

Credit Card Fraud Detection

The project is to recognize fraudulent credit card transactions so that the customers of credit card companies are not charged for items that they did not purchase.

1. Importing all the necessary Libraries

```
In [1]: import numpy as np
import pandas as pd
import matplotlib.pyplot as plt
import seaborn as sns
from matplotlib import gridspec
```

2. Loading the Data

You can download the dataset from kaggle. dataset link- <https://www.kaggle.com/mlg-ulb/creditcardfraud/download> You only need to put dataset the model will detect the frauds.

```
In [2]: data = pd.read_csv("https://media.githubusercontent.com/media/yashwantaditya009/Fintech/master/creditcard.csv")
```

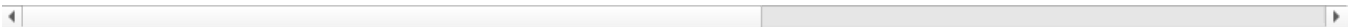
3. Understanding the Data

```
In [3]: data.head()
```

```
Out[3]:
```

| | Time | V1 | V2 | V3 | V4 | V5 | V6 | V7 | V8 | V9 | ... | V21 | V22 |
|---|------|-----------|-----------|----------|-----------|-----------|-----------|-----------|-----------|-----------|-----|-----------|-----------|
| 0 | 0.0 | -1.359807 | -0.072781 | 2.536347 | 1.378155 | -0.338321 | 0.462388 | 0.239599 | 0.098698 | 0.363787 | ... | -0.018307 | 0.277838 |
| 1 | 0.0 | 1.191857 | 0.266151 | 0.166480 | 0.448154 | 0.060018 | -0.082361 | -0.078803 | 0.085102 | -0.255425 | ... | -0.225775 | -0.638672 |
| 2 | 1.0 | -1.358354 | -1.340163 | 1.773209 | 0.379780 | -0.503198 | 1.800499 | 0.791461 | 0.247676 | -1.514654 | ... | 0.247998 | 0.771679 |
| 3 | 1.0 | -0.966272 | -0.185226 | 1.792993 | -0.863291 | -0.010309 | 1.247203 | 0.237609 | 0.377436 | -1.387024 | ... | -0.108300 | 0.005274 |
| 4 | 2.0 | -1.158233 | 0.877737 | 1.548718 | 0.403034 | -0.407193 | 0.095921 | 0.592941 | -0.270533 | 0.817739 | ... | -0.009431 | 0.798278 |

5 rows × 31 columns



4. Describing the Data

```
In [4]: print(data.shape)
print(data.describe())
```

```
(284807, 31)
```

| | Time | V1 | V2 | V3 | V4 | \ |
|-------|---------------|---------------|---------------|---------------|---------------|---|
| count | 284807.000000 | 2.848070e+05 | 2.848070e+05 | 2.848070e+05 | 2.848070e+05 | |
| mean | 94813.859575 | 1.165980e-15 | 3.416908e-16 | -1.373150e-15 | 2.086869e-15 | |
| std | 47488.145955 | 1.958696e+00 | 1.651309e+00 | 1.516255e+00 | 1.415869e+00 | |
| min | 0.000000 | -5.640751e+01 | -7.271573e+01 | -4.832559e+01 | -5.683171e+00 | |
| 25% | 54201.500000 | -9.203734e-01 | -5.985499e-01 | -8.903648e-01 | -8.486401e-01 | |
| 50% | 84692.000000 | 1.810880e-02 | 6.548556e-02 | 1.798463e-01 | -1.984653e-02 | |
| 75% | 139320.500000 | 1.315642e+00 | 8.037239e-01 | 1.027196e+00 | 7.433413e-01 | |
| max | 172792.000000 | 2.454930e+00 | 2.205773e+01 | 9.382558e+00 | 1.687534e+01 | |

| | V5 | V6 | V7 | V8 | V9 | \ |
|-------|---------------|---------------|---------------|---------------|---------------|---|
| count | 2.848070e+05 | 2.848070e+05 | 2.848070e+05 | 2.848070e+05 | 2.848070e+05 | |
| mean | 9.604066e-16 | 1.490107e-15 | -5.556467e-16 | 1.177556e-16 | -2.406455e-15 | |
| std | 1.380247e+00 | 1.332271e+00 | 1.237094e+00 | 1.194353e+00 | 1.098632e+00 | |
| min | -1.137433e+02 | -2.616051e+01 | -4.355724e+01 | -7.321672e+01 | -1.343407e+01 | |
| 25% | -6.915971e-01 | -7.682956e-01 | -5.540759e-01 | -2.086297e-01 | -6.430976e-01 | |
| 50% | -5.433583e-02 | -2.741871e-01 | 4.010308e-02 | 2.235804e-02 | -5.142873e-02 | |
| 75% | 6.119264e-01 | 3.985649e-01 | 5.704361e-01 | 3.273459e-01 | 5.971390e-01 | |
| max | 3.480167e+01 | 7.330163e+01 | 1.205895e+02 | 2.000721e+01 | 1.559499e+01 | |

| | ... | V21 | V22 | V23 | V24 | \ |
|-------|-----|---------------|---------------|---------------|---------------|---|
| count | ... | 2.848070e+05 | 2.848070e+05 | 2.848070e+05 | 2.848070e+05 | |
| mean | ... | 1.656562e-16 | -3.444850e-16 | 2.578648e-16 | 4.471968e-15 | |
| std | ... | 7.345240e-01 | 7.257016e-01 | 6.244603e-01 | 6.056471e-01 | |
| min | ... | -3.483038e+01 | -1.093314e+01 | -4.480774e+01 | -2.836627e+00 | |
| 25% | ... | -2.283949e-01 | -5.423504e-01 | -1.618463e-01 | -3.545861e-01 | |
| 50% | ... | -2.945017e-02 | 6.781943e-03 | -1.119293e-02 | 4.097606e-02 | |
| 75% | ... | 1.863772e-01 | 5.285536e-01 | 1.476421e-01 | 4.395266e-01 | |
| max | ... | 2.720284e+01 | 1.050309e+01 | 2.252841e+01 | 4.584549e+00 | |

| | V25 | V26 | V27 | V28 | Amount | \ |
|-------|---------------|---------------|---------------|---------------|---------------|---|
| count | 2.848070e+05 | 2.848070e+05 | 2.848070e+05 | 2.848070e+05 | 284807.000000 | |
| mean | 5.340915e-16 | 1.687098e-15 | -3.666453e-16 | -1.220404e-16 | 88.349619 | |
| std | 5.212781e-01 | 4.822270e-01 | 4.036325e-01 | 3.300833e-01 | 250.120109 | |
| min | -1.029540e+01 | -2.604551e+00 | -2.256568e+01 | -1.543008e+01 | 0.000000 | |
| 25% | -3.171451e-01 | -3.269839e-01 | -7.083953e-02 | -5.295979e-02 | 5.600000 | |
| 50% | 1.659350e-02 | -5.213911e-02 | 1.342146e-03 | 1.124383e-02 | 22.000000 | |
| 75% | 3.507156e-01 | 2.409522e-01 | 9.104512e-02 | 7.827995e-02 | 77.165000 | |
| max | 7.519589e+00 | 3.517346e+00 | 3.161220e+01 | 3.384781e+01 | 25691.160000 | |

| | Class |
|-------|---------------|
| count | 284807.000000 |
| mean | 0.001727 |
| std | 0.041527 |
| min | 0.000000 |
| 25% | 0.000000 |
| 50% | 0.000000 |
| 75% | 0.000000 |
| max | 1.000000 |

[8 rows x 31 columns]

5. Imbalance in the data

```
In [5]: fraud = data[data['Class'] == 1]
valid = data[data['Class'] == 0]
outlierFraction = len(fraud)/float(len(valid))
print(outlierFraction)
print('Fraud Cases: {}'.format(len(data[data['Class'] == 1])))
print('Valid Transactions: {}'.format(len(data[data['Class'] == 0])))
```

0.0017304750013189597

Fraud Cases: 492

Valid Transactions: 284315

Only 0.17% fraudulent transaction out all the transactions. The data is highly Unbalanced. Lets first apply the models without balancing it and if we don't get a good accuracy then we can find a way to balance this dataset.

6. Print the amount details for Fraudulent Transaction

```
In [6]: print("Amount details of the fraudulent transaction")
fraud.Amount.describe()
```

Amount details of the fraudulent transaction

```
Out[6]: count      492.000000
       mean       122.211321
       std        256.683288
       min         0.000000

       25%         1.000000
       50%         9.250000
       75%        105.890000
       max        2125.870000
       Name: Amount, dtype: float64
```

Here we can clearly see from this, the average Money transaction for the fraudulent ones is more. This makes this problem crucial to deal with.

7. Print the amount details for Normal Transaction

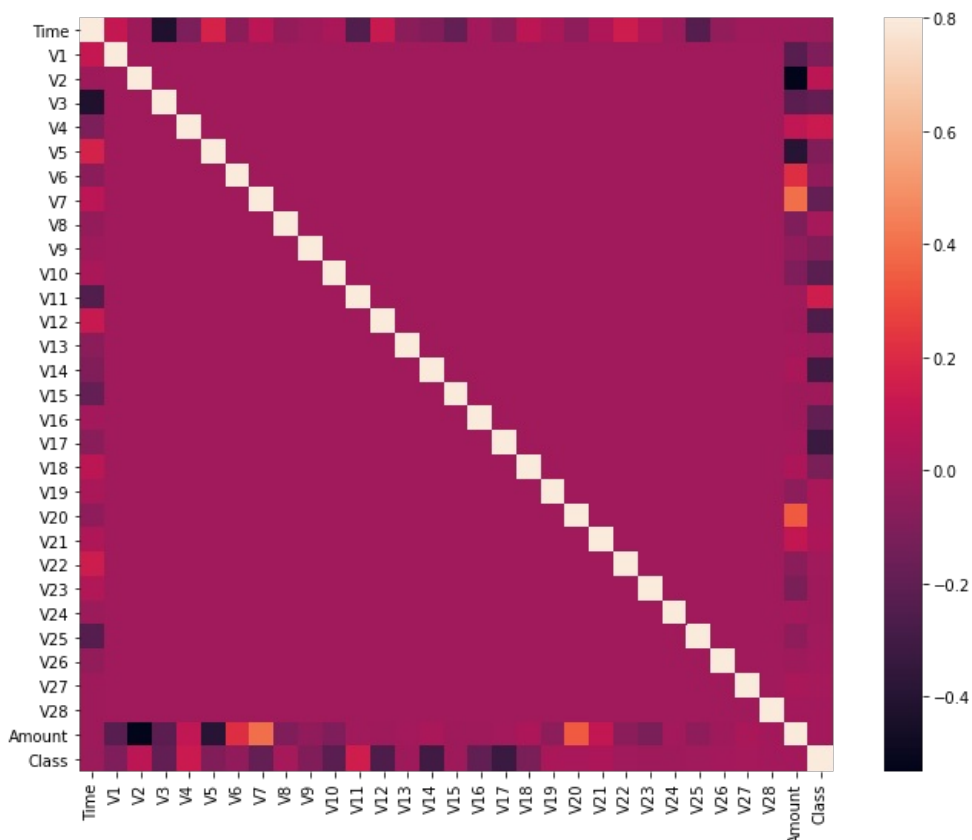
```
In [7]: print("details of valid transaction")
       valid.Amount.describe()
```

details of valid transaction

```
Out[7]: count      284315.000000
       mean         88.291022
       std         250.105092
       min          0.000000
       25%          5.650000
       50%          22.000000
       75%          77.050000
       max        25691.160000
       Name: Amount, dtype: float64
```

8. Plotting the Correlation Matrix

```
In [8]: corrmat = data.corr()
       fig = plt.figure(figsize = (12, 9))
       sns.heatmap(corrmat, vmax = .8, square = True)
       plt.show()
```



In the HeatMap we can clearly see that most of the features do not correlate to other features but there are some features that either has a positive or a negative correlation with each other. For example, V2 and V5 are highly negatively correlated with the feature called Amount. We also see some correlation with V20 and Amount. This gives us a deeper understanding of the Data available to us.

9. Separating the X and the Y values

Dividing the data into inputs parameters and outputs value format

```
In [9]: X = data.drop(['Class'], axis = 1)
Y = data["Class"]
print(X.shape)
print(Y.shape)
xData = X.values
yData = Y.values
```

```
(284807, 30)
(284807,)
```

10. Training and Testing Data Bifurcation

We will be dividing the dataset into two main groups. One for training the model and the other for Testing our trained model's performance.

```
In [10]: from sklearn.model_selection import train_test_split
xTrain, xTest, yTrain, yTest = train_test_split(
    xData, yData, test_size = 0.2, random_state = 42)
```

```
In [17]: from sklearn.ensemble import RandomForestClassifier
```

```
In [18]: rfc = RandomForestClassifier()
rfc.fit(xTrain, yTrain)
yPred = rfc.predict(xTest)
```

11. Building all kinds of evaluating parameters

```
In [19]: from sklearn.metrics import classification_report, accuracy_score
from sklearn.metrics import precision_score, recall_score
from sklearn.metrics import f1_score, matthews_corrcoef
from sklearn.metrics import confusion_matrix

n_outliers = len(fraud)
n_errors = (yPred != yTest).sum()
print("The model used is Random Forest classifier")

acc = accuracy_score(yTest, yPred)
print("The accuracy is {}".format(acc))

prec = precision_score(yTest, yPred)
print("The precision is {}".format(prec))

rec = recall_score(yTest, yPred)
print("The recall is {}".format(rec))

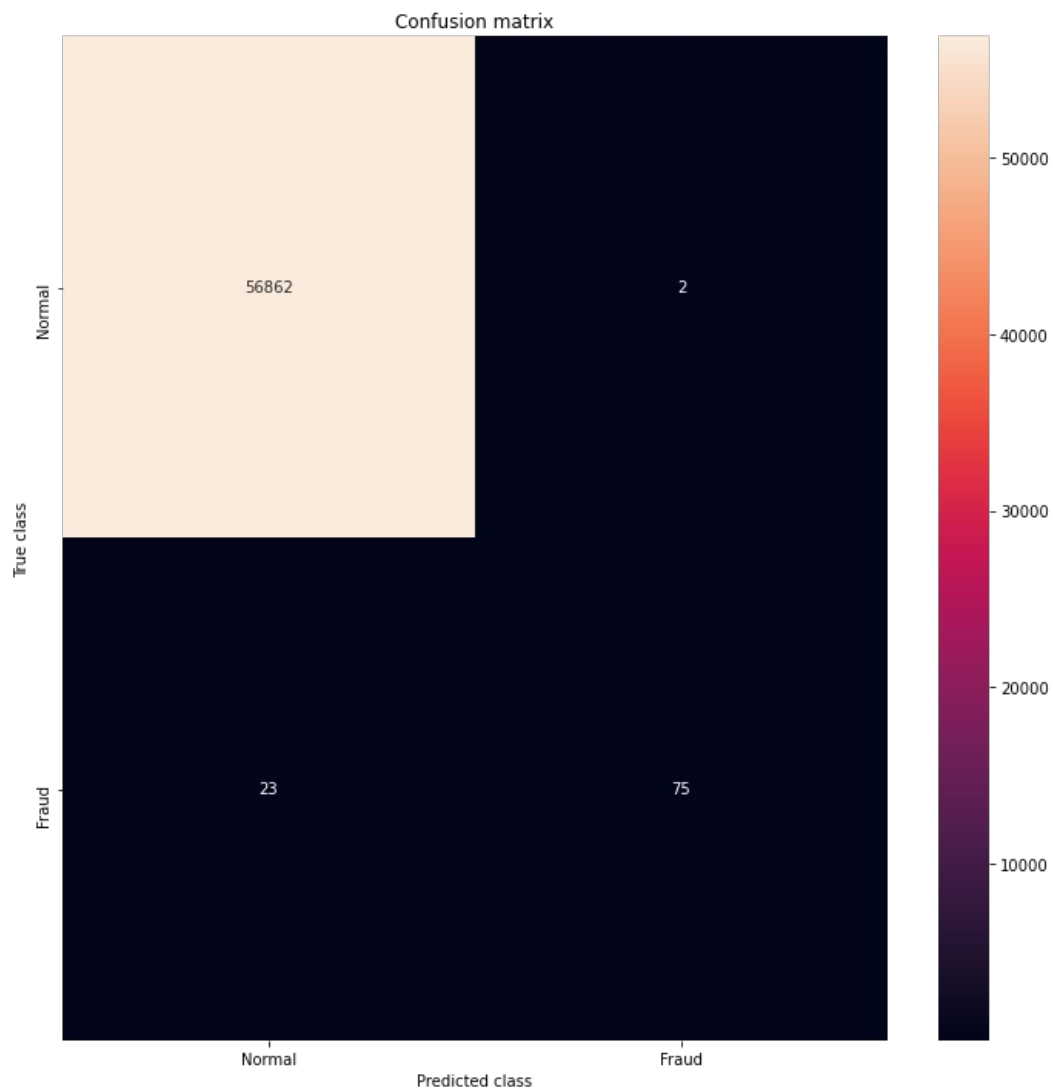
f1 = f1_score(yTest, yPred)
print("The F1-Score is {}".format(f1))

MCC = matthews_corrcoef(yTest, yPred)
print("The Matthews correlation coefficient is{}".format(MCC))
```

```
The model used is Random Forest classifier
The accuracy is 0.9995611109160493
The precision is 0.974025974025974
The recall is 0.7653061224489796
The F1-Score is 0.8571428571428571
The Matthews correlation coefficient is0.8631826952924256
```

12. Visualizing the Confusion Matrix

```
In [20]: LABELS = ['Normal', 'Fraud']
conf_matrix = confusion_matrix(yTest, yPred)
plt.figure(figsize =(12, 12))
sns.heatmap(conf_matrix, xticklabels = LABELS,
            yticklabels = LABELS, annot = True, fmt ="d");
plt.title("Confusion matrix")
plt.ylabel('True class')
plt.xlabel('Predicted class')
plt.show()
```



As we can see with our Random Forest Model we are getting a better result even for the recall which is the most tricky part.

In []:

In []:

Loading [MathJax]/jax/output/CommonHTML/fonts/TeX/fontdata.js