# **Assignment -2**

Data Visualization and Pre-processing

Assignment Date	21 September 2022
Student Name	B.RAJESH
Student Roll Number	510119205013
Maximum Marks	2 Marks

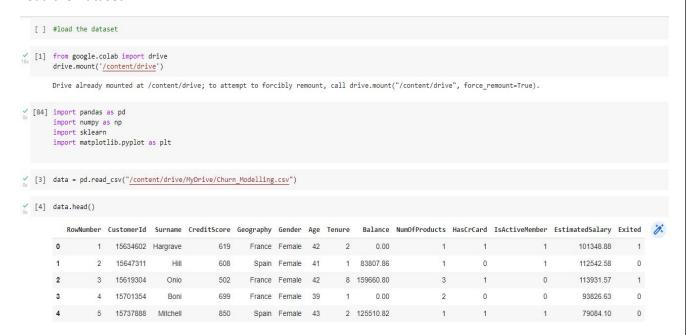
# Question-1:

#### Download the dataset:

The dataset "Churn\_Modelling.csv" was downloaded Successfully

# Question-2:

# Load the Dataset:



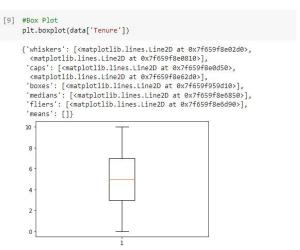
#### Question-3:

**Perform Below Visualization:** 

**Univariate Analysis** 

[5] #Univariate Analysis for Numerical data

# [6] #Histogram data['Age'].plot(kind='hist') <matplotlib.axes.\_subplots.AxesSubplot at 0x7f65a0462590> 3500 2500 2000 1000 500 20 30 40 50 60 70 80 90

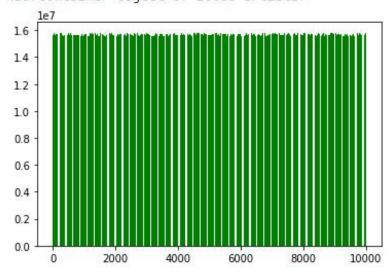


# [ ] #Univariate Analysis for Categorical Data

```
[14] #Bar Chart
    df = pd.DataFrame(data)

X = list(df.iloc[:, 0])
    Y = list(df.iloc[:, 1])
    plt.bar(X, Y, color='g')
```

<BarContainer object of 10000 artists>



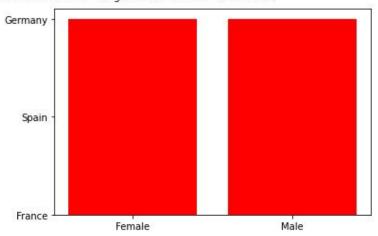
# **Bivariate Analysis**

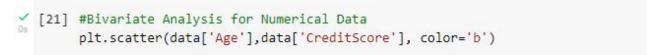
```
[23] #Bivariate Analysis for Categorical Data

#Stacked Bar chart

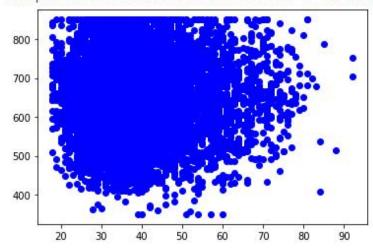
plt.bar(data['Gender'], data['Geography'], color='r')
```

<BarContainer object of 10000 artists>





<matplotlib.collections.PathCollection at 0x7f6589f606d0>



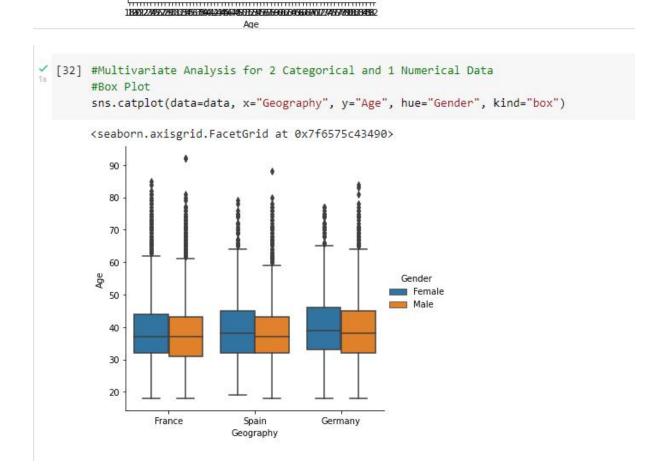
# **Multivariate Analysis**

```
#Multivariate Analysis for 2 Numerical and 1 Categorical Data
#Scatter Plot
import seaborn as sns
sns.catplot(data=data, x="Age", y="CreditScore", hue="Gender")

C> <seaborn.axisgrid.FacetGrid at 0x7f657aab5d90>

#Multivariate Analysis for 2 Numerical and 1 Categorical Data
#Scatter Plot
import seaborn as sns
sns.catplot(data=data, x="Age", y="CreditScore", hue="Gender")

Gender
Female
Male
```



#### Question-4:

(10000, 14)

#### Perform Descriptive Statistics on the dataset:

```
[ ] #Perform Descriptive Statistics on the Dataset
        data.mean()
   /usr/local/lib/python3.7/dist-packages/ipykernel_launcher.py:1:
          """Entry point for launching an IPython kernel.
                        5.000500e+03
1.569094e+07
        RowNumber
        CustomerId
                            6.505288e+02
       CreditScore
       Age
                             3.892180e+01
        Tenure
                             5.012800e+00
                             7.648589e+04
        Balance
        NumOfProducts
                              1.530200e+00
        HasCrCard
                               7.055000e-01
        IsActiveMember
                               5.151000e-01
        EstimatedSalary 1.000902e+05
        Exited
                              2,037000e-01
        dtype: float64
[34] data.median()
        /usr/local/lib/python3.7/dist-packages/ipykernel_launcher.py:1:
          """Entry point for launching an IPython kernel.
        RowNumber
                              5.000500e+03
                           1.569074e+07
        CustomerId
        CreditScore
                            6.520000e+02
                              3.700000e+01
        Age
        Tenure
                              5.000000e+00
        Balance
                               9.719854e+04
        NumOfProducts
                             1.000000e+00
        HasCrCard
                              1.000000e+00
        IsActiveMember
                              1.000000e+00
        EstimatedSalary 1.001939e+05
                              0.000000e+00
        Exited
        dtype: float64
[36] data.describe()
                                               Tenure
           RowNumber CustomerId CreditScore
                                                        Balance NumOfProducts HasCrCard IsActiveMember EstimatedSalary
                                                                                                           Exited 🎢
      count 10000.00000 1.000000e+04 10000.000000 10000.000000 10000.000000 10000.000000 10000.000000 10000.000000 10000.000000 10000.000000 10000.000000 10000.000000 10000.000000
      mean 5000.50000 1.569094e+07 650.528800 38.921800 5.012800 76485.889288
                                                                                     0.515100 100090.239881
                                                                 1.530200 0.70550
      std 2886.89568 7.193619e+04 96.653299 10.487806 2.892174 62397.405202 0.581654 0.45584 0.499797 57510.492818 0.402769
            1.00000 1.556570e+07 350.000000 18.000000
                                              0.000000
                                                        0.000000
                                                                   1.000000
                                                                            0.00000
                                                                                      0.000000
                                                                                                11.580000
                                                       0.000000 1.000000 0.00000 0.000000 51002.110000 0.000000
      25% 2500.75000 1.562853e+07 584.000000 32.000000 3.000000
          5000.50000 1.569074e+07 652.000000
                                     37.000000
                                               5.000000 97198.540000
                                                                   1.000000
                                                                            1.00000
                                                                                      1.000000
                                                                                             100193.915000
                                     44.000000 7.000000 127644.240000 2.000000
                                                                            1.00000 1.000000 149388.247500
      75% 7500.25000 1.575323e+07 718.000000
      max 10000.00000 1.581569e+07 850.000000 92.000000
                                              10.000000 250898.090000
                                                                   4 000000
                                                                            1.00000
                                                                                      1.000000 199992.480000
                                                                                                          1.000000
/ [38] data.shape
```

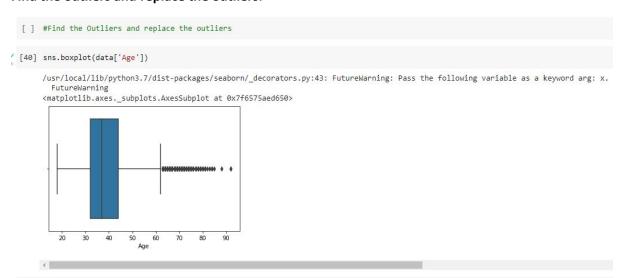
# Question-5:

# Handle the Missing values:

```
\frac{\checkmark}{O_{0}} [39] #Handling the missing values
        data.isnull().sum()
        RowNumber
                             0
        CustomerId
                             0
        Surname
                             0
        CreditScore
                             0
        Geography
                             0
        Gender
                             0
                             0
        Age
        Tenure
        Balance
                             0
        NumOfProducts
                             0
        HasCrCard
        IsActiveMember
                             0
        EstimatedSalary
                             0
        Exited
        dtype: int64
```

# Question-6:

# Find the outliers and replace the outliers:



```
(41] qnt=data.quantile(q=[0.25,0.75])
         RowNumber CustomerId CreditScore Age Tenure Balance NumOfProducts HasCrCard IsActiveMember EstimatedSalary Exited 🎉
     0.25 2500.75 15628528.25 584.0 32.0 3.0
                                               1.0 0.0 0.0
                                           0.00
                                                                             51002 1100
                                                                                      0.0
     0.75 7500.25 15753233.75
                          718.0 44.0 7.0 127644.24
                                                                            149388.2475
(42) IQR = qnt.loc[0.75] - qnt.loc[0.25]
     IQR
     RowNumber
                  4999.5000
              124705.5000
     CustomerId
               134.0000
     CreditScore
                   12.0000
     Age
     Tenure
    4.0000
127644.2400
NumOfProducts
HasCnC--'
                    4.0000
                1.0000
     HasCrCard
IsActiveMember
                    1.0000
     EstimatedSalary 98386.1375
     Exited
                   0.0000
     dtype: float64
   [43] upper_extreme = qnt.loc[0.75]+1.5*IQR
          upper extreme
                                1.499950e+04
         RowNumber
         CustomerId
                                1.594029e+07
         CreditScore
                                9.190000e+02
                                6.200000e+01
         Age
         Tenure
                                1.300000e+01
         Balance
                                3.191106e+05
         NumOfProducts
                               3.500000e+00
         HasCrCard
                                2.500000e+00
         IsActiveMember
                                2.500000e+00
         EstimatedSalary
                                2.969675e+05
         Exited
                                0.000000e+00
         dtype: float64
   [44] lower_extreme = qnt.loc[0.25]-1.5*IQR
          lower extreme
         RowNumber
                               -4.998500e+03
         CustomerId
                               1.544147e+07
         CreditScore
                                3.830000e+02
                                1.400000e+01
         Age
         Tenure
                               -3.000000e+00
         Balance
                               -1.914664e+05
         NumOfProducts
                               -5.000000e-01
         HasCrCard
                               -1.500000e+00
         IsActiveMember
                               -1.500000e+00
         EstimatedSalary
                               -9.657710e+04
         Exited
                                0.000000e+00
         dtype: float64
```

```
[51] df2 = data[(data['Age']<upper_extreme['Age']) & (data['Age']>lower_extreme['Age'])]

[50] data.shape
(10000, 14)

[49] df2.shape
(9589, 14)

[52] sns.boxplot(df2['Age'])

/usr/local/lib/python3.7/dist-packages/seaborn/_decorators.py:43: FutureWarning: Pass the following variable as a keyword arg: x.
FutureWarning
<matplotlib.axes__subplots.AxesSubplot at 0x7f6573caad10>

[53] data.shape
(9589, 14)

[54] df2.shape
(9589, 14)

[55] sns.boxplot(df2['Age'])

/usr/local/lib/python3.7/dist-packages/seaborn/_decorators.py:43: FutureWarning: Pass the following variable as a keyword arg: x.
```

#### Question-7:

#### **Check for Categorical columns and perform Encoding:**

```
[53] #Check for Categorical columns and perform encoding
    #Categorical are Geography and Gender
    from sklearn.preprocessing import LabelEncoder

[75] le=LabelEncoder()
    df2['Geography'] = le.fit_transform(df2['Geography'])
    df2['Gender'] = le.fit_transform(df2['Gender'])
```

[76] df2.head() RowNumber CustomerId Surname CreditScore Geography Gender Age Tenure Balance NumOfProducts HasCrCard IsActiveMember EstimatedSalary Exited 15634602 Hargrave 0 42 0.00 101348.88 15647311 0 41 1 83807.86 112542.58 0 Hill 608 0 42 8 159660.80 113931.57 15619304 Onio 502 15701354 Boni 0 39 93826.63 0 15737888 Mitchell 850 0 43 2 125510.82 79084.10

#### Question-8:

#### Split the data into dependent and independent variables:

```
[77] #Split the data into dependent and independent variables.
    y=df2['EstimatedSalary']
    x=df2.drop(columns=['EstimatedSalary'],axis=1)
```



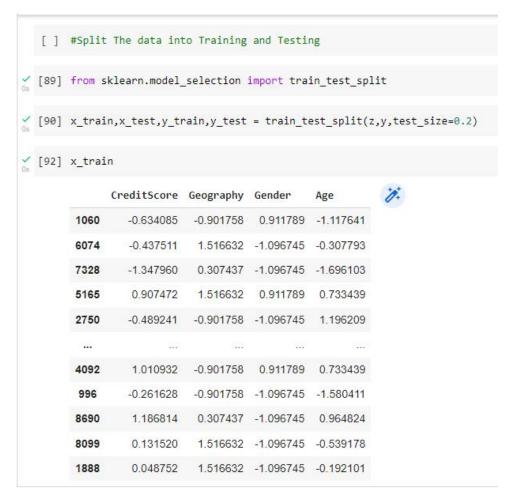
#### Question-9:

#### Scale the independent variables:



#### Question-10:

#### Split the data into training and testing:



#### y\_train C 1104 151645.96 6334 143463.28 7638 37577.66 5392 43018.82 100478.60 2851 ... 4269 2048.55 1037 180969.55 9056 166896.01 8440 36864.05 1960 86013.96

Name: EstimatedSalary, Length: 7671, dtype: float64

# / [94] x\_test

	CreditScore	Geography	Gender	Age
962	0.772974	0.307437	0.911789	0.154977
5257	1.248890	1.516632	-1.096745	0.386361
7515	-0.841005	0.307437	-1.096745	-0.654871
6844	0.959202	-0.901758	-1.096745	-0.886256
4102	-0.996196	1.516632	-1.096745	0.386361
	***	3753	1320	
60	0.379825	0.307437	-1.096745	-1.233333
5555	0.503977	-0.901 <mark>75</mark> 8	0.911789	-0.076408
5112	1.704115	1.516632	-1.096745	2.237441
138	0.131520	-0.901758	0.911789	-0.423486
4973	0.328095	-0.901758	-1.096745	2.353134

1918 rows × 4 columns

```
// [95] y_test
       1002
              184023.54
       5486
              92914.67
       7838
              132038.65
       7133
              138780.89
       4281
               36242.19
              126494.82
       61
       5797
               83263.04
               38941.44
       5337
              180427.24
       141
       5191
                 706.50
       Name: EstimatedSalary, Length: 1918, dtype: float64
```