

```
In [1]: import pandas as pd
import numpy as np
import matplotlib.pyplot as plt
import seaborn as sns
%matplotlib inline
pd.set_option('display.max_columns', None)
```

```
In [2]: df=pd.read_csv("creditcard.csv")
```

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

```
Out[3]:
```

	Time	V1	V2	V3	V4	V5	V6	V7	V8	
0	0.0	-1.359807	-0.072781	2.536347	1.378155	-0.338321	0.462388	0.239599	0.098698	0.3
1	0.0	1.191857	0.266151	0.166480	0.448154	0.060018	-0.082361	-0.078803	0.085102	-0.2
2	1.0	-1.358354	-1.340163	1.773209	0.379780	-0.503198	1.800499	0.791461	0.247676	-1.5
3	1.0	-0.966272	-0.185226	1.792993	-0.863291	-0.010309	1.247203	0.237609	0.377436	-1.3
4	2.0	-1.158233	0.877737	1.548718	0.403034	-0.407193	0.095921	0.592941	-0.270533	0.8

```
In [4]: df.dtypes
```

```
Out[4]:
```

Time	float64
V1	float64
V2	float64
V3	float64
V4	float64
V5	float64
V6	float64
V7	float64
V8	float64
V9	float64
V10	float64
V11	float64
V12	float64
V13	float64
V14	float64
V15	float64
V16	float64
V17	float64
V18	float64
V19	float64
V20	float64
V21	float64
V22	float64
V23	float64
V24	float64
V25	float64
V26	float64
V27	float64
V28	float64
Amount	float64
Class	int64
dtype:	object

```
In [5]: df.Class.value_counts()
```

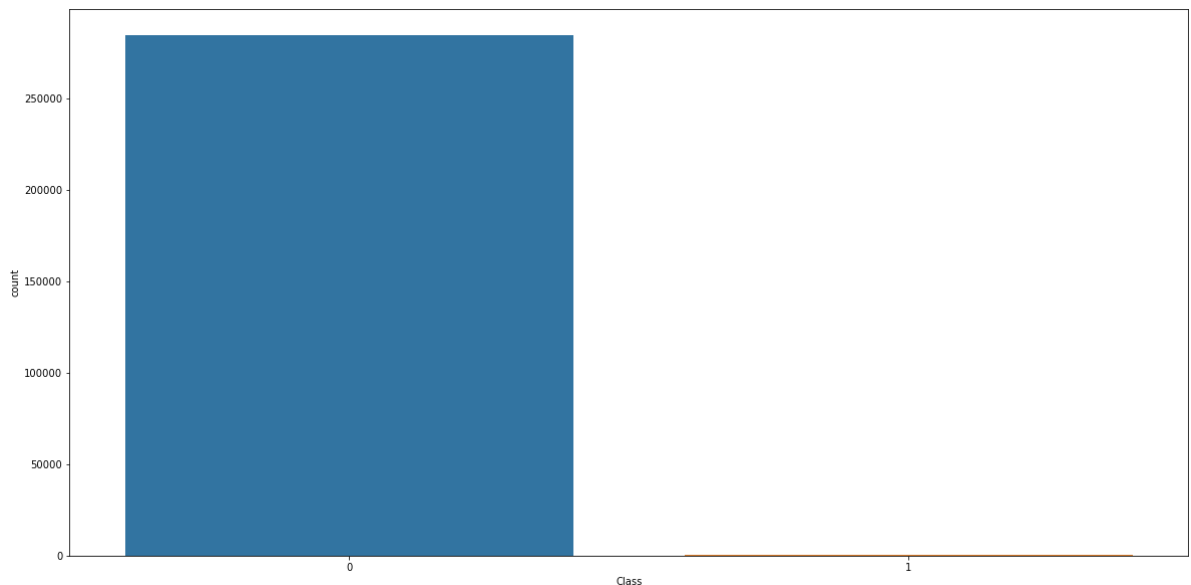
```
Out[5]: 0    284315
        1      492
        Name: Class, dtype: int64
```

0 Means not froud data 1 froud data

```
In [6]: plt.figure(figsize=(20,10))
        sns.countplot('Class',data=df)
```

C:\ProgramData\Anaconda3\lib\site-packages\seaborn_decorators.py:36: FutureWarning: Pass the following variable as a keyword arg: x. From version 0.12, the only valid positional argument will be `data`, and passing other arguments without an explicit keyword will result in an error or misinterpretation.

```
Out[6]: warnings.warn(
        <AxesSubplot:xlabel='Class', ylabel='count'>
```



```
In [7]: df.Class.value_counts()/df.Time.count()*100
```

```
Out[7]: 0    99.827251
        1     0.172749
        Name: Class, dtype: float64
```

We have less than 1% froud data sample to train the moder to predect the froud.

ML Model Selection

We have categorical dependent variable or Feature as 'Class' so i have decided to select the classification ML model.

Example

1. SVM (Support vector Machine)

```
In [8]: fraud = df[df['Class'] == 1]
        valid = df[df['Class'] == 0]
        outlierFraction = len(fraud)/float(len(valid))
        print(outlierFraction)
        print('Fraud Cases: {}'.format(len(df[df['Class'] == 1])))
        print('Valid Transactions: {}'.format(len(df[df['Class'] == 0])))
```

0.0017304750013189597
Fraud Cases: 492
Valid Transactions: 284315

```
In [9]: print('Amount details of the fraudulent transaction')  
        fraud.Amount.describe()
```

Amount details of the fraudulent transaction

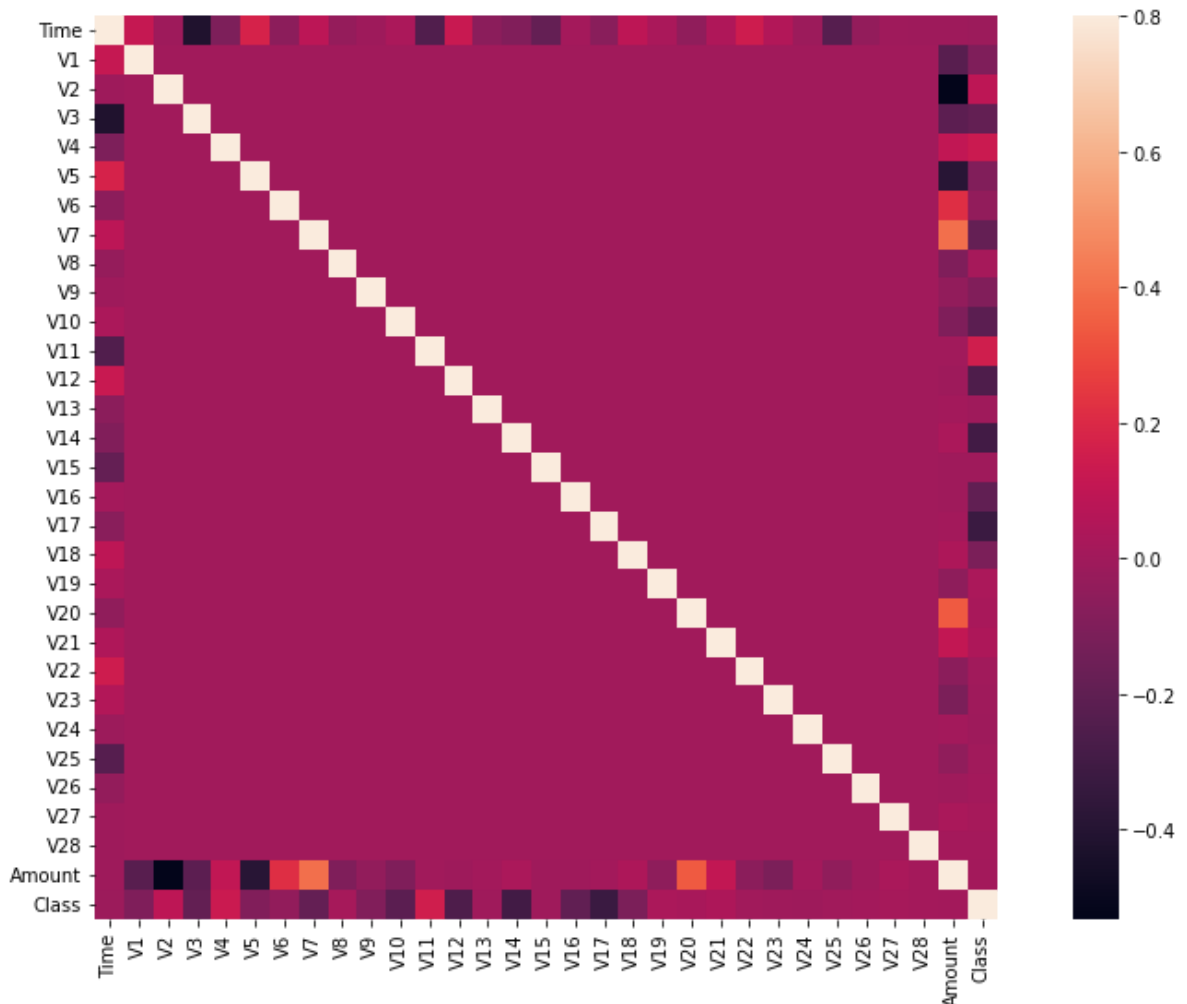
```
Out[9]: count      492.000000  
        mean       122.211321  
        std        256.683288  
        min         0.000000  
        25%         1.000000  
        50%         9.250000  
        75%        105.890000  
        max        2125.870000  
        Name: Amount, dtype: float64
```

```
In [10]: print('Amount details of the valid transaction')  
         valid.Amount.describe()
```

Amount details of the valid transaction

```
Out[10]: 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 [11]: # Correlation matrix  
         corrmatrix = df.corr()  
         fig = plt.figure(figsize = (15, 9))  
         sns.heatmap(corrmatrix, 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.

```
In [12]: # dividing the X and the Y from the dataset
X = df.drop(['Class'], axis = 1)
Y = df["Class"]
print(X.shape)
print(Y.shape)
```

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

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 [13]: # Using Scikit-Learn to split data into training and testing sets
from sklearn.model_selection import train_test_split

xTrain, xTest, yTrain, yTest = train_test_split(
    X, Y, test_size = 0.2, random_state = 42)
```

```
In [14]: # Building the Random Forest Classifier (RANDOM FOREST)
from sklearn.ensemble import RandomForestClassifier
```

```
# random forest model creation
rfc = RandomForestClassifier()
rfc.fit(xTrain, yTrain)
# predictions
yPred = rfc.predict(xTest)
```

```
In [15]: # Evaluating the classifier
# printing every score of the classifier
# scoring in anything
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 0.9995435553526912
The precision is 0.9615384615384616
The recall is 0.7653061224489796
The F1-Score is 0.8522727272727273
The Matthews correlation coefficient is0.8576194535617819
```

```
In [ ]:
```