

**Date: 17-10-2023**

**Project Title: Credit Card Fraudlent Detection**

**Team ID: 3890**

## Importing required libraries

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

## Set the jupyter notebook to show maximum number of columns

```
In [2]: pd.options.display.max_columns = None
```

## Loading the datasets

```
In [3]: ccfd = pd.read_excel("C:\\Users\\Mohamed Safthar\\OneDrive\\Documents\\IBM\\creditcard.csv\\CreditCardFraudDataset.xlsx")
```

## Displaying top 5 rows

```
In [4]: ccfd.head()
```

```
Out[4]:
```

	Time	V1	V2	V3	V4	V5	V6	V7	V8	V9	V10	V11	V12	V13
0	0	-1.359807	-0.072781	2.536347	1.378155	-0.338321	0.462388	0.239599	0.098698	0.363787	0.090794	-0.551600	-0.617801	-0.991390
1	0	1.191857	0.266151	0.166480	0.448154	0.060018	-0.082361	-0.078803	0.085102	-0.255425	-0.166974	1.612727	1.065235	0.489095
2	1	-1.358354	-1.340163	1.773209	0.379780	-0.503198	1.800499	0.791461	0.247676	-1.514654	0.207643	0.624501	0.066084	0.717293
3	1	-0.966272	-0.185226	1.792993	-0.863291	-0.010309	1.247203	0.237609	0.377436	-1.387024	-0.054952	-0.226487	0.178228	0.507757
4	2	-1.158233	0.877737	1.548718	0.403034	-0.407193	0.095921	0.592941	-0.270533	0.817739	0.753074	-0.822843	0.538196	1.345852

## Displaying bottom 5 rows

```
In [5]: ccfd.tail()
```

```
Out[5]:
```

	Time	V1	V2	V3	V4	V5	V6	V7	V8	V9	V10	V11	V12
284802	172786	-11.881118	10.071785	-9.834783	-2.066656	-5.364473	-2.606837	-4.918215	7.305334	1.914428	4.356170	-1.593105	2.711941
284803	172787	-0.732789	-0.055080	2.035030	-0.738589	0.868229	1.058415	0.024330	0.294869	0.584800	-0.975926	-0.150189	0.915802
284804	172788	1.919565	-0.301254	-3.249640	-0.557828	2.630515	3.031260	-0.296827	0.708417	0.432454	-0.484782	0.411614	0.063119
284805	172788	-0.240440	0.530483	0.702510	0.689799	-0.377961	0.623708	-0.686180	0.679145	0.392087	-0.399126	-1.933849	-0.962886
284806	172792	-0.533413	-0.189733	0.703337	-0.506271	-0.012546	-0.649617	1.577006	-0.414650	0.486180	-0.915427	-1.040458	-0.031513

## Shows number of rows and columns

```
In [7]: print("Number of rows in given dataset ",ccfd.shape[0])  
        print("Number of columns in the given dataset ",ccfd.shape[1])
```

```
Number of rows in given dataset  284807  
Number of columns in the given dataset  31
```

## Getting basis information

In [8]:

```
ccfd.info()
```

```

<class 'pandas.core.frame.DataFrame'>
RangeIndex: 284807 entries, 0 to 284806
Data columns (total 31 columns):
#   Column  Non-Null Count  Dtype
---  -
0   Time    284807 non-null  int64
1   V1      284807 non-null  float64
2   V2      284807 non-null  float64
3   V3      284807 non-null  float64
4   V4      284807 non-null  float64
5   V5      284807 non-null  float64
6   V6      284807 non-null  float64
7   V7      284807 non-null  float64
8   V8      284807 non-null  float64
9   V9      284807 non-null  float64
10  V10     284807 non-null  float64
11  V11     284807 non-null  float64
12  V12     284807 non-null  float64
13  V13     284807 non-null  float64
14  V14     284807 non-null  float64
15  V15     284807 non-null  float64
16  V16     284807 non-null  float64
17  V17     284807 non-null  float64
18  V18     284807 non-null  float64
19  V19     284807 non-null  float64
20  V20     284807 non-null  float64
21  V21     284807 non-null  float64
22  V22     284807 non-null  float64
23  V23     284807 non-null  float64
24  V24     284807 non-null  float64
25  V25     284807 non-null  float64
26  V26     284807 non-null  float64
27  V27     284807 non-null  float64
28  V28     284807 non-null  float64
29  Amount  284807 non-null  float64
30  Class   284807 non-null  int64
dtypes: float64(29), int64(2)
memory usage: 67.4 MB

```

## Checking null values in the given data

```
In [9]: ccfd.isnull().sum()
```

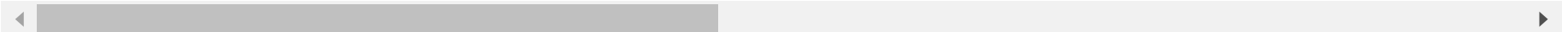
```
Out[9]: Time      0  
V1      0  
V2      0  
V3      0  
V4      0  
V5      0  
V6      0  
V7      0  
V8      0  
V9      0  
V10     0  
V11     0  
V12     0  
V13     0  
V14     0  
V15     0  
V16     0  
V17     0  
V18     0  
V19     0  
V20     0  
V21     0  
V22     0  
V23     0  
V24     0  
V25     0  
V26     0  
V27     0  
V28     0  
Amount   0  
Class    0  
dtype: int64
```

## Scaling the Amount features, removing the independent columns

```
In [10]: #removing the column name Time, it is unnecessary to our training purposes
ccfd.head(2)
```

```
Out[10]:
```

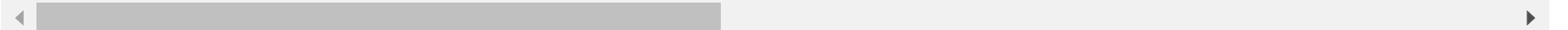
	Time	V1	V2	V3	V4	V5	V6	V7	V8	V9	V10	V11	V12	V13
0	0	-1.359807	-0.072781	2.536347	1.378155	-0.338321	0.462388	0.239599	0.098698	0.363787	0.090794	-0.551600	-0.617801	-0.991390
1	0	1.191857	0.266151	0.166480	0.448154	0.060018	-0.082361	-0.078803	0.085102	-0.255425	-0.166974	1.612727	1.065235	0.489095



```
In [13]: #time features is unnecessary here
ccfd.drop('Time',axis=1,inplace=True).head()
```

```
Out[13]:
```

	V1	V2	V3	V4	V5	V6	V7	V8	V9	V10	V11	V12	V13	
0	-1.359807	-0.072781	2.536347	1.378155	-0.338321	0.462388	0.239599	0.098698	0.363787	0.090794	-0.551600	-0.617801	-0.991390	-0.31
1	1.191857	0.266151	0.166480	0.448154	0.060018	-0.082361	-0.078803	0.085102	-0.255425	-0.166974	1.612727	1.065235	0.489095	-0.14
2	-1.358354	-1.340163	1.773209	0.379780	-0.503198	1.800499	0.791461	0.247676	-1.514654	0.207643	0.624501	0.066084	0.717293	-0.16
3	-0.966272	-0.185226	1.792993	-0.863291	-0.010309	1.247203	0.237609	0.377436	-1.387024	-0.054952	-0.226487	0.178228	0.507757	-0.28
4	-1.158233	0.877737	1.548718	0.403034	-0.407193	0.095921	0.592941	-0.270533	0.817739	0.753074	-0.822843	0.538196	1.345852	-1.11



```
In [17]: ccfdf.head()
```

```
Out[17]:
```

	V1	V2	V3	V4	V5	V6	V7	V8	V9	V10	V11	V12	V13	
0	-1.359807	-0.072781	2.536347	1.378155	-0.338321	0.462388	0.239599	0.098698	0.363787	0.090794	-0.551600	-0.617801	-0.991390	-0.31
1	1.191857	0.266151	0.166480	0.448154	0.060018	-0.082361	-0.078803	0.085102	-0.255425	-0.166974	1.612727	1.065235	0.489095	-0.14
2	-1.358354	-1.340163	1.773209	0.379780	-0.503198	1.800499	0.791461	0.247676	-1.514654	0.207643	0.624501	0.066084	0.717293	-0.16
3	-0.966272	-0.185226	1.792993	-0.863291	-0.010309	1.247203	0.237609	0.377436	-1.387024	-0.054952	-0.226487	0.178228	0.507757	-0.28
4	-1.158233	0.877737	1.548718	0.403034	-0.407193	0.095921	0.592941	-0.270533	0.817739	0.753074	-0.822843	0.538196	1.345852	-1.11

## Scaling the Amount column data

```
In [18]: from sklearn.preprocessing import StandardScaler
```

```
In [19]: ss = StandardScaler()
```

```
In [23]: ccfdf['Amounts'] = ss.fit_transform(pd.DataFrame(ccfdf['Amount']))
```

```
In [24]: ccfdf.head()
```

```
Out[24]:
```

	V1	V2	V3	V4	V5	V6	V7	V8	V9	V10	V11	V12	V13	
0	-1.359807	-0.072781	2.536347	1.378155	-0.338321	0.462388	0.239599	0.098698	0.363787	0.090794	-0.551600	-0.617801	-0.991390	-0.31
1	1.191857	0.266151	0.166480	0.448154	0.060018	-0.082361	-0.078803	0.085102	-0.255425	-0.166974	1.612727	1.065235	0.489095	-0.14
2	-1.358354	-1.340163	1.773209	0.379780	-0.503198	1.800499	0.791461	0.247676	-1.514654	0.207643	0.624501	0.066084	0.717293	-0.16
3	-0.966272	-0.185226	1.792993	-0.863291	-0.010309	1.247203	0.237609	0.377436	-1.387024	-0.054952	-0.226487	0.178228	0.507757	-0.28
4	-1.158233	0.877737	1.548718	0.403034	-0.407193	0.095921	0.592941	-0.270533	0.817739	0.753074	-0.822843	0.538196	1.345852	-1.11



```
In [25]: ccfd.shape
```

```
Out[25]: (284807, 31)
```

```
In [26]: ccfd.drop('Amount',axis=1,inplace=True)
```

```
In [27]: ccfd.shape
```

```
Out[27]: (284807, 30)
```

## Dropping the duplicate records

```
In [31]: ccfd.duplicated().any()
```

```
Out[31]: True
```

```
In [34]: ccfd.drop_duplicates(inplace=True)
```

```
In [35]: ccfd.shape
```

```
Out[35]: (275663, 30)
```

```
In [36]: 284807 - 275663
```

```
Out[36]: 9144
```

## Exploring Class columns

```
In [38]: ccfd['Class'].unique()
```

```
Out[38]: array([0, 1], dtype=int64)
```

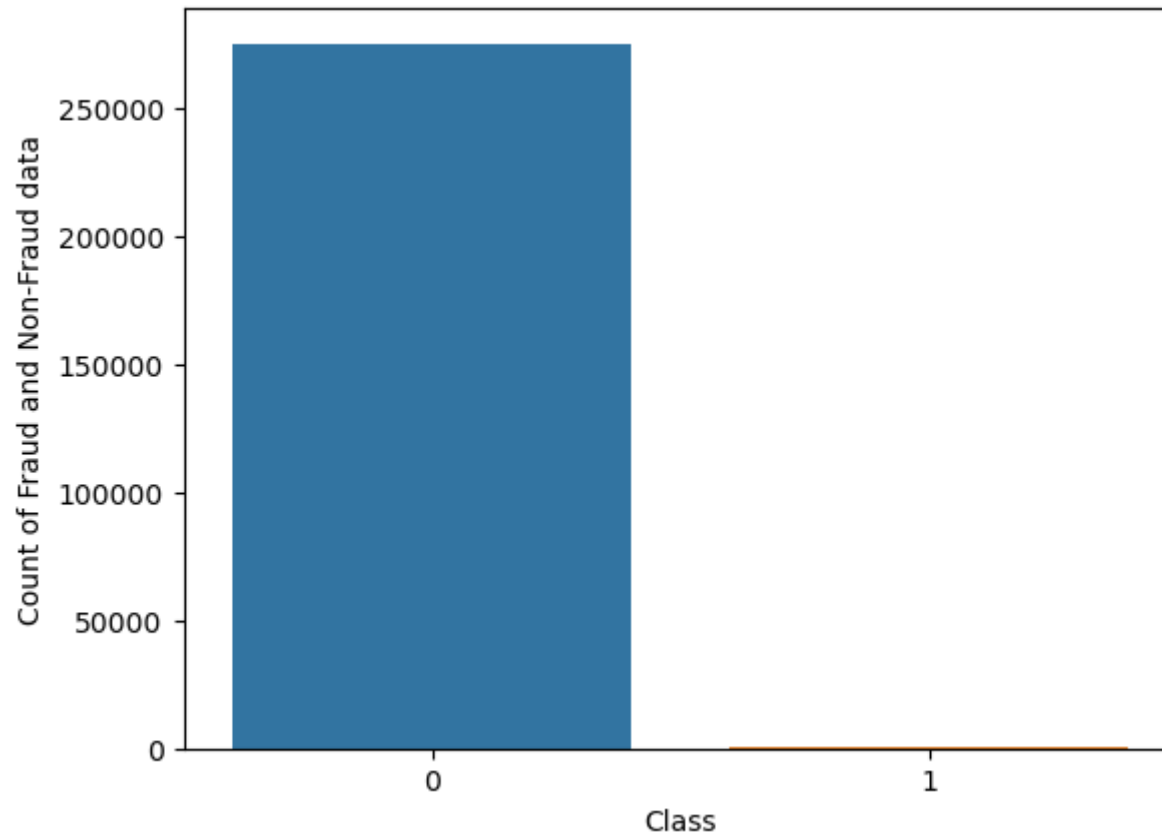
```
In [39]: ccfd['Class'].nunique()
```

```
Out[39]: 2
```

```
In [40]: ccfd['Class'].value_counts()
```

```
Out[40]: Class
0      275190
1         473
Name: count, dtype: int64
```

```
In [46]: #visualizing the distribution of 0 and 1 using seaborn countplot
sns.countplot(ccfd,x = ccfd['Class'])
plt.xlabel('Class')
plt.ylabel('Count of Fraud and Non-Fraud data')
plt.show()
```



From the above information, We can say that our data is high imbalanced, so need to apply oversampling and undersampling technique to train our model

```
In [ ]:
```