

```
In [1]: import pandas as pd
import numpy as np
import matplotlib.pyplot as plt
import seaborn as sns
from sklearn.preprocessing import scale
```

Load the Dataset

```
In [3]: 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 | Huigrave | 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 | 15616304 | Ono | 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 |

```
In [4]: df.shape
Out[4]: (10000, 14)
```

```
In [5]: df.isnull().any()
Out[5]: RowNumber      False
CustomerId      False
Surname         False
CreditScore     False
Geography       False
Gender          False
Age             False
Tenure          False
Balance         False
NumOfProducts  False
HasCrCard       False
IsActiveMember  False
EstimatedSalary False
Exited          False
dtype: bool
```

```
In [6]: df.info()
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 10000 entries, 0 to 9999
Data columns (total 14 columns):
#   Column                Non-Null Count  Dtype
---  -
0   RowNumber             10000 non-null  int64
1   CustomerId            10000 non-null  int64
2   Surname                10000 non-null  object
3   CreditScore            10000 non-null  int64
4   Geography              10000 non-null  object
5   Gender                 10000 non-null  object
6   Age                   10000 non-null  int64
7   Tenure                 10000 non-null  int64
8   Balance                10000 non-null  float64
9   NumOfProducts          10000 non-null  int64
10  HasCrCard              10000 non-null  int64
11  IsActiveMember         10000 non-null  int64
12  EstimatedSalary         10000 non-null  float64
13  Exited                  10000 non-null  int64
dtypes: float64(2), int64(9), object(3)
memory usage: 1.1+ MB
```

```
In [7]: df.describe()
Out[7]:
```

| | RowNumber | CustomerId | CreditScore | Age | Tenure | 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 |
| mean | 5000.50000 | 1.569094e+07 | 650.528800 | 38.921800 | 5.012800 | 76485.889288 | 1.530200 | 0.70550 | 0.515100 | 100090.239881 | 0.203700 |
| std | 2886.89568 | 7.193619e+04 | 96.653299 | 10.487806 | 2.892174 | 62397.405202 | 0.581654 | 0.45584 | 0.499797 | 57510.492818 | 0.402769 |
| min | 1.00000 | 1.565670e+07 | 350.000000 | 18.000000 | 0.000000 | 0.000000 | 1.000000 | 0.000000 | 0.000000 | 11.580000 | 0.000000 |
| 25% | 2500.75000 | 1.562853e+07 | 584.000000 | 32.000000 | 3.000000 | 0.000000 | 1.000000 | 0.000000 | 0.000000 | 51002.110000 | 0.000000 |
| 50% | 5000.50000 | 1.569074e+07 | 652.000000 | 37.000000 | 5.000000 | 97198.540000 | 1.000000 | 1.000000 | 1.000000 | 100193.915000 | 0.000000 |
| 75% | 7500.25000 | 1.575323e+07 | 718.000000 | 44.000000 | 7.000000 | 127644.240000 | 2.000000 | 1.000000 | 1.000000 | 149388.247500 | 0.000000 |
| max | 10000.00000 | 1.581569e+07 | 850.000000 | 92.000000 | 10.000000 | 250896.090000 | 4.000000 | 1.000000 | 1.000000 | 199992.480000 | 1.000000 |

```
In [8]: df.Geography.value_counts()
Out[8]: France      5614
Germany     2599
Spain       2477
Name: Geography, dtype: int64
```

```
In [9]: df.Surname.value_counts()
Out[9]: Smith      32
Scott      29
Martin     29
Walker     28
Brown      28
..
Izmailov   1
Boile      1
Bonhain    1
Poninski   1
Burnside   1
Name: Surname, Length: 2932, dtype: int64
```

Visualizations.

Univariate Analysis

```
In [10]: sns.displot(df.Tenure)
Out[10]: <seaborn.axisgrid.FacetGrid at 0x239e957b010>
```



```
In [23]: plt.pie(df.Geography.value_counts(), [0, 0, 0.1], autopct='%1.1f%%', labels=['France', 'Germany', 'Spain'], shadow=False, colors=['blue', 'red', 'green'])
plt.show()
```

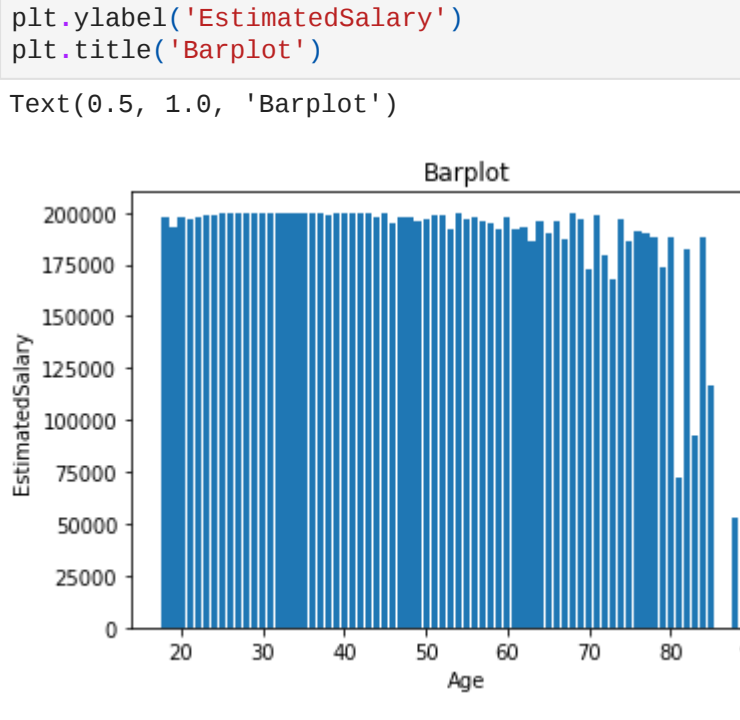


```
In [18]: sns.histplot(df.Age)
Out[18]: <AxesSubplot: xlabel='Age', ylabel='Count'>
```

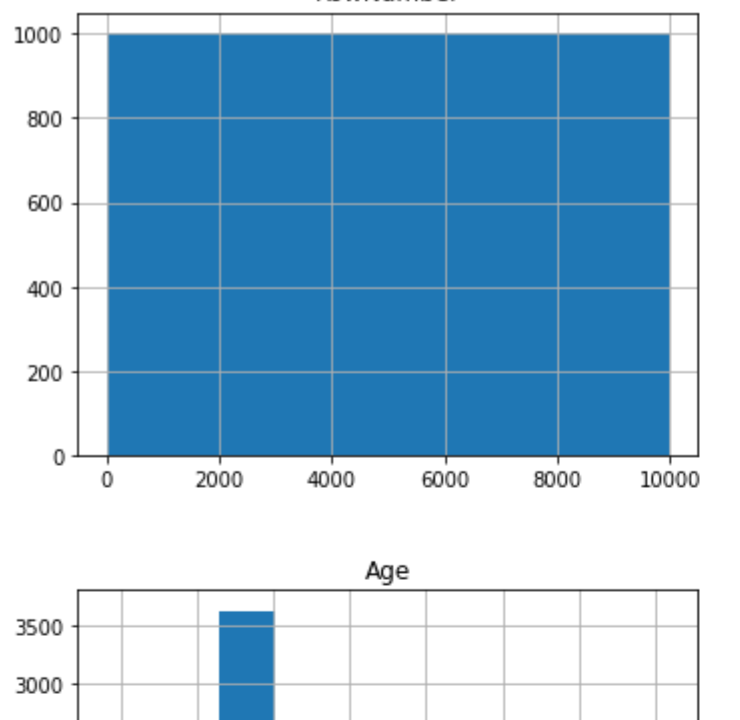


Bi - Variate Analysis

```
In [26]: plt.bar(df.Tenure, df.CreditScore)
Out[26]: <BarContainer object of 10000 artists>
```



```
In [29]: plt.bar(df.Age, df.EstimatedSalary)
plt.xlabel('Age')
plt.ylabel('EstimatedSalary')
plt.title('Barplot')
Text(0.5, 1.0, 'Barplot')
```



Multi - Variate Analysis

```
In [32]: df.hist(figsize=[20,20])
Out[32]: array([[<AxesSubplot: title='{center': 'RowNumber'}>,
<AxesSubplot: title='{center': 'CustomerId'}>,
<AxesSubplot: title='{center': 'CreditScore'}>],
[<AxesSubplot: title='{center': 'Age'}>,
<AxesSubplot: title='{center': 'Tenure'}>,
<AxesSubplot: title='{center': 'Balance'}>],
[<AxesSubplot: title='{center': 'NumOfProducts'}>,
<AxesSubplot: title='{center': 'HasCrCard'}>,
<AxesSubplot: title='{center': 'IsActiveMember'}>],
[<AxesSubplot: title='{center': 'EstimatedSalary'}>,
<AxesSubplot: title='{center': 'Exited'}>],
dtype=object])
```



Descriptive statistics on the dataset.

```
In [33]: df.describe()
Out[33]:
```

| | RowNumber | CustomerId | CreditScore | Age | Tenure | 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 |
| mean | 5000.50000 | 1.569094e+07 | 650.528800 | 38.921800 | 5.012800 | 76485.889288 | 1.530200 | 0.70550 | 0.515100 | 100090.239881 | 0.203700 |
| std | 2886.89568 | 7.193619e+04 | 96.653299 | 10.487806 | 2.892174 | 62397.405202 | 0.581654 | 0.45584 | 0.499797 | 57510.492818 | 0.402769 |
| min | 1.00000 | 1.565670e+07 | 350.000000 | 18.000000 | 0.000000 | 0.000000 | 1.000000 | 0.000000 | 0.000000 | 11.580000 | 0.000000 |
| 25% | 2500.75000 | 1.562853e+07 | 584.000000 | 32.000000 | 3.000000 | 0.000000 | 1.000000 | 0.000000 | 0.000000 | 51002.110000 | 0.000000 |
| 50% | 5000.50000 | 1.569074e+07 | 652.000000 | 37.000000 | 5.000000 | 97198.540000 | 1.000000 | 1.000000 | 1.000000 | 100193.915000 | 0.000000 |
| 75% | 7500.25000 | 1.575323e+07 | 718.000000 | 44.000000 | 7.000000 | 127644.240000 | 2.000000 | 1.000000 | 1.000000 | 149388.247500 | 0.000000 |
| max | 10000.00000 | 1.581569e+07 | 850.000000 | 92.000000 | 10.000000 | 250896.090000 | 4.000000 | 1.000000 | 1.000000 | 199992.480000 | 1.000000 |

```
In [34]: df.corr()
C:\Users\LECO\AppData\Local\Temp\ipykernel_6756\1134722465.py:1: FutureWarning: The default value of numeric_only in DataFrame.corr is deprecated. In a future version, it will default to False. Select only valid columns or specify the value of numeric_only to silence this warning.
df.corr()
```

| | RowNumber | CustomerId | CreditScore | Age | Tenure | Balance | NumOfProducts | HasCrCard | IsActiveMember | EstimatedSalary | Exited |
|-----------------|-----------|------------|-------------|-----------|-----------|-----------|---------------|-----------|----------------|-----------------|-----------|
| RowNumber | 1.000000 | 0.004202 | 0.000000 | 0.005840 | 0.000783 | -0.006495 | -0.009067 | 0.007246 | 0.000599 | 0.012044 | -0.005988 |
| CustomerId | 0.004202 | 1.000000 | 0.005308 | 0.009497 | -0.014883 | -0.012419 | 0.016972 | -0.014025 | 0.001665 | -0.015271 | -0.006671 |
| CreditScore | 0.000000 | 0.005308 | 1.000000 | -0.003965 | 0.000842 | 0.006268 | 0.012238 | -0.005458 | 0.025651 | -0.001384 | -0.027094 |
| Age | 0.000783 | 0.009497 | -0.003965 | 1.000000 | -0.009997 | 0.028308 | -0.030680 | -0.011721 | 0.085472 | -0.007201 | 0.285323 |
| Tenure | -0.006495 | -0.012419 | 0.000842 | -0.009997 | 1.000000 | -0.012254 | 0.013444 | -0.022583 | -0.028362 | 0.007784 | -0.014001 |
| Balance | -0.009067 | -0.012419 | 0.006268 | 0.028308 | -0.012254 | 1.000000 | -0.304180 | -0.014858 | -0.010084 | 0.012797 | -0.118533 |
| NumOfProducts | 0.007246 | 0.016972 | 0.012238 | -0.030680 | 0.013444 | -0.304180 | 1.000000 | 0.003183 | 0.009612 | -0.014204 | -0.047180 |
| HasCrCard | 0.000599 | -0.014025 | -0.005458 | -0.011721 | 0.022583 | -0.014858 | 0.003183 | 1.000000 | -0.011866 | -0.009933 | -0.007320 |
| IsActiveMember | 0.012044 | 0.001665 | 0.025651 | 0.085472 | -0.028362 | -0.010084 | 0.009612 | -0.011866 | 1.000000 | -0.011421 | -0.156128 |
| EstimatedSalary | -0.005988 | -0.015271 | -0.001384 | -0.007201 | 0.007784 | 0.012797 | 0.014204 | -0.009933 | -0.011421 | 1.000000 | 0.012097 |
| Exited | -0.016571 | -0.006248 | -0.027094 | 0.285323 | -0.014001 | 0.118533 | -0.047180 | -0.007138 | -0.156128 | 0.012097 | 1.000000 |

```
In [35]: df.CreditScore.mean()
Out[35]: 650.5288
```

```
In [38]: df.median()
C:\Users\LECO\AppData\Local\Temp\ipykernel_6756\530051474.py:1: FutureWarning: The default value of numeric_only in DataFrame.median is deprecated. In a future version, it will default to False. In addition, specifying 'numeric_only=None' is deprecated. Select only valid columns or specify the value of numeric_only to silence this warning.
df.median()
```

| | RowNumber | CustomerId | CreditScore | Age | Tenure | Balance | NumOfProducts | HasCrCard | IsActiveMember | EstimatedSalary | Exited |
|-----------------|--------------|------------|-------------|-----|--------|---------|---------------|-----------|----------------|-----------------|--------|
| RowNumber | 5.000500e+03 | | | | | | | | | | |
| CustomerId | 1.569074e+07 | | | | | | | | | | |
| CreditScore | 6.520000e+02 | | | | | | | | | | |
| Age | 3.789000e+01 | | | | | | | | | | |
| Tenure | 5.000000e+00 | | | | | | | | | | |
| Balance | 9.719854e+04 | | | | | | | | | | |
| NumOfProducts | 1.000000e+00 | | | | | | | | | | |
| HasCrCard | 1.000000e+00 | | | | | | | | | | |
| IsActiveMember | 1.000000e+00 | | | | | | | | | | |
| EstimatedSalary | 1.001939e+05 | | | | | | | | | | |
| Exited | 0.000000e+00 | | | | | | | | | | |
| dtype: | float64 | | | | | | | | | | |

Handle the Missing values

```
In [37]: df.isnull().any()
Out[37]: RowNumber      False
CustomerId      False
Surname         False
CreditScore     False
Geography       False
Gender          False
Age             False
Tenure          False
Balance         False
NumOfProducts  False
HasCrCard       False
IsActiveMember  False
EstimatedSalary False
Exited          False
dtype: bool
```

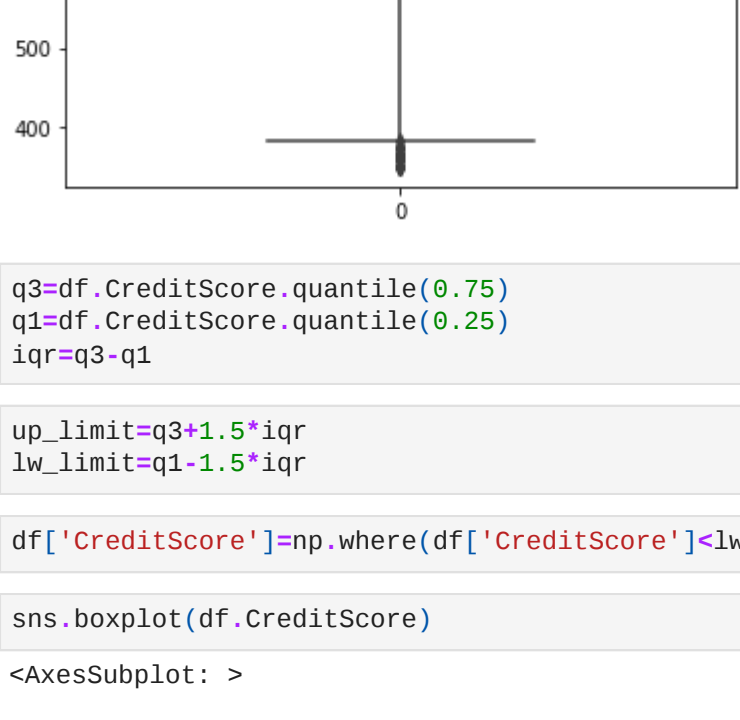
Outliers and replace the outliers

```
In [38]: df.shape
Out[38]: (10000, 14)
```

```
In [39]: df.median()
C:\Users\LECO\AppData\Local\Temp\ipykernel_6756\530051474.py:1: FutureWarning: The default value of numeric_only in DataFrame.median is deprecated. In a future version, it will default to False. In addition, specifying 'numeric_only=None' is deprecated. Select only valid columns or specify the value of numeric_only to silence this warning.
df.median()
```

| | RowNumber | CustomerId | CreditScore | Age | Tenure | Balance | NumOfProducts | HasCrCard | IsActiveMember | EstimatedSalary | Exited |
|-----------------|--------------|------------|-------------|-----|--------|---------|---------------|-----------|----------------|-----------------|--------|
| RowNumber | 5.000500e+03 | | | | | | | | | | |
| CustomerId | 1.569074e+07 | | | | | | | | | | |
| CreditScore | 6.520000e+02 | | | | | | | | | | |
| Age | 3.789000e+01 | | | | | | | | | | |
| Tenure | 5.000000e+00 | | | | | | | | | | |
| Balance | 9.719854e+04 | | | | | | | | | | |
| NumOfProducts | 1.000000e+00 | | | | | | | | | | |
| HasCrCard | 1.000000e+00 | | | | | | | | | | |
| IsActiveMember | 1.000000e+00 | | | | | | | | | | |
| EstimatedSalary | 1.001939e+05 | | | | | | | | | | |
| Exited | 0.000000e+00 | | | | | | | | | | |
| dtype: | float64 | | | | | | | | | | |

```
In [40]: sns.boxplot(df.CreditScore)
Out[40]: <AxesSubplot: >
```

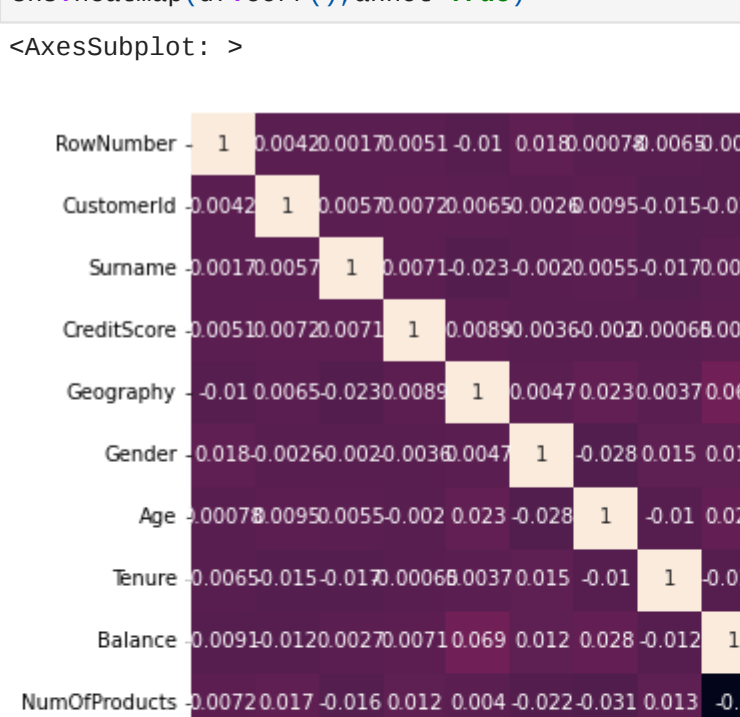


```
In [41]: q3=df.CreditScore.quantile(0.75)
q1=df.CreditScore.quantile(0.25)
iqr=q3-q1
```

```
In [42]: up_limit=q3+1.5*iqr
lw_limit=q1-1.5*iqr
```

```
In [43]: df['CreditScore']=np.where(df['CreditScore']<lw_limit,652,df['CreditScore'])
```

```
In [44]: sns.boxplot(df.CreditScore)
Out[44]: <AxesSubplot: >
```



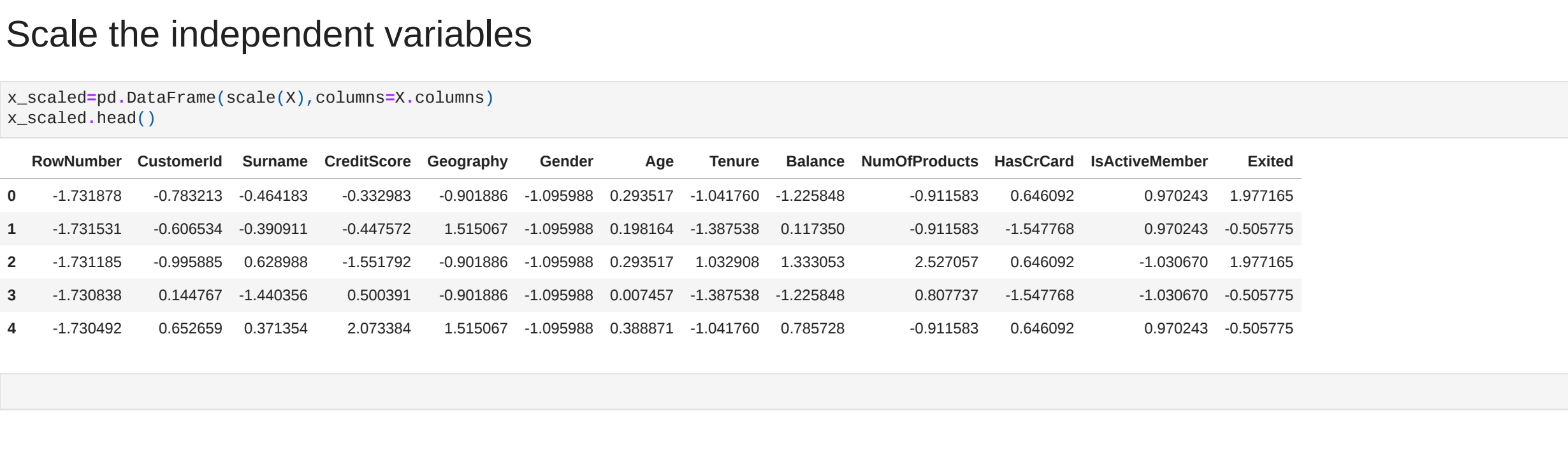
Categorical columns and perform encoding.

```
In [45]: from sklearn.preprocessing import LabelEncoder
le=LabelEncoder()
```

```
In [46]: df.Surname=le.fit_transform(df.Surname)
df.Geography=le.fit_transform(df.Geography)
df.Gender=le.fit_transform(df.Gender)
df.head()
```

| | RowNumber | CustomerId | Surname | CreditScore | Geography | Gender | Age | Tenure | Balance | NumOfProducts | HasCrCard | IsActiveMember | EstimatedSalary | Exited |
|---|-----------|------------|---------|-------------|-----------|--------|-----|--------|-----------|---------------|-----------|----------------|-----------------|--------|
| 0 | 1 | 15634602 | 1117 | 619 | 0 | 0 | 42 | 2 | 0.00 | 1 | 1 | 1 | 101348.88 | 1 |
| 1 | 2 | 15647311 | 1175 | 608 | 2 | 0 | 41 | 1 | 83807.86 | 1 | 0 | 1 | 112542.58 | 0 |
| 2 | 3 | 15616304 | 2040 | 502 | 0 | 0 | 42 | 8 | 159660.80 | 3 | 1 | 0 | 113931.57 | 1 |
| 3 | 4 | 15701354 | 289 | 699 | 0 | 0 | 39 | 1 | 0.00 | 2 | 0 | 0 | 93826.63 | 0 |
| 4 | 5 | 15737888 | 1822 | 850 | 2 | 0 | 43 | 2 | 125510.82 | 1 | 1 | 1 | 79084.10 | 0 |

```
In [47]: plt.figure(figsize=[10,8])
sns.heatmap(df.corr(),annot=True)
```



Split the data into dependent and independent variables.

```
In [48]: y=df['EstimatedSalary']
print(y)
```

| | |
|------|-----------|
| 0 | 101348.88 |
| 1 | 112542.58 |
| 2 | 113931.57 |
| 3 | 93826.63 |
| 4 | 79084.10 |
| ... | ... |
| 9995 | 96279.64 |
| 9996 | 101699.77 |
| 9997 | 42065.58 |
| 9998 | 92888.52 |
| 9999 | 38190.78 |

```
Name: EstimatedSalary, Length: 10000, dtype: float64
```

```
In [49]: X=df.drop(columns=['EstimatedSalary'],axis=1)
X.head()
```

| | RowNumber | CustomerId | Surname | CreditScore | Geography | Gender | Age | Tenure | Balance | NumOfProducts | HasCrCard | IsActiveMember | Exited |
|---|-----------|------------|---------|-------------|-----------|--------|-----|--------|-----------|---------------|-----------|----------------|-----------|
| 0 | 1 | 15634602 | 1117 | 619 | 0 | 0 | 42 | 2 | 0.00 | 1 | 1 | 1 | 101348.88 |
| 1 | 2 | 15647311 | 1175 | 608 | 2 | 0 | 41 | 1 | 83807.86 | 1 | 0 | 1 | 112542.58 |
| 2 | 3 | 15616304 | 2040 | 502 | 0 | 0 | 42 | 8 | 159660.80 | 3 | 1 | 0 | 113931.57 |
| 3 | 4 | 15701354 | 289 | 699 | 0 | 0 | 39 | 1 | 0.00 | 2 | 0 | 0 | 93826.63 |
| 4 | 5 | 15737888 | 1822 | 850 | 2 | 0 | 43 | 2 | 125510.82 | 1 | 1 | | |