# Assignment 2

| Assignment Date | 21.09.2022 |
|---|---|
| Student Name | V.Srividhya |
| Student Roll Number | 2019115106 |
| Maximum Marks | 2 Marks |

**Question-1:**

Download Dataset

**Question-2:**

Loading the dataset

**Solution:**
```
import pandas as pd
df=pd.read_csv('Churn_Modelling.csv')
```

**Screenshot:**

```
[ ] import pandas as pd
    df=pd.read_csv('Churn_Modelling.csv')
```

**Question-3.1:**

Univariate Analysis

**Solution:**
```
df.dtypes
df.head()
df.shape
import matplotlib.pyplot as plt
import seaborn as sns
plt.scatter(df.index,df['Geography'])
plt.scatter(df.index,df['NumOfProducts'])
plt.scatter(df.index,df['Age'])
import numpy as np
print('Number of users who have a credit card ',np.sum(df['HasCrCard']==1))
print('Number of users who dont have a credit card ',np.sum(df['HasCrCard']==0))
plt.hist(df['HasCrCard'])
```
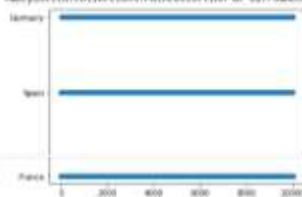
**Screenshot:**



```
df.dtypes
```
```
RowNumber          int64
CustomerId         int64
Surname            object
CreditScore        int64
Geography          object
Gender             object
Age                int64
Tenure             int64
Balance            float64
NumOfProducts      int64
HasCrCard          int64
IsActiveMember     int64
EstimatedSalary    float64
Exited             int64
dtype: object
```

```
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 |

```
df.shape
```
```
(10000, 14)
```

```
df.isnull().value_counts()
```
```
RowNumber  CustomerId  Surname  CreditScore  Geography  Gender  Age  Tenure  Balance  NumOfProducts  HasCrCard  IsActiveMember  EstimatedSalary  Exited
False      False       False    False        False      False   False False   False    False          False      False           False            False    10000
dtype: int64
```

```
import matplotlib.pyplot as plt
import seaborn as sns
```

```
plt.scatter(df.index,df['Geography'])
```
```
<matplotlib.collections.PathCollection at 0x7f8d84ee91b6>
```



```
plt.scatter(df.index,df['NumOfProducts'])
```
```
<matplotlib.collections.PathCollection at 0x7f9d84001456>
```



```
plt.scatter(df.index,df['Age'])
```
```
<matplotlib.collections.PathCollection at 0x7f8d84edbfd6>
```



● Os  completed at 10:01 PM

,



```
import numpy as np
print('Number of users who have a credit card ',np.sum(df['HasCrCard']==1))
print('Number of users who dont have a credit card ',np.sum(df['HasCrCard']==0))
```
```
Number of users who have a credit card  7055
Number of users who dont have a credit card  2945
```

```
plt.hist(df['HasCrCard'])
```
```
(array([2945.,    0.,    0.,    0.,    0.,    0.,    0.,    0.,    0.,
        7055.]),
 array([0. , 0.1, 0.2, 0.3, 0.4, 0.5, 0.6, 0.7, 0.8, 0.9, 1. ]),
 <a list of 10 Patch objects>)
```
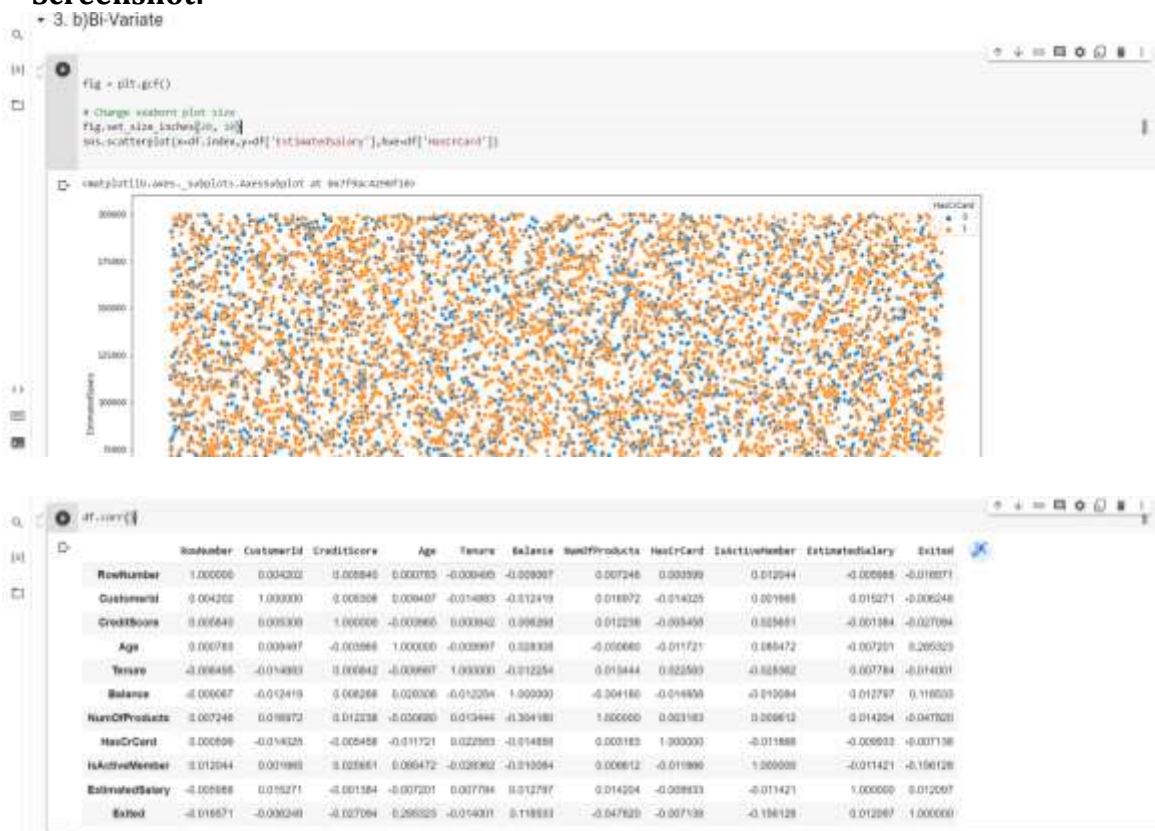
## Question-3.2:

Bi-Variate

### Solution:

```
fig = plt.gcf()

# Change seaborn plot size
fig.set_size_inches(20, 10)
sns.scatterplot(x=df.index,y=df['EstimatedSalary'],hue=df['HasCrCard'])

df.corr()
```
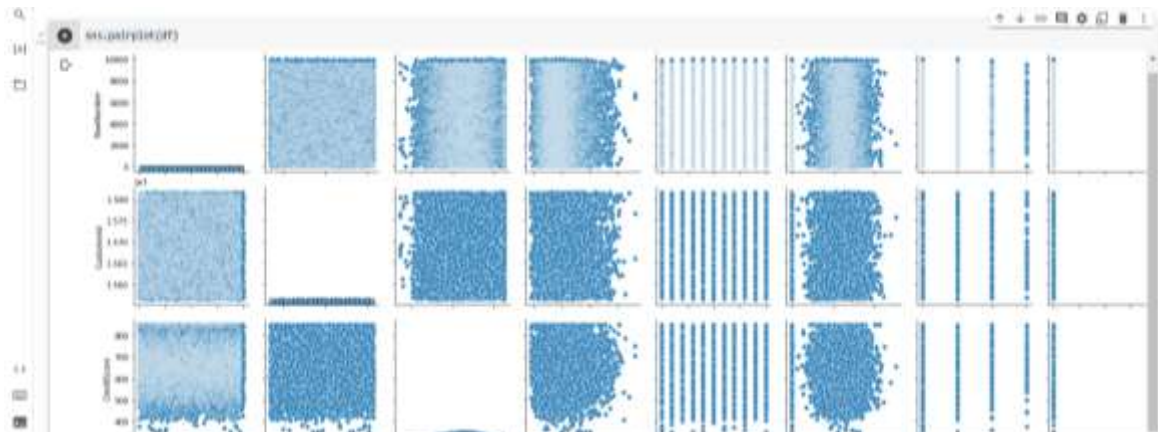
### Screenshot:



## Question-3.3:

Multi Variate

### Solution:

```
sns.pairplot(df)
```

**Screenshot:**



### Question-4:

Perform descriptive statistics on the dataset.

**Solution:**
df.describe()

**Screenshot:**



### Question-5:

Handle the Missing values.

**Solution:**
df.isnull().value_counts()
df.notnull().value_counts()

**Screenshot:**

**Question-6:**

Find the outliers and replace the outliers

**Solution:**

df.columns

df.boxplot('HasCrCard')

numeric_col = ['Balance','CreditScore','Tenure','Age','IsActiveMember']

df.boxplot('Balance')

df.boxplot('CreditScore')
df.boxplot('Tenure')
df.boxplot('Age')
df.boxplot('IsActiveMember')
for x in ['CreditScore']:
  q75,q25 = np.percentile(df.loc[:,x],[75,25])
  intr_qr = q75-q25

  max = q75+(1.5*intr_qr)
  min = q25-(1.5*intr_qr)

  df.loc[df[x] < min,x] = np.nan
  df.loc[df[x] > max,x] = np.nan
for x in ['Age']:
  q75,q25 = np.percentile(df.loc[:,x],[75,25])
  intr_qr = q75-q25

  max = q75+(1.5*intr_qr)
  min = q25-(1.5*intr_qr)

  df.loc[df[x] < min,x] = np.nan
  df.loc[df[x] > max,x] = np.nan
df.boxplot('CreditScore')
df.boxplot('Age')

**Screenshot:**

Replacing outliers in CreditScore and Age

```
[16] for x in ['CreditScore']:
         q75,q25 = np.percentile(df.loc[:,x],[75,25])
         intr_qr = q75-q25

         max = q75+(1.5*intr_qr)
         min = q25-(1.5*intr_qr)

         df.loc[df[x] < min,x] = np.nan
         df.loc[df[x] > max,x] = np.nan

[17] for x in ['Age']:
         q75,q25 = np.percentile(df.loc[:,x],[75,25])
         intr_qr = q75-q25

         max = q75+(1.5*intr_qr)
         min = q25-(1.5*intr_qr)

         df.loc[df[x] < min,x] = np.nan
         df.loc[df[x] > max,x] = np.nan
```

df.boxplot('CreditScore')

<matplotlib.axes._subplots.AxesSubplot at 0x7f9abb040000>

○ 0s    completed at 10:01 PM



df.boxplot('CreditScore')

<matplotlib.axes._subplots.AxesSubplot at 0x7f9abb040000>

[18] df.boxplot('Age')

<matplotlib.axes._subplots.AxesSubplot at 0x7f9abb305250>

○ 0s    completed at 10:01 PM

## Question-7:

Check for Categorical columns and perform encoding.

### Solution:
df.dtypes
df.head()

df_categorical=df[['Geography','Gender']]

using OneHotEncoding
from sklearn.preprocessing import OneHotEncoder
df = pd.get_dummies(df, columns = ['Geography','Gender'])
df.head()

**Screenshot:**

Geography and Gender — dtype = object

using OneHotEncoding



## Question-8:

Split the data into dependent and independent variables

### Solution:
```
df.shape
df['CustomerId'].nunique()
df['RowNumber'].nunique()
df.columns
df_independent=df[['CreditScore','Age','Tenure','Balance','NumOfProducts','IsActive
Member','EstimatedSalary','Exited', 'Geography_France', 'Geography_Germany',
    'Geography_Spain','Gender_Female', 'Gender_Male']]
df_dependent=df['HasCrCard']
```

**Screenshot:**



```
2023-09-25 00:00:00
2023-09-26 00:00:00
2023-09-27 00:00:00
2023-09-28 00:00:00
2023-09-29 00:00:00
2023-09-30 00:00:00
2023-10-01 00:00:00
2023-10-02 00:00:00
```

### Question-9:

Scale the independent variables

#### Solution:
```python
from sklearn.preprocessing import MinMaxScaler

scaler = MinMaxScaler()

df_independent = scaler.fit_transform(df_independent)
print(df_independent)
```

**Screenshot:**

**Question-10:**

Split the data into training and testing

**Solution:**
```
import numpy as np
from sklearn.model_selection import train_test_split
X_train, X_test, y_train, y_test = train_test_split(df_independent, df_dependent, test_size=0.20, random_state=42)
```

**Screenshot:**