

DATE	28/09/2022
TEAM ID	PNT2022TMID27851
STUDENT NAME	SWAATHI CM
DOMAIN NAME	HEALTH CARE
PROJECT NAME	EARLY DETECTION OF CHRONIC KIDNEY DISEASE USING MACHINE LEARNING
MAXIMUM MARKS	2 MARKS

LOAD THE DATASET

IMPORT THE REQUIRED LIBRARIES

```
[1] import pandas as pd
import numpy as np
import matplotlib.pyplot as plt
import seaborn as sns
```

LOAD THE DATASET

```
[2] df = pd.read_csv('Churn_Modelling.csv')
```

```
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
4	5	15737888	Mitchell	850	Spain	Female	43	2	125510.82	1	1	1	79084.10	0

PERFORM VISUALIZATIONS

Removing Unwanted Values

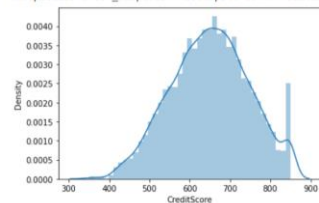
```
[4] df = df.drop(columns=['RowNumber', 'CustomerId', 'Surname'])
```

● UNIVARIATE

```
[5] ## FOR NUMERICAL COLUMN
sns.distplot(df.CreditScore)
```

/usr/local/lib/python3.7/dist-packages/seaborn/distributions.py:2619: FutureWarning: "distplot" is a deprecated function and will be removed in a future version. Please adapt your code to warnings.warn(msg, FutureWarning)

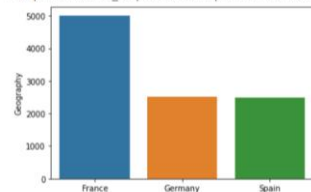
<matplotlib.axes._subplots.AxesSubplot at 0x7f52ac13d190>



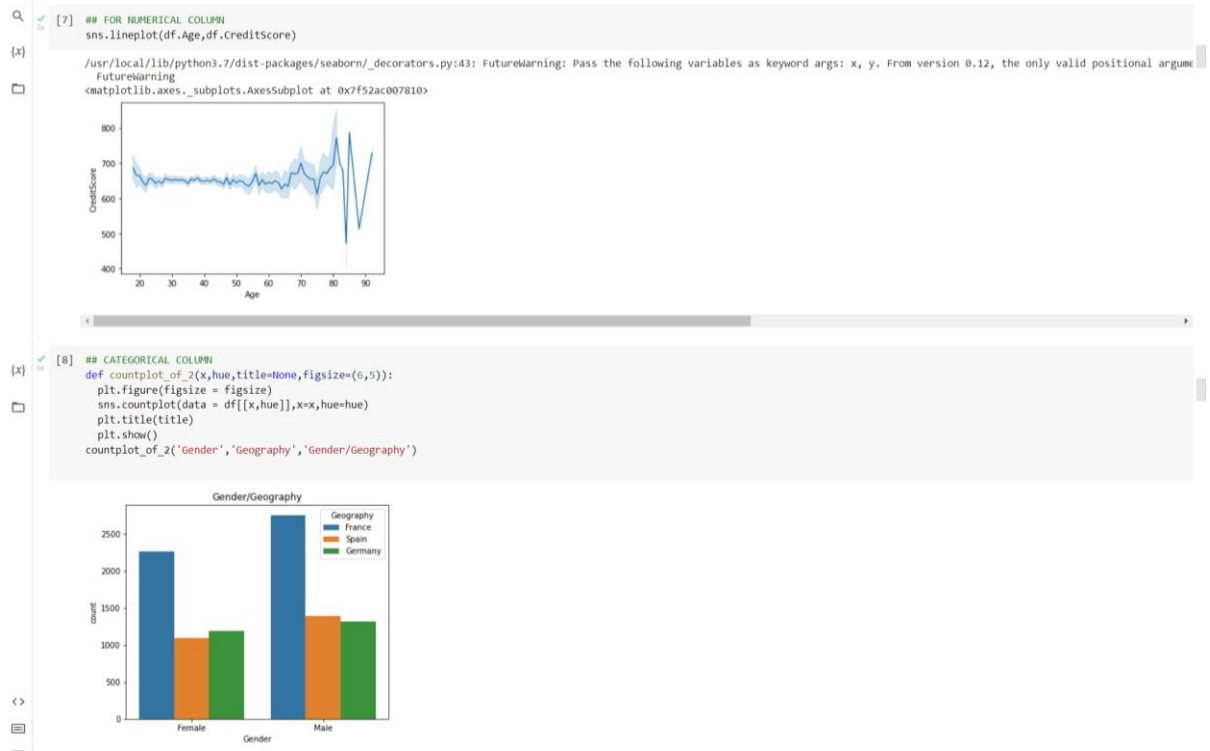
```
[6] ## FOR CATEGORICAL COLUMN
sns.barplot(df.Geography.value_counts().index, df.Geography.value_counts())
```

/usr/local/lib/python3.7/dist-packages/seaborn/_decorators.py:43: FutureWarning: Pass the following variables as keyword args: x, y. From version 0.12, the only valid positional argument is FutureWarning

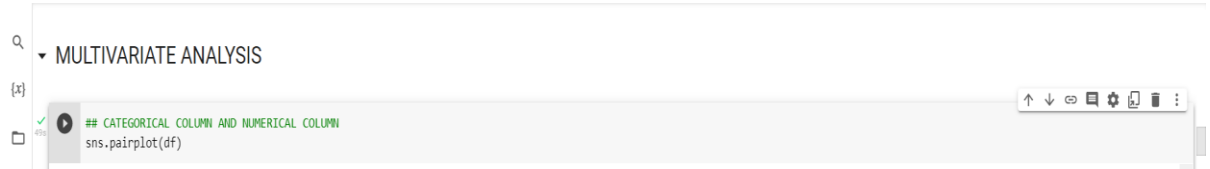
<matplotlib.axes._subplots.AxesSubplot at 0x7f52ac007bd0>

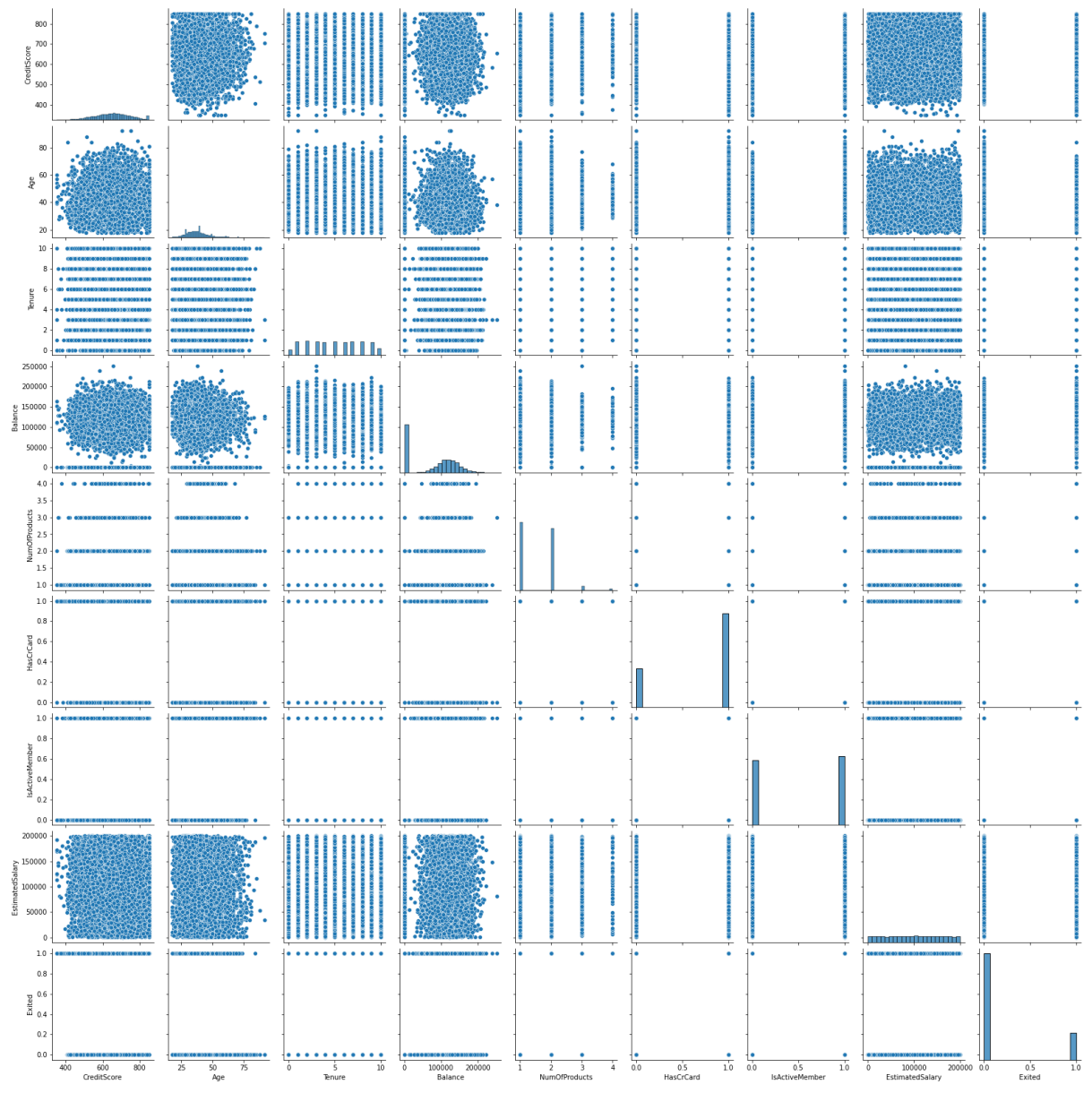


● BIVARIATE



● MULTIVARIATE

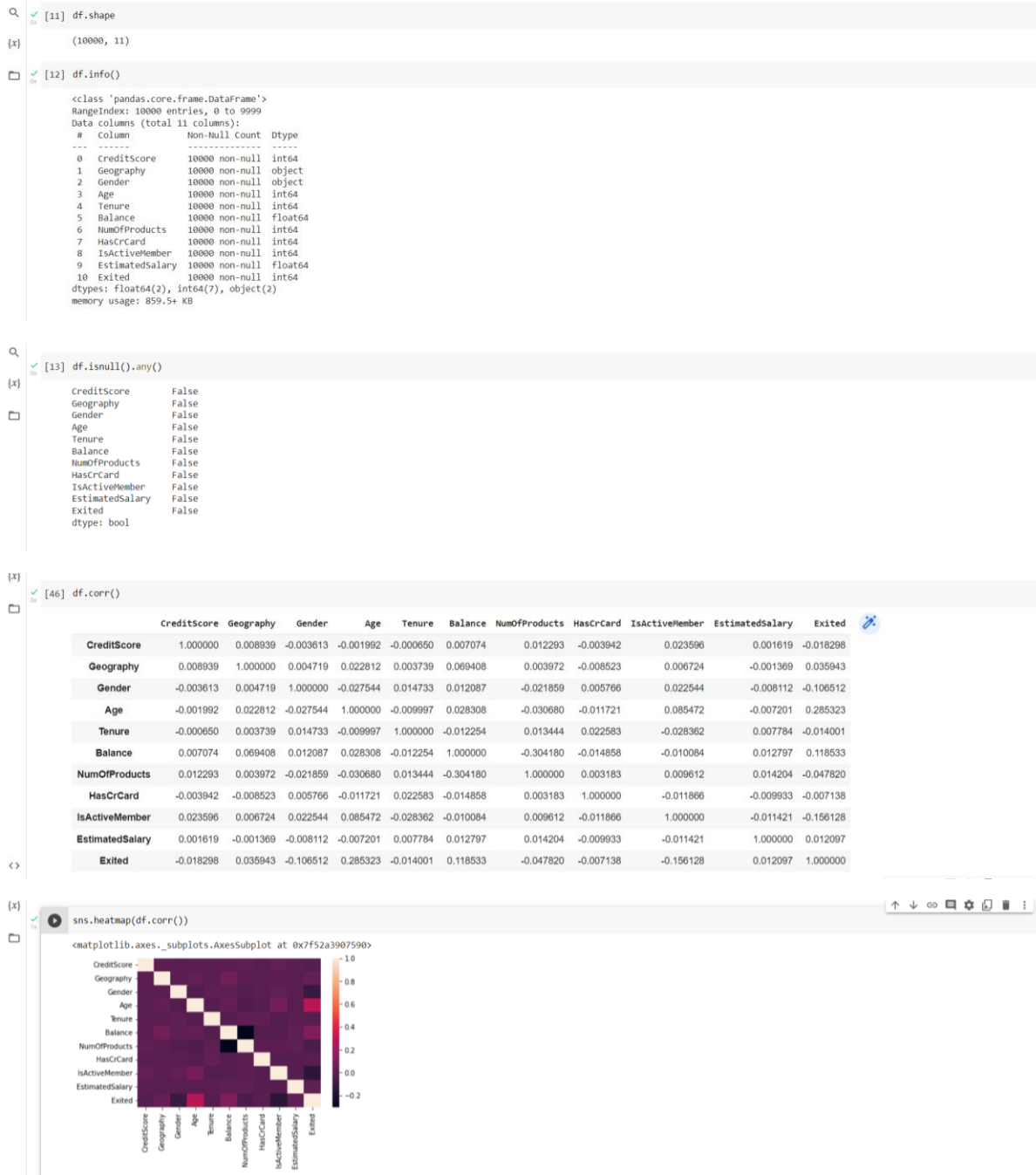




PERFORM DESCRIPTIVE ANALYSIS ON THE DATASET

```
[10] df.describe()
```

	CreditScore	Age	Tenure	Balance	NumOfProducts	HasCrCard	IsActiveMember	EstimatedSalary	Exited
count	10000.000000	10000.000000	10000.000000	10000.000000	10000.000000	10000.000000	10000.000000	10000.000000	10000.000000
mean	650.528800	38.921800	5.012800	76485.889288	1.530200	0.705500	0.515100	100090.239881	0.203700
std	96.653299	10.487806	2.892174	62397.405202	0.581654	0.455840	0.499797	57510.492818	0.402769
min	350.000000	18.000000	0.000000	0.000000	1.000000	0.000000	0.000000	11.580000	0.000000
25%	584.000000	32.000000	3.000000	0.000000	1.000000	0.000000	0.000000	51002.110000	0.000000
50%	652.000000	37.000000	5.000000	97198.540000	1.000000	1.000000	1.000000	100193.915000	0.000000
75%	718.000000	44.000000	7.000000	127644.240000	2.000000	1.000000	1.000000	149388.247500	0.000000
max	850.000000	92.000000	10.000000	250898.090000	4.000000	1.000000	1.000000	199992.480000	1.000000



HANDLE THE MISSING VALUES

```
[14] ## NUMERICAL COLUMN

[15] sns.distplot(df.CreditScore)

/usr/local/lib/python3.7/dist-packages/seaborn/distributions.py:2619: FutureWarning: `distplot` is a deprecated function and will be removed in a future version. Please adapt your code
warnings.warn(msg, FutureWarning)
<matplotlib.axes._subplots.AxesSubplot at 0x7f52a6fe1910>

Density
0.0040
0.0035
0.0030
0.0025
0.0020
0.0015
0.0010
0.0005
0.0000
300 400 500 600 700 800 900
CreditScore

[16] df['CreditScore'].fillna(df['CreditScore'].mean(),inplace=True)

[17] sns.distplot(df.Age)

/usr/local/lib/python3.7/dist-packages/seaborn/distributions.py:2619: FutureWarning: `distplot` is a deprecated function and will be removed in a future version. Please adapt your code
warnings.warn(msg, FutureWarning)
<matplotlib.axes._subplots.AxesSubplot at 0x7f52a5693d90>

Density
0.06
0.05
0.04
0.03
0.02
0.01
0.00
20 40 60 80 100
Age

[18] df['Age'].fillna(df['Age'].median(),inplace=True)

[19] sns.distplot(df.Tenure)

/usr/local/lib/python3.7/dist-packages/seaborn/distributions.py:2619: FutureWarning: `distplot` is a deprecated function and will be removed in a future version. Please adapt your code
warnings.warn(msg, FutureWarning)
<matplotlib.axes._subplots.AxesSubplot at 0x7f52a55978d0>

Density
0.25
0.20
0.15
0.10
0.05
0.00
-2 0 2 4 6 8 10 12
Tenure

[20] df['Tenure'].fillna(df['Tenure'].median(),inplace = True)

[21] sns.distplot(df.Balance)

/usr/local/lib/python3.7/dist-packages/seaborn/distributions.py:2619: FutureWarning: `distplot` is a deprecated function and will be removed in a future version. Please adapt your code
warnings.warn(msg, FutureWarning)
<matplotlib.axes._subplots.AxesSubplot at 0x7f52a559ef90>

Density
3e-5
3.0
2.5
2.0
1.5
1.0
0.5
0.0
0 50000 100000 150000 200000 250000
Balance

[22] df['Balance'].fillna(df['Balance'].median(),inplace = True)
```

```
[x] [23] sns.distplot(df.NumOfProducts)

/usr/local/lib/python3.7/dist-packages/seaborn/distributions.py:2619: FutureWarning: 'distplot' is a deprecated function and will be removed in a future version. Please adapt your code
warnings.warn(msg, FutureWarning)
<matplotlib.axes._subplots.AxesSubplot at 0x7f52a544e690>


A density plot for the variable 'NumOfProducts'. The x-axis ranges from 0 to 40, and the y-axis (Density) ranges from 0 to 5. The plot shows two distinct peaks: one at approximately 1.0 with a density of about 2.2, and another at approximately 2.0 with a density of about 2.0. There are very small peaks at higher values like 3.0 and 4.0.

[24] df['NumOfProducts'].fillna(df['NumOfProducts'].median(),inplace = True)

[x] [25] sns.distplot(df.HasCrCard)

/usr/local/lib/python3.7/dist-packages/seaborn/distributions.py:2619: FutureWarning: 'distplot' is a deprecated function and will be removed in a future version. Please adapt your code
warnings.warn(msg, FutureWarning)
<matplotlib.axes._subplots.AxesSubplot at 0x7f52a6fbc810>


A density plot for the variable 'HasCrCard'. The x-axis ranges from -0.2 to 1.2, and the y-axis (Density) ranges from 0 to 8. The plot shows two main peaks: one at 0.0 with a density of about 3.5, and a much larger one at 1.0 with a density of about 7.5. There are also very small peaks at approximately 0.1 and 0.9.

[26] df['HasCrCard'].fillna(df['HasCrCard'].median(),inplace = True)

[x] [27] sns.distplot(df.IsActiveMember)

/usr/local/lib/python3.7/dist-packages/seaborn/distributions.py:2619: FutureWarning: 'distplot' is a deprecated function and will be removed in a future version. Please adapt your code
warnings.warn(msg, FutureWarning)
<matplotlib.axes._subplots.AxesSubplot at 0x7f52a531b450>


A density plot for the variable 'IsActiveMember'. The x-axis ranges from -0.2 to 1.2, and the y-axis (Density) ranges from 0 to 5. The plot shows two main peaks: one at 0.0 with a density of about 2.5, and another at 1.0 with a density of about 2.5. There are also very small peaks at approximately 0.1 and 0.9.

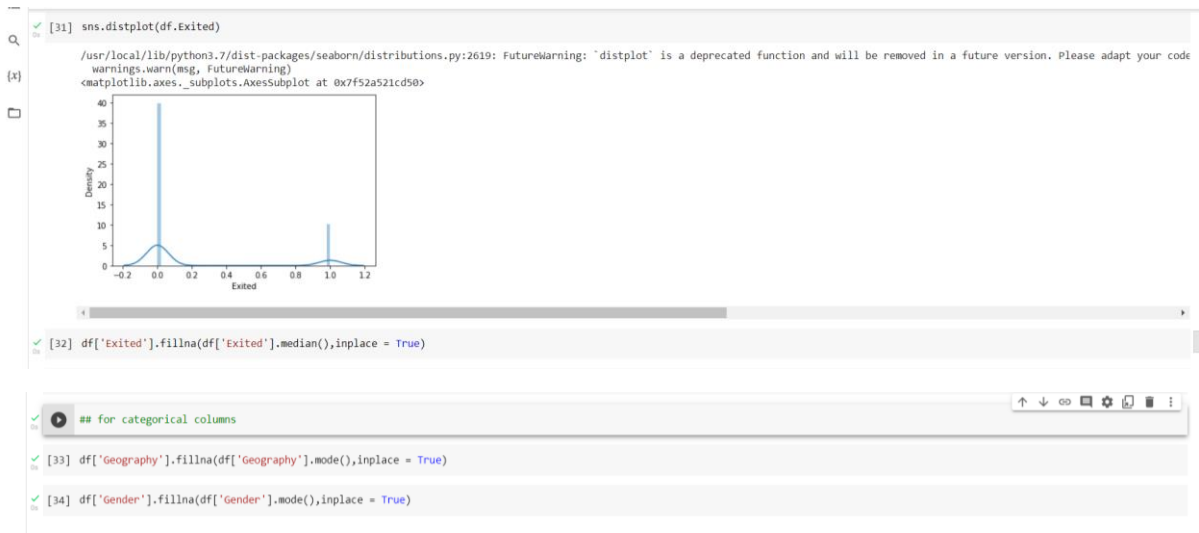
[28] df['IsActiveMember'].fillna(df['IsActiveMember'].median(),inplace = True)

[x] [29] sns.distplot(df.EstimatedSalary)

/usr/local/lib/python3.7/dist-packages/seaborn/distributions.py:2619: FutureWarning: 'distplot' is a deprecated function and will be removed in a future version. Please adapt your code
warnings.warn(msg, FutureWarning)
<matplotlib.axes._subplots.AxesSubplot at 0x7f52a5291e90>

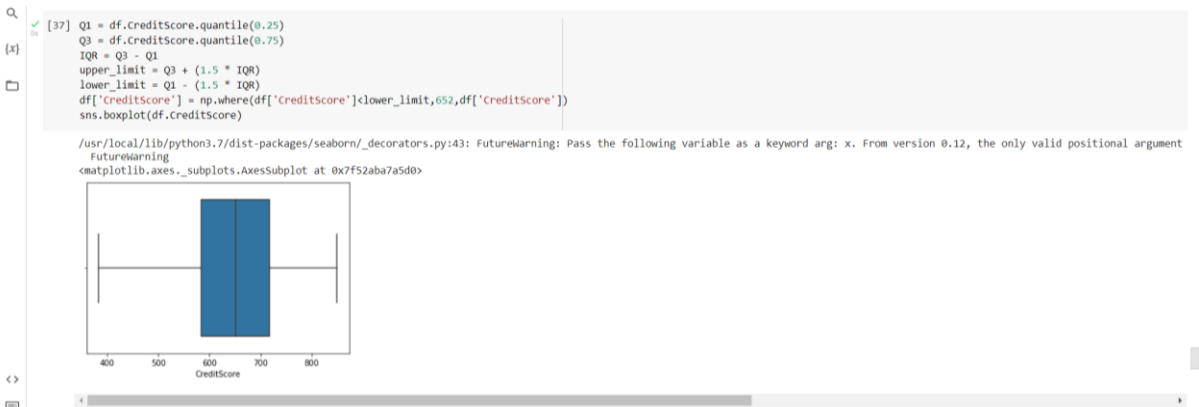
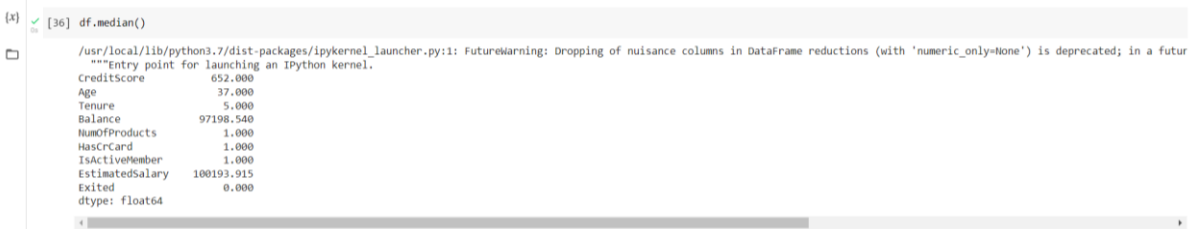
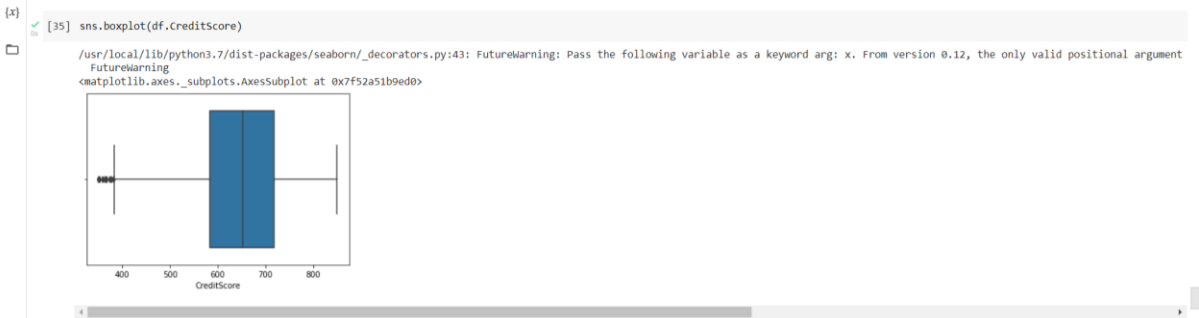

A density plot for the variable 'EstimatedSalary'. The x-axis ranges from 0 to 200,000, and the y-axis (Density) ranges from 0 to 5. The plot shows a broad, flat distribution with a slight peak around 100,000. The density is relatively constant across the range, with a slight increase towards the center.

[30] df['EstimatedSalary'].fillna(df['EstimatedSalary'].median(),inplace = True)
```



FIND THE OUTLIERS AND REPLACE THE OUTLIERS

▼ FINDING AND REPLACING OUTLIERS



CHECK FOR CATEGORICAL COLUMNS AND PERFORM ENCODING

▼ CHECK FOR CATEGORICAL COLUMN AND PERFORM ENCODING

```
[38] from sklearn.preprocessing import LabelEncoder
le = LabelEncoder()
df.Gender = le.fit_transform(df.Gender)
df.head()
```

	CreditScore	Geography	Gender	Age	Tenure	Balance	NumOfProducts	HasCrCard	IsActiveMember	EstimatedSalary	Exited
0	619	France	0	42	2	0.00	1	1	1	101348.88	1
1	608	Spain	0	41	1	83807.86	1	0	1	112542.58	0
2	502	France	0	42	8	159660.80	3	1	0	113931.57	1
3	699	France	0	39	1	0.00	2	0	0	93826.63	0
4	850	Spain	0	43	2	125510.82	1	1	1	79084.10	0

```
[39] df.Geography = le.fit_transform(df.Geography)
df.head()
```

	CreditScore	Geography	Gender	Age	Tenure	Balance	NumOfProducts	HasCrCard	IsActiveMember	EstimatedSalary	Exited
0	619	0	0	42	2	0.00	1	1	1	101348.88	1
1	608	2	0	41	1	83807.86	1	0	1	112542.58	0
2	502	0	0	42	8	159660.80	3	1	0	113931.57	1
3	699	0	0	39	1	0.00	2	0	0	93826.63	0
4	850	2	0	43	2	125510.82	1	1	1	79084.10	0

SPLIT THE DATA INTO DEPENDENT AND INDEPENDENT VARIABLES

▼ SPLIT THE DATA SET INTO DEPENDENT AND INDEPENDENT VARIABLES

```
[40] X = df.drop(columns = ['Exited'], axis = 1)
X.head()
```

	CreditScore	Geography	Gender	Age	Tenure	Balance	NumOfProducts	HasCrCard	IsActiveMember	EstimatedSalary
0	619	0	0	42	2	0.00	1	1	1	101348.88
1	608	2	0	41	1	83807.86	1	0	1	112542.58
2	502	0	0	42	8	159660.80	3	1	0	113931.57
3	699	0	0	39	1	0.00	2	0	0	93826.63
4	850	2	0	43	2	125510.82	1	1	1	79084.10

```
[41] y = df['Exited']
y.head()
```

```
0    1
1    0
2    1
3    0
4    0
Name: Exited, dtype: int64
```

SCALE THE INDEPENDENT VARIABLES

▼ SCALE THE INDEPENDENT COLUMNS

```
[42] from sklearn.preprocessing import MinMaxScaler
scale = MinMaxScaler()
X_scaled = pd.DataFrame(scale.fit_transform(X), columns = X.columns)
X_scaled.head()
```

	CreditScore	Geography	Gender	Age	Tenure	Balance	NumOfProducts	HasCrCard	IsActiveMember	EstimatedSalary
0	0.505353	0.0	0.0	0.324324	0.2	0.000000	0.000000	1.0	1.0	0.506735
1	0.481799	1.0	0.0	0.310811	0.1	0.334031	0.000000	0.0	1.0	0.562709
2	0.254818	0.0	0.0	0.324324	0.8	0.636357	0.666667	1.0	0.0	0.569654
3	0.676660	0.0	0.0	0.283784	0.1	0.000000	0.333333	0.0	0.0	0.469120
4	1.000000	1.0	0.0	0.337838	0.2	0.500246	0.000000	1.0	1.0	0.395400

SPLIT THE DATA INTO TRAINING AND TESTING

[x] ▾ SPLIT THE DATASET INTO TRAINING AND TESTING



```
[43] ## SPLITTING DATASET INTO TRAIN TEST
      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)
```

```
[44] x_train.shape
```

```
(8000, 10)
```

```
[45] x_test.shape
```

```
(2000, 10)
```