

ADFA-LD - Logistic Regression

In [1]:

```
1 import seaborn as sns
2 import matplotlib.pyplot as plt
3 import pandas as pd
4 from IPython.display import display
5 pd.options.display.max_columns = None
6 from sklearn.ensemble import RandomForestRegressor, RandomForestClassifier
7 from IPython.display import display
8 from sklearn import metrics
9 from sklearn.model_selection import train_test_split
10 import statistics
11 import numpy as np
12 from sklearn import metrics
13 from sklearn.preprocessing import MinMaxScaler, StandardScaler, LabelEncoder
14 from sklearn.feature_selection import SelectKBest
15 from sklearn.pipeline import Pipeline
16 from sklearn.model_selection import train_test_split, GridSearchCV, RandomizedSearchCV
```

In [2]:

```

1  import glob
2  import math
3  from collections import Counter
4  import csv
5
6  import numpy as np
7
8
9  def plot_confusion_matrix(cm,
10                          target_names,
11                          title='Confusion matrix',
12                          cmap=None,
13                          normalize=True):
14      import matplotlib.pyplot as plt
15      import numpy as np
16      import itertools
17
18      accuracy = np.trace(cm) / float(np.sum(cm))
19      misclass = 1 - accuracy
20
21      if cmap is None:
22          cmap = plt.get_cmap('Blues')
23
24      plt.figure(figsize=(8, 6))
25      plt.imshow(cm, interpolation='nearest', cmap=cmap)
26      plt.title(title)
27      plt.colorbar()
28
29      if target_names is not None:
30          tick_marks = np.arange(len(target_names))
31          plt.xticks(tick_marks, target_names, rotation=45)
32          plt.yticks(tick_marks, target_names)
33
34      if normalize:
35          cm = cm.astype('float') / cm.sum(axis=1)[:, np.newaxis]
36
37
38      thresh = cm.max() / 1.5 if normalize else cm.max() / 2
39      for i, j in itertools.product(range(cm.shape[0]), range(cm.shape[1])):
40          if normalize:
41              plt.text(j, i, "{:0.4f}".format(cm[i, j]),
42                      horizontalalignment="center",
43                      color="white" if cm[i, j] > thresh else "black")
44          else:
45              plt.text(j, i, "{:,}".format(cm[i, j]),
46                      horizontalalignment="center",
47                      color="white" if cm[i, j] > thresh else "black")
48      plt.tight_layout()
49      plt.ylabel('True label')
50      plt.xlabel('Predicted label\naccuracy={:0.4f}; misclass={:0.4f}'.format(accuracy,
51      plt.show()
52
53  # returns a dictionary of n-grams frequency for any list
54  def ngrams_freq(listname, n):
55      counts = dict()
56      # make n-grams as string iteratively
57      grams = [' '.join(listname[i:i+n]) for i in range(len(listname)-n)]
58      for gram in grams:
59          if gram not in counts:

```

```

60         counts[gram] = 1
61     else:
62         counts[gram] += 1
63     return counts
64
65 # returns the values of features for any list
66 def feature_freq(listname,n,features):
67     counts = dict()
68     # make n-grams as string iteratively
69     grams = [' '.join(listname[i:i+n]) for i in range(len(listname)-n)]
70     for gram in grams:
71         counts[gram] = 0
72     for gram in grams:
73         if gram in features:
74             counts[gram] += 1
75     return counts
76
77 # values of n for finding n-grams
78 n_values = [1]
79
80 # Base address for attack data files
81 add = "ADFA-LD/ADFA-LD/Attack_Data_Master/"
82 # list of attacks
83 attack = ['Adduser','Hydra_FTP','Hydra_SSH','Java_Meterpreter','Meterpreter','Web_Shell']
84
85 # initializing dictionary for n-grams from all files
86 tra indict = {}
87
88 Attack_list_new = []
89 print("Generating Training Data .....")
90 for term in attack:
91     print(" Training data from " + term)
92     globals()['%s_list' % term] = []
93     in_address = add+term
94     k = 1
95     # finding list of data from all files
96     for i in range (1,11):
97         read_files = glob.glob(in_address+"_"+str(i)+"/*.txt")
98         for f in read_files:
99             with open(f, "r") as infile:
100                 globals()['%s_list_array' % term+str(k)] = ALine =infile.read()
101                 #ALine = ALine[:820]
102                 Attack_list_new.append(term + ',' + str(ALine))
103                 globals()['%s_list' % term].extend(globals()['%s_list_array' % term+str(k)])
104                 k += 1
105     # number of lists for distinct files
106     globals()['%s_size' % term] = k-1
107     # combined list of all files
108     listname = globals()['%s_list' % term]
109     # finding n-grams and extracting top 30%
110     for n in n_values:
111         #print("      Extracting top 30% "+str(n)+"-grams from "+term+".....")
112         dictname = ngrams_freq(listname,n)
113         top = math.ceil(0.3*len(dictname))
114         dictname = Counter(dictname)
115         for k, v in dictname.most_common(top):
116             tra indict.update({k : v})
117
118 # finding training data for Normal file
119 print(" Training data from Normal")
120 Normal_list = []

```

```

121 Normal_list_new = []
122 in_address = "ADFA-LD/ADFA-LD/Training_Data_Master/"
123 k = 1
124 read_files = glob.glob(in_address+"/*.txt")
125 for f in read_files:
126     with open(f, "r") as infile:
127         globals()['Normal%s_list_array' % str(k)] = Line = infile.read()
128         Normal_list_new.append('Normal,'+ str(Line))
129         Normal_list.extend(globals()['Normal%s_list_array' % str(k)])
130         k += 1
131
132 # number of lists for distinct files
133 Normal_list_size = k-1
134 # combined list of all files
135 listname = Normal_list
136
137
138 print("\nnew_train.csv created.....\n")
139

```

Generating Training Data

```

    Training data from Adduser
    Training data from Hydra_FTP
    Training data from Hydra_SSH
    Training data from Java_Meterpreter
    Training data from Meterpreter
    Training data from Web_Shell
    Training data from Normal

```

new_train.csv created.....

In [3]:

```

1 new_train_list = []
2 new_train_list = Normal_list_new + Attack_list_new
3 #new_train_list[1]
4 #Attack_list_new[1]
5

```

In [4]:

```

1 new_train_list = []
2 new_train_list = Normal_list_new + Attack_list_new
3
4
5 with open('new_train.csv', 'w') as f:
6     for item in new_train_list:
7         f.write("%s\n" % item)

```

In [5]:

```

1 train = pd.read_csv("./new_train.csv", sep=',', error_bad_lines=False, header=None, name
2 train.head(5)
3 train.shape
4 #train.info()
5
6 #train.describe(include = 'all')
7 train_df = train.copy()
8 train['Label'] = train['Label'].astype('category')
9 train['CallTrace'] = train['CallTrace'].astype('category')
10
11 train['Label'].value_counts()
12 #train['CallTrace'].value_counts()

```

Out[5]:

```

Normal      833
Hydra_SSH   176
Hydra_FTP   162
Java_Meterpreter 124
Web_Shell   118
Adduser     91
Meterpreter 75
Name: Label, dtype: int64

```

In [6]:

```

1 train['Label_Codes'] = train['Label'].cat.codes
2 train['CallTrace_Codes'] = train['CallTrace'].cat.codes
3 train['Label_Codes'].value_counts()

```

Out[6]:

```

5      833
2      176
1      162
3      124
6      118
0       91
4       75
Name: Label_Codes, dtype: int64

```

In [7]:

```

1 train.head()

```

Out[7]:

	Label	CallTrace	Label_Codes	CallTrace_Codes
0	Normal	6 6 63 6 42 120 6 195 120 6 6 114 114 1 1 252 ...	5	1407
1	Normal	54 175 120 175 175 3 175 175 120 175 120 175 1...	5	1239
2	Normal	6 11 45 33 192 33 5 197 192 6 33 5 3 197 192 1...	5	1286
3	Normal	7 174 174 5 197 197 6 13 195 4 4 118 6 91 38 5...	5	1465
4	Normal	11 45 33 192 33 5 197 192 6 33 5 3 197 192 192...	5	93

Multinomial Logistic Regression

In [8]:

```
1 import warnings
2 warnings.filterwarnings("ignore")
3
4 # split the dataset in train and test
5 X = train.iloc[:, [3]].values
6 y = train.iloc[:, 2].values
7
8
9 # Splitting the dataset into the Training set and Test set
10 from sklearn.model_selection import train_test_split
11 X_train, X_test, y_train, y_test = train_test_split(X, y, test_size = 0.2, random_state=0)
12
13 # Feature Scaling
14 from sklearn.preprocessing import StandardScaler
15 sc = StandardScaler()
16 X_train = sc.fit_transform(X_train)
17 X_test = sc.transform(X_test)
18
19 # Fitting Logistic Regression to the Training set
20 from sklearn.linear_model import LogisticRegression
21 classifier = LogisticRegression(multi_class='ovr', solver = 'lbfgs')
22 classifier.fit(X_train, y_train)
23
24 # Predicting the Test set results
25 y_pred = classifier.predict(X_test)
26
27 # How did our model perform?
28 from sklearn import metrics
29 count_misclassified = (y_test != y_pred).sum()
30 print('Misclassified samples: {}'.format(count_misclassified))
31 accuracy = metrics.accuracy_score(y_test, y_pred)
32 print('Accuracy: {:.2f}'.format(accuracy))
33
34
35
```

Misclassified samples: 145

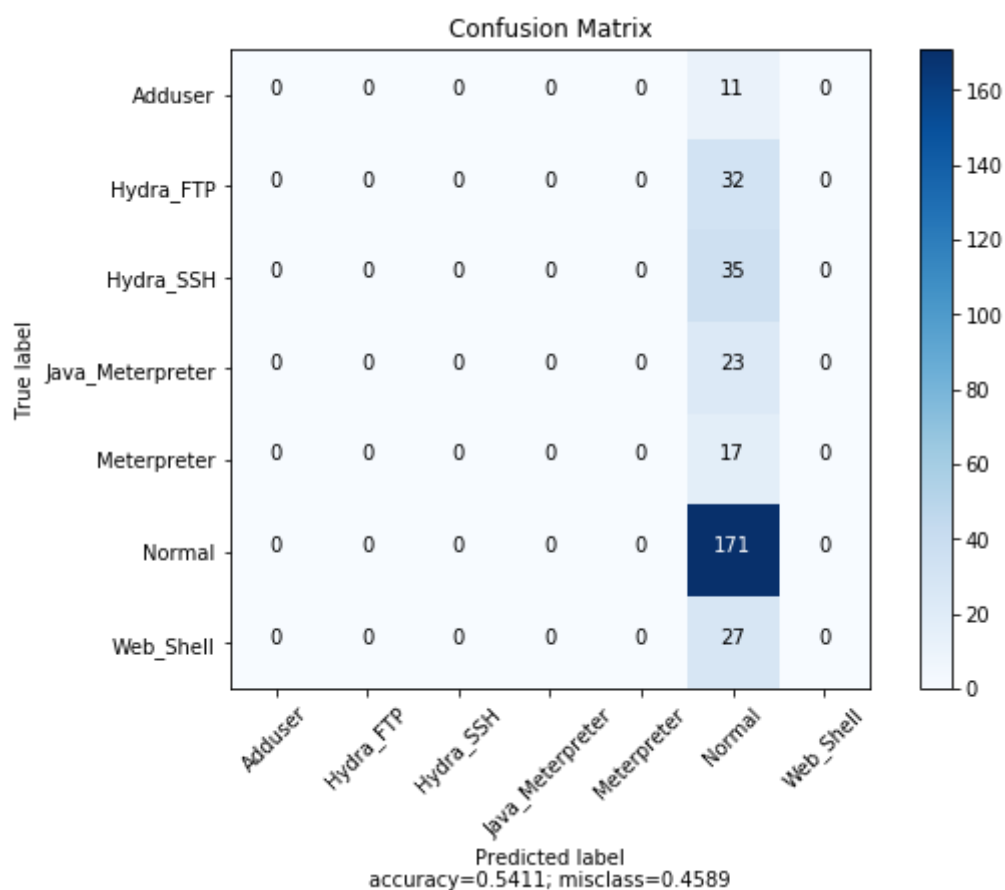
Accuracy: 0.54

In [9]:

```

1 #classifier.predict_proba(X_test)
2
3 # Making the Confusion Matrix
4 from sklearn.metrics import confusion_matrix
5 cm = confusion_matrix(y_test, y_pred)
6 plot_confusion_matrix(cm,
7                        normalize      = False,
8                        target_names  = ['Adduser', 'Hydra_FTP', 'Hydra_SSH', 'Java_Meterpreter', 'Meterpreter', 'Normal', 'Web_Shell'],
9                        title         = "Confusion Matrix")

```



Logistic Regression Binary Classification

In [10]:

```
1 train.loc[train.Label != 'Normal','Label_Binary']= 1
2 train.loc[train.Label == 'Normal','Label_Binary']= 0
3 train['Label_Binary'].value_counts()
4 #train.head()
```

Out[10]:

0.0 833
1.0 746
Name: Label_Binary, dtype: int64

In [11]:

```
1 train.head()
```

Out[11]:

	Label	CallTrace	Label_Codes	CallTrace_Codes	Label_Binary
0	Normal	6 6 63 6 42 120 6 195 120 6 6 114 114 1 1 252 ...	5	1407	0.0
1	Normal	54 175 120 175 175 3 175 175 120 175 120 175 1...	5	1239	0.0
2	Normal	6 11 45 33 192 33 5 197 192 6 33 5 3 197 192 1...	5	1286	0.0
3	Normal	7 174 174 5 197 197 6 13 195 4 4 118 6 91 38 5...	5	1465	0.0
4	Normal	11 45 33 192 33 5 197 192 6 33 5 3 197 192 192...	5	93	0.0

In [12]:

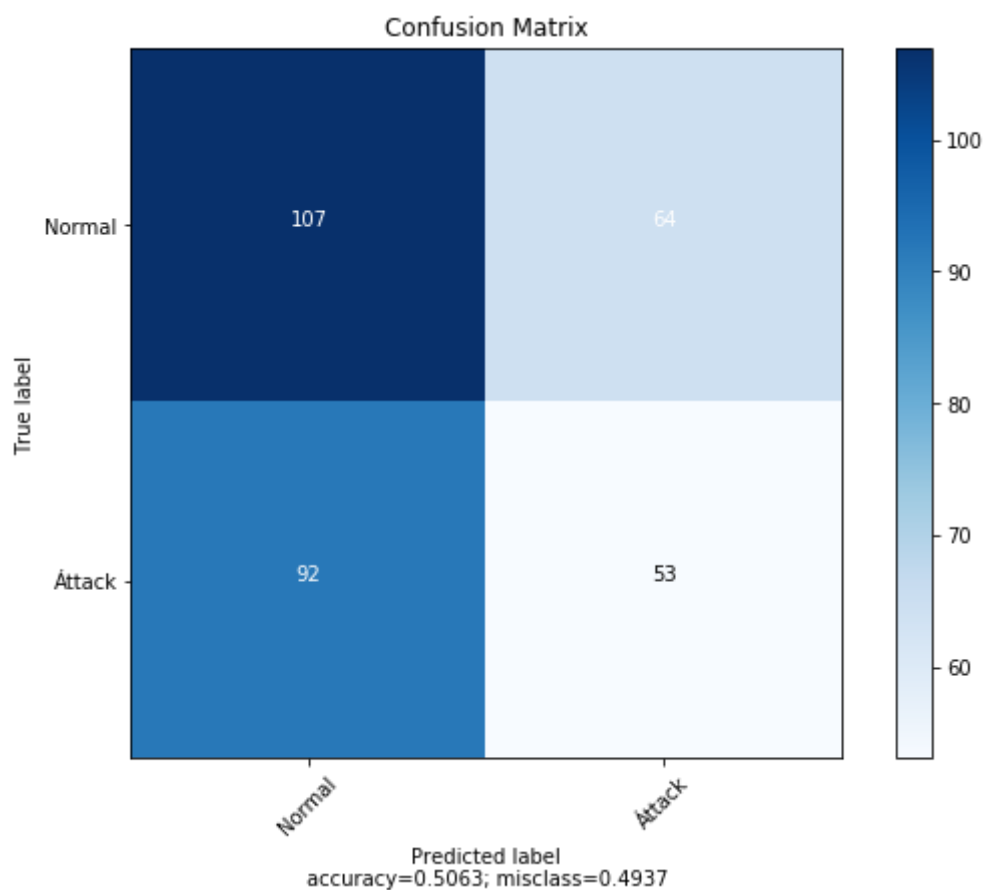
```
1 import warnings
2 warnings.filterwarnings("ignore")
3
4 # split the dataset in train and test
5 X = train.iloc[:, [3]].values
6 y = train.iloc[:, 4].values
7
8
9 # Splitting the dataset into the Training set and Test set
10 from sklearn.model_selection import train_test_split
11 X_train, X_test, y_train, y_test = train_test_split(X, y, test_size = 0.2, random_state=0)
12
13 # Feature Scaling
14 from sklearn.preprocessing import StandardScaler
15 sc = StandardScaler()
16 X_train = sc.fit_transform(X_train)
17 X_test = sc.transform(X_test)
18
19 # Fitting Logistic Regression to the Training set
20 from sklearn.linear_model import LogisticRegression
21 #classifier = LogisticRegression(multi_class='ovr', solver = 'lbfgs')
22 classifier = LogisticRegression()
23 classifier.fit(X_train, y_train)
24
25 # Predicting the Test set results
26 y_pred = classifier.predict(X_test)
27
28 # How did our model perform?
29 from sklearn import metrics
30 count_misclassified = (y_test != y_pred).sum()
31 print('Misclassified samples: {}'.format(count_misclassified))
32 accuracy = metrics.accuracy_score(y_test, y_pred)
33 print('Accuracy: {:.2f}'.format(accuracy))
34
35
36
```

Misclassified samples: 156

Accuracy: 0.51

In [13]:

```
1 # Making the Confusion Matrix
2 from sklearn.metrics import confusion_matrix
3 cm = confusion_matrix(y_test, y_pred)
4 plot_confusion_matrix(cm,
5                       normalize = False,
6                       target_names = ['Normal', 'Attack'],
7                       title = "Confusion Matrix")
```



In [14]:

```
1 print(metrics.classification_report(y_pred, y_test))
```

	precision	recall	f1-score	support
0.0	0.63	0.54	0.58	199
1.0	0.37	0.45	0.40	117
micro avg	0.51	0.51	0.51	316
macro avg	0.50	0.50	0.49	316
weighted avg	0.53	0.51	0.51	316

OneHotEncoding for LogisticRegression

In [15]:

```
1 # Split into predictor and response dataframes.
2 train_df_enc = train_df.copy()
3 X_df = train_df_enc.drop('Label', axis=1)
4 y = train_df_enc['Label']
5
6 X_df.shape,y.shape
```

Out[15]:

((1579, 1), (1579,))

In [16]:

```
1 X_df.head()
```

Out[16]:

CallTrace														
0	6	6	63	6	42	120	6	195	120	6	6	114	114	1 1 252 ...
1	54	175	120	175	175	3	175	175	120	175	120	175	1...	
2	6	11	45	33	192	33	5	197	192	6	33	5	3	197 192 1...
3	7	174	174	5	197	197	6	13	195	4	4	118	6	91 38 5...
4	11	45	33	192	33	5	197	192	6	33	5	3	197	192 192...

In [17]:

```
1 train_df.head()
```

Out[17]:

	Label	CallTrace
0	Normal	6 6 63 6 42 120 6 195 120 6 6 114 114 1 1 252 ...
1	Normal	54 175 120 175 175 3 175 175 120 175 120 175 1...
2	Normal	6 11 45 33 192 33 5 197 192 6 33 5 3 197 192 1...
3	Normal	7 174 174 5 197 197 6 13 195 4 4 118 6 91 38 5...
4	Normal	11 45 33 192 33 5 197 192 6 33 5 3 197 192 192...

In [18]:

```
1 # Map response variable to integers 0,1.
2 y = pd.Series(np.where(y.values != 'Normal',1,0), y.index)
3 y.value_counts()
```

Out[18]:

```
0    833
1    746
dtype: int64
```

In [19]:

```
1 # Label Encode instead of dummy variables
2
3 mappings = []
4
5 from sklearn.preprocessing import LabelEncoder
6
7 label_encoder = LabelEncoder()
8
9 label_df = train.drop('Label', axis=1)
10 label_df = train.drop('Label_Binary', axis=1)
11 label_df = train.drop('Label_Codes', axis=1)
12 label_df['CallTrace'] = label_df['CallTrace_Codes']
13 label_df = X_df.copy()
14 for i, col in enumerate(label_df):
15     if label_df[col].dtype == 'object':
16         label_df[col] = label_encoder.fit_transform(np.array(label_df[col].astype(str)))
17         mappings.append(dict(zip(label_encoder.classes_, range(1, len(label_encoder.cl
```

In [20]:

```
1 label_df.head()
```

Out[20]:

	CallTrace
0	1407
1	1239
2	1286
3	1465
4	93

In [21]:

```
1 from sklearn.preprocessing import OneHotEncoder
2
3
4 onehot_encoder = OneHotEncoder()
5 for i, col in enumerate(label_df):
6     if label_df[col].dtype == 'object':
7         label_df[col] = onehot_encoder.fit_transform(np.array(label_df[col].astype(str)
8         mappings.append(dict(zip(onehot_encoder.classes_, range(1, len(onehot_encoder.c
```

In [22]:

```
1 X_train, X_test, y_train, y_test = train_test_split(label_df, y, test_size = 0.2, rand
2 X_train.shape, X_test.shape, y_train.shape, y_test.shape
```

Out[22]:

```
((1263, 1), (316, 1), (1263,), (316,))
```

In [23]:

```
1 clf = LogisticRegression()
2 model_mix = clf.fit(X_train, y_train)
3 # y_pred = model_norm.predict(X_test)
4 print("Model accuracy is", model_mix.score(X_test, y_test))
```

```
Model accuracy is 0.5569620253164557
```

In [24]:

```
1 model_mix
```

Out[24]:

```
LogisticRegression(C=1.0, class_weight=None, dual=False, fit_intercept=True,
    intercept_scaling=1, max_iter=100, multi_class='warn',
    n_jobs=None, penalty='l2', random_state=None, solver='warn',
    tol=0.0001, verbose=0, warm_start=False)
```

In [25]:

```

1 # logit_roc_auc = roc_auc_score(y_test, model_norm.predict(X_test))
2 # fpr, tpr, thresholds = roc_curve(y_test, model_norm.predict_proba(X_test)[:,1])
3
4 classes = model_mix.predict(X_test)
5 probs = model_mix.predict_proba(X_test)
6 preds = probs[:,1]
7 #preds

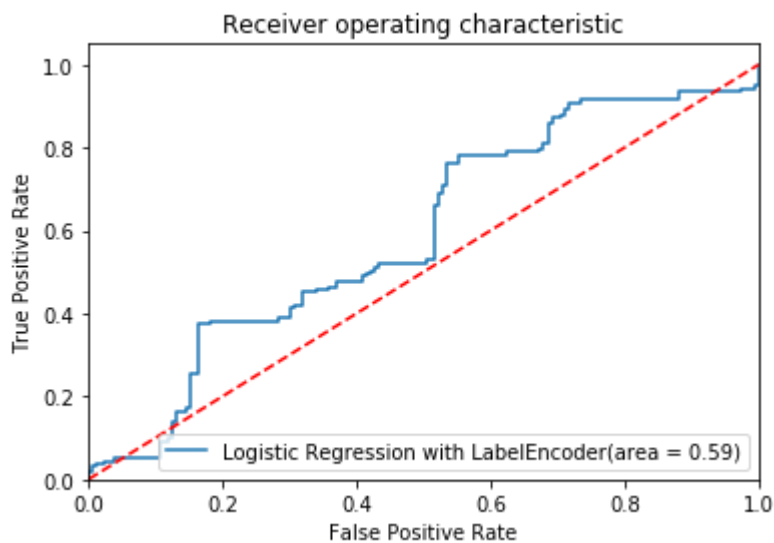
```

In [26]:

```

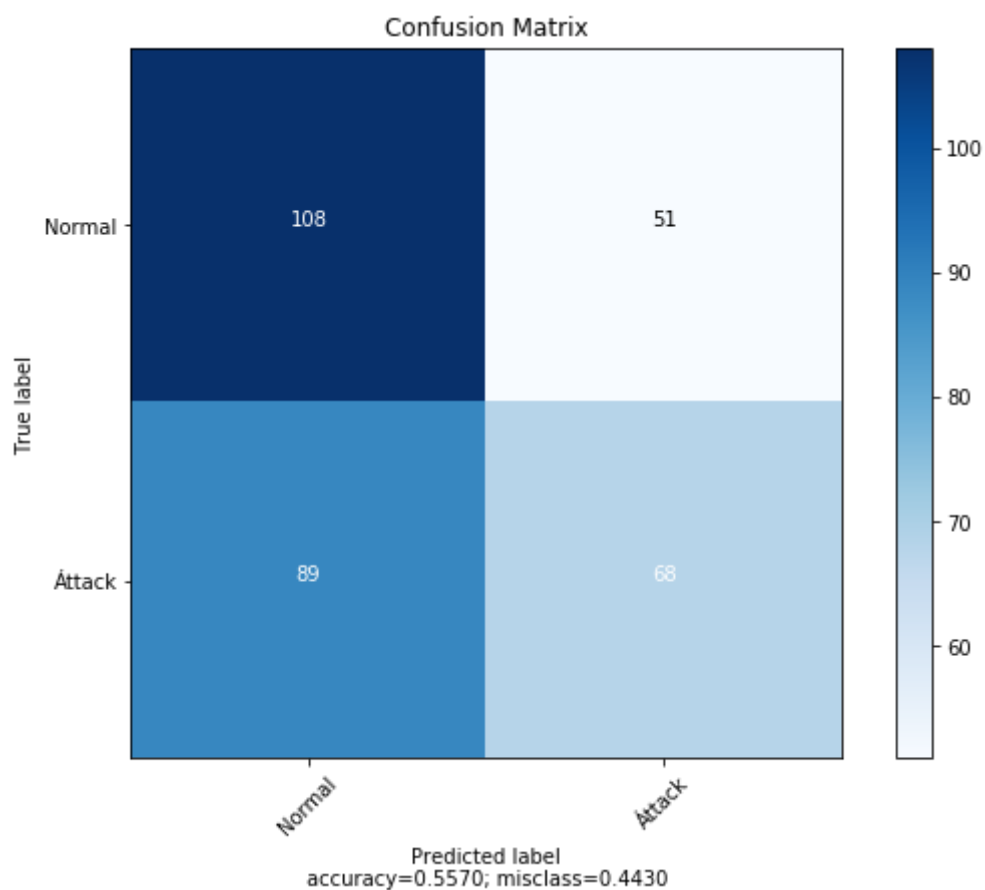
1 labelfpr, labeltpr, labelthreshold = metrics.roc_curve(y_test, preds)
2 label_roc_auc = metrics.auc(labelfpr, labeltpr)
3
4 plt.figure()
5 plt.plot(labelfpr, labeltpr, label='Logistic Regression with LabelEncoder(area = %0.2f'
6 plt.plot([0, 1], [0, 1], 'r--')
7 plt.xlim([0.0, 1.0])
8 plt.ylim([0.0, 1.05])
9 plt.xlabel('False Positive Rate')
10 plt.ylabel('True Positive Rate')
11 plt.title('Receiver operating characteristic')
12 plt.legend(loc="lower right")
13 plt.savefig('Log_ROC')
14 plt.show()

```



In [27]:

```
1 # Making the Confusion Matrix
2 from sklearn.metrics import confusion_matrix
3 cm = confusion_matrix(y_test, classes)
4 plot_confusion_matrix(cm,
5                       normalize = False,
6                       target_names = ['Normal', 'Attack'],
7                       title = "Confusion Matrix")
```



In [28]:

```
1 X_train.shape, X_test.shape, y_train.shape, y_test.shape
```

Out[28]:

```
((1263, 1), (316, 1), (1263,), (316,))
```

In [29]:

```
1 print(metrics.classification_report(classes, y_test))
```

	precision	recall	f1-score	support
0	0.68	0.55	0.61	197
1	0.43	0.57	0.49	119
micro avg	0.56	0.56	0.56	316
macro avg	0.56	0.56	0.55	316
weighted avg	0.59	0.56	0.56	316

RandomForest Classification

In [30]:

```
1 # Normalize using MinMaxScaler to constrain values to between 0 and 1.
2 from sklearn.preprocessing import MinMaxScaler, StandardScaler
3
4 scaler = MinMaxScaler(feature_range = (0,1))
5
6 scaler.fit(X_train)
7 X_train = scaler.transform(X_train)
8 X_test = scaler.transform(X_test)
```

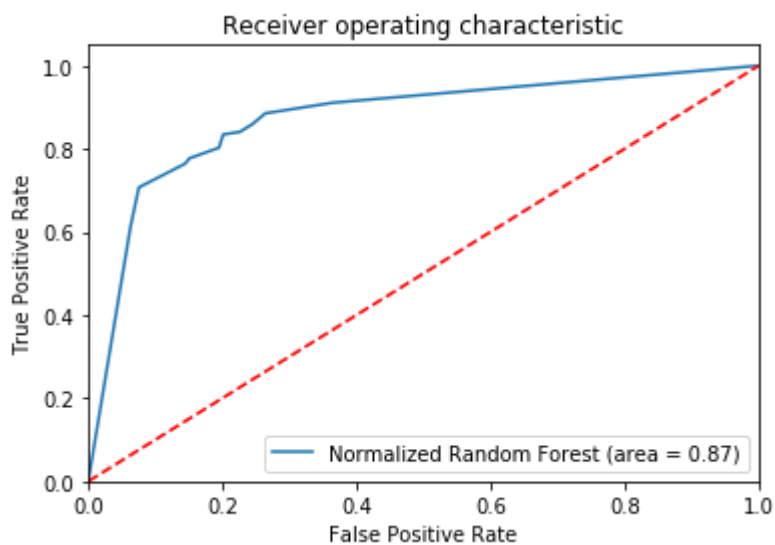
In [31]:

```
1 clf = RandomForestClassifier(n_jobs=-1)
2 model_rf = clf.fit(X_train, y_train)
3 print('Model accuracy is', model_rf.score(X_test, y_test))
```

Model accuracy is 0.8037974683544303

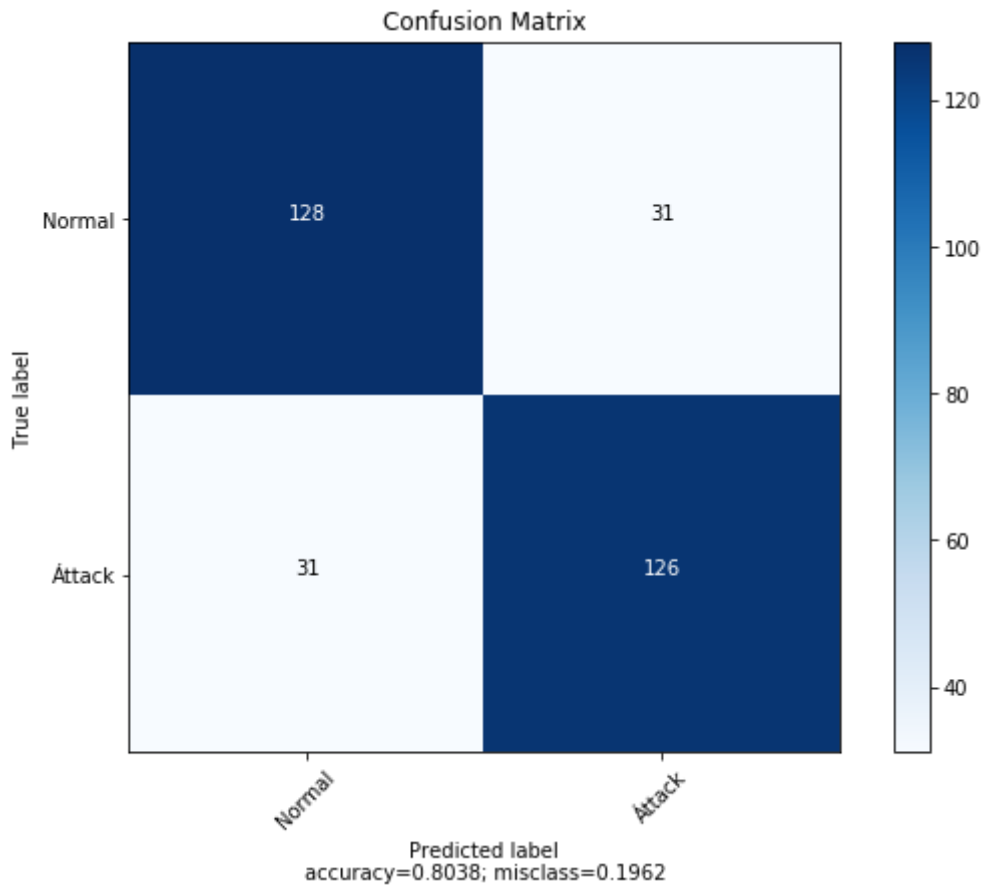
In [32]:

```
1 probs = model_rf.predict_proba(X_test)
2 preds = probs[:,1]
3 rffpr, rftpr, rfthreshold = metrics.roc_curve(y_test, preds)
4 rf_roc_auc = metrics.auc(rffpr, rftpr)
5
6 plt.figure()
7 plt.plot(rffpr, rftpr, label='Normalized Random Forest (area = %0.2f)' % rf_roc_auc)
8 plt.plot([0, 1], [0, 1], 'r--')
9 plt.xlim([0.0, 1.0])
10 plt.ylim([0.0, 1.05])
11 plt.xlabel('False Positive Rate')
12 plt.ylabel('True Positive Rate')
13 plt.title('Receiver operating characteristic')
14 plt.legend(loc="lower right")
15 plt.savefig('Log_ROC')
16 plt.show()
```



In [33]:

```
1 # Making the Confusion Matrix
2 from sklearn.metrics import confusion_matrix
3 classes = model_rf.predict(X_test)
4 cm = confusion_matrix(y_test, classes)
5 plot_confusion_matrix(cm,
6                       normalize = False,
7                       target_names = ['Normal', 'Attack'],
8                       title = "Confusion Matrix")
```



In [34]:

```
1 print(metrics.classification_report(classes, y_test))
```

	precision	recall	f1-score	support
0	0.81	0.81	0.81	159
1	0.80	0.80	0.80	157
micro avg	0.80	0.80	0.80	316
macro avg	0.80	0.80	0.80	316
weighted avg	0.80	0.80	0.80	316

Train Data with ngrams

In [35]:

```
1 from sklearn.datasets import make_classification
2 from sklearn.model_selection import train_test_split
3 from sklearn.linear_model import LogisticRegression
4 from sklearn.metrics import accuracy_score, confusion_matrix, recall_score, roc_auc_score
5
6 X, y = make_classification(
7     n_classes=2, class_sep=1.5, weights=[0.1, 0.9],
8     n_features=20, n_samples=1000, random_state=10
9 )
10
11 #X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.33, random_state=10)
12
13 clf = LogisticRegression(class_weight="balanced")
14 clf.fit(X_train, y_train)
15 THRESHOLD = 0.5
16 preds = np.where(clf.predict_proba(X_test)[:,:1] > THRESHOLD, 1, 0)
17
18 pd.DataFrame(data=[accuracy_score(y_test, preds), recall_score(y_test, preds),
19                     precision_score(y_test, preds), roc_auc_score(y_test, preds)],
20             index=["accuracy", "recall", "precision", "roc_auc_score"])
```

Out[35]:

	0
accuracy	0.531646
recall	0.579618
precision	0.526012
roc_auc_score	0.531947

In [36]:

```

1  from sklearn import model_selection, preprocessing, linear_model, naive_bayes, metrics
2  from sklearn.feature_extraction.text import TfidfVectorizer, CountVectorizer
3  from sklearn import decomposition, ensemble
4
5  import pandas, xgboost, numpy, textblob, string
6  from keras.preprocessing import text, sequence
7  from keras import layers, models, optimizers
8
9  def train_model(classifier, feature_vector_train, label, feature_vector_valid, is_neural_net):
10     # fit the training dataset on the classifier
11     classifier.fit(feature_vector_train, label)
12
13     # predict the labels on validation dataset
14     predictions = classifier.predict(feature_vector_valid)
15
16     if is_neural_net:
17         predictions = predictions.argmax(axis=-1)
18
19     return metrics.accuracy_score(predictions, valid_y)
20
21 # Load the dataset
22 #data = open('data/corpus').read()
23 #labels, texts = [], []
24 #for i, line in enumerate(data.split("\n")):
25 #     content = line.split()
26 #     labels.append(content[0])
27 #     texts.append(" ".join(content[1:]))
28
29 # create a dataframe using texts and labels
30 #trainDF = pandas.DataFrame()
31 #trainDF['text'] = texts
32 #trainDF['label'] = labels

```

Using TensorFlow backend.

In [37]:

```
1 X_df.head()
```

Out[37]:

CallTrace

0	6 6 63 6 42 120 6 195 120 6 6 114 114 1 1 252 ...
1	54 175 120 175 175 3 175 175 120 175 120 175 1...
2	6 11 45 33 192 33 5 197 192 6 33 5 3 197 192 1...
3	7 174 174 5 197 197 6 13 195 4 4 118 6 91 38 5...
4	11 45 33 192 33 5 197 192 6 33 5 3 197 192 192...

In [38]:

```

1  # create a dataframe using texts and lables
2  trainDF = train_df.copy()
3
4  trainDF['CallTrace_T'] = trainDF.CallTrace.str.split(' ').str.join(',').astype(str)
5  #X_df = trainDF.drop('Label', axis=1)
6  X_df = trainDF.drop(['Label', 'CallTrace'], axis=1)
7  y = trainDF['Label']
8
9  # split the dataset into training and validation datasets
10 train_x, valid_x, train_y, valid_y = model_selection.train_test_split(X_df, y)
11
12 # label encode the target variable
13 encoder = preprocessing.LabelEncoder()
14 train_y = encoder.fit_transform(train_y)
15 valid_y = encoder.fit_transform(valid_y)
16
17 X_df.head()

```

Out[38]:

	CallTrace_T
0	6,6,63,6,42,120,6,195,120,6,6,114,114,1,1,252,...
1	54,175,120,175,175,3,175,175,120,175,120,175,1...
2	6,11,45,33,192,33,5,197,192,6,33,5,3,197,192,1...
3	7,174,174,5,197,197,6,13,195,4,4,118,6,91,38,5...
4	11,45,33,192,33,5,197,192,6,33,5,3,197,192,192...

In [39]:

```
1 train_x.shape, valid_x.shape, train_y.shape, valid_y.shape
```

Out[39]:

```
((1184, 1), (395, 1), (1184,), (395,))
```

Feature Engineering

2.1 Count Vectors as features

2.2 TF-IDF Vectors as features

Word level

N-Gram level

Character level

2.3 Word Embeddings as features

2.4 Text / NLP based features

2.5 Topic Models as features

```

1  # create a count vectorizer object
2  count_vect = CountVectorizer(analyzer='word', token_pattern=r'\w{1,}')
```

```

3  count_vect.fit(trainDF['CallTrace_T'])
4
5  # transform the training and validation data using count vectorizer object
6  xtrain_count = count_vect.transform(train_x)
7  xvalid_count = count_vect.transform(valid_x)
8
9  print(xtrain_count.toarray())
10 print(count_vect.vocabulary_)
11 print(xtrain_count.shape)
12 #xt_count =pd.DataFrame(xtrain_count.A, columns=count_vect.get_feature_names())
13 #xt_count.shape
14 #print(pd.DataFrame(xtrain_count.A, columns=count_vect.get_feature_names()).to_string(
```

(1, 153)

In [41]:

```

1 # Linear Classifier on Count Vectors
2 accuracy = train_model(linear_model.LogisticRegression(), xtrain_count, train_y, xvalid_count)
3 print("LR, Count Vectors: ", accuracy)

```

ValueError

Traceback (most recent call last)

<ipython-input-41-3924f370ab03> in <module>

```

1 # Linear Classifier on Count Vectors
----> 2 accuracy = train_model(linear_model.LogisticRegression(), xtrain_count, train_y, xvalid_count)
3 print("LR, Count Vectors: ", accuracy)

```

```

<ipython-input-36-c766096d1c1b> in train_model(classifier, feature_vector_train, label, feature_vector_valid, is_neural_net)
9 def train_model(classifier, feature_vector_train, label, feature_vector_valid, is_neural_net=False):
10     # fit the training dataset on the classifier
--> 11     classifier.fit(feature_vector_train, label)
12
13     # predict the labels on validation dataset

```

```

~\AppData\Local\Continuum\anaconda3\lib\site-packages\sklearn\linear_model\logistic.py in fit(self, X, y, sample_weight)
1283
1284     X, y = check_X_y(X, y, accept_sparse='csr', dtype=_dtype, order="C",
-> 1285                     accept_large_sparse=solver != 'liblinear')
1286     check_classification_targets(y)
1287     self.classes_ = np.unique(y)

```

```

~\AppData\Local\Continuum\anaconda3\lib\site-packages\sklearn\utils\validation.py in check_X_y(X, y, accept_sparse, accept_large_sparse, dtype, order, copy, force_all_finite, ensure_2d, allow_nd, multi_output, ensure_min_samples, ensure_min_features, y_numeric, warn_on_dtype, estimator)
764     y = y.astype(np.float64)
765
--> 766     check_consistent_length(X, y)
767
768     return X, y

```

```

~\AppData\Local\Continuum\anaconda3\lib\site-packages\sklearn\utils\validation.py in check_consistent_length(*arrays)
233     if len(uniques) > 1:
234         raise ValueError("Found input variables with inconsistent numbers of"
-> 235                          " samples: %r" % [int(1) for 1 in lengths])
236
237

```

```

~\AppData\Local\Continuum\anaconda3\lib\site-packages\sklearn\utils\validation.py in check_consistent_length(*arrays)
233     if len(uniques) > 1:
234         raise ValueError("Found input variables with inconsistent numbers of"
-> 235                          " samples: %r" % [int(1) for 1 in lengths])
236
237

```

```

~\AppData\Local\Continuum\anaconda3\lib\site-packages\sklearn\utils\validation.py in check_X_y(X, y, accept_sparse, accept_large_sparse, dtype, order, copy, force_all_finite, ensure_2d, allow_nd, multi_output, ensure_min_samples, ensure_min_features, y_numeric, warn_on_dtype, estimator)
764     y = y.astype(np.float64)
765
--> 766     check_consistent_length(X, y)
767
768     return X, y

```

```

~\AppData\Local\Continuum\anaconda3\lib\site-packages\sklearn\utils\validation.py in check_X_y(X, y, accept_sparse, accept_large_sparse, dtype, order, copy, force_all_finite, ensure_2d, allow_nd, multi_output, ensure_min_samples, ensure_min_features, y_numeric, warn_on_dtype, estimator)
764     y = y.astype(np.float64)
765
--> 766     check_consistent_length(X, y)
767
768     return X, y

```

```

~\AppData\Local\Continuum\anaconda3\lib\site-packages\sklearn\utils\validation.py in check_X_y(X, y, accept_sparse, accept_large_sparse, dtype, order, copy, force_all_finite, ensure_2d, allow_nd, multi_output, ensure_min_samples, ensure_min_features, y_numeric, warn_on_dtype, estimator)
764     y = y.astype(np.float64)
765
--> 766     check_consistent_length(X, y)
767
768     return X, y

```

```

~\AppData\Local\Continuum\anaconda3\lib\site-packages\sklearn\utils\validation.py in check_X_y(X, y, accept_sparse, accept_large_sparse, dtype, order, copy, force_all_finite, ensure_2d, allow_nd, multi_output, ensure_min_samples, ensure_min_features, y_numeric, warn_on_dtype, estimator)
764     y = y.astype(np.float64)
765
--> 766     check_consistent_length(X, y)
767
768     return X, y

```

```

ValueError: Found input variables with inconsistent numbers of samples: [1, 1184]

```

In []:

```
1 # word level tf-idf
2 #tfidf_vect = TfidfVectorizer(analyzer='word', token_pattern=r'\w{1,}', max_features=5000)
3 #tfidf_vect.fit(trainDF['CallTrace'])
4 #xtrain_tfidf = tfidf_vect.transform(train_x)
5 #xvalid_tfidf = tfidf_vect.transform(valid_x)
6 #print(xtrain_tfidf.toarray())
7 #print(tfidf_vect.vocabulary_)
8 #print(xtrain_tfidf.shape)
```

In []:

```
1 # Linear Classifier on Count Vectors
2 #accuracy = train_model(linear_model.LogisticRegression(), xtrain_count, train_y, xvalid_count)
3 #print("LR, Count Vectors: ", accuracy)
```

In []:

```
1 # ngram level tf-idf
2 #tfidf_vect_ngram = TfidfVectorizer(analyzer='word', token_pattern=r'\w{1,}', ngram_range=(1,2))
3 #tfidf_vect_ngram.fit(trainDF['CallTrace_T'])
4 #xtrain_tfidf_ngram = tfidf_vect_ngram.transform(train_x)
5 #xvalid_tfidf_ngram = tfidf_vect_ngram.transform(valid_x)
```

In []:

```
1 # Linear Classifier on Ngram Level TF IDF Vectors
2 #accuracy = train_model(linear_model.LogisticRegression(), xtrain_tfidf_ngram, train_y, xvalid_tfidf_ngram)
3 #print("LR, N-Gram Vectors: ", accuracy)
```

In []:

```
1 # characters level tf-idf
2 #tfidf_vect_ngram_chars = TfidfVectorizer(analyzer='char', token_pattern=r'\w{1,}', ngram_range=(1,2))
3 #tfidf_vect_ngram_chars.fit(trainDF['CallTrace_T'])
4 #xtrain_tfidf_ngram_chars = tfidf_vect_ngram_chars.transform(train_x)
5 #xvalid_tfidf_ngram_chars = tfidf_vect_ngram_chars.transform(valid_x)
```

In []:

```
1 # Linear Classifier on Character Level TF IDF Vectors
2 #accuracy = train_model(linear_model.LogisticRegression(), xtrain_tfidf_ngram_chars, train_y, xvalid_tfidf_ngram_chars)
3 #print("LR, CharLevel Vectors: ", accuracy)
```