Assignment -2

Data Visualization and Pre-Processing

Date: 26 September 2022

Student Name:R.Jeyapriya

Student Roll Number : 9517201906018
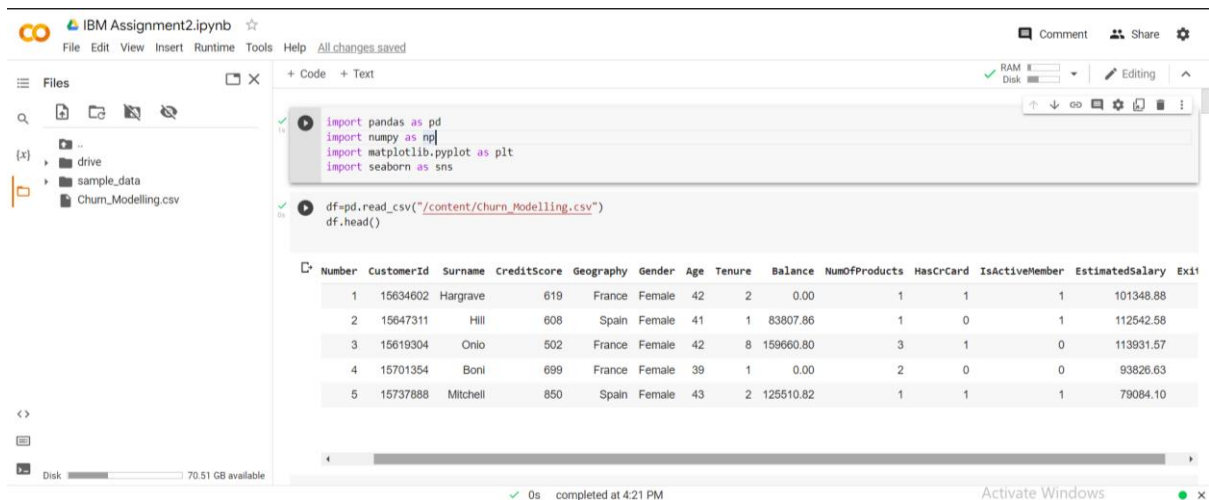
**Question 1 - Load the dataset.**

SOLUTION:

import pandas as pd

import numpy as np

import matplotlib.pyplot as plt

import seaborn as sns

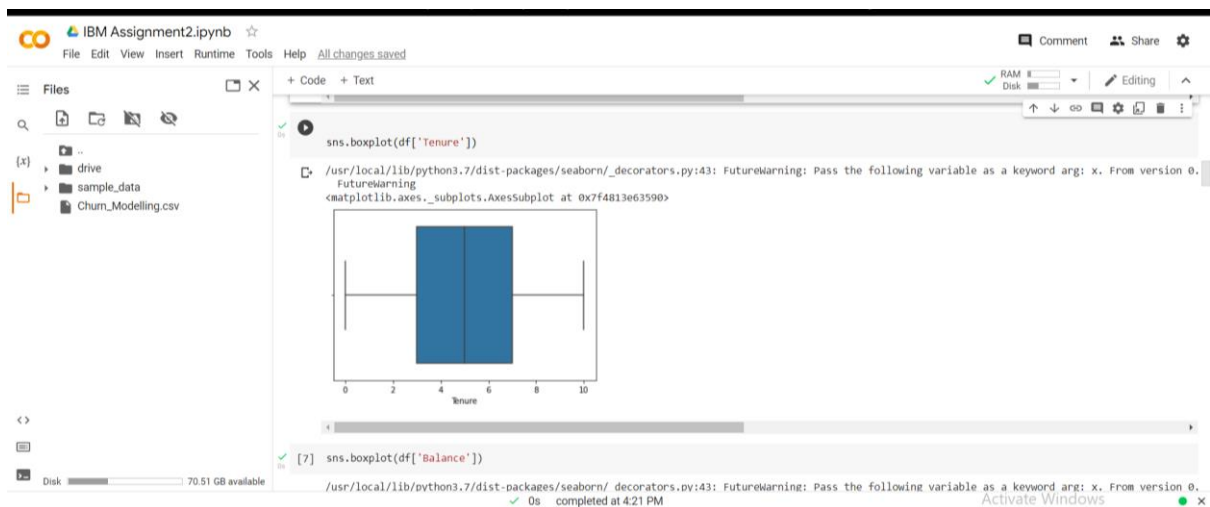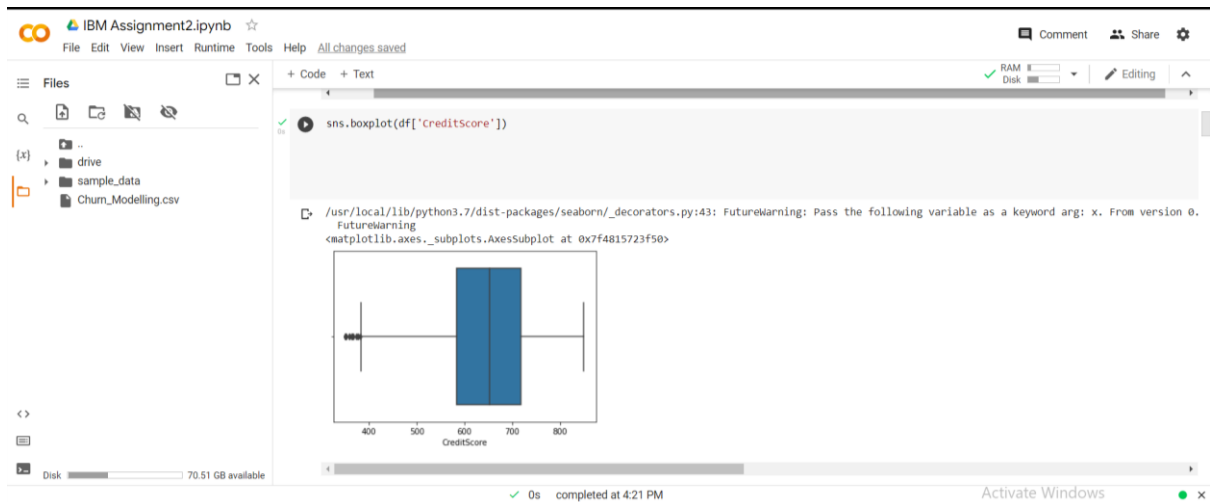 df=pd.read_csv("/content/Churn_Modelling.csv")

df.head()

**OUTPUT:**



**Question 2 - Perform Univariate, Bivariate and Multivariate Analysis**

 SOLUTION:

sns.boxplot(df['CreditScore'])

 sns.boxplot(df['Age'])

 sns.boxplot(df['Tenure'])

sns.boxplot(df['Balance'])

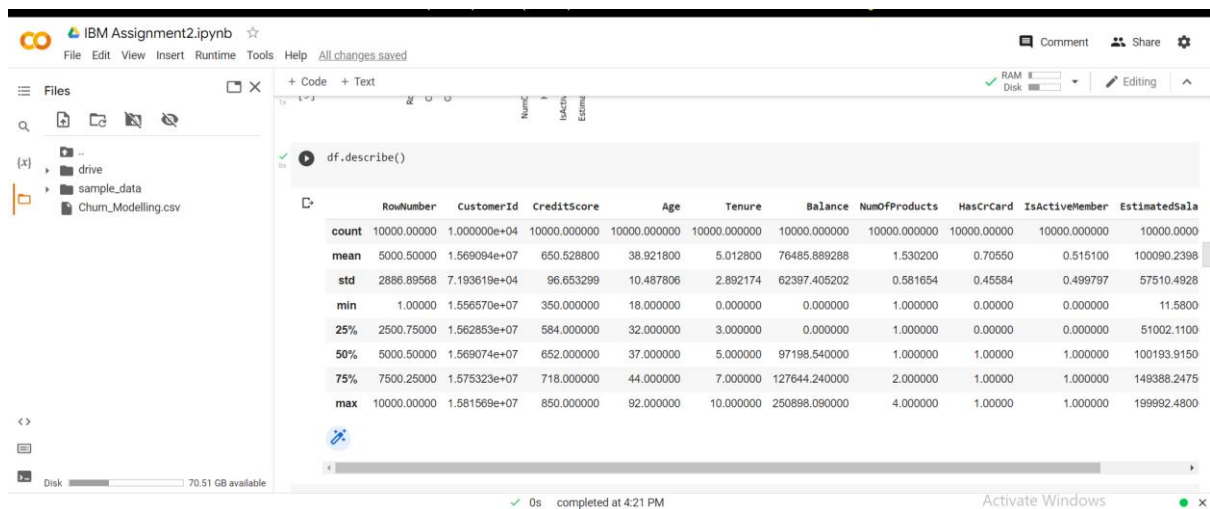sns.boxplot(df['EstimatedSalary'])

sns.heatmap(df.corr(), annot=True)

```python
sns.boxplot(df['Balance'])
```

/usr/local/lib/python3.7/dist-packages/seaborn/_decorators.py:43: FutureWarning: Pass the following variable as a keyword arg: x. From version 0.
FutureWarning
<matplotlib.axes._subplots.AxesSubplot at 0x7f4813dc7f10>

```python
sns.boxplot(df['EstimatedSalary'])
```

/usr/local/lib/python3.7/dist-packages/seaborn/_decorators.py:43: FutureWarning: Pass the following variable as a keyword arg: x. From version 0.
FutureWarning
<matplotlib.axes._subplots.AxesSubplot at 0x7f4813d2ca50>

```python
sns.heatmap(df.corr(), annot=True)
```

<matplotlib.axes._subplots.AxesSubplot at 0x7f4813e56750>

# Question 3 - Perform descriptive statistics on the dataset.

SOLUTION:

df.describe()

OUTPUT:



# Question 4 – Handle the missing values

SOLUTION:

 df.duplicated().sum()

df.nunique()

df.info()

OUTPUT:

## Question 5 - Find the outliers and replace the outliers

SOLUTION:

out = df.drop(columns=['Gender', 'Tenure', 'HasCrCard', 'IsActiveMember', 'NumOfProducts', 'Exi ted']).quantile(q=[0.25, 0.50])

out

Q1 = out.iloc[0]

Q3 = out.iloc[1]

iqr = Q3 - Q1

iqr

```
dtype: int64

[17] out = df.drop(columns=['Gender', 'Tenure', 'HasCrCard', 'IsActiveMember', 'NumOfProducts', 'Exited']).quantile(q=[0.25, 0.50])
     out
```

|       | RowNumber | CustomerId  | CreditScore | Age  | Balance  | EstimatedSalary |
|-------|-----------|-------------|-------------|------|----------|-----------------|
| 0.25  | 2500.75   | 15628528.25 | 584.0       | 32.0 | 0.00     | 51002.110       |
| 0.50  | 5000.50   | 15690738.00 | 652.0       | 37.0 | 97198.54 | 100193.915      |

```
[18] Q1 = out.iloc[0]
     Q3 = out.iloc[1]
     iqr = Q3 - Q1
     iqr

     RowNumber          2499.750
     CustomerId        62209.750
     CreditScore          68.000
     Age                   5.000
     Balance           97198.540
     EstimatedSalary   49191.805
     dtype: float64

[19] upper = out.iloc[1] + 1.5*iqr
```

upper = out.iloc[1] + 1.5*iqr

 upper

```
[18] RowNumber          2499.750
     CustomerId        62209.750
     CreditScore          68.000
     Age                   5.000
     Balance           97198.540
     EstimatedSalary   49191.805
     dtype: float64

[19] upper = out.iloc[1] + 1.5*iqr
     upper

     RowNumber         8.750125e+03
     CustomerId        1.578405e+07
     CreditScore       7.540000e+02
     Age               4.450000e+01
     Balance           2.429964e+05
     EstimatedSalary   1.739816e+05
     dtype: float64

[20] lower = out.iloc[0] - 1.5*iqr
     lower

     RowNumber        -1.248875e+03
     CustomerId        1.553521e+07
     CreditScore       4.820000e+02
     Age               2.450000e+01
     Balance          -1.457978e+05
```

lower = out.iloc[0] - 1.5*iqr

 lower

## Replace outliers

SOLUTION:

df['CreditScore'] = np.where(df['CreditScore']>756, 650.5288, df['CreditScore'])

df['Age'] = np.where(df['Age']>62, 38.9218, df['Age'])

## Question 6 - Check for Categorical columns and perform encoding.

SOLUTION:

df['Gender'].replace({'Male': 1, 'Female': 0}, inplace=True)

df.head(10)

OUTPUT:

# Question 7 – Split the data into dependent and independent variables.

SOLUTION:

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

df.head()



x = df.iloc[:, :-1]

x.head()

y = df.iloc[:, -1]

y.head()

## Question 8 – Scale the independent variables

SOLUTION:

from sklearn.preprocessing import StandardScaler

ss = StandardScaler()

x = ss.fit_transform(x)

x

OUTPUT:



## Question 9 - Split the data into training and testing

SOLUTION:

from sklearn.model_selection import train_test_split

x_train,x_text,y_train, y_test = train_test_split(x, y, test_size=0.2, random_state=0)

print(x_train.shape)

print(x_test.shape)

 print(y_train.shape)

print(y_test.shape)

 OUTPUT:

Files

drive
sample_data
Churn_Modelling.csv

```
              0.97024255, -1.00864308],
       [ 0.29416906,  0.91241915,  0.48205148, ...,  0.64609167,
        -1.03067011, -0.12523071],
       [ 0.29416906, -1.09598752, -1.13723705, ...,  0.64609167,
        -1.03067011, -1.07636976]])
```

```
[29] from sklearn.model_selection import train_test_split
     x_train,x_test,y_train, y_test = train_test_split(x, y, test_size=0.2, random_state=0)
```

```
print(x_train.shape)
print(x_test.shape)
print(y_train.shape)
print(y_test.shape)
```

```
(8000, 9)
(2000, 9)
(8000,)
(2000,)
```