.pern	android.pern	android.pern	android.pern	android.pern	android.pern	android.pern	android.pern	android.pern	android.inter	android.inter	android.inter	MD5				
2	1	1	1	1	0	0	0	0	0	0	1	1	1	0	1	0	0	0	0	0	0	0	0	0	0	1c4e357a8e...	1	
3	1	1	1	1	0	0	0	0	0	0	1	1	1	0	1	0	0	0	0	0	0	0	0	0	0	1d15765ffec...	1	
4	1	1	1	1	0	0	0	0	0	0	0	0	1	0	1	0	0	0	0	0	0	0	0	0	0	31c657b777x...	1	
5	1	1	1	1	1	0	0	0	0	0	1	1	1	0	1	0	0	0	0	0	0	0	0	0	0	37b993b5f5...	1	
6	1	1	1	1	0	0	0	0	0	0	0	0	1	0	1	0	0	0	0	0	0	0	0	0	0	3e3072644a...	1	
7	1	1	1	1	0	0	0	0	0	1	1	1	1	0	1	0	0	0	0	0	0	0	0	0	0	41172215c...	1	
8	1	1	1	1	1	0	1	0	1	0	1	0	0	0	1	0	0	0	0	0	0	0	0	0	0	005465c7b...	1	
9	1	1	1	1	1	0	0	0	0	1	0	0	1	0	1	0	0	0	0	0	0	0	0	0	0	03e8e08061...	1	
10	1	1	1	1	1	0	1	0	1	0	0	0	0	1	1	0	0	0	0	0	0	0	0	0	0	1de012d8de...	1	
11	1	1	1	1	1	0	1	0	1	0	1	0	0	1	1	0	0	0	0	0	0	0	0	0	0	1e6801a75...	1	
12	1	1	1	1	1	0	1	0	1	0	0	0	0	1	1	0	0	0	0	0	0	0	0	0	0	1e87d0abe9...	1	
13	1	1	1	1	1	0	1	0	1	0	0	0	0	1	1	0	0	0	0	0	0	0	0	0	0	2bcf8dd123...	1	
14	1	1	1	1	1	0	0	0	0	0	1	1	1	1	1	1	1	1	0	0	0	0	0	0	0	04c12809d3...	1	
15	1	1	1	1	1	0	0	0	0	0	0	1	1	1	1	1	1	1	0	0	0	0	0	0	0	0	40b9b74525...	1
16	1	1	1	1	1	0	0	0	0	0	1	1	1	1	1	1	1	1	0	0	0	0	0	0	0	0	5d184d033...	1
17	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	0	0	0	0	0	0	0	0	a8fe9d03651...	1
18	1	1	1	1	0	0	0	0	0	0	1	0	0	1	0	0	0	0	0	0	0	0	0	0	0	0	b2bb7d5f5e...	1
19	1	1	1	1	0	0	1	0	0	0	1	0	0	1	1	0	0	0	0	0	0	0	0	0	0	0	b3c7575f23...	1
20	1	1	1	1	0	0	0	0	0	0	1	0	1	0	1	0	0	0	0	0	0	0	0	0	0	0	b464a695df...	1
21	1	1	1	1	0	0	0	0	0	0	1	1	1	0	1	0	0	0	0	0	0	0	0	0	0	0	b51837f70c...	1
22	1	1	1	1	0	0	0	0	0	0	1	0	1	0	0	0	0	0	0	0	0	0	0	0	0	0	b6d0bb5f7e...	1
23	1	1	1	1	1	0	0	0	0	0	1	1	1	0	1	0	1	0	0	0	0	0	0	0	0	0	030423f268...	1
24	1	1	1	1	1	0	0	0	0	0	1	1	1	0	1	1	1	0	0	0	0	0	0	0	0	0	1634b1fb3b...	1
25	1	1	1	1	1	0	0	0	0	0	0	0	1	0	1	0	1	0	0	0	0	0	0	0	0	0	1652df0226...	1
26	1	1	1	1	1	1	1	1	0	1	1	0	1	0	1	0	1	0	1	0	0	0	0	0	0	0	7739442c33...	1
27	1	1	1	1	1	0	0	0	0	0	0	0	1	0	1	0	1	0	0	0	0	0	0	0	0	0	87206282b5...	1
28	1	1	1	1	1	0	0	0	0	0	1	1	1	0	1	0	1	0	1	0	0	0	0	0	0	0	9324376e27...	1
29	1	1	1	1	1	1	1	1	1	1	1	1	1	0	1	1	1	1	0	0	0	0	0	0	0	0	a3b065085e...	1
30	1	1	1	1	1	0	0	0	0	0	1	1	1	0	1	0	1	0	0	0	0	0	0	0	0	0	a4fcdf1992d...	1
31	1	1	1	1	1	0	0	0	0	0	1	1	1	0	1	0	0	0	0	0	0	0	0	0	0	0	0c67d0919e...	1
32	1	1	1	1	1	0	0	0	0	0	1	1	1	0	1	0	1	0	0	0	0	0	0	0	0	0	162cb09e2e...	1
33	1	1	1	1	1	1	1	1	1	1	1	1	1	0	1	0	1	0	1	0	0	0	0	0	0	0	1be29a6e22...	1
34	1	1	1	1	1	0	0	0	0	0	1	1	1	0	1	0	0	0	0	0	0	0	0	0	0	0	6a66635b6b...	1
35	1	1	1	1	1	0	0	0	0	0	1	1	1	0	1	0	1	0	1	1	1	1	1	1	1	0	6fc29ab75df...	1
36	1	1	1	1	1	1	1	1	1	1	1	1	1	0	1	1	1	1	1	1	1	1	1	1	1	0	7cd86d83d9...	1

Figure 2: Dataset CSV

3.3 Database:

We used a sqlite database to store the android application's unique identifier "Package-name" or MD5 hash key-value. Database helped us in making our system more efficient by returning responses in very less time. Database contains a table named "Apps" which stores all the app's information. We have data of almost 600 Normal and botnet applications in a json format. This json format data was inserted in our database for later use.

```

4
5
6 def create_connect_db():
7     try:
8         sqlite_connection = sqlite3.connect('SQLite_Python.db')
9         return sqlite_connection
10
11     except sqlite3.Error as error:
12         print("Error while connecting to sqlite", error)
13
14
15 def create_table():
16     db_conn = create_connect_db()
17     cursor = db_conn.cursor()
18     # Create a table.
19     query = """CREATE TABLE IF NOT EXISTS Apps (ID INTEGER PRIMARY KEY, PACKAGE_NAME VARCHAR(100), IS_BOTNET INTEGER); """
20     cursor.execute(query)
21     db_conn.commit()
22     cursor.close()
23
24
25 # Returns true if table already contains apps data else returns false
26 def table_contain_data():
27     db_conn = create_connect_db()
28     query = "SELECT * FROM Apps"
29     cursor = db_conn.execute(query)
30     data = cursor.fetchall()
31     cursor.close()
32     db_conn.close()
33
34     if len(data) > 0:

```

Figure 3: Database and Table Creation

```

82
83 # method checks if app does not exists in db then returns 0,
84 # if exists and app is normal then -1,
85 # if exists and app is botnet then 1.
86 def check_app_db(package_name):
87     db_conn = create_connect_db()
88     query = "SELECT * FROM Apps WHERE PACKAGE_NAME LIKE '" + package_name + "'"
89     cursor = db_conn.execute(query)
90     data = cursor.fetchall()
91     cursor.close()
92     db_conn.close()
93     if len(data) > 0:
94         return data[0][2]
95     else:
96         return 0
97
98

```

Figure 4: Querying Database

3.4 Data Pre-processing:

We selected 18 total features (permission and intent), which are most important in distinguishing botnet applications, for training our ML models. In preprocessing part, we removed the noise from data by deleting rows which are incomplete. Also, we dropped the columns that are not necessary for our training process which are “family”, “category” and “MD5”. After filtering process, data was spitted into 2 parts with 70:30 ratio, 70% for training and 30% for testing.

```

12
13 def get_dataset():
14     dataset = pd.read_csv("datasets/data.csv")
15     # Preprocess Data (removing extra columns)
16     dataset = dataset.drop('family', 1)
17     dataset = dataset.drop('category', 1)
18     dataset = dataset.drop('MD5', 1)
19     x = dataset.iloc[:, :-1].values
20     y = dataset.iloc[:, -1:].values
21     x_train, x_test, y_train, y_test = train_test_split(x, y, test_size=0.30, random_state=1)
22     return x_train, x_test, y_train, y_test
23

```

Figure 5: Dataset preprocessing and splitting

3.5 Classifier Training:

We have trained 5 different classifiers which includes: Logistic Regression, Random Forest, Support Vector Machine, Decision Tree and Naive Bayes, for our machine learning predictor. Models were trained on 70% of dataset. In training process, first the model is created and trained by “model.fit()” method. “Model.predict()” is used for getting test results out of the model and returns a confusion matrix which is used for finding the accuracy of the trained model. We extracted 3 different performance measures for our models which are Accuracy, Precision and Recall. At last, after training process, all trained models were stored in “pkl” file format.

Source code for training the Logistic Regression in mentioned in figure 6.

```
24
25 def train_lr():
26     data_train, data_test, label_train, label_test = get_dataset()
27     model = LogisticRegression(random_state=0)
28     model.fit(data_train, label_train.ravel())
29
30     # predicting the tests set result
31     pred = model.predict(data_test)
32     create_confusion_matrix(label_test, pred)
33
34     # storing trained models
35     store_model(model, "trained_models/lr.pkl")
36
37
```

Figure 6: Logistic Regression Model

Figure 7 contains the code for training Random forest Model with best and appropriate parameters.

```
37
38 def train_rf():
39     data_train, data_test, label_train, label_test = get_dataset()
40     model = RandomForestClassifier(n_estimators=100, criterion="gini", max_features='log2', random_state=0)
41     model.fit(data_train, label_train.ravel())
42
43     # predicting the tests set result
44     pred = model.predict(data_test)
45     create_confusion_matrix(label_test, pred)
46
47     # storing trained models
48     store_model(model, "trained_models/rf.pkl")
49
```

Figure 7: Random Forest Model

Source code for training the SVM in mentioned in figure 6.

```
50
51 def train_svm():
52     data_train, data_test, label_train, label_test = get_dataset()
53     model = SVC(C=1000, kernel='rbf', gamma=0.2, random_state=0)
54     # classifier.fit(x_train, y_train.ravel())
55     model.fit(data_train, label_train.ravel())
56
57     # predicting the tests set result
58     pred = model.predict(data_test)
59     create_confusion_matrix(label_test, pred)
60
61     # storing trained models
62     store_model(model, "trained_models/svm.pkl")
63
```

Figure 8: Support Vector Machine Model

Following figure 9 contains the code of Naïve Bayes model training and storing in the “pkl” file format.

```

79 # Naive Bayes
80 def train_nb():
81     data_train, data_test, label_train, label_test = get_dataset()
82     model = GaussianNB()
83     # classifier.fit(x_train, y_train.ravel())
84     model.fit(data_train, label_train.ravel())
85
86     # predicting the tests set result
87     pred = model.predict(data_test)
88     create_confusion_matrix(label_test, pred)
89
90     # storing trained models
91     store_model(model, "trained_models/nb.pkl")
92
93

```

Figure 9: Naive Bayes Model

Source code for training the Decision Tree in mentioned in figure 10.

```

64
65 def train_dt():
66     data_train, data_test, label_train, label_test = get_dataset()
67     model = DecisionTreeClassifier(random_state=0)
68     # classifier.fit(x_train, y_train.ravel())
69     model.fit(data_train, label_train.ravel())
70
71     # predicting the tests set result
72     pred = model.predict(data_test)
73     create_confusion_matrix(label_test, pred)
74
75     # storing trained models
76     store_model(model, "trained_models/dt.pkl")
77

```

Figure 10: Decision Tree Model

Figure 11 is the code for creating confusion matrix , calculating accuracy, precision and recall values.


```

93
94 def create_confusion_matrix(testset, predictions):
95     # confusion matrix
96     cm = confusion_matrix(testset, predictions)
97     print("Confusion Matrix: ", cm)
98
99     # Accuracy
100    accuracy = (cm[0][0] + cm[1][1]) / (cm[0][0] + cm[0][1] + cm[1][0] + cm[1][1])
101    print("Accuracy:", accuracy)
102
103    # Precision
104    precision = cm[0][0]/(cm[0][0]+cm[1][0])
105    print("Precision:", precision)
106
107    # Recall
108    recall = cm[0][0] / (cm[0][0] + cm[0][1])
109    print("Recall:", recall)
110

```

Figure 11: Confusion Matrix Creation

3.6 Confusion Matrix

Data trained to the five above mentioned algorithms is analyzed during the training session with predicted results against the actual results.

1. Logistic Regression

The predicted true and predicted false results for the Logistic regression algorithm in comparison with actual true and actual false values are given in Table 2 below.

Logistic regression	Predicted True	Predicted False
Actual True	328	10
Actual False	45	61

Table 2: LR Confusion Matrix

The accuracy, precision and recall values on the basis of recorded results for Logistic regression has been identified as given in Table 3 below.

Logistic Regression		
1	Accuracy	87.61%
2	Precision	87.93%

3	Recall	97.04%
---	--------	--------

Table 3

2. Support Vector Machine (SVM)

The predicted results for the Support Vector Machine (SVM) algorithm in comparison with actual true and actual false values are given in Table 4 below.

SVM	Predicted True	Predicted False
Actual True	336	2
Actual False	18	88

Table 4: SVM Confusion Matrix

The recorded accuracy, precision and recall values for Support Vector Machine (SVM) algorithm is recorded as given in Table 5 below.

Support Vector Machine (SVM)		
1	Accuracy	95.49%
2	Precision	94.91%
3	Recall	99.40%

Table 5

3. Random Forest

The predicted results for the Random Forest algorithm with actual and predicted values as recorded in Table 6 given below.

Random Forest	Predicted True	Predicted False
Actual True	332	6
Actual False	19	87

Table 6: RF Confusion Matrix

The calculated values of precision, accuracy and recall have been decided on the basis of the total number of predicted (True, False) and the total number of actual (True, False) values given in table 7 below.

Random Forest		
1	Accuracy	94.36%
2	Precision	94.58%
3	Recall	98.22%

Table 7

4. Decision Tree

Results table recording the predicted and actually true, false values for the Decision Tree algorithm during the testing phase is shown in Table 8 below.

Decision Tree	Predicted True	Predicted False
Actual True	333	5
Actual False	18	88

Table 8: DT Confusion Matrix

Accuracy, precision and recall values calculated on the basis of results recorded in table 9 are listed below in table 9.

Decision Tree		
1	Accuracy	94.81%
2	Precision	94.87%
3	Recall	98.22%

Table 9

5. Naïve Bayes

Training datasets to the Naïve Bayes algorithm are further tested on the basis of 20% of the overall dataset and the results (true, false) have been recorded as listed in table 10 below.

Naïve Bayes	Predicted True	Predicted False
Actual True	624	16
Actual False	61	45

Table 10: NB Confusion Matrix

The accuracy, precision and recall values recorded from the predicted and actual result calculations are given in Table 11 below.

Naïve Bayes		
1	Accuracy	83.10%
2	Precision	84.15%
3	Recall	95.85%

Table 11

3.7 Analysis and Model selection

The recorded results are compared based on the values of the overall percentage of accuracy, precision and recall. The top three models are then further shortlisted with higher accuracy values i.e., Support Vector Machine (SVM), Decision Tree and Random Forest.

Table 1.10 given below gives a better understanding of the top 3 selected models.

ID	Parameters	SVM	Random Forest	Decision Tree
1	Accuracy	95.49%	94.36%	94.81%
2	Precision	94.91%	94.58%	94.87%
3	Recall	99.40%	98.22%	98.22%

3.8 Voting-Based Ensemble Method:

We used ensemble methods in our predictor. Ensemble is a technique which combines different models to produce improved results. Ensembles can produce more precise and accurate outputs. We have created voting-based ensemble model which includes 3 models: Random Forest, SVM and Decision Tree. 3 models out of total 5 models were chosen on the basis of accuracy, so only top 3 model with best accuracy results were used in creation of ensembles. Source code written for ensemble is mentioned in figure 12.

```

4 classifier_rf = None
5 classifier_dt = None
6 classifier_svm = None
7
8
9 # load the pickle file
10 def load_classifiers():
11     global classifier_rf, classifier_dt, classifier_svm
12     classifier_rf = joblib.load('models/trained_models/rf.pkl')
13     classifier_dt = joblib.load('models/trained_models/dt.pkl')
14     classifier_svm = joblib.load('models/trained_models/svm.pkl')
15
16
17 def check_app(vec):
18     # checking and predicting
19     feature_vec = vec
20     prediction_rf = classifier_rf.predict(feature_vec)
21     prediction_dt = classifier_dt.predict(feature_vec)
22     prediction_svm = classifier_svm.predict(feature_vec)
23
24     # Voting Mechanism
25     prediction_result = prediction_rf + prediction_dt + prediction_svm
26
27     if prediction_result > 0:
28         return 1
29     else:
30         return -1
31

```

Figure 12: Ensemble Source code

References