

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

loading the dataset

```
In [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	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 [3]: df.shape
Out[3]: (10000, 14)
```

```
In [4]: 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 [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
```

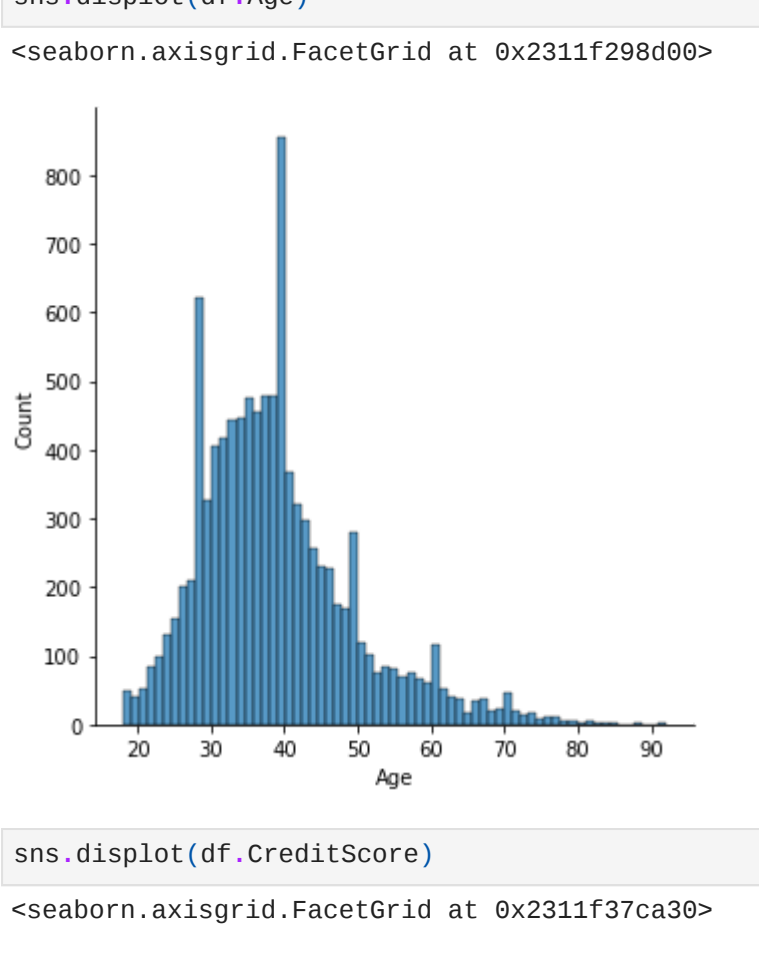
```
In [6]: df.describe()
```

	RowNumber	CustomerId	CreditScore	Age	Tenure	Balance	NumOfProducts	HasCrCard	IsActiveMember	EstimatedSalary	Exited
count	10000.000000	10000000+04	10000.000000	10000.000000	10000.000000	10000.000000	10000.000000	10000.000000	10000.000000	10000.000000	10000.000000
mean	5000.500000	1.569094e+07	650.528890	38.921800	5.012800	76485.889288	1.530200	0.70550	0.515100	100090.239881	0.203700
std	2886.89568	7.193619e+04	86.653299	10.487806	2.892174	62397.405202	0.581654	0.45584	0.499797	57510.492818	0.402769
min	1.00000	1.565570e+07	350.000000	18.000000	0.000000	0.000000	1.000000	0.00000	0.000000	11.580000	0.000000
25%	2500.75000	1.562836e+07	584.000000	32.000000	3.000000	0.000000	1.000000	0.00000	0.000000	51002.110000	0.000000
50%	5000.50000	1.569074e+07	652.000000	37.000000	5.000000	97198.540000	1.000000	1.00000	1.000000	100193.915000	0.000000
75%	7500.25000	1.575323e+07	718.000000	44.000000	7.000000	127644.240000	2.000000	1.00000	1.000000	149388.247500	0.000000
max	10000.00000	1.581509e+07	850.000000	92.000000	10.000000	250898.090000	4.000000	1.00000	1.000000	199992.480000	1.000000

Univariate analysis

```
In [7]: sns.displot(df.Age)
```

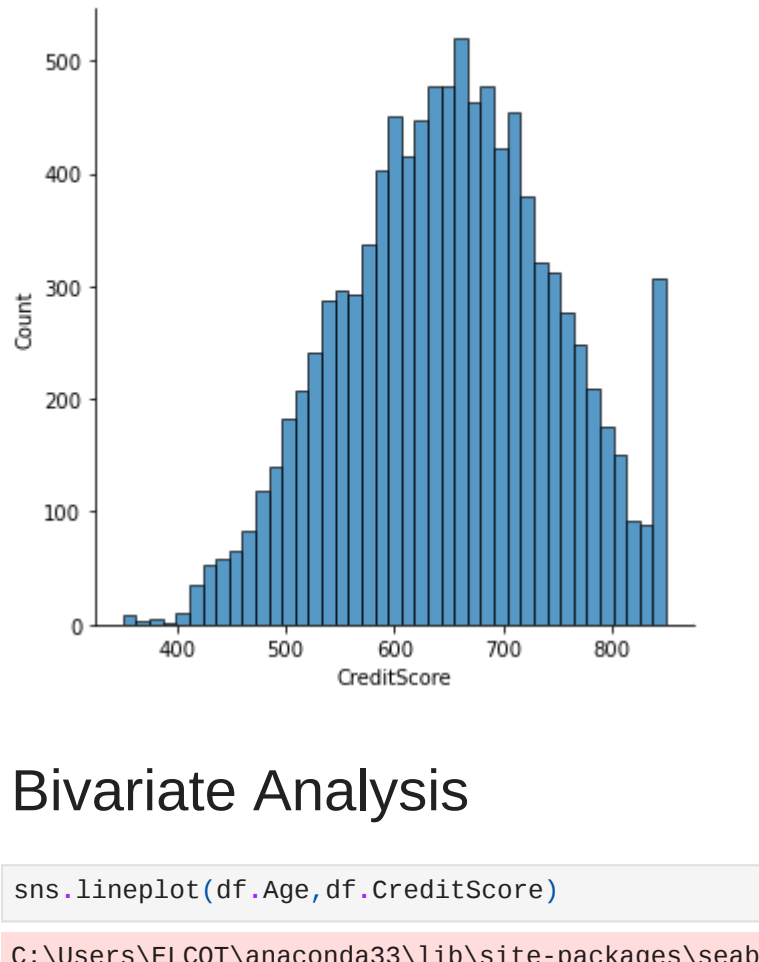
Out[7]: <seaborn.axisgrid.FacetGrid at 0x2311f298d00>



A histogram showing the distribution of Age. The x-axis is labeled 'Age' and ranges from 20 to 90. The y-axis is labeled 'Count' and ranges from 0 to 800. The distribution is unimodal and slightly right-skewed, with a peak count of approximately 800 around age 40.

```
In [8]: sns.displot(df.CreditScore)
```

Out[8]: <seaborn.axisgrid.FacetGrid at 0x2311f37ca30>



A histogram showing the distribution of CreditScore. The x-axis is labeled 'CreditScore' and ranges from 400 to 800. The y-axis is labeled 'Count' and ranges from 0 to 500. The distribution is unimodal and slightly left-skewed, with a peak count of approximately 500 around a credit score of 650.

Bivariate Analysis

```
In [9]: sns.lineplot(df.Age,df.CreditScore)
```

C:\Users\ELCOT\anaconda3\lib\site-packages\seaborn_decorators.py:36: FutureWarning: Pass the following variables as keyword args: x, y. From version 0.12, the only valid positional argument will be 'data', and passing other arguments without an explicit keyword will result in an error or misinterpretation.

warnings.warn(
<AxesSubplot: xlabel='Age', ylabel='CreditScore'>



A line plot showing the relationship between Age and CreditScore. The x-axis is labeled 'Age' and ranges from 20 to 90. The y-axis is labeled 'CreditScore' and ranges from 400 to 800. The plot shows a noisy trend line with a light blue shaded area representing the confidence interval. The credit score generally increases with age, with some fluctuations.

```
In [10]: sns.scatterplot(df.Age,df.CreditScore)
```

C:\Users\ELCOT\anaconda3\lib\site-packages\seaborn_decorators.py:36: FutureWarning: Pass the following variables as keyword args: x, y. From version 0.12, the only valid positional argument will be 'data', and passing other arguments without an explicit keyword will result in an error or misinterpretation.

warnings.warn(
<AxesSubplot: xlabel='Age', ylabel='CreditScore'>



A scatter plot showing the relationship between Age and CreditScore. The x-axis is labeled 'Age' and ranges from 20 to 90. The y-axis is labeled 'CreditScore' and ranges from 400 to 800. The plot shows a dense cloud of blue points, indicating a positive correlation between age and credit score.

```
In [11]: df.Gender.value_counts()
```

Out[11]: Male 5457
Female 4543
Name: Gender, dtype: int64

```
In [12]: df.Geography.value_counts()
```

Out[12]: France 5614
Germany 2599
Spain 2477
Name: Geography, dtype: int64

```
In [13]: plt.pie(df.Gender.value_counts(),labels=['Male','Female'])
```

Out[13]: ([<matplotlib.patches.Wedge at 0x2311f8e2250>,
<matplotlib.patches.Wedge at 0x2311f5e2730>],
[Text(-0.1573859183753977, 1.088682539906438, 'Male'),
Text(0.1573859183753977, -1.088682539906438, 'Female')])

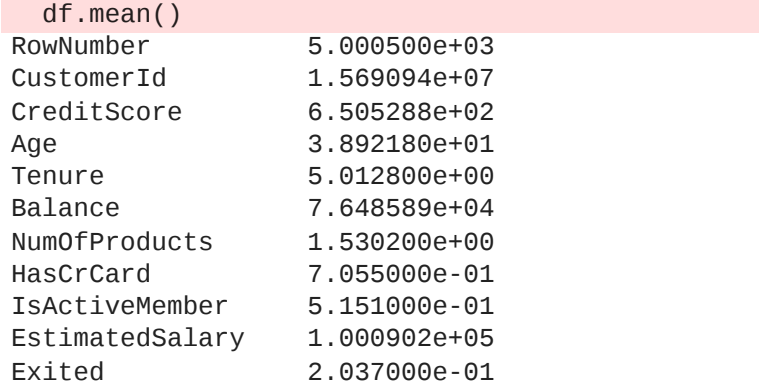


A pie chart showing the distribution of Gender. The chart is divided into two segments: a larger blue segment for 'Male' (approximately 54.6%) and a smaller orange segment for 'Female' (approximately 45.4%).

```
In [14]: sns.barpplot(df.Geography.value_counts().index,df.Geography.value_counts())
```

C:\Users\ELCOT\anaconda3\lib\site-packages\seaborn_decorators.py:36: FutureWarning: Pass the following variables as keyword args: x, y. From version 0.12, the only valid positional argument will be 'data', and passing other arguments without an explicit keyword will result in an error or misinterpretation.

warnings.warn(
<AxesSubplot: ylabel='Geography'>



A bar plot showing the distribution of Geography. The x-axis is labeled with the categories 'France', 'Germany', and 'Spain'. The y-axis is labeled 'Geography' and ranges from 0 to 5000. The bars are colored blue for France (approx. 5600), orange for Germany (approx. 2600), and green for Spain (approx. 2500).

Descriptive Statistics

```
In [15]: 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	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 [16]: df.mean()
```

C:\Users\ELCOT\AppData\Local\Temp\ipykernel_7404\3698961737.py:1: FutureWarning: Dropping of nuisance columns in DataFrame reductions (with 'numeric_only=None') is deprecated; in a future version this will raise TypeError. Select only valid columns before calling the reduction.

df.mean()

	RowNumber	CustomerId	CreditScore	Age	Tenure	Balance	NumOfProducts	HasCrCard	IsActiveMember	EstimatedSalary	Exited
RowNumber	5.000500e+03	1.569094e+07	6.505289e+02	3.892180e+01	5.012800e+00	7.648589e+04	1.530200e+00	7.055000e-01	5.151000e-01	1.000902e+05	2.037000e-01
CustomerId	1.569074e+07	6.520800e+02	3.708080e+01	5.009090e+00	9.719854e+04	1.000000e+00	1.000000e+00	1.000000e+00	1.000000e+00	1.001939e+05	0.000000e+00
CreditScore	6.520800e+02	3.708080e+01	5.009090e+00	9.719854e+04	1.000000e+00	1.000000e+00	1.000000e+00	1.000000e+00	1.000000e+00	1.001939e+05	0.000000e+00
Age	3.892180e+01	5.012800e+00	7.648589e+04	1.530200e+00	7.055000e-01	5.151000e-01	1.000000e+00	1.000000e+00	1.000000e+00	1.001939e+05	2.037000e-01
Tenure	5.012800e+00	9.719854e+04	1.000000e+00	1.000000e+00	1.000000e+00	1.000000e+00	1.000000e+00	1.000000e+00	1.000000e+00	1.001939e+05	2.037000e-01
Balance	7.648589e+04	1.530200e+00	7.055000e-01	5.151000e-01	1.000000e+00	1.000000e+00	1.000000e+00	1.000000e+00	1.000000e+00	1.001939e+05	2.037000e-01
NumOfProducts	1.530200e+00	7.055000e-01	5.151000e-01	1.000000e+00	1.000000e+00	1.000000e+00	1.000000e+00	1.000000e+00	1.000000e+00	1.001939e+05	2.037000e-01
HasCrCard	7.055000e-01	5.151000e-01	1.000000e+00	1.000000e+00	1.000000e+00	1.000000e+00	1.000000e+00	1.000000e+00	1.000000e+00	1.001939e+05	2.037000e-01
IsActiveMember	5.151000e-01	1.000000e+00	1.000000e+00	1.000000e+00	1.000000e+00	1.000000e+00	1.000000e+00	1.000000e+00	1.000000e+00	1.001939e+05	2.037000e-01
EstimatedSalary	1.000902e+05	2.037000e-01	0.000000e+00	0.000000e+00	0.000000e+00	0.000000e+00	0.000000e+00	0.000000e+00	0.000000e+00	0.000000e+00	0.000000e+00
Exited	2.037000e-01	0.000000e+00	0.000000e+00	0.000000e+00	0.000000e+00	0.000000e+00	0.000000e+00	0.000000e+00	0.000000e+00	0.000000e+00	0.000000e+00
dtype:	Float64										

```
In [17]: df.median()
```

C:\Users\ELCOT\AppData\Local\Temp\ipykernel_7404\530951474.py:1: FutureWarning: Dropping of nuisance columns in DataFrame reductions (with 'numeric_only=None') is deprecated; in a future version this will raise TypeError. Select only valid columns before calling the reduction.

df.median()

	RowNumber	CustomerId	CreditScore	Age	Tenure	Balance	NumOfProducts	HasCrCard	IsActiveMember	EstimatedSalary	Exited
RowNumber	5.000500e+03	1.569074e+07	6.520800e+02	3.892180e+01	5.012800e+00	7.648589e+04	1.530200e+00	7.055000e-01	5.151000e-01	1.000902e+05	2.037000e-01
CustomerId	1.569074e+07	6.520800e+02	3.708080e+01	5.009090e+00	9.719854e+04	1.000000e+00	1.000000e+00	1.000000e+00	1.000000e+00	1.001939e+05	0.000000e+00
CreditScore	6.520800e+02	3.708080e+01	5.009090e+00	9.719854e+04	1.000000e+00	1.000000e+00	1.000000e+00	1.000000e+00	1.000000e+00	1.001939e+05	0.000000e+00
Age	3.892180e+01	5.012800e+00	7.648589e+04	1.530200e+00	7.055000e-01	5.151000e-01	1.000000e+00	1.000000e+00	1.000000e+00	1.001939e+05	2.037000e-01
Tenure	5.012800e+00	9.719854e+04	1.000000e+00	1.000000e+00	1.000000e+00	1.000000e+00	1.000000e+00	1.000000e+00	1.000000e+00	1.001939e+05	2.037000e-01
Balance	7.648589e+04	1.530200e+00	7.055000e-01	5.151000e-01	1.000000e+00	1.000000e+00	1.000000e+00	1.000000e+00	1.000000e+00	1.001939e+05	2.037000e-01
NumOfProducts	1.530200e+00	7.055000e-01	5.151000e-01	1.000000e+00	1.000000e+00	1.000000e+00	1.000000e+00	1.000000e+00	1.000000e+00	1.001939e+05	2.037000e-01
HasCrCard	7.055000e-01	5.151000e-01	1.000000e+00	1.000000e+00	1.000000e+00	1.000000e+00	1.000000e+00	1.000000e+00	1.000000e+00	1.001939e+05	2.037000e-01
IsActiveMember	5.151000e-01	1.000000e+00	1.000000e+00	1.000000e+00	1.000000e+00	1.000000e+00	1.000000e+00	1.000000e+00	1.000000e+00	1.001939e+05	2.037000e-01
EstimatedSalary	1.000902e+05	2.037000e-01	0.000000e+00	0.000000e+00	0.000000e+00	0.000000e+00	0.000000e+00	0.000000e+00	0.000000e+00	0.000000e+00	0.000000e+00
Exited	2.037000e-01	0.000000e+00	0.000000e+00	0.000000e+00	0.000000e+00	0.000000e+00	0.000000e+00	0.000000e+00	0.000000e+00	0.000000e+00	0.000000e+00
dtype:	Float64										

```
In [18]: df.mode()
```

	RowNumber	CustomerId	Surname	CreditScore	Geography	Gender	Age	Tenure	Balance	NumOfProducts	HasCrCard	IsActiveMember	EstimatedSalary	Exited
0	1	15655701	Smith	850.0	France	Male	37.0	2.0	0.0	1.0	1.0	1.0	24924.92	0.0
1	2	15655706	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN
2	3	15655714	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN
3	4	15655779	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN
4	5	15655796	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN
...
9995	9996	15815628	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN
9996	9997	15815645	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN
9997	9998	15815656	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN
9998	9999	15815660	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN
9999	10000	15815690	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN

10000 rows x 14 columns

Handle missing values

```
In [19]: df2=df.fillna(value=0)
```

```
In [20]: df2
```

	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	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
...
9995	9996	15606229	Obijaku	771	France	Male	39	5	0.00	2	1	0	96270.64	0
9996	9997	15668292	Johnstone											