

Credit Card Fraud Detection Done by:Samyuktha Rajkumaran 190701182 CSE-D

In [1]:

```
import numpy as np
import pandas as pd
import matplotlib.pyplot as plt
import seaborn as sns
from matplotlib import gridspec
```

In [2]:

```
data=pd.read_csv(r"C:\Users\Rajkumar\Downloads\archive (2)\creditcard.csv")
```

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	0.0
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	0.0
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	0.0
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	0.0
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	0.0

5 rows x 31 columns

We'll clean the data. Here the only column which doesn't have an impact on our analysis is Time. Therefore, we can drop that column.

In [4]:

```
data = data.drop(['Time'], axis=1)
data.head()
```

Out[4]:

	V1	V2	V3	V4	V5	V6	V7	V8	V9	V10	...	V21	V22	
0	1.359807	0.072781	2.536347	1.378155	0.338321	0.462388	0.239599	0.098698	0.363787	0.090794	...	0.018307	0.277838	0.0
1	1.191857	0.266151	0.166480	0.448154	0.060018	0.082361	0.078803	0.085102	0.255425	0.166974	...	0.225775	0.638672	0.0
2	1.358354	1.340163	1.773209	0.379780	0.503198	1.800499	0.791461	0.247676	1.514654	0.207643	...	0.247998	0.771679	0.0
3	0.966272	0.185226	1.792993	0.863291	0.010309	1.247203	0.237609	0.377436	1.387024	0.054952	...	0.108300	0.005274	0.0
4	1.158233	0.877737	1.548718	0.403034	0.407193	0.095921	0.592941	0.270533	0.817739	0.753074	...	0.009431	0.798278	0.0

5 rows x 30 columns

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

It's evident that the data is quite unbalanced since only 0.17% of the cases are fraudulent. However we'll apply our models without balancing it and if we don't get a good accuracy we can find a way to balance the dataset.

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
```

In [7]:

```
print('Amount details of valid transaction')
valid.Amount.describe()
```

Amount details of valid transaction

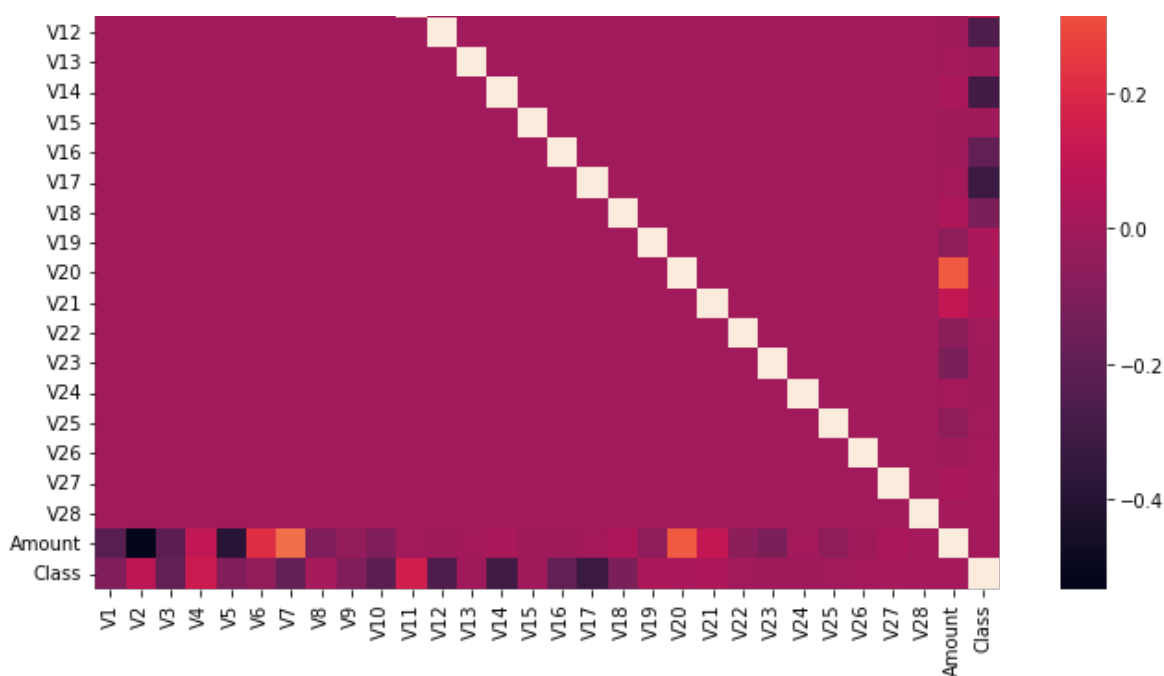
Out[7]:

```
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
```

In [8]:

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





In [9]:

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

```
(284807, 29)
(284807,)
```

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)
from sklearn.ensemble import RandomForestClassifier
rfc = RandomForestClassifier()
rfc.fit(xTrain, yTrain)
yPred = rfc.predict(xTest)
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 0.9995962220427653

The precision is 0.9746835443037974

The recall is 0.7857142857142857

The Matthews correlation coefficient is 0.8808846153846154

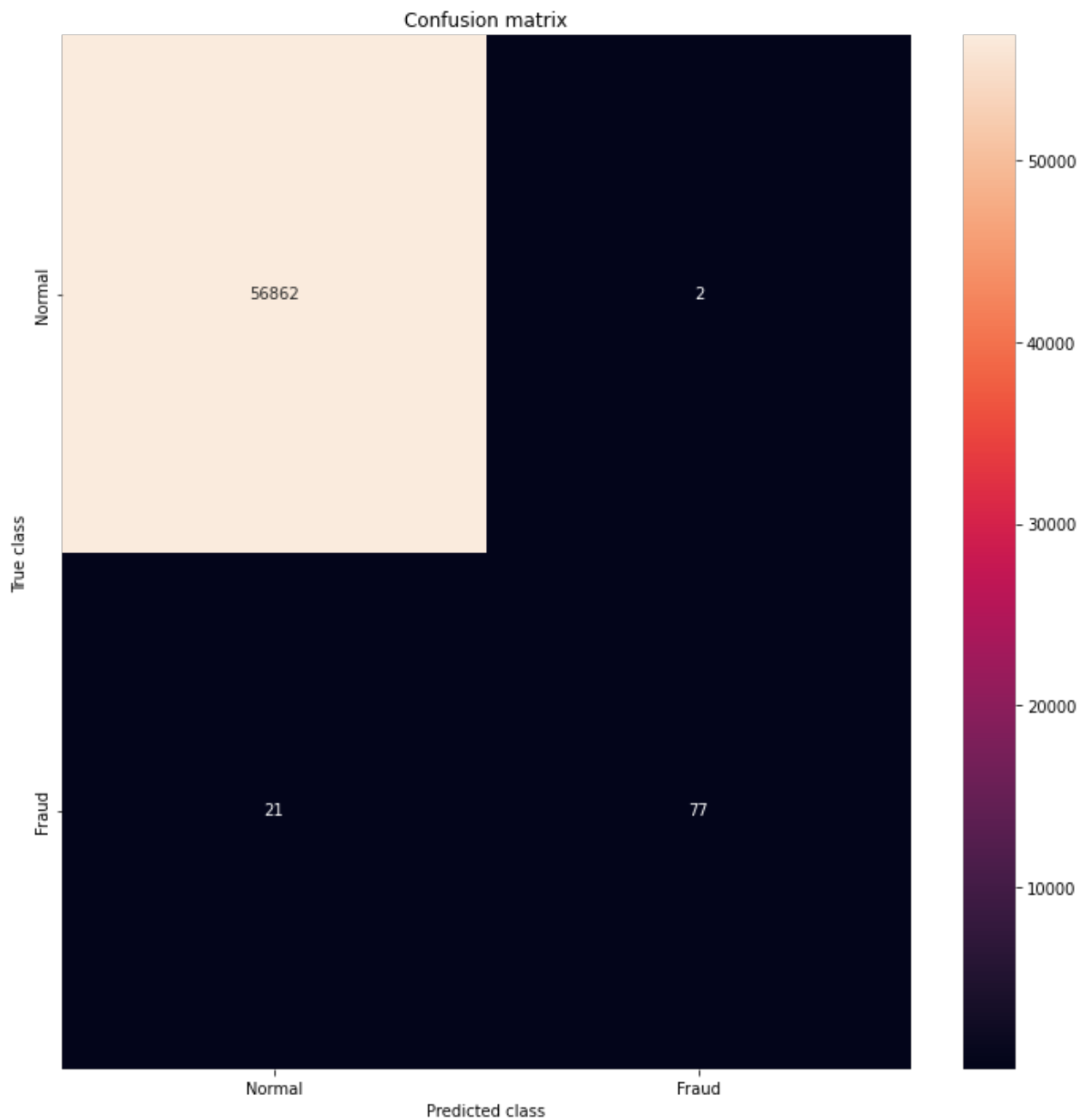
The F1-Score is 0.8700564971751412

The Matthews correlation coefficient is 0.8749276812909632

As the accuracy is high we don't need to balance our dataset.

In [11]:

```
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()
```



In []: