

ASSIGNMENT-02

DATA VISUALIZATION AND PRE PROCESSING

Assignment Date: 22 September 2022
Student Name: BHOO MIHA.M
Student Roll Number: 113219071003
Maximum Marks: 2 Marks

1. Download the dataset: Dataset
Dataset downloaded in csv form.

2. Load the dataset.

```
import pandas as pd
df = pd.read_csv("/content/drive/MyDrive/IBM Assignments/Churn_Modelling.csv")
```

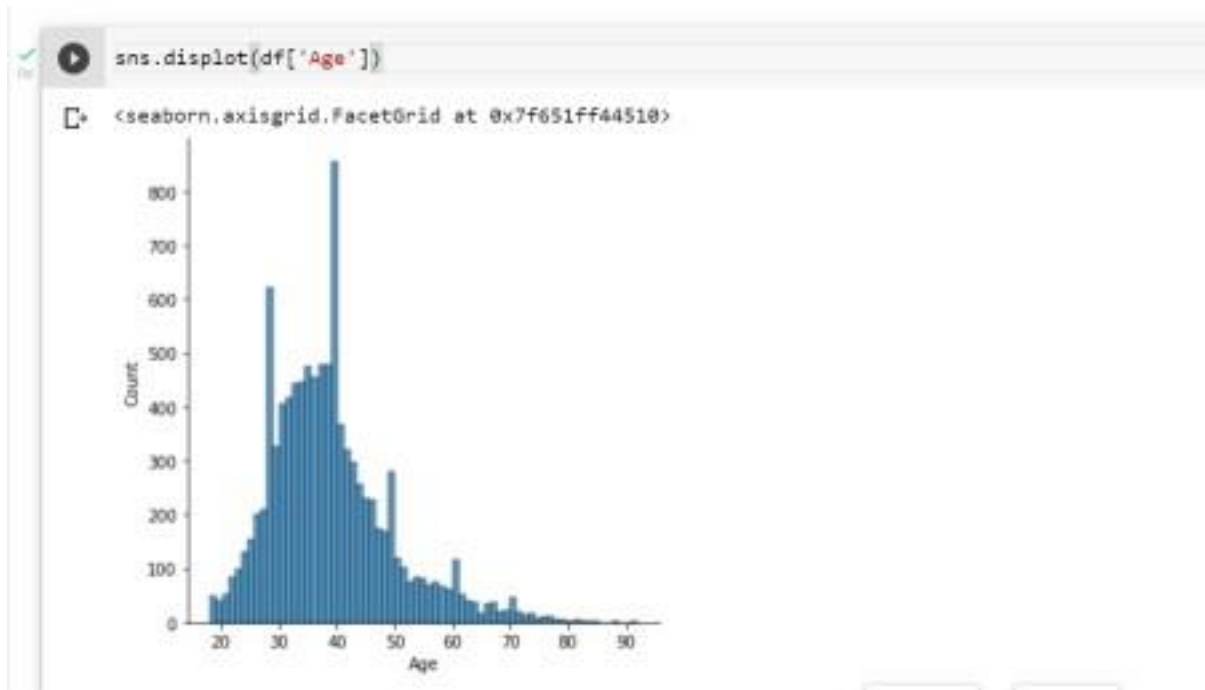
```
import pandas as pd
df = pd.read_csv("/content/drive/MyDrive/IBM Assignments/Churn_Modelling.csv")
```

3. Perform Below Visualizations.

- Univariate Analysis

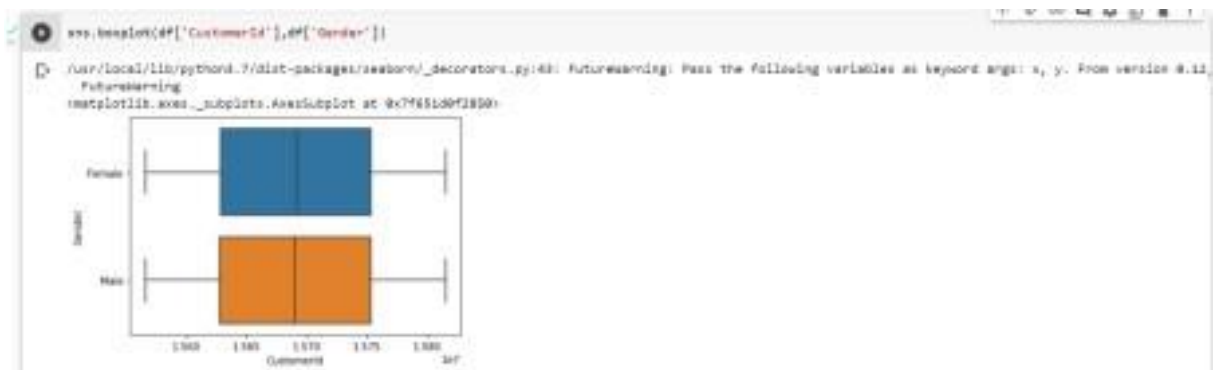
```
sns.displot(df['Age'])
```

```
[2] import matplotlib.pyplot as plt
    %matplotlib inline
    import seaborn as sns
```

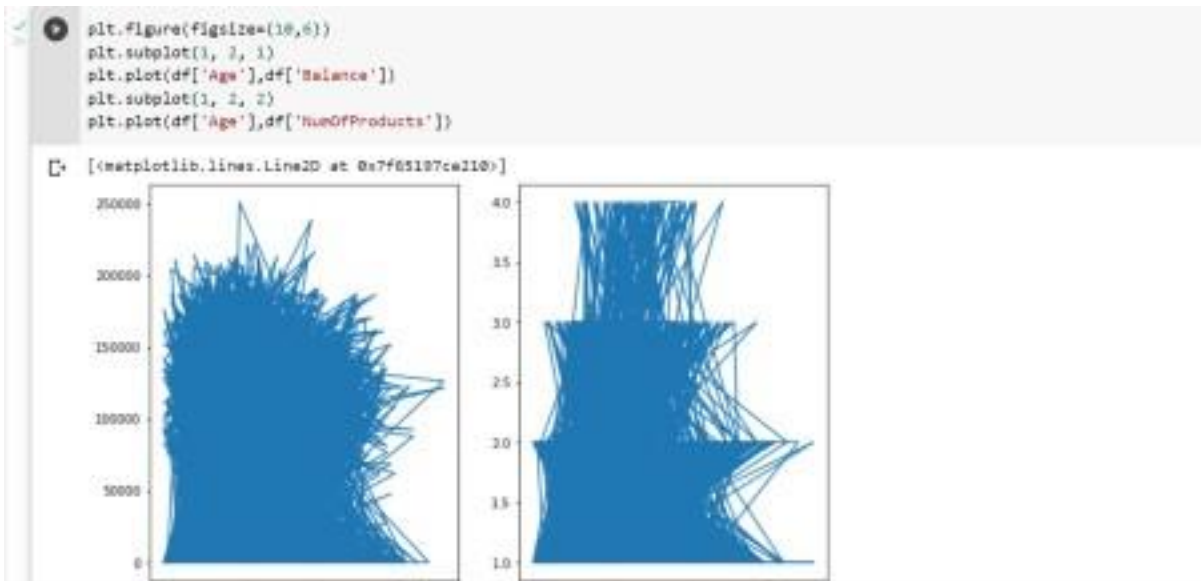


- Bi - Variate Analysis

```
sns.boxplot(df['CustomerId'],df['Gender'])
```



- Multi - Variate Analysis



4. Perform descriptive statistics on the dataset.

```
df.describe()
```

	RowNumber	CustomerId	CreditScore	Age	Tenure	Balance	NumOfProducts	HasCrCard	IsActiveMember	EstimatedSalary
count	10000.00000	1.000000e+04	10000.000000	10000.000000	10000.000000	10000.000000	10000.000000	10000.000000	10000.000000	10
mean	5000.50000	1.500000e+03	650.526600	38.021800	5.812800	76485.860288	1.530200	0.70580	0.515100	100066.236661
std	2886.89000	7.100019e+04	95.953288	10.467800	2.882174	82307.405202	0.581054	0.40584	0.499707	57510.482818
min	1.00000	1.000000e+03	350.000000	18.000000	0.000000	0.000000	1.000000	0.000000	0.000000	11.580000
25%	2500.75000	1.902853e+07	594.000000	32.000000	3.000000	0.000000	1.000000	0.000000	0.000000	81002.110000
50%	5000.50000	1.500074e+07	652.000000	37.000000	5.000000	97196.540000	1.000000	1.000000	1.000000	100193.915000
75%	7500.25000	1.575323e+07	716.000000	44.000000	7.000000	127944.240000	2.000000	1.000000	1.000000	140388.247500
max	10000.00000	1.581589e+07	850.000000	50.000000	10.000000	253996.090000	4.000000	1.000000	1.000000	199062.480000

Mean:

```
df.mean()
```

```
/usr/local/lib/python3.7/dist-packages/ipykernel_launcher.py:1: FutureWarning: Dropping of nuisance columns in DataFrame reductions (with 'numeric_only')
  """Entry point for launching an IPython kernel.
```

RowNumber	5.000000e+03
CustomerId	1.500000e+03
CreditScore	6.505266e+02
Age	3.802180e+01
Tenure	5.812800e+00
Balance	7.648586e+04
NumOfProducts	1.530200e+00
HasCrCard	7.058000e-01
IsActiveMember	5.151000e-01
EstimatedSalary	1.000662e+05
Exited	2.037000e-01
dtype: float64	

5. Handle the Missing values.

```
df.isnull().sum()

RowNumber      0
CustomerId     0
Surname        0
CreditScore    0
Geography      0
Gender         0
Age            0
Tenure         0
Balance        0
NumOfProducts 0
HasCrCard      0
IsActiveMember 0
EstimatedSalary 0
Exited         0
dtype: int64
```

6. Find the outliers and replace the outliers

Finding Outliers:

Using Boxplot



Using method

```
[85] qnt = df.quantile(q=[0.25,0.75])
qnt
```

	RowNumber	CustomerId	CreditScore	Age	Tenure	Balance	NumOfProducts	HasCrCard	IsActiveMember	EstimatedSalary	Exited
0.25	2500.75	10028528.25	584.0	32.0	3.0	0.00	1.0	0.0	0.0	51002.1100	0.0
0.75	7500.25	15753233.75	718.0	44.0	7.0	127644.24	2.0	1.0	1.0	149388.2475	0.0

```

iqr = qnt.loc[0.75]-qnt.loc[0.25]
iqr
```

	RowNumber	CustomerId	CreditScore	Age	Tenure	Balance	NumOfProducts	HasCrCard	IsActiveMember	EstimatedSalary	Exited
	4999.5000	124705.5000	134.0000	12.0000	4.0000	127644.2400	1.0000	1.0000	1.0000	98386.1375	0.0000

dtype: float64

```

lower = qnt.loc[0.25]-1.5*Iqr
print("Lower bound:",lower)
upper = qnt.loc[0.75]+1.5*Iqr
print("Upper bound:",upper)

Lower bound: RowNumber      -2.888888e+00
CustomerId      1.534147e+07
CreditScore      5.838888e+02
Age      1.488888e+01
Tenure      -3.888888e+00
Balance      -1.914654e+05
NumOfProducts      5.000000e-01
HasCrCard      -1.500000e+00
IsActiveMember      -1.500000e+00
EstimatedSalary      -9.657710e+04
Exited      0.000000e+00
dtype: float64

Upper bound: RowNumber      1.488888e+04
CustomerId      1.594029e+07
CreditScore      9.190000e+02
Age      6.100000e+01
Tenure      1.300000e+01
Balance      5.191106e+05
NumOfProducts      5.000000e+00
HasCrCard      2.500000e+00
IsActiveMember      2.500000e+00
EstimatedSalary      2.909675e+05
Exited      0.000000e+00
dtype: float64

```

Replacing Outliers:

```

''' replacing outliers '''
df['Balance'] = np.where(df['Balance']>127644,0.00,df['Balance'])

```

7. Check for Categorical columns and perform encoding.

Categorical columns: Geography,Gender

```

[98]: from sklearn.preprocessing import LabelEncoder
LabelEncoder_df = LabelEncoder()
df['Geography'] = LabelEncoder_df.fit_transform(df['Geography'])

df.head()

```

	RowNumber	CustomerId	Surname	CreditScore	Geography	Gender	Age	Tenure	Balance	NumOfProducts	HasCrCard	IsActiveMember	EstimatedSalary
0	1	15634602	Hargrave	819	0	Female	42	2	0.00	1	1	1	101348.88
1	2	15647311	Hill	688	2	Female	41	1	83807.86	1	0	1	112542.58
2	3	15618304	Oric	502	0	Female	42	8	0.00	3	1	0	113631.57
3	4	15701354	Bori	688	0	Female	39	1	0.00	2	0	0	90829.63
4	5	15737888	Michael	890	2	Female	43	2	125510.82	1	1	1	79064.10

```

df['Gender'] = LabelEncoder_df.fit_transform(df['Gender'])

df.head(7)

```

	RowNumber	CustomerId	Surname	CreditScore	Geography	Gender	Age	Tenure	Balance	NumOfProducts	HasCrCard	IsActiveMember	EstimatedSalary
0	1	15634602	Hargrave	819	0	0	42	2	0.00	1	1	1	101348.88
1	2	15647311	Hill	688	2	0	41	1	83807.86	1	0	1	112542.58
2	3	15618304	Oric	502	0	0	42	8	0.00	3	1	0	113631.57
3	4	15701354	Bori	688	0	0	39	1	0.00	2	0	0	90829.63
4	5	15737888	Michael	890	2	0	43	2	125510.82	1	1	1	79064.10
5	6	15574612	Chu	845	2	1	44	6	113755.78	2	1	0	146756.71
6	7	15592631	Barbott	622	0	1	39	7	0.00	2	1	1	10062.80

8. Split the data into dependent and independent variables.



9. Scale the independent variables



10. Split the data into training and testing

