

Configuration Manual

MSc Research Project
Data Analytics

Shibi Murugesan
Student ID: X18170871

School of Computing
National College of Ireland

Supervisor: Dr. Pierpaolo Dondio

National College of Ireland
Project Submission Sheet
School of Computing



Student Name:	Shibi Murugesan
Student ID:	X18170871
Programme:	Data Analytics
Year:	2019
Module:	MSc Research Project
Supervisor:	Dr.Pierpaolo Dondio
Submission Due Date:	12/12/2019
Project Title:	Configuration Manual
Word Count:	512
Page Count:	11

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.

Signature:	
Date:	12th December 2019

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

Shibi Murugesan
X18170871

1 Introduction

This configuration manual explains most of the system setup, software & hardware requirements and finally programming code required to implement the research project: *"Application of Machine Learning Models for Network Intrusion Detection Systems Based on Feature Selection Approach"*

2 Systems Configurations

2.1 Hardware Requirements

- **Processor:** Intel(R) Core(TM) i7-8750H CPU @ 2.20GHz
- **RAM:** 16.0 GB
- **System Type:** 64-bit Operating System, x64-based processor

2.2 Software Requirements

- **Pycharm Community Edition:** Much of the data aggregation, merging, cleaning, etc. was done in python.
- **Packages Required:** sklearn, yellowbrick, matplotlib, pandas and zipfile.

3 Dataset Preparation

The *NDL-KDD* dataset is downloaded from <https://www.unb.ca/cic/datasets/nsl.html>. The downloaded file would be zipped, hence we need to use *zipfile* python package to unzip it and convert the required set to the pandas dataframe format.

3.1 *NSL-KDD* unzip and converting to dataframe

3.1.1 Extraction

Extract the downloaded *NSL-KDD.zip* using *zipfile* python package in given location.

```
1 """
2 Extract required dataset from zipfile
3 """
4 import zipfile
5 import os
```

```

6 import pandas as pd
7
8 PATH_TO_ZIP_FILE = 'NSL-KDD.zip'
9 PATH_TO_ARFF_FILE = 'NSL-KDD/'
10 KDD_TRAIN_PATH = 'KDDTrain+.arff'
11 KDD_TEST_PATH = 'KDDTest+.arff'
12 IMAGE_DIR = 'images/'
13 LATEX_DIR = 'latex/'
14 CM_DIR = IMAGE_DIR + 'cm/'
15 ROC_DIR = IMAGE_DIR + 'roc/'
16 REPORT_DIR = IMAGE_DIR + 'report/'
17
18 def extractzipfile():
19     """
20     Unzip the zip file to a folder
21     Also make images directory if not exist.
22     """
23     # Create the Image directory to save any plots
24     if not os.path.exists(IMAGE_DIR):
25         os.makedirs(IMAGE_DIR)
26
27     # Extract the zipfile
28     if os.path.exists(PATH_TO_ZIP_FILE):
29         with zipfile.ZipFile(PATH_TO_ZIP_FILE, 'r') as zip_ref:
30             zip_ref.extractall(PATH_TO_ARFF_FILE)
31         return True
32
33     # If ZipFile is not found
34     return False

```

After unzipping, we might see below files in extracted folder

1. NSL-KDD/index.html
2. *NSL-KDD/KDDTest+.arff – > **We will use this as an Input to our algorithm**
3. NSL-KDD/KDDTest+.txt
4. NSL-KDD/KDDTest1.jpg
5. NSL-KDD/KDDTest-21.arff
6. NSL-KDD/KDDTest-21.txt
7. *NSL-KDD/KDDTrain+.arff – > **We will use this as an Input to our algorithm**
8. NSL-KDD/KDDTrain+.txt
9. NSL-KDD/KDDTrain+_20Percent.arff
10. NSL-KDD/KDDTrain+_20Percent.txt
11. NSL-KDD/KDDTrain1.jpg

3.1.2 Conversion to Dataframe

Convert the extracted **NSL-KDD/KDDTrain+.arff** and **NSL-KDD/KDDTest+.arff** into corresponding dataframes

```

1 def loaddataframe(filename, display=False):
2     """
3     Extract the required arff file into a dataframe.
4     :return: <pd.DataFrame> Dataframe related with KDD data
5     """
6
7     dataset = pd.read_table(os.path.join(PATH_TO_ARFF_FILE, filename), sep=',',
8                             index_col=False, header=None, names=COLUMNS)
9
10    if display:
11        print("Size of %s: " % filename)
12        print("\tSamples Count: %d" % dataset.shape[0])
13        print("\tAttributes Count: %d" % dataset.shape[1])
14    return dataset
15
16 # Load the Dataset
17 traindata = extract.loaddataframe('KDDTrain+.txt')
18 testdata = extract.loaddataframe('KDDTest+.txt')

```

4 Dataset Pre-processing

Dataset Preprocessing follows below 3 steps:

1. Aggregating different attacks into 5 attack types
2. Label Encoding and Resampling minority Classes
3. Normalisation

4.1 Aggregating different attacks into 5 attack types

Here we will create a new column **class_group** which maps different attacks into different classes.

```

1     """
2     Aggregate different types of attacks into different groups.
3     """
4
5     def aggregateintogroups(value):
6         """
7         We will map different classes into 5 different groups.
8         :return: <String> Return class group type
9         """
10        mapping = {'Probe': ['ipsweep', 'satan', 'nmap', 'portsweep', 'saint', 'mscan',
11                            ''],
12                  'DoS': ['teardrop', 'pod', 'land', 'back', 'neptune', 'smurf', 'mailbomb',
13                          'udpstorm', 'apache2', 'processtable'],
14                  'U2R': ['perl', 'loadmodule', 'rootkit', 'buffer_overflow', 'xterm',
15                          'ps',
16                          'sqlattack', 'httptunnel'],
17                  'R2L': ['ftp_write', 'phf', 'guess_passwd', 'warezmaster', 'warezclient',
18                          'imap', 'spy', 'multihop', 'named', 'snmpguess', 'worm', 'snmpgetattack',
19                          'xsnoop', 'xlock', 'sendmail'],
20                  'normal': ['normal']}
21
22        # return the group type
23        for grp, names in mapping.items():
24            for name in names:
25                if name == value:
26                    return grp
27
28        return None

```

4.2 Label Encoding and Re-sampling minority Classes

We will convert the categorical columns into a **Label Encoded** format. This is then followed by up-sampling the *U2R* and *R2L* samples. And then followed by down-sampling the *Normal* samples.

```
1 def createfeatureslabelset(dataframe, train=True):
2     """
3     Label encodes the 'object' type columns
4     :param dataframe: <pd.DataFrame> Dataframe having NSL-KDD records
5     """
6     # Drop the class column
7     dataframe = dataframe.drop('class', axis=1)
8
9     encoder = LabelEncoder()
10    # extract categorical attributes from dataframe
11    # encode the categorical attributes
12    categorical = dataframe.select_dtypes(include=['object']).apply(encoder.fit_transform)
13    # Now let us merge both scaled and encoded feature dataframes
14    dataframe = pd.concat([dataframe.drop(['protocol_type', 'service',
15                                         'flag', 'class_group'],
16                                         axis=1), categorical],
17                          axis=1)
18
19    if train:
20        # separate minority and majority classes
21        DoSsection = dataframe[dataframe['class_group'] == 0]
22        Probesection = dataframe[dataframe['class_group'] == 1]
23        R2Lsection = dataframe[dataframe['class_group'] == 2]
24        U2Rsection = dataframe[dataframe['class_group'] == 3]
25        Normalsection = dataframe[dataframe['class_group'] == 4]
26
27        # upsample minority
28        R2L_upsampled = resample(R2Lsection,
29                                # sample with replacement
30                                replace=True,
31                                # match number in majority class
32                                n_samples=len(Probesection),
33                                # reproducible results
34                                random_state=27)
35
36        U2R_upsampled = resample(U2Rsection,
37                                replace=True,
38                                n_samples=len(Probesection),
39                                random_state=27)
40
41        # downsample majority
42        Normal_downsample = resample(Normalsection,
43                                    replace=False,
44                                    n_samples=len(Probesection),
45                                    random_state=27)
46
47        # combine majority and upsampled minority
48        dataframe = pd.concat([DoSsection, Probesection,
49                               R2L_upsampled,
50                               U2R_upsampled,
51                               Normal_downsample])
52
53    features = dataframe.drop('class_group', axis=1)
54    labels = pd.DataFrame()
55    labels['labels'] = dataframe['class_group']
56
57    return features, labels
```

4.3 Normalisation

We will use *MinMaxScalar* from sklearn to adjust all continuous values from dataset between 0 and 1. Also, some features with binary results like either 0 or 1 should not be considered for normalisation. Those features with binary values are *'land'*, *'logged_in'*, *'root_shell'*, *'su_attempted'*, *'is_host_login'*, *'is_guest_login'*

```
1 def normalisation(dataframe):
2     """
3     Apply Minmaxscalar on all int64 and float64 columns from dataframe.
4     """
5     # Extract if any binary columns present in dataframe
6     binarycolumns = ['land', 'logged_in',
7                     'root_shell', 'su_attempted',
8                     'is_host_login', 'is_guest_login']
9
10    bc = []
11    for col in dataframe.columns:
12        if col in binarycolumns:
13            bc.append(col)
14    binarydf = dataframe[bc]
15    dataframe = dataframe.drop(bc, 1)
16    dataframe = MinMaxScaler().fit(dataframe)
17
18    categoricaldf = dataframe[dataframe.select_dtypes(include=['int32',
19                                                         'int64',
20                                                         'float64']).
21                                                         columns]
22
23    result = pd.concat([binarydf, categoricaldf], axis=1)
24    return result
```

5 Feature Selection Techniques

5.1 Embedded Technique

A two different feature selection techniques were applied i.e. *LASSO* and *Elastic-Net*

```
1 """
2 This program implements different embedded/wrapper techniques for feature
3 selection
4 """
5 from sklearn.linear_model import Lasso, ElasticNet
6 import numpy as np
7 import pandas as pd
8 import seaborn as sns
9 import matplotlib.pyplot as plt
10
11 class LassoSelection:
12     """
13     Class to apply LASSO feature selection Technique.
14     """
15     def __init__(self, X, y):
16         """
17         Initialize feature columns and label column
18         :param X: <DataFrame> Features from Sensor Dataframe
19         :param y: <DataFrame> Label column from Sensor Dataframe
20         """
21         self.X = X
22         self.y = y
23         self.coef = None
24
25     def applyLasso(self):
26         lasso = Lasso(alpha=0.1)
27         lasso.fit(self.X, self.y)
```

```

28         self.coef = lasso.coef_
29         return np.abs(self.coef)
30
31
32 class Elasticnet(LassoSelection):
33     """
34     Class to apply Elastic-Net feature selection Technique. It is hybrid
35     of LASSO and Ridge
36     """
37     def applyElasticNet(self):
38         regr = ElasticNet(alpha=0.01)
39         regr.fit(self.X, self.y)
40         return np.abs(regr.coef_)

```

5.1.1 Creating Feature Ranking table

A ranks table is created which arranges the features in ranked order of their feature importance on the class column.

```

1  def rank_to_dict(ranks, names=trainX.columns, order=1):
2      """
3      Apply minmax scaler to normalise score or rank values of features
4      """
5      minmax = MinMaxScaler()
6      ranks = minmax.fit_transform(order*np.array([ranks]).T).T[0]
7      ranks = map(lambda x: round(x, 2), ranks)
8      return dict(zip(names, ranks))
9
10
11 def create_ranks_dataframe(ranks):
12     """
13     This function does the following manipulations:
14     1. Converts dictionary of ranks to a dataframe 'ranks_dataframe'.
15     2. Creates a new column 'Mean' in 'ranks_dataframe' which specifies the
16     mean along the rows.
17     3. Sorts the 'ranks_dataframe' in descending order of 'Mean' Value.
18     4. Creates a another column 'Rank' based on 'Mean' Column.
19
20     :param ranks: <Dict> Dictionary of feature importance scores from all feature \
                ranking
21     methods
22     :return: <DataFrame> Converted ranks to a dataframe format
23     """
24     mean_ranks = {}
25
26     for name in trainX.columns:
27         mean_ranks[name] = np.mean([ranks[method][name] for method in ranks.keys \
                ()])
28
29     methods = sorted(ranks.keys())
30     ranks["Mean"] = mean_ranks
31     methods.append("Mean")
32
33     ranks_dataframe = pd.DataFrame(ranks).sort_values(by="Mean", ascending=False)
34
35     # Consider index of 'sensorX' inside column name
36     ranks_dataframe.reset_index(level=0, inplace=True)
37
38     # Create new column called as Rank and sequentially increment its value
39     ranks_dataframe["Rank"] = ranks_dataframe.index + 1
40     return ranks_dataframe
41
42
43 *****
44 # APPLY EMBEDDED METHOD FOR FEATURE RANKING #
45 *****
46 def callembeddedmethod(X=trainX, y=trainy):
47     """
48     In this function, various embedded feature rankings/selection methods are \
                implemented

```



```

49     side-by-side and the feature importance values are returned to a output dict ↘
        for each
50     of those methods.
51
52     :param X: <Numpy> Array of features from training set
53     :param y: <Numpy> Array of labels from training set
54     :return: <Dict> Dictionary of feature importance scores from all feature ↘
        ranking methods
55     """
56     # Create objects for respective embedded methods
57     lasso = featureselection.LassoSelection(X, y)
58     net = featureselection.Elasticnet(X, y)
59
60     dict_rank = {"Lasso": rank_to_dict(lasso.applyLasso()),
61                 "Elastic-Net": rank_to_dict(net.applyElasticNet())}
62     return dict_rank

```

5.2 Filter-based Technique

Pearson Correlation coefficient matrix is used to remove highly correlated features.

```

1  class Correlation(LassoSelection):
2      """
3      Brute Force Method to find Correlation between features
4      """
5      def correlation(self, threshold=0.8):
6          # Set of all names of correlated columns
7          data = pd.concat([self.X, self.y], axis=1)
8          data = data.drop(['num_outbound_cmds',
9                           'is_host_login'], axis=1)
10         col_corr = set()
11         corr_mat = data.corr()
12         for i in range(len(corr_mat.columns)):
13             for j in range(i):
14                 if abs(corr_mat.iloc[i, j]) > threshold:
15                     colname = corr_mat.columns[i]
16                     col_corr.add(colname)
17         # plot the heatmap
18         CorrMtx(corr_mat)
19         return col_corr
20
21     #*****
22     #          APPLY FILTER METHOD FOR FEATURE SELECTION          #
23     #*****
24     def callfiltermethod(X=trainX, y=trainy):
25         """
26         In this function, pearson correlation coefficient used to eliminate features
27         """
28         # Create objects for respective filter method
29         corr = featureselection.Correlation(X, y)
30         droppedfeatures = corr.correlation()
31         return droppedfeatures

```

5.2.1 To Plot heatmap of Correlation matrix

```

1  def CorrMtx(df, dropDuplicates = True):
2
3      # Your dataset is already a correlation matrix.
4      # If you have a dataset where you need to include the calculation
5      # of a correlation matrix, just uncomment the line below:
6      # df = df.corr()
7
8      # Exclude duplicate correlations by masking upper right values
9      if dropDuplicates:
10         mask = np.zeros_like(df, dtype=np.bool)
11         mask[np.triu_indices_from(mask)] = True
12
13     # Set background color / chart style
14     sns.set_style(style = 'white')

```

```

15
16 # Set up matplotlib figure
17 f, ax = plt.subplots(figsize=(11, 9))
18
19 # Add diverging colormap from red to blue
20 cmap = sns.diverging_palette(250, 10, as_cmap=True)
21
22 # Draw correlation plot with or without duplicates
23 if dropDuplicates:
24     map = sns.heatmap(df, mask=mask, cmap=cmap,
25                       square=True,
26                       linewidth=.5, cbar_kws={"shrink": .8}, ax=ax)
27 else:
28     map = sns.heatmap(df, cmap=cmap,
29                       square=True,
30                       linewidth=.5, cbar_kws={"shrink": .8}, ax=ax)
31 map.set_yticklabels(map.get_yticklabels(),
32                    horizontalalignment='right',
33                    fontweight='light',
34                    fontsize='small')
35
36 map.set_xticklabels(map.get_xticklabels(), rotation=60,
37                    horizontalalignment='right',
38                    fontweight='light',
39                    fontsize='small')
40 plt.show(block=True)

```

6 Implementation of Classical Machine Learning Algorithms

A series of different machine learning algorithms are applied to evaluate performance on validation and test set.

- Logistic Regression
- K-Nearest Neighbours
- Decision Tree
- Random Forest
- Linear-SVM
- PCA with Linear-SVM

```

1 def fitmodels(dataset, train_label, test_label, val_label, model,
2               model_name='logistic_regression'):
3     """
4     This model implements each algorithm
5     """
6     Xtrain, Ytrain, Xval, Yval, Xtest, Ytest = dataset
7     test_accuracy = []
8     valid_accuracy = []
9     classes = ['DoS', 'Probe', 'R2L', 'U2R', 'normal']
10    precision = {'DoS': None, 'Probe': None,
11               'R2L': None, 'U2R': None, 'normal': None}
12    recall = {'DoS': None, 'Probe': None,
13             'R2L': None, 'U2R': None, 'normal': None}
14    f1score = {'DoS': None, 'Probe': None,
15              'R2L': None, 'U2R': None, 'normal': None}
16
17    for attack_type in range(5):

```

```

18     model.fit(Xtrain, train_label[attack_type])
19     y_pred = model.predict(Xtest)
20     stacked_train_df[attack_type][model_name] = model.predict(Xtrain)
21     stacked_test_df[attack_type][model_name] = y_pred
22     y_pred_val = model.predict(Xval)
23     stacked_val_df[attack_type][model_name] = y_pred_val
24     valid_accuracy += [metrics.accuracy_score(val_label[attack_type],
25                                               y_pred_val)]
26     test_accuracy += [metrics.accuracy_score(test_label[attack_type],
27                                             y_pred)]
28
29     precision[classess[attack_type]] = \
30         np.round(metrics.average_precision_score(test_label[attack_type],
31                                                 y_pred)*100, 2)
32
33     recall[classess[attack_type]] = \
34         np.round(metrics.recall_score(test_label[attack_type],
35                                     y_pred,
36                                     average='weighted')*100, 2)
37
38     f1score[classess[attack_type]] = \
39         np.round(metrics.f1_score(test_label[attack_type],
40                                 y_pred,
41                                 average='weighted')*100, 2)
42     cm = metrics.confusion_matrix(test_label[attack_type], y_pred)
43     plot_confusion_matrix(cm, labels=[classes[attack_type], 'rest'],
44                          title=extract.CM_DIR +
45                              model_name + '-{}.jpeg'.format(classes[attack_type]))
46
47     mean_test_accuracy = np.mean(test_accuracy)
48     mean_valid_accuracy = np.mean(valid_accuracy)
49     print("Mean Accuracy score "
50           "on Validation Set:\t\t{} %".format(np.round(mean_valid_accuracy*100,
51                                                       2)))
52     print("Mean Accuracy score "
53           "on Test Set:\t\t{} %".format(np.round(mean_test_accuracy*100,
54                                                  2)))
55     print("Average Precision score "
56           "on Test Set:\t\t{} ==> {} %".format(precision,
57                                                np.round(np.array([precision[k]
58                                                                for k in
59                                                                precision]).mean(),
60                                                         2)))
61     print("Average Recall score "
62           "on Test Set:\t\t{} ==> {} %".format(recall,
63                                                np.round(np.array([recall[k]
64                                                                for k in
65                                                                recall]).mean(),
66                                                         2)))
67     print("Average f1-score "
68           "on Test Set:\t\t\t{} ==> {} %".format(f1score,
69                                                np.round(np.array([f1score[k]
70                                                                for k in
71                                                                f1score]).mean(),
72                                                         2)))
73
74 def applyModel(dataset, selected_columns=columns):
75     """
76     Apply different model on dataset
77     """
78     Xtrain, Ytrain, Xval, Yval, Xtest, Ytest = dataset
79     Xtrain = splitdataset.normalisation(Xtrain[selected_columns])
80     Xtest = splitdataset.normalisation(Xtest[selected_columns])
81     Xval = splitdataset.normalisation(Xval[selected_columns])
82     scalar = MinMaxScaler()
83     scalar.fit(Xtrain)
84     Xtrain = scalar.transform(Xtrain)
85     Xtest = scalar.transform(Xtest)
86     Xval = scalar.transform(Xval)

```

```

87
88 dataset = Xtrain, Ytrain, Xval, Yval, Xtest, Ytest
89 train_label = pd.DataFrame()
90 test_label = pd.DataFrame()
91 val_label = pd.DataFrame()
92
93 # One hot encode the class column
94 for attack_type in range(5):
95     train_label[attack_type] = Ytrain['labels'].apply(lambda x: int(x ==
96     attack_type))
97     test_label[attack_type] = Ytest['labels'].apply(lambda x: int(x ==
98     attack_type))
99     val_label[attack_type] = Yval['labels'].apply(lambda x: int(x ==
100     attack_type))
101
102 print("\n***** LogisticRegression *****")
103 fitmodels(dataset, train_label, test_label, val_label,
104           model=LogisticRegression(C=0.2, solver='liblinear'))
105
106 print("\n***** SGDClassifier *****")
107 fitmodels(dataset, train_label, test_label, val_label,
108           model=SGDClassifier(), model_name='SGDClassifier')
109
110 print("\n***** KNearestNeighbours *****")
111 fitmodels(dataset, train_label, test_label, val_label,
112           model=KNeighborsClassifier(n_neighbors=30),
113           model_name='KNearestNeighbours')
114
115 print("\n***** DecisionTree *****")
116 fitmodels(dataset, train_label, test_label, val_label,
117           model=DecisionTreeClassifier(max_depth=8),
118           model_name='DecisionTree')
119
120 print("\n***** RandomForest *****")
121 fitmodels(dataset, train_label, test_label, val_label,
122           model=RandomForestClassifier(max_features='sqrt',
123           max_depth=8,
124           n_estimators=900),
125           model_name='RandomForest')
126
127 print("\n***** Support Vector Classifier
128 *****")
129 fitmodels(dataset, train_label, test_label, val_label,
130           model=LinearSVC(C=100),
131           model_name='SupportVectorClassifier')
132
133 print("\n***** Principal Component Analysis
134 *****")
135 pca = PCA(n_components=10)
136 pca.fit(Xtrain)
137 Xtrain = pca.transform(Xtrain)
138 Xtest = pca.transform(Xtest)
139 Xval = pca.transform(Xval)
140 print("Shape of PCA Components: {}".format(pca.components_.shape))
141 dataset = Xtrain, Ytrain, Xval, Yval, Xtest, Ytest
142 fitmodels(dataset, train_label, test_label, val_label,
143           model=LinearSVC(C=100),
144           model_name='PCAwithSVC')
145
146 print("\n***** XGBClassifier *****")
147 fitmodels(dataset, train_label, test_label, val_label,
148           model=XGBClassifier(max_depth=6, learning_rate=0.01, n_estimators
149           =300, nthread=4,
150           min_child_weight=4, subsample=0.7,
151           colsample_bytree=0.7),
152           model_name='XGBClassifier')
153
154 print("\n***** StackedAccuracy *****")
155 stacked_accuracy = []
156
157 for attack_type in range(5):
158     clf = KNeighborsClassifier(n_neighbors=30)
159     clf.fit(stacked_train_df[attack_type], train_label[attack_type])

```

```
153     y_pred = clf.predict(stacked_test_df[attack_type])
154     stacked_accuracy += [metrics.accuracy_score(test_label[attack_type],
155                                                y_pred)]
156
157 mean_stacked_accuracy = np.mean(stacked_accuracy)
158 print("Mean Accuracy score on Test Set: {}".format(mean_stacked_accuracy))
```