

```
import pandas as pd
import numpy as np
from sklearn.preprocessing import MinMaxScaler
from sklearn.metrics import confusion_matrix, accuracy_score
import matplotlib.pyplot as plt
import seaborn as sns
```

```
In [106]: #import the dataset
dataset = pd.read_csv('D:\IBM\Churn_Modelling.csv')
dataset.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	83007.86	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	Michell	850	Spain	Female	43	2	125510.82	1	1	1	79084.10	0

```
In [107]: dataset.tail()
```

9999	9999	15682395	Sabbatini	772	Germany	Male	42	3	75075.31	2	1	0	92888.52	1
9999	10000	15628319	Walker	792	France	Female	28	4	120142.79	1	1	0	38190.78	0

```
In [108]: dataset.describe()
```

```
Out[108]:
```

	RowNumber	CustomerId	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	10000.000000	10000.000000

```
In [108]: dataset.describe()
```

	RowNumber	CustomerId	CreditScore	Age	Tenure	Balance	NumOfProducts	HasCrCard	IsActiveMember	EstimatedSalary	Exited
count	10000	10000	1.0	95000.0e+04	10000	1000000	10000	1000000	10000	1000000	10000
mean	5000	50000	1.56398e+07	650.528800	38.921800	5.012800	76485.880388	1.530000	0.705800	0.515100	100060.239881
std	2086.95568	7.13961e+04	96.563200	10.487006	2.882174	62387.405202	1.581654	0.458400	0.499797	57510.628818	0.402769
min	1	10000	1.56398e+07	350.000000	18.000000	0.000000	0.000000	0.000000	0.000000	0.000000	51002.110000
50%	5000	50000	1.56285e+07	584.000000	32.000000	3.000000	0.000000	1.000000	0.000000	51002.110000	0.000000
75%	7500	75000	1.56907e+07	652.000000	37.000000	5.000000	97188.540000	1.000000	1.000000	1.000000	149388.915000
90%	9000	90000	1.57532e+07	718.000000	44.000000	7.000000	127644.240000	2.000000	1.000000	1.000000	149388.915000
max	10000	100000	1.56156e+07	850.000000	92.000000	10.000000	250898.000000	4.000000	1.000000	1.000000	199982.450000

## univariate analysis

```
In [109]: #histogram
sns.distplot(dataset['CustomerId'], kde=False)
```

C:\Users\vinhai\AppData\Local\Temp\ipykernel\_2619\7619: FutureWarning: 'distplot' is a deprecated function and will be removed in a future version. Please adapt your code to use either 'displot' (a figure-level function with similar flexibility) or 'histplot' (an axes-level function for histograms).

```
Out[109]: <AxesSubplot: xlabel='CustomerId', ylabel='count'>
```



```
In [110]: #count plot
sns.countplot(x='CreditScore', data=dataset)
```

```
Out[110]: <AxesSubplot: xlabel='CreditScore', ylabel='count'>
```



## Bi-Variate Analysis

```
In [141]: #bar chart
sns.barplot(x='Surname', y='CreditScore', data=dataset)
```

```
Out[141]: <AxesSubplot: xlabel='Surname', ylabel='CreditScore'>
```



```
In [142]: #box plot
sns.boxplot(data=dataset, x='Surname', y='CreditScore')
```

```
Out[112]: <AxesSubplot: xlabel='Surname', ylabel='CreditScore'>
```



```
In [139]: #regression plot
sns.lmplot(x='EstimatedSalary', y='Balance', data=dataset)
```

```
Out[138]: <seaborn.axisgrid.FacetGrid at 0x261061f6678>
```



## Multivariate Analysis

```
In [37]: sns.pairplot(dataset)
```

```
Out[97]: <seaborn.axisgrid.PairGrid at 0x1013c78980>
```



## descriptive statistics on the dataset

```
In [51]: dataset = pd.read_csv('D:\IBM\Churn_Modelling.csv')
dataset.head()
```

```
Out[51]:
```

	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	83007.86	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	Michell	850	Spain	Female	43	2	125510.82	1	1	1	79084.10	0

```
In [55]: dataset.tail()
```

```
Out[55]:
```

```
in [10]: #box plot
sns.boxplot(data=dataset, x='Surname', y='CreditScore')
<AxesSubplot:xlabel='Surname', ylabel='CreditScore'>
```

```
In [57]:
```

```
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 [62]: #mean
dataset.mean()
```

C:\Users\vinhai\AppData\Local\Temp\ipykernel\_2888\165591842.py:3: FutureWarning: Dropping of nuisance columns in DataFrame reductions (with 'numeric\_only=no') is deprecated; in a future version this will raise TypeError. Select only valid columns before calling the reduction.

```
Out[62]:
```

RowNumber	0.005000e+03
CustomerId	1.56994e+07
CreditScore	6.58528e+02
Age	3.88218e+01
Tenure	5.01289e+00
Balance	5.04585e+04
NumOfProducts	1.53020e+00
HasCrCard	7.05509e+01
IsActiveMember	5.15390e+01
EstimatedSalary	1.00992e+05
Exited	2.03780e+01
dtype:	float64

```
In [63]: #median
dataset.median()
```

C:\Users\vinhai\AppData\Local\Temp\ipykernel\_2888\3979661894.py:3: FutureWarning: Dropping of nuisance columns in DataFrame reductions (with 'numeric\_only=no') is deprecated; in a future version this will raise TypeError. Select only valid columns before calling the reduction.

```
Out[63]:
```

RowNumber	0.005000e+03
CustomerId	3.56997e+07
CreditScore	6.52059e+02
Age	3.78009e+01
Tenure	5.00909e+00
Balance	9.71985e+04
NumOfProducts	1.00000e+00
HasCrCard	1.00000e+00
IsActiveMember	1.00000e+00
EstimatedSalary	1.00123e+05
Exited	0.00000e+00
dtype:	float64

```
In [64]: #mode
dataset.mode()
```

```
Out[64]:
```

	RowNumber	CustomerId	Surname	CreditScore	Geography	Gender	Age	Tenure	Balance	NumOfProducts	HasCrCard	IsActiveMember	EstimatedSalary	Exited
0	1	15635701	Smith	850.0	France	Male	37.0	2.0	0.0	1.0	1.0	1.0	24924.92	0.0
1	2	15665706	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN
2	3	15665714	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN
3	4	15665779	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN
4	5	15657768	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN

10000 rows x 14 columns

## handle the missing values

```
In [65]: dataset.isna()
```

```
Out[65]:
```

	RowNumber	CustomerId	Surname	CreditScore	Geography	Gender	Age	Tenure	Balance	NumOfProducts	HasCrCard	IsActiveMember	EstimatedSalary	Exited
0	False	False	False	False	False	False	False	False	False	False	False	False	False	False
1	False	False	False	False	False	False	False	False	False	False	False	False	False	False
2	False	False	False	False	False	False	False	False	False	False	False	False	False	False
3	False	False	False	False	False	False	False	False	False	False	False	False	False	False
4	False	False	False	False	False	False	False	False	False	False	False	False	False	False

10000 rows x 14 columns

```
In [66]: dataset.isna().any()
```

```
Out[66]:
```

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 [70]: #skewness
dataset.skew()
```

C:\Users\vinhai\AppData\Local\Temp\ipykernel\_2888\2006622695.py:3: FutureWarning: Dropping of nuisance columns in DataFrame reductions (with 'numeric\_only=no') is deprecated; in a future version this will raise TypeError. Select only valid columns before calling the reduction.

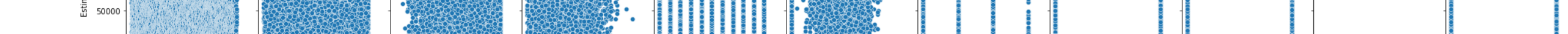
```
Out[70]:
```

RowNumber	0.000000
CustomerId	0.081149
CreditScore	-0.073807
Age	0.011329
Tenure	0.021691
Balance	-0.141199
NumOfProducts	0.582568
HasCrCard	-0.981812
IsActiveMember	-0.060437
EstimatedSalary	0.080425
Exited	1.473611
dtype:	float64

```
In [60]: print(stats.distplot(dataset['Age']))
```

C:\Users\vinhai\AppData\Local\Temp\ipykernel\_2888\1784932129.py:1: FutureWarning: Dropping of nuisance columns in DataFrame reductions (with 'numeric\_only=no') is deprecated; in a future version this will raise TypeError. Select only valid columns before calling the reduction.

```
Out[60]:
```



```
In [71]: dataset.kurt()
```

C:\Users\vinhai\AppData\Local\Temp\ipykernel\_2888\113884202.py:1: FutureWarning: Dropping of nuisance columns in DataFrame reductions (with 'numeric\_only=no') is deprecated; in a future version this will raise TypeError. Select only valid columns before calling the reduction.

```
Out[71]:
```

RowNumber	-1.209990
CustomerId	-1.196113
CreditScore	-0.425728
Age	1.395547
Tenure	-1.165225
Balance	-1.486412
NumOfProducts	0.582981
HasCrCard	-1.186972
IsActiveMember	-1.395290
EstimatedSalary	-1.181518
Exited	0.155671
dtype:	float64

```
In [72]: dataset.var()
```

C:\Users\vinhai\AppData\Local\Temp\ipykernel\_2888\2458428938.py:1: FutureWarning: Dropping of nuisance columns in DataFrame reductions (with 'numeric\_only=no') is deprecated; in a future version this will raise TypeError. Select only valid columns before calling the reduction.

```
Out[72]:
```

RowNumber	8.334167e+06
CustomerId	5.174351e+09
CreditScore	9.343588e+02
Age	1.099941e+02
Tenure	0.364673e+00
Balance	3.893436e+09
NumOfProducts	3.382238e+01
HasCrCard	2.077805e+01
IsActiveMember	2.497970e+01
EstimatedSalary	3.387437e+05
Exited	1.622225e+01
dtype:	float64

```
In [73]: dataset.std()
```

C:\Users\vinhai\AppData\Local\Temp\ipykernel\_2888\1784932129.py:1: FutureWarning: Dropping of nuisance columns in DataFrame reductions (with 'numeric\_only=no') is deprecated; in a future version this will raise TypeError. Select only valid columns before calling the reduction.

```
Out[73]:
```

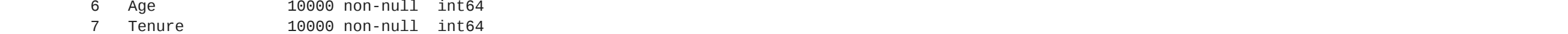
RowNumber	2886.895680
CustomerId	71895.186123
CreditScore	96.653299
Age	2.074673e+00
Tenure	2.892174
Balance	62397.485202
NumOfProducts	0.581654
HasCrCard	0.645840
IsActiveMember	0.499797
EstimatedSalary	57510.628818
Exited	0.402769
dtype:	float64

## find the outlier andreplace the outliers

```
In [80]: sns.boxplot(dataset['Age'])
```

C:\Users\vinhai\AppData\Local\Temp\ipykernel\_2619\15909084.py:36: FutureWarning: Pass the following variable as a keyword arg: x. 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.

```
Out[80]: <AxesSubplot: xlabel='Age'>
```



```
In [82]: qnt=dataset.quantile(q=[0.39,0.45])
qnt
```

	RowNumber	CustomerId	CreditScore	Age	Tenure	Balance	NumOfProducts	HasCrCard	IsActiveMember	EstimatedSalary	Exited
0.39	5000	70541363.90	598.7	33.0	3.0	0.000	1.0	1.0	0.0	60736.079	0.0
0.45	4500	15616399.85	639.0	36.0	5.0	87					