

## Assignment - 2

Assignment Date	21 September 2022
Student Name	MANCHINELLA KUMUD BANDHAV
Student Roll Number	111519104075
Maximum Marks	2 Marks

1. Download the dataset: Dataset

2. Load the dataset.

```
[1] import pandas as pd
import numpy as np
```

```
[2] file=pd.read_csv("/content/Churn_Modelling.csv")
df=pd.DataFrame(file)
df.head()
```

	RowNumber	CustomerId	Surname	CreditScore	Geography	Gender	Age	Tenure	Balance	NumOfProducts	HasCrCard	IsActiveMember	EstimatedSalary	Exited
0	1	15634602	Hargrave	619	France	Female	42	2	0.00	1	1	1	101348.88	1
1	2	15647311	Hill	608	Spain	Female	41	1	83807.86	1	0	1	112542.58	0
2	3	15619304	Onio	502	France	Female	42	8	159660.80	3	1	0	113931.57	1
3	4	15701354	Boni	699	France	Female	39	1	0.00	2	0	0	93826.63	0
4	5	15737888	Mitchell	850	Spain	Female	43	2	125510.82	1	1	1	79084.10	0

```
[3] df['HasCrCard'] = df['HasCrCard'].astype('category')
```

```
[4] df['IsActiveMember'] = df['IsActiveMember'].astype('category')
df['Exited'] = df['Exited'].astype('category')
```

```
[5] df = df.drop(columns=['RowNumber', 'CustomerId', 'Surname'])
```

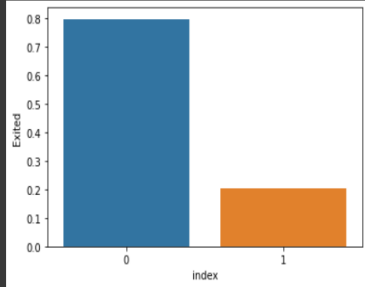
```
[6] df.head()
```

	CreditScore	Geography	Gender	Age	Tenure	Balance	NumOfProducts	HasCrCard	IsActiveMember	EstimatedSalary	Exited
0	619	France	Female	42	2	0.00	1	1	1	101348.88	1
1	608	Spain	Female	41	1	83807.86	1	0	1	112542.58	0
2	502	France	Female	42	8	159660.80	3	1	0	113931.57	1
3	699	France	Female	39	1	0.00	2	0	0	93826.63	0
4	850	Spain	Female	43	2	125510.82	1	1	1	79084.10	0

3. Perform Below Visualizations: Univariate Analysis, Bi - Variate Analysis, Multi - Variate Analysis

```
[7] import seaborn as sns
density = df['Exited'].value_counts(normalize=True).reset_index()
sns.barplot(data=density, x='index', y='Exited', );
density
```

	index	Exited
0	0	0.7963
1	1	0.2037



the data is significantly imbalanced

```
[8] import matplotlib.pyplot as plt
```

```
[9] categorical = df.drop(columns=['CreditScore', 'Age', 'Tenure', 'Balance', 'EstimatedSalary'])
rows = int(np.ceil(categorical.shape[1] / 2)) - 1

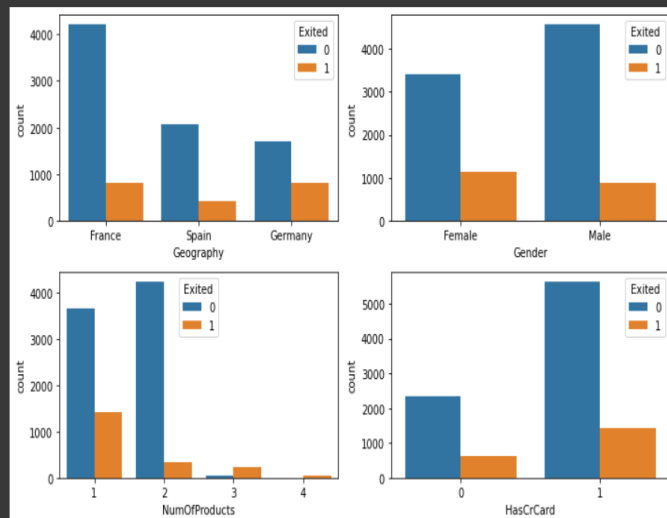
# create sub-plots and title them
fig, axes = plt.subplots(nrows=rows, ncols=2, figsize=(10,6))
```

+ Code + Text

```
[9] for col in range(cols):
    col_name = categorical.columns[2 * row + col]
    ax = axes[row*2 + col]

    sns.countplot(data=categorical, x=col_name, hue="Exited", ax=ax);

plt.tight_layout()
```



[10] df.info()

<class 'pandas.core.frame.DataFrame'>  
RangeIndex: 10000 entries, 0 to 9999  
Data columns (total 11 columns):  
# Column Non-Null Count Dtype  
---  
0 CreditScore 10000 non-null int64  
1 Geography 10000 non-null object  
2 Gender 10000 non-null object  
3 Age 10000 non-null int64  
4 Tenure 10000 non-null int64  
5 Balance 10000 non-null float64  
6 NumOfProducts 10000 non-null int64  
7 HasCrCard 10000 non-null category  
8 IsActiveMember 10000 non-null category  
9 EstimatedSalary 10000 non-null float64  
10 Exited 10000 non-null category  
dtypes: category(3), float64(2), int64(4), object(2)  
memory usage: 654.8+ KB

[11] df.describe()

	CreditScore	Age	Tenure	Balance	NumOfProducts	EstimatedSalary
count	10000.000000	10000.000000	10000.000000	10000.000000	10000.000000	10000.000000
mean	650.528800	38.921800	5.012800	76485.889288	1.530200	100090.239881
std	96.653299	10.487806	2.892174	62397.405202	0.581654	57510.492818
min	350.000000	18.000000	0.000000	0.000000	1.000000	11.580000
max	850.000000	92.000000	10.000000	250898.090000	4.000000	199992.480000

[11]

min	350.000000	18.000000	0.000000	0.000000	1.000000	11.580000
25%	584.000000	32.000000	3.000000	0.000000	1.000000	51002.110000
50%	652.000000	37.000000	5.000000	97198.540000	1.000000	100193.915000
75%	718.000000	44.000000	7.000000	127644.240000	2.000000	149388.247500
max	850.000000	92.000000	10.000000	250898.090000	4.000000	199992.480000

5. Handle the Missing values.

[12] df.isna().sum()

CreditScore 0  
Geography 0  
Gender 0  
Age 0  
Tenure 0  
Balance 0  
NumOfProducts 0  
HasCrCard 0  
IsActiveMember 0  
EstimatedSalary 0  
Exited 0  
dtype: int64

there is no missing values in dataset