

▼ CREDIT CARD FRAUD DETECTION

IMPORTING REQUIRED LIBRARIES

```
1 import pandas as pd
2 import numpy as np
3 import itertools
4 import time
5
6 from sklearn.preprocessing import StandardScaler, RobustScaler
7 from sklearn.model_selection import train_test_split, cross_val_score, KFold, StratifiedKFold
8 from sklearn.metrics import confusion_matrix, precision_score, recall_score, f1_score,
9
10 from sklearn.linear_model import LogisticRegression
11 from sklearn.ensemble import RandomForestClassifier, AdaBoostClassifier
12 from sklearn.tree import DecisionTreeClassifier
13 from sklearn.neural_network import MLPClassifier
14 from sklearn.neighbors import KNeighborsClassifier
15 from sklearn.svm import SVC
16 from xgboost import XGBClassifier
17
18 import seaborn as sns
19 import matplotlib.pyplot as plt
20 %matplotlib inline
21
22 import warnings
23 warnings.simplefilter('ignore')
```

➞ /usr/local/lib/python3.6/dist-packages/statsmodels/tools/_testing.py:19: FutureWarning:
import pandas.util.testing as tm

LOAD IN THE DATASET

```
1 data = pd.read_csv('creditcard.csv')
2 data.head()
```

➞

	Time	V1	V2	V3	V4	V5	V6	V7	
0	0	-1.359807	-0.072781	2.536347	1.378155	-0.338321	0.462388	0.239599	0.098

EXPLORATORY DATA ANALYSIS

0	1	0.066272	0.485226	1.702002	0.862204	0.040200	1.247202	0.227600	0.277
---	---	----------	----------	----------	----------	----------	----------	----------	-------

```
1 # Find the rows and columns size
2 data.shape
```

```
(93181, 31)
```

```
1 # Columns present for the dataset
2 data.columns
```

```
Index(['Time', 'V1', 'V2', 'V3', 'V4', 'V5', 'V6', 'V7', 'V8', 'V9', 'V10',
      'V11', 'V12', 'V13', 'V14', 'V15', 'V16', 'V17', 'V18', 'V19', 'V20',
      'V21', 'V22', 'V23', 'V24', 'V25', 'V26', 'V27', 'V28', 'Amount',
      'Class'],
      dtype='object')
```

```
1 # Checking the datatypes of attributes
2 data.dtypes
```

Time int64

```
1 # Summary of the attributes
2 data.describe(include='all')
```

↗

	Time	V1	V2	V3	V4	
count	93181.000000	93181.000000	93181.000000	93181.000000	93181.000000	93181.0000
mean	40720.460405	-0.263002	-0.040780	0.676306	0.163064	-0.2800
std	16392.444505	1.869011	1.664308	1.341575	1.355105	1.3686
min	0.000000	-56.407510	-72.715728	-33.680984	-5.172595	-42.1478
25%	32678.000000	-1.028016	-0.605706	0.179018	-0.716783	-0.8987
50%	42848.000000	-0.258489	0.072818	0.756532	0.189117	-0.3149
75%	53622.000000	1.153021	0.728656	1.381324	1.034477	0.2509
max	64283.000000	1.960497	18.902453	4.226108	16.715537	34.8016

V22 T10dL04

```
1 # Checking for null values
2 data.isnull().sum()
```

↗

Time	0
V1	0
V2	0
V3	0
V4	0
V5	0
V6	0
V7	0
V8	0
V9	0
V10	0
V11	0
V12	0
V13	1
V14	1
V15	1
V16	1
V17	1
V18	1
V19	1
V20	1
V21	1
V22	1
V23	1
V24	1
V25	1
V26	1
V27	1
V28	1
Amount	1
Class	1
dtype:	int64

```
1 # Null values are present
2 # Replacing null values with mean
3 data.fillna(data.mean(), inplace=True)
```

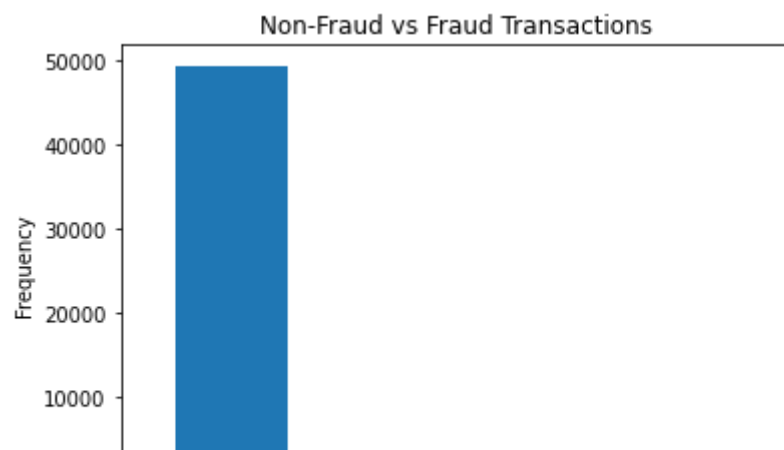
```
1 # No null values
2 data.isnull().sum()
```

```
Time      0
V1        0
V2        0
V3        0
V4        0
V5        0
V6        0
V7        0
V8        0
V9        0
V10       0
V11       0
V12       0
V13       0
V14       0
V15       0
V16       0
V17       0
V18       0
V19       0
V20       0
V21       0
V22       0
V23       0
V24       0
V25       0
V26       0
V27       0
V28       0
Amount    0
Class     0
dtype: int64
```

```
1 # 'Class' is the target variable
2 # It has two values counts - Fraud and Non-fraud
3 # Plot showing the counts of both counts
4 # Most of them are Non-fraud
5 # This is an unbalanced dataset
6
7 count_classes = pd.value_counts(data['Class'], sort=True)
8 count_classes.plot(kind='bar', rot=0)
9 plt.title('Non-Fraud vs Fraud Transactions')
10 plt.xticks(range(2), ['Non-Fraud', 'Fraud'])
11 plt.xlabel('Class')
12 plt.ylabel('Frequency')
```

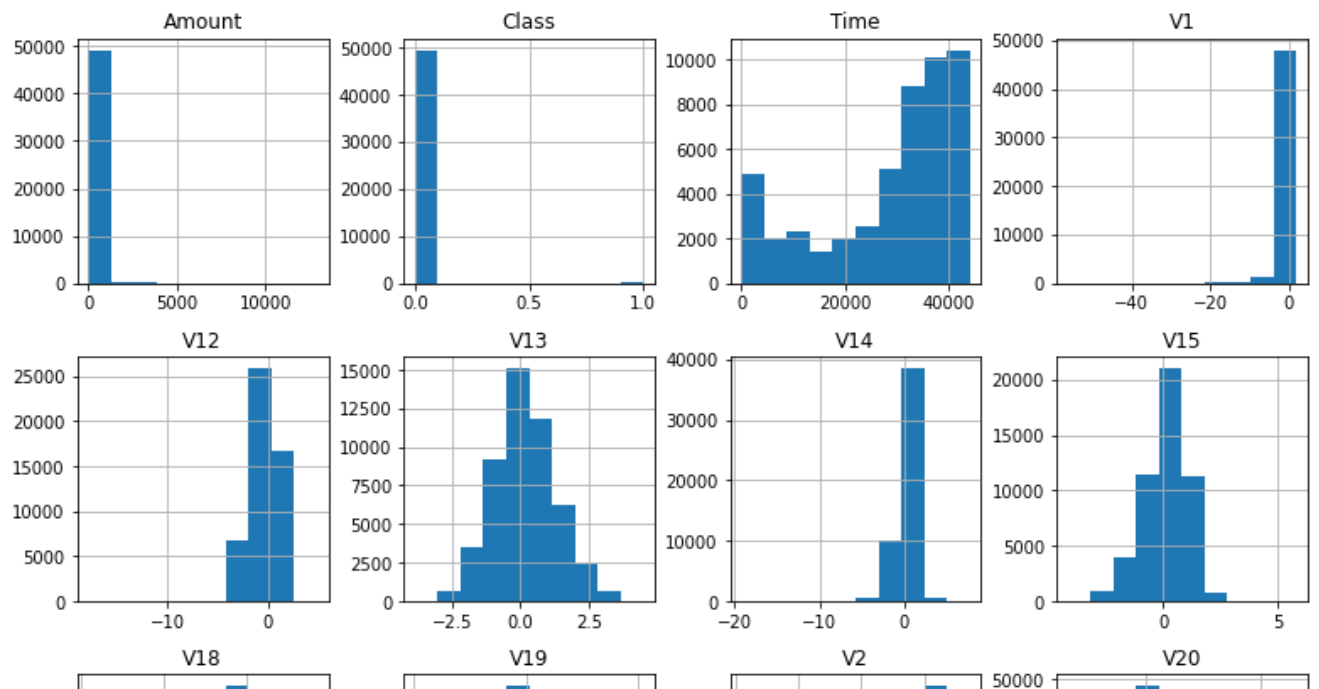
```
↳
```

```
Text(0, 0.5, 'Frequency')
```



```
1 # Plotting a histogram for each attribute  
2 data.hist(figsize=(20,20));
```





```

1 # Finding exact number of Fraud and Non-fraud values
2 # Fraud percentage is less than 1% of total transactions
3
4 fraud = data[data['Class'] == 1]
5 normal = data[data['Class'] == 0]
6
7 print('Fraud shape: {}'.format(fraud.shape))
8 print('Non-Fraud shape: {}'.format(normal.shape))
9 print('Percentage of fraud transcation: {:.4f}%'.format(len(fraud)/len(normal) * 100))

```

```

[ ]> Fraud shape: (148, 31)
Non-Fraud shape: (49461, 31)
Percentage of fraud transcation: 0.2992%

```

```

1 # Summary of Fraud transactions
2 fraud.Amount.describe()

```

```

[ ]> count      148.000000
mean        100.170676
std         233.347471
min           0.000000
25%           1.000000
50%           9.560000
75%          99.990000
max        1809.680000
Name: Amount, dtype: float64

```

```

1 # Summary of Normal transactions
2 normal.Amount.describe()

```

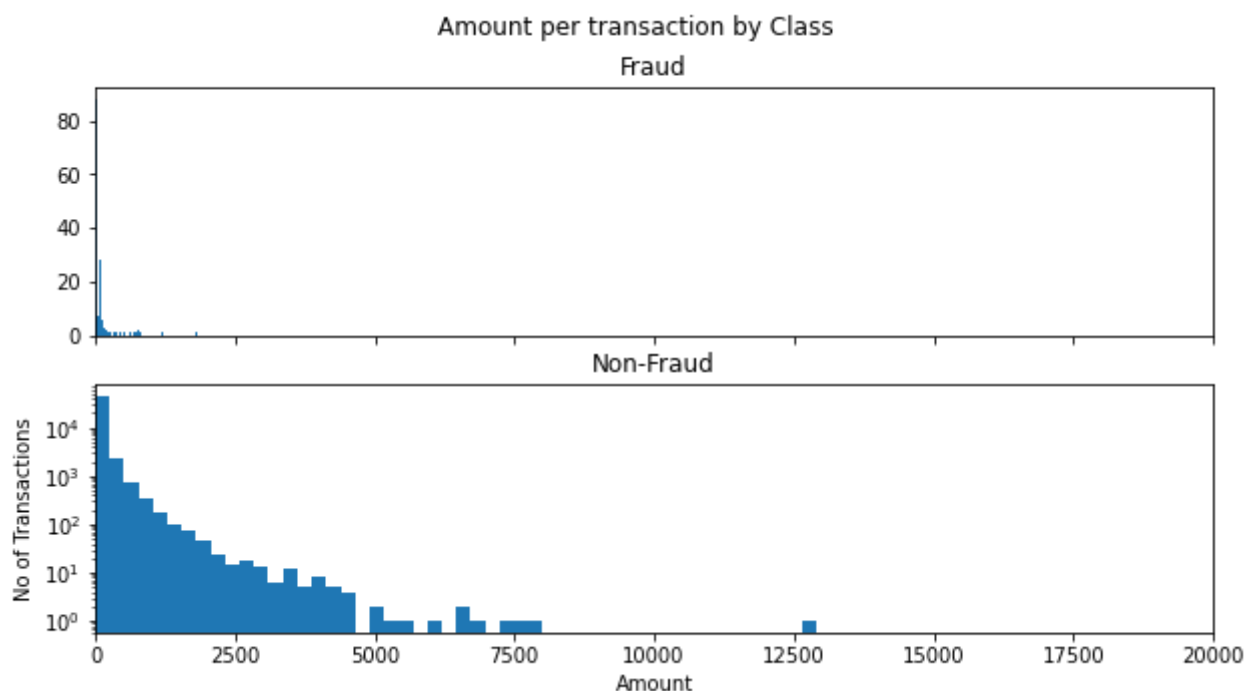
```

[ ]>

```

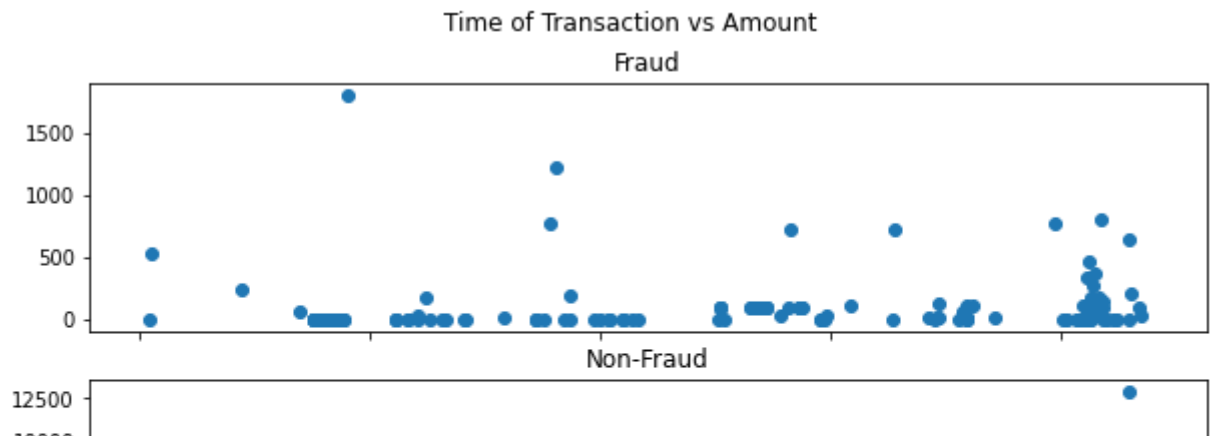
```
count    49461.000000
mean      93.099593
std       253.325102
min        0.000000
```

```
1 # Plot for showing the amount for both classes
2 fig, (ax1, ax2) = plt.subplots(2, 1, sharex=True, figsize=(10,5))
3 fig.suptitle('Amount per transaction by Class')
4 ax1.hist(fraud.Amount, bins=50)
5 ax1.set_title('Fraud')
6 ax2.hist(normal.Amount, bins=50)
7 ax2.set_title('Non-Fraud')
8 plt.xlabel('Amount')
9 plt.ylabel('No of Transactions')
10 plt.xlim(0, 20000)
11 plt.yscale('log')
```



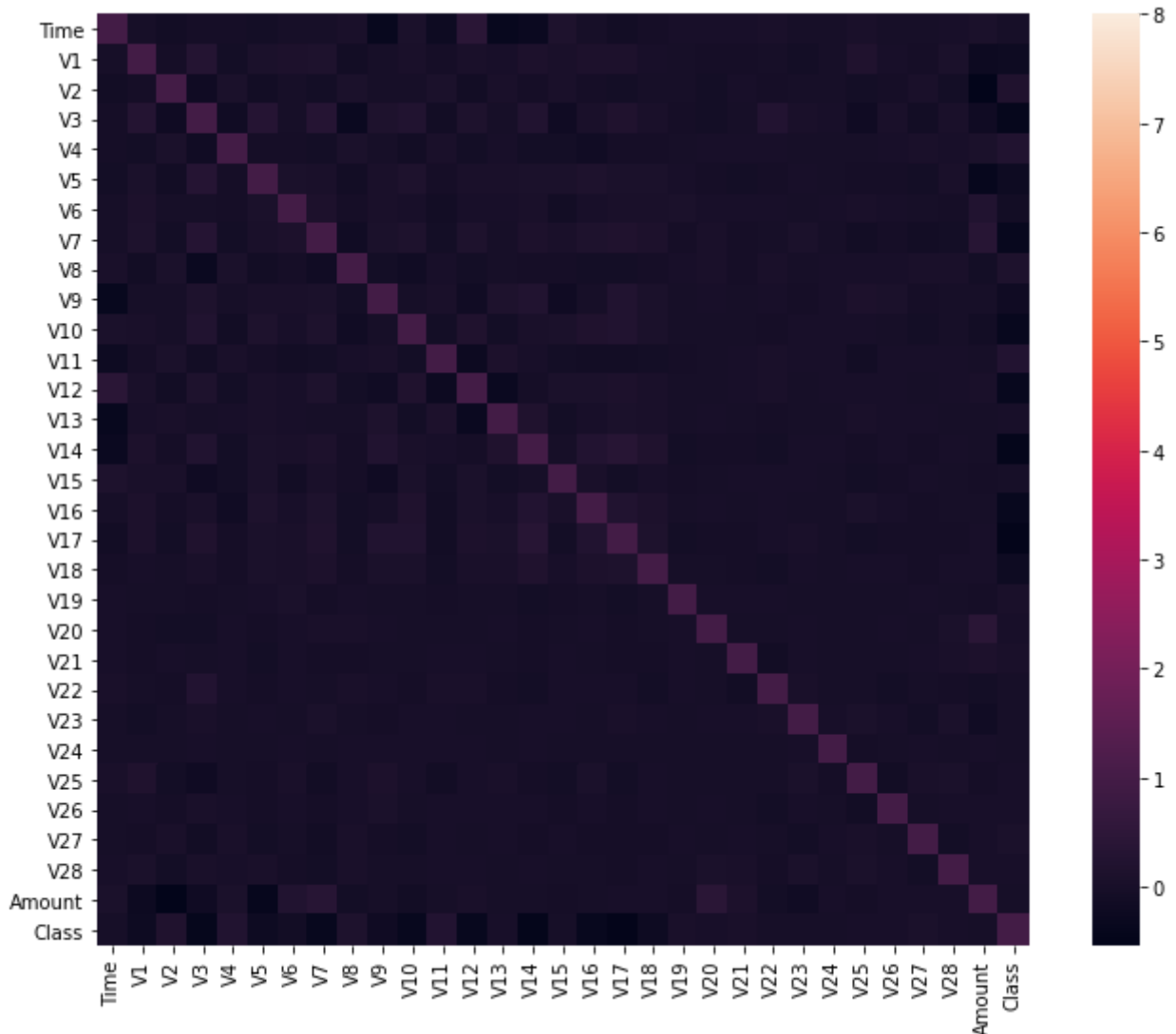
```
1 # Plot for Fraud and Normal transactions over time intervals
2 fig, (ax1, ax2) = plt.subplots(2, 1, sharex=True, figsize=(10,5))
3 fig.suptitle('Time of Transaction vs Amount')
4 ax1.scatter(fraud.Time, fraud.Amount)
5 ax1.set_title('Fraud')
6 ax2.scatter(normal.Time, normal.Amount)
7 ax2.set_title('Non-Fraud')
8 plt.xlabel('Time in seconds')
9 plt.ylabel('Amount')
10 plt.show()
```





```
1 # Heatmap to find any correlation between attributes
2 corr_matrix = data.corr()
3 fig = plt.figure(figsize=(12,9))
4 sns.heatmap(corr_matrix, vmax=8, square=True)
```

→ <matplotlib.axes._subplots.AxesSubplot at 0x7f2468830e48>



APPLYING TRANSFORMATION ON DATA

```
1 # The columns from 'V1' to 'V28' are already normalized
```



```

2 # Applying scaling on the 'Time' and 'Amount' attributes
3 # Using RobustScaler as its less prone to outliers
4 # and we need to detect outliers
5
6 rob_scaler = RobustScaler()
7
8 data['Scaled_time'] = rob_scaler.fit_transform(data['Time'].values.reshape(-1, 1))
9 data['Scaled_amount'] = rob_scaler.fit_transform(data['Amount'].values.reshape(-1, 1))
10
11 data.drop(['Time', 'Amount'], axis=1, inplace=True)

```

```

1 # Adding and removing columns
2
3 scaled_amount = data['Scaled_amount']
4 scaled_time = data['Scaled_time']
5
6 data.drop(['Scaled_amount', 'Scaled_time'], axis=1, inplace=True)
7 data.insert(0, 'Scaled_amount', scaled_amount)
8 data.insert(1, 'Scaled_time', scaled_time)
9
10 # Amount and Time are Scaled!
11 data.head()

```

```

↳

```

	Scaled_amount	Scaled_time	V1	V2	V3	V4	V5	
0	1.488894	-2.045837	-1.359807	-0.072781	2.536347	1.378155	-0.338321	0.40
1	-0.294453	-2.045837	1.191857	0.266151	0.166480	0.448154	0.060018	-0.00
2	4.268843	-2.045789	-1.358354	-1.340163	1.773209	0.379780	-0.503198	1.80
3	1.171866	-2.045789	-0.966272	-0.185226	1.792993	-0.863291	-0.010309	1.20
4	0.522393	-2.045741	-1.158233	0.877737	1.548718	0.403034	-0.407193	0.00

```

1 # Splitting of data for model
2 # X - Predictor variable
3 # y - Target variable
4
5 X = data.iloc[:, :-1]
6 y = data[['Class']]

```

```

1 # Splitting into train and test values
2
3 X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=.2, random_state=42)
4 print('X_train shape: {}'.format(X_train.shape))
5 print('X_test shape: {}'.format(X_test.shape))

```

```

↳ X_train shape: (74544, 30)
   X_test shape: (18637, 30)

```

TRAINING THE MODELS

```
1 # This is a classification task
2 # We will use the following models
3
4 classifiers = {
5     'LogisticRegression': LogisticRegression(),
6     'RandomForestClassifier': RandomForestClassifier(),
7     'AdaBoostClassifier': AdaBoostClassifier(),
8     'DecisionTreeClassifier': DecisionTreeClassifier(),
9     'SVC': SVC(),
10    'MLPClassifier': MLPClassifier(),
11    'KNeighborsClassifier': KNeighborsClassifier(),
12    'XGBClassifier': XGBClassifier()
13 }
```

```
1 # Function to train model and do predictions
2 # We want a model with high recall as to detect outliers - Fraud transactions
3
4 precision = []
5 recall = []
6
7 for name, clf in classifiers.items():
8     start = time.time()
9     name = clf.fit(X_train, y_train)
10    end = time.time()
11    y_pred = name.predict(X_test)
12
13    print('*****')
14    print('\nModel: {}'.format(name))
15    print('\nTime taken: {:.2f}min'.format((end-start)/60))
16    print('\nTrainig Accuracy: {:.2f}%'.format(name.score(X_train, y_train)*100))
17    print('\nTest Accuracy: {:.2f}%'.format(name.score(X_test, y_test)*100))
18    #print('\nConfusion Matrix: \n')
19    #print(confusion_matrix(y_test, y_pred))
20    print('\nPrecision Score: {:.3f}'.format(precision_score(y_test, y_pred)))
21    print('\nRecall Score: {:.3f}'.format(recall_score(y_test, y_pred)))
22    precision.append(precision_score(y_test, y_pred))
23    recall.append(recall_score(y_test, y_pred))
24    print('\n')
25    print('*****')
26
27 # i = recall.index(max(recall))
28 # print('\nModel with best recall: {}'.format(classifiers[i]))
29 # print('\nRecall: {}, Precision: {}'.format(recall[i], precision[i]))
```



```
Model: LogisticRegression(C=1.0, class_weight=None, dual=False, fit_intercept=True,
                           intercept_scaling=1, l1_ratio=None, max_iter=100,
                           multi_class='auto', n_jobs=None, penalty='l2',
                           random_state=None, solver='lbfgs', tol=0.0001, verbose=0,
                           warm_start=False)
```

Time taken: 0.01min

Trainig Accuracy: 99.87%

Test Accuracy: 99.82%

Precision Score: 0.692

Recall Score: 0.692

```
Model: RandomForestClassifier(bootstrap=True, ccp_alpha=0.0, class_weight=None,
                              criterion='gini', max_depth=None, max_features='auto',
                              max_leaf_nodes=None, max_samples=None,
                              min_impurity_decrease=0.0, min_impurity_split=None,
                              min_samples_leaf=1, min_samples_split=2,
                              min_weight_fraction_leaf=0.0, n_estimators=100,
                              n_jobs=None, oob_score=False, random_state=None,
                              verbose=0, warm_start=False)
```

Time taken: 0.09min

Trainig Accuracy: 100.00%

Test Accuracy: 99.93%

Precision Score: 0.917

Recall Score: 0.846

```
Model: AdaBoostClassifier(algorithm='SAMME.R', base_estimator=None, learning_rate=1.0
                           n_estimators=50, random_state=None)
```

Time taken: 0.07min

Trainig Accuracy: 100.00%

Test Accuracy: 99.89%

Precision Score: 0.900

Recall Score: 0.692

```
Model: DecisionTreeClassifier(ccp_alpha=0.0, class_weight=None, criterion='gini',
                             max_depth=None, max_features=None, max_leaf_nodes=None,
                             min_impurity_decrease=0.0, min_impurity_split=None,
                             min_samples_leaf=1, min_samples_split=2,
                             min_weight_fraction_leaf=0.0, presort='deprecated',
                             random_state=None, splitter='best')
```

Time taken: 0.01min

Trainig Accuracy: 100.00%

Test Accuracy: 99.89%

Precision Score: 0.833

Recall Score: 0.769

```
*****
*****
```

```
Model: SVC(C=1.0, break_ties=False, cache_size=200, class_weight=None, coef0=0.0,
           decision_function_shape='ovr', degree=3, gamma='scale', kernel='rbf',
           max_iter=-1, probability=False, random_state=None, shrinking=True,
           tol=0.001, verbose=False)
```

Time taken: 0.01min

Trainig Accuracy: 99.95%

Test Accuracy: 99.86%

Precision Score: 0.889

Recall Score: 0.615

```
*****
*****
```

```
Model: MLPClassifier(activation='relu', alpha=0.0001, batch_size='auto', beta_1=0.9,
                    beta_2=0.999, early_stopping=False, epsilon=1e-08,
                    hidden_layer_sizes=(100,), learning_rate='constant',
                    learning_rate_init=0.001, max_fun=15000, max_iter=200,
                    momentum=0.9, n_iter_no_change=10, nesterovs_momentum=True,
                    power_t=0.5, random_state=None, shuffle=True, solver='adam',
                    tol=0.0001, validation_fraction=0.1, verbose=False,
                    warm_start=False)
```

Time taken: 0.11min

Trainig Accuracy: 99.97%

Test Accuracy: 99.91%

Precision Score: 0.846

Recall Score: 0.846

```
*****
```

```
*****
```

```
Model: KNeighborsClassifier(algorithm='auto', leaf_size=30, metric='minkowski',
                           metric_params=None, n_jobs=None, n_neighbors=5, p=2,
                           weights='uniform')
```

```
Time taken: 0.00min
```

```
Trainig Accuracy: 99.91%
```

```
Test Accuracy: 99.89%
```

```
Precision Score: 0.833
```

```
Recall Score: 0.769
```

```
*****
*****
```

```
Model: XGBClassifier(base_score=0.5, booster='gbtree', colsample_bylevel=1,
                    colsample_bynode=1, colsample_bytree=1, gamma=0,
                    learning_rate=0.1, max_delta_step=0, max_depth=3,
                    min_child_weight=1, missing=None, n_estimators=100, n_jobs=1,
                    nthread=None, objective='binary:logistic', random_state=0,
                    reg_alpha=0, reg_lambda=1, scale_pos_weight=1, seed=None,
                    silent=None, subsample=1, verbosity=1)
```

```
Time taken: 0.05min
```

```
Trainig Accuracy: 100.00%
```

```
Test Accuracy: 99.91%
```

```
Precision Score: 0.833
```

```
1 # Wow our scores are getting even high scores even when applying cross validation.
2
3 for key, classifier in classifiers.items():
4     classifier.fit(X_train, y_train)
5     training_score = cross_val_score(classifier, X_train, y_train, cv=5)
6     print(classifier.__class__.__name__, ": ", round(training_score.mean(), 2) * 100, "%")
```

```
LogisticRegression : 100.0 % accuracy score
RandomForestClassifier : 100.0 % accuracy score
AdaBoostClassifier : 100.0 % accuracy score
DecisionTreeClassifier : 100.0 % accuracy score
SVC : 100.0 % accuracy score
MLPClassifier : 100.0 % accuracy score
KNeighborsClassifier : 100.0 % accuracy score
XGBClassifier : 100.0 % accuracy score
```

```
1 # Use GridSearchCV to find the best parameters.
2
3
4 # Logistic Regression
5 log_reg_params = {"penalty": ['l1', 'l2'], 'C': [0.001, 0.01, 0.1, 1, 10, 100, 1000]}
6 grid_log_reg = GridSearchCV(LogisticRegression(), log_reg_params)
7 grid_log_reg.fit(X_train, y_train)
8 # We automatically get the logistic regression with the best parameters
```

```
8 # We automatically get the logistic regression with the best parameters.
9 log_reg = grid_log_reg.best_estimator_
10
11
12 # K-Nears
13 knears_params = {"n_neighbors": list(range(2,5,1)), 'algorithm': ['auto', 'ball_tree',
14 grid_knears = GridSearchCV(KNeighborsClassifier(), knears_params)
15 grid_knears.fit(X_train, y_train)
16 # KNears best estimator
17 knears_neighbors = grid_knears.best_estimator_
18
19
20 # Support Vector Classifier
21 svc_params = {'C': [0.5, 0.7, 0.9, 1], 'kernel': ['rbf', 'poly', 'sigmoid', 'linear']}
22 grid_svc = GridSearchCV(SVC(), svc_params)
23 grid_svc.fit(X_train, y_train)
24 # SVC best estimator
25 svc = grid_svc.best_estimator_
26
27
28 # DecisionTree Classifier
29 tree_params = {"criterion": ["gini", "entropy"], "max_depth": list(range(2,4,1)),
30                 "min_samples_leaf": list(range(5,7,1))}
31 grid_tree = GridSearchCV(DecisionTreeClassifier(), tree_params)
32 grid_tree.fit(X_train, y_train)
33 # tree best estimator
34 tree_clf = grid_tree.best_estimator_
35
36 # Random
37 random_params = {
38     'n_estimators': [200, 500],
39     'max_features': ['auto', 'sqrt', 'log2'],
40     'max_depth' : [4,5,6,7,8],
41     'criterion' :['gini', 'entropy']}
42 grid_random = GridSearchCV(RandomForestClassifier(), random_params)
43 grid_random.fit(X_train, y_train)
44 # random best estimator
45 random_clf = grid_random.best_estimator_
46
47 # MLP
48 mlp_params = {'solver': ['lbfgs'],
49               'max_iter': [1000,1100,1200,1300,1400,1500,1600,1700,1800,1900,2000],
50               'alpha': 10.0 ** -np.arange(1, 10),
51               'hidden_layer_sizes':np.arange(10, 15)}
52 grid_mlp = GridSearchCV(MLPClassifier(), mlp_params)
53 grid_mlp.fit(X_train, y_train)
54 # random best estimator
55 mlp_clf = grid_mlp.best_estimator_
56
57 # Ada
58 ada_params = {'base_estimator__max_depth':[1,50],
59               'base_estimator':[DecisionTreeClassifier(max_features=2), DecisionTreeClassif
60 grid_ada = GridSearchCV(AdaBoostClassifier(base_estimator=DecisionTreeClassifier()), ad
61 grid_ada.fit(X_train, y_train)
62 # best
63 ada_clf = grid_ada.best_estimator_
```

```
ada_clf = grid_ada.best_estimator_
```

```
1 # Now we will train models again with the parameters obtained from
2 # GridSearchCV
3 # Lets see if we can get higher recall values
4
5 estimators = [log_reg, random_clf, ada_clf ,tree_clf, svc, mlp_clf, knears_neighbors]
6
7 names = ["Logistic Regression", "Nearest Neighbors",
8          "Decision Tree", "Random Forest", "Neural Net", "AdaBoost",
9          "Naive Bayes", "QDA"]
10
11 classifiers = [
12     LogisticRegression(),
13     RandomForestClassifier(),
14     AdaBoostClassifier(),
15     DecisionTreeClassifier(),
16     SVC(),
17     MLPClassifier(),
18     KNeighborsClassifier()]
```

```
1 # Function to train model and do predictions, with best parameters
2 precision = []
3 recall = []
4
5 for name, clf, est in zip(names, classifiers, estimators):
6     start = time.time()
7     name = est.fit(X_train, y_train)
8     end = time.time()
9     y_pred = name.predict(X_test)
10
11     print('*****')
12     print('\nModel: {}'.format(name))
13     print('\nTime taken: {:.2f}min'.format((end-start)/60))
14     print('\nTrainig Accuracy: {:.2f}%'.format(name.score(X_train, y_train)*100))
15     print('\nTest Accuracy: {:.2f}%'.format(name.score(X_test, y_test)*100))
16     #print('\nConfusion Matrix: \n')
17     #print(confusion_matrix(y_test, y_pred))
18     print('\nPrecision Score: {:.3f}'.format(precision_score(y_test, y_pred)))
19     print('\nRecall Score: {:.3f}'.format(recall_score(y_test, y_pred)))
20     precision.append(precision_score(y_test, y_pred))
21     recall.append(recall_score(y_test, y_pred))
22     print('\n')
23     print('*****')
24
25 i = recall.index(max(recall))
26 print('\nModel with best recall: {}'.format(names[i]))
27 print('\nRecall: {}, Precision: {}'.format(recall[i], precision[i]))
```



```
Model: LogisticRegression(C=10, class_weight=None, dual=False, fit_intercept=True,
                           intercept_scaling=1, l1_ratio=None, max_iter=100,
                           multi_class='auto', n_jobs=None, penalty='l2',
                           random_state=None, solver='lbfgs', tol=0.0001, verbose=0,
                           warm_start=False)
```

Time taken: 0.01min

Trainig Accuracy: 99.89%

Test Accuracy: 99.84%

Precision Score: 0.750

Recall Score: 0.692

```
Model: RandomForestClassifier(bootstrap=True, ccp_alpha=0.0, class_weight=None,
                              criterion='entropy', max_depth=7, max_features='sqrt',
                              max_leaf_nodes=None, max_samples=None,
                              min_impurity_decrease=0.0, min_impurity_split=None,
                              min_samples_leaf=1, min_samples_split=2,
                              min_weight_fraction_leaf=0.0, n_estimators=200,
                              n_jobs=None, oob_score=False, random_state=None,
                              verbose=0, warm_start=False)
```

Time taken: 0.16min

Trainig Accuracy: 99.98%

Test Accuracy: 99.89%

Precision Score: 0.900

Recall Score: 0.692

```
Model: AdaBoostClassifier(algorithm='SAMME.R',
                          base_estimator=DecisionTreeClassifier(ccp_alpha=0.0,
                                                                    class_weight=None,
                                                                    criterion='gini',
                                                                    max_depth=1,
                                                                    max_features=2,
                                                                    max_leaf_nodes=None,
                                                                    min_impurity_decrease=0.0,
                                                                    min_impurity_split=None,
                                                                    min_samples_leaf=1,
                                                                    min_samples_split=2,
                                                                    min_weight_fraction_leaf=0.0,
                                                                    presort='deprecated',
                                                                    random_state=None,
                                                                    splitter='best'),
                          learning_rate=1.0, n_estimators=50, random_state=None)
```


Time taken: 0.01min

Trainig Accuracy: 100.00%

Test Accuracy: 99.89%

Precision Score: 0.900

Recall Score: 0.692


```
Model: DecisionTreeClassifier(ccp_alpha=0.0, class_weight=None, criterion='gini',
                             max_depth=2, max_features=None, max_leaf_nodes=None,
                             min_impurity_decrease=0.0, min_impurity_split=None,
                             min_samples_leaf=6, min_samples_split=2,
                             min_weight_fraction_leaf=0.0, presort='deprecated',
                             random_state=None, splitter='best')
```

Time taken: 0.00min

Trainig Accuracy: 99.90%

Test Accuracy: 99.89%

Precision Score: 0.900

Recall Score: 0.692


```
Model: SVC(C=0.5, break_ties=False, cache_size=200, class_weight=None, coef0=0.0,
           decision_function_shape='ovr', degree=3, gamma='scale', kernel='poly',
           max_iter=-1, probability=False, random_state=None, shrinking=True,
           tol=0.001, verbose=False)
```

Time taken: 0.00min

Trainig Accuracy: 99.97%

Test Accuracy: 99.93%

Precision Score: 0.917

Recall Score: 0.846


```
Model: MLPClassifier(activation='relu', alpha=0.001, batch_size='auto', beta_1=0.9,
                    beta_2=0.999, early_stopping=False, epsilon=1e-08,
                    hidden_layer_sizes=10, learning_rate='constant',
                    learning_rate_init=0.001, max_fun=15000, max_iter=1600,
                    momentum=0.9, n_iter_no_change=10, nesterovs_momentum=True,
                    power_t=0.5, random_state=None, shuffle=True, solver='lbfgs',
                    tol=0.0001, validation_fraction=0.1, verbose=False,
```

```

    tol=0.0001, validation_fraction=0.1, verbose=False,
    warm_start=False)

```

Time taken: 0.01min

Trainig Accuracy: 100.00%

Test Accuracy: 99.93%

Precision Score: 0.917

Recall Score: 0.846

```

*****
*****

```

```

Model: KNeighborsClassifier(algorithm='auto', leaf_size=30, metric='minkowski',
                             metric_params=None, n_jobs=None, n_neighbors=3, p=2,
                             weights='uniform')

```

Time taken: 0.00min

Trainig Accuracy: 99.92%

Test Accuracy: 99.89%

Precision Score: 0.900

Recall Score: 0.692

```

*****

```

Model with best recall: Neural Net

```

1 # As Neural net and KNN has same score we will test using both
2 # The better one will be chosen
3
4 mlp = MLPClassifier(activation='relu', alpha=0.001, batch_size='auto', beta_1=0.9,
5                     beta_2=0.999, early_stopping=False, epsilon=1e-08,
6                     hidden_layer_sizes=10, learning_rate='constant',
7                     learning_rate_init=0.001, max_fun=15000, max_iter=1600,
8                     momentum=0.9, n_iter_no_change=10, nesterovs_momentum=True,
9                     power_t=0.5, random_state=None, shuffle=True, solver='lbfgs',
10                    tol=0.0001, validation_fraction=0.1, verbose=False,
11                    warm_start=False)
12 mlp.fit(X_train, y_train)
13 y_pred_mlp = mlp.predict(X_test)
14
15 print('Confusion Matrix:\n',confusion_matrix(y_test.round(), y_pred_mlp.round()))
16 print('Recall score: ',recall_score(y_test.round(),y_pred_mlp.round()))
17 print('Precision score:',precision_score(y_test.round(),y_pred_mlp.round()))

```



Confusion Matrix:

```

[[18586   8]
 [   12  31]]

```

Recall score: 0.7209302325581395

Precision score: 0.7948717948717948

```
1 # KNN has lower recall than MLP on test set
2
3 knn = KNeighborsClassifier(algorithm='auto', leaf_size=30, metric='minkowski',
4                             metric_params=None, n_jobs=None, n_neighbors=3, p=2,
5                             weights='uniform')
6 knn.fit(X_train, y_train)
7 y_pred_knn = knn.predict(X_test)
8
9 print('Confusion Matrix: \n',confusion_matrix(y_test.round(), y_pred_knn.round()))
10 print('Recall score: ',recall_score(y_test.round(),y_pred_knn.round()))
11 print('Precision score: ',precision_score(y_test.round(),y_pred_knn.round()))
```

↳ Confusion Matrix:
[[18592 2]
[13 30]]
Recall score: 0.6976744186046512
Precision score: 0.9375

```
1 # As Neural net has higher recall we will use that and
2 # Prepare submission file
3
4 output = pd.DataFrame({'Time': X_test.Scaled_time, 'Amount': X_test.Scaled_amount, 'Pre
5 output.to_csv('submission.csv', index=False)
```

```
1
```