| Assignment Date | 28 September 2022 |
|---|---|
| Student Name | M Hari Haran |
| Student Roll Number | 913119104028 |
| Maximum Marks | 2 Marks |

ASSIGNMENT 2

**Question-1:**

**Download the dataset**

**Question-2:**

**Load the dataset :**



**Question-3:**

**Perform Below Visualizations.**

**Univariate Analysis**

**Bi – Variate Analysis**

**Mult – Variate Analysis**

Untitled0.ipynb
File  Edit  View  Insert  Runtime  Tools  Help  All changes saved

Comment    Share    ⚙

Files

[6]

```python
df.boxplot("CreditScore")
```

<matplotlib.axes._subplots.AxesSubplot at 0x7f435da76a50>



0s    completed at 6:49 AM

Untitled0.ipynb
File  Edit  View  Insert  Runtime  Tools  Help  Saving...

Comment    Share    ⚙

Files

```python
df.boxplot("CreditScore","Tenure")
```

/usr/local/lib/python3.7/dist-packages/matplotlib/cbook/__init__.py:1376: VisibleDeprecationWarning: Creating an ndarray fr
    X = np.atleast_1d(X.T if isinstance(X, np.ndarray) else np.asarray(X))
<matplotlib.axes._subplots.AxesSubplot at 0x7f4373abb090>



0s    completed at 6:50 AM
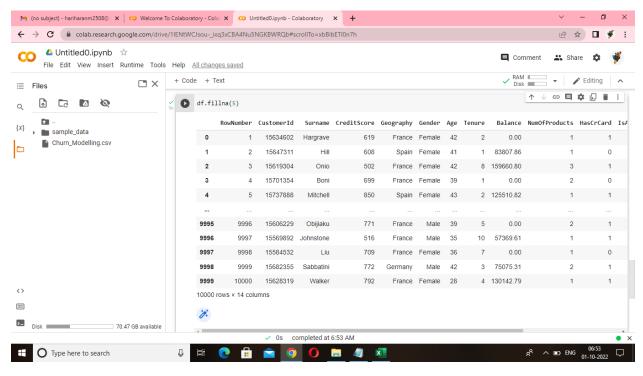
**Question – 4**

**Perform descriptive statistics on the dataset**

## Question-5

### Handle the Missing Values



## Question-6:

### Find the outline & replace the outliners

```python
Q1 = out.iloc[0]
Q2=out.iloc[1]
iqr=Q2-Q1
iqr
```

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

0s    completed at 6:56 AM

```python
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
```

0s    completed at 6:57 AM

CO **Untitled0.ipynb** ☆

File Edit View Insert Runtime Tools Help

Comment    Share ⚙

+ Code   + Text         RAM / Disk   Editing ∧

```
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
EstimatedSalary  -2.278560e+04
dtype: float64
```

Disk   70.47 GB available

✓ 0s   completed at 6:57 AM

CO **Untitled0.ipynb** ☆

File Edit View Insert Runtime Tools Help   All changes saved

Comment    Share ⚙

+ Code   + Text         RAM / Disk   Editing ∧

```
[18] 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
EstimatedSalary  -2.278560e+04
dtype: float64
```

```
df['CreditScore']= np.where(df['CreditScore']>756, 650.5288,df['CreditScore'])
df['Age']=np.where(df['Age']>62, 38.9218, df['Age'])
```

Disk   70.47 GB available

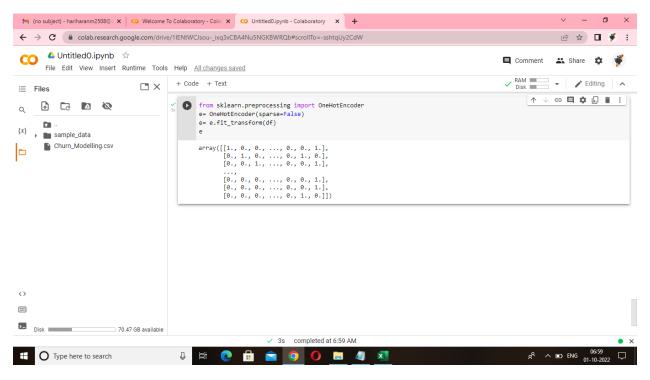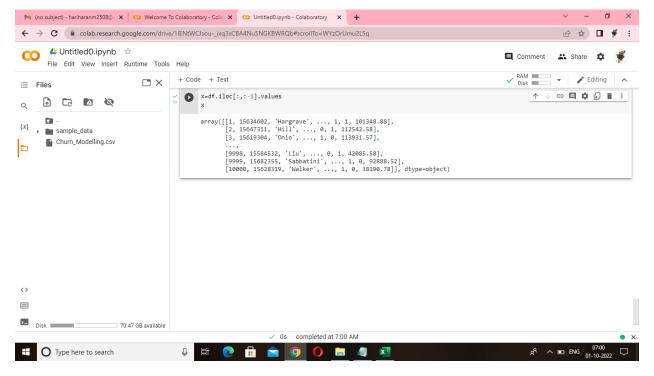✓ 0s   completed at 6:58 AM

## Question-7:

**Check for Categorical columns and Perform encoding**

## Question-8

**Split the data into depentant and independent Variables.**



## Question-9

**Scale the independent Variables**

```
from sklearn.preprocessing import StandardScaler

df.head()
```

| | RowNumber | CustomerId | Surname | CreditScore | Geography | Gender | Age | Tenure | Balance | NumOfProducts | HasCrCard | IsActi |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 0 | 1 | 15634602 | Hargrave | 619.0000 | France | Female | 42.0 | 2 | 0.00 | 1 | 1 | |
| 1 | 2 | 15647311 | Hill | 608.0000 | Spain | Female | 41.0 | 1 | 83807.86 | 1 | 0 | |
| 2 | 3 | 15619304 | Onio | 502.0000 | France | Female | 42.0 | 8 | 159660.80 | 3 | 1 | |
| 3 | 4 | 15701354 | Boni | 699.0000 | France | Female | 39.0 | 1 | 0.00 | 2 | 0 | |
| 4 | 5 | 15737888 | Mitchell | 650.5288 | Spain | Female | 43.0 | 2 | 125510.82 | 1 | 1 | |

## Question-10

## Split the data into training & testing



```
from sklearn.model_selection import train_test_split
x_train,x_test,y_train,y_test = train_test_split(x,y, random_state=0, train_size=.75)

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

(7500, 13)
(2500, 13)
(7500,)
(2500,)
```