

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

In [2]:

```
df = pd.read_csv('/Churn_Modelling.csv')
df
```

Out[2]:

	Row Num ber	Cust omer Id	Sur na me	Cred itSco re	Geo grap hy	Ge nd er	A ge	Te nu re	Bal anc e	NumO fProdu cts	Has CrC ard	IsActiv eMem ber	Estima tedSal ary	Ex ite d
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
...
9995	9996	15606229	Obijaku	771	France	Male	39	5	0.00	2	1	0	96270.64	0
9996	9997	15569892	Johnstone	516	France	Male	35	10	57369.61	1	1	1	101699.77	0
9997	9998	15584532	Liu	709	France	Female	36	7	0.00	1	0	1	42085.58	1

	Row Num ber	Cust omer Id	Sur na me	Cred itSco re	Geo grap hy	Ge nd er	A ge	Te nu re	Bal anc e	NumO fProdu cts	Has CrC ard	IsActiv eMem ber	Estima tedSal ary	Ex ite d
9998	9999	15682355	Sab bati ni	772	Ger man y	Ma le	42	3	75075.31	2	1	0	92888.52	1
9999	10000	15628319	Wal ker	792	Fran ce	Fe ma le	28	4	130142.79	1	1	0	38190.78	0

10000 rows × 14 columns

1)Univariate Analysis

```
In [3]:
plt.scatter(df.index,df['Age'])

Out[3]:
<matplotlib.collections.PathCollection at 0x7fcded5146d0>

In [4]:
df['Age'].value_counts().head(12).plot.bar()

Out[4]:
<matplotlib.axes._subplots.AxesSubplot at 0x7fcdecfa96d0>

In [5]:
plt.hist(df['Age'])

Out[5]:
(array([ 611., 2179., 3629., 1871.,  828.,  523.,  208.,  127.,   20.,
         4.]),
 array([18. , 25.4, 32.8, 40.2, 47.6, 55. , 62.4, 69.8, 77.2, 84.6, 92. ]),
 <a list of 10 Patch objects>)
```

2)Bi-variate Analysis

```
In [6]:
sns.scatterplot(x='CreditScore',y='Balance',data=df,hue='Gender')
plt.show()

In [7]:
df.plot.line()

Out[7]:
<matplotlib.axes._subplots.AxesSubplot at 0x7fcdecd93350>
```

3)Multi-Variate Analysis

```
In [8]:
sns.pairplot(df)
plt.show()
```

4)Descriptive Statistics

In [9]:
df.describe()

Out[9]:

	RowN umbe r	Custo merId	Credit Score	Age	Tenur e	Balanc e	NumOf Product s	HasC rCar d	IsActive Membe r	Estimat edSalar y	Exited
co un t	10000 .0000 0	1.0000 00e+0 4	10000. 00000 0	10000. 00000 0	10000. 00000 0	10000. 000000	10000.0 00000	10000 .0000 0	10000.0 00000	10000.0 00000	10000. 00000 0
m ea n	5000. 50000	1.5690 94e+0 7	650.52 8800	38.921 800	5.0128 00	76485. 889288	1.53020 0	0.705 50	0.51510 0	100090. 239881	0.2037 00
st d	2886. 89568	7.1936 19e+0 4	96.653 299	10.487 806	2.8921 74	62397. 405202	0.58165 4	0.455 84	0.49979 7	57510.4 92818	0.4027 69
mi n	1.000 00	1.5565 70e+0 7	350.00 0000	18.000 000	0.0000 00	0.0000 00	1.00000 0	0.000 00	0.00000 0	11.5800 00	0.0000 00
25 %	2500. 75000	1.5628 53e+0 7	584.00 0000	32.000 000	3.0000 00	0.0000 00	1.00000 0	0.000 00	0.00000 0	51002.1 10000	0.0000 00
50 %	5000. 50000	1.5690 74e+0 7	652.00 0000	37.000 000	5.0000 00	97198. 540000	1.00000 0	1.000 00	1.00000 0	100193. 915000	0.0000 00
75 %	7500. 25000	1.5753 23e+0 7	718.00 0000	44.000 000	7.0000 00	127644 .24000 0	2.00000 0	1.000 00	1.00000 0	149388. 247500	0.0000 00
m ax	10000 .0000 0	1.5815 69e+0 7	850.00 0000	92.000 000	10.000 000	250898 .09000 0	4.00000 0	1.000 00	1.00000 0	199992. 480000	1.0000 00

In [11]:
df.sum()

Out[11]:

RowNumber	50005000
CustomerId	156909405694
Surname	HargraveHillOnioBoniMitchellChuBartlettObinnaH...
CreditScore	6505288
Geography	FranceSpainFranceFranceSpainSpainFranceGermany...
Gender	FemaleFemaleFemaleFemaleFemaleMaleMaleFemaleMa...

```

Age                                     389218
Tenure                                 50128
Balance                               764858892.88
NumOfProducts                         15302
HasCrCard                             7055
IsActiveMember                        5151
EstimatedSalary                       1000902398.81
Exited                                2037
dtype: object

```

In [12]:

```
df.min()
```

Out[12]:

```

RowNumber      1
CustomerId    15565701
Surname        Abazu
CreditScore    350
Geography      France
Gender         Female
Age            18
Tenure         0
Balance        0.0
NumOfProducts  1
HasCrCard      0
IsActiveMember 0
EstimatedSalary 11.58
Exited         0
dtype: object

```

In [13]:

```
df.max()
```

Out[13]:

```

RowNumber      10000
CustomerId    15815690
Surname        Zuyeva
CreditScore    850
Geography      Spain
Gender         Male
Age            92
Tenure         10
Balance        250898.09
NumOfProducts  4
HasCrCard      1
IsActiveMember 1
EstimatedSalary 199992.48
Exited         1
dtype: object

```

5)Handling Missing the Values

In [14]:

```
df.isnull()
```

Out[14]:

10000 rows × 14 columns

df.dropna()

In [15]:

Out[15]:

	Row Number	Customer Id	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	Michell	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	15569892	Johnstone	516	France	Male	35	10	57369.61	1	1	1	101699.77	0
9997	9998	15584532	Liu	709	France	Female	36	7	0.00	1	0	1	42085.58	1
9999	9999	15682355	Sabatin	772	Germany	Male	42	3	75075.31	2	1	0	92888.52	1

Row Num ber	Cust omer Id	Sur na me	Cred itSco re	Geo grap hy	Ge nd er	A g e	Te nu re	Bal anc e	NumO fProdu cts	Has CrC ard	IsActiv eMem ber	Estima tedSal ary	Ex ite d
9													
8													
9													
9	1000	1562	Wal		Fe	2		130				38190.	
9	0	8319	ker	792	ma	8	4	142.	1	1	0	78	0
9					le			79					

10000 rows × 14 columns

6)Find the outliers and replace the outliers

In [16]:

```
sns.boxplot(df['Age'])

/usr/local/lib/python3.7/dist-packages/seaborn/_decorators.py:43: 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
FutureWarning
```

Out[16]:

```
<matplotlib.axes._subplots.AxesSubplot at 0x7fcde5a6cf50>
```

In [17]:

```
print(np.where(df['Age']>48))

(array([ 6, 16, 41, ..., 9975, 9979, 9991]),)
```

In [18]:

```
missing_values=df.isnull().sum()
missing_values[missing_values>0]/len(df)*100
```

Out[18]:

```
Series([], dtype: float64)
```

In [19]:

```
Q1 = df['Age'].quantile(0.25)
Q3 = df['Age'].quantile(0.75)
IQR = Q3 - Q1
whisker_width = 1.5
lower_whisker = Q1 - (whisker_width*IQR)
upper_whisker = Q3 + (whisker_width*IQR)
df['Age']=np.where(df['Age']>upper_whisker,upper_whisker,np.where(df['Age']<lower_whisker,lower_whisker,df['Age']))
```

In [20]:

```
sns.boxplot(df['Age'],data=df)

/usr/local/lib/python3.7/dist-packages/seaborn/_decorators.py:43: 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
FutureWarning
```

Out[20]:

```
<matplotlib.axes._subplots.AxesSubplot at 0x7fcde59c8410>
```

7) Check for Categorical columns and perform encoding

In [21]:

```
df_categorical = df[['Geography', 'Gender', 'CustomerId', 'Surname']]
df_categorical.head()
```

Out[21]:

	Geography	Gender	CustomerId	Surname
0	France	Female	15634602	Hargrave
1	Spain	Female	15647311	Hill
2	France	Female	15619304	Onio
3	France	Female	15701354	Boni
4	Spain	Female	15737888	Mitchell

In [22]:

```
from sklearn.preprocessing import LabelEncoder
```

```
Gender_encoder = LabelEncoder()
Gender_encoder.fit(df_categorical['Gender'])
```

Out[22]:

```
LabelEncoder()
```

In [23]:

```
gender_values = Gender_encoder.transform(df_categorical['Gender'])
print("Before Encoding:", list(df_categorical['Gender'][-10:]))
print("After Encoding:", gender_values[-10:])
print("The inverse from the encoding result:",
      Gender_encoder.inverse_transform(gender_values[-10:]))

Before Encoding: ['Male', 'Female', 'Male', 'Male', 'Female', 'Male', 'Male', 'Male', 'Female', 'Male']
After Encoding: [1 0 1 1 0 1 1 0 1 0]
The inverse from the encoding result: ['Male' 'Female' 'Male' 'Male' 'Female' 'Male' 'Male' 'Male' 'Female' 'Male']
```

8) Split the data into dependent and independent variables

In [24]:

```
#independent
X = df.iloc[:, :-1].values
print(X)

[[1 15634602 'Hargrave' ... 1 1 101348.88]
 [2 15647311 'Hill' ... 0 1 112542.58]
 [3 15619304 'Onio' ... 1 0 113931.57]
 ...]
```



```
[9998 15584532 'Liu' ... 0 1 42085.58]
[9999 15682355 'Sabbatini' ... 1 0 92888.52]
[10000 15628319 'Walker' ... 1 0 38190.78]]
```

In [25]:

```
#dependent
Y = df.iloc[:, -1].values
print(Y)

[1 0 1 ... 1 1 0]
9)Scale the independent variables
```

In [26]:

```
from sklearn.preprocessing import StandardScaler
from sklearn.preprocessing import MinMaxScaler
scaler = MinMaxScaler()
df[["CustomerId"]] = scaler.fit_transform(df[["CustomerId"]])
print(df)
```

	RowNumber	CustomerId	Surname	CreditScore	Geography	Gender	Age
\							
0	1	0.275616	Hargrave	619	France	Female	42.0
1	2	0.326454	Hill	608	Spain	Female	41.0
2	3	0.214421	Onio	502	France	Female	42.0
3	4	0.542636	Boni	699	France	Female	39.0
4	5	0.688778	Mitchell	850	Spain	Female	43.0
...
9995	9996	0.162119	Obijiaku	771	France	Male	39.0
9996	9997	0.016765	Johnstone	516	France	Male	35.0
9997	9998	0.075327	Liu	709	France	Female	36.0
9998	9999	0.466637	Sabbatini	772	Germany	Male	42.0
9999	10000	0.250483	Walker	792	France	Female	28.0

	Tenure	Balance	NumOfProducts	HasCrCard	IsActiveMember	\
0	2	0.00	1	1		1
1	1	83807.86	1	0		1
2	8	159660.80	3	1		0
3	1	0.00	2	0		0
4	2	125510.82	1	1		1
...
9995	5	0.00	2	1		0
9996	10	57369.61	1	1		1
9997	7	0.00	1	0		1
9998	3	75075.31	2	1		0
9999	4	130142.79	1	1		0

	EstimatedSalary	Exited
0	101348.88	1
1	112542.58	0
2	113931.57	1
3	93826.63	0
4	79084.10	0
...
9995	96270.64	0
9996	101699.77	0
9997	42085.58	1
9998	92888.52	1
9999	38190.78	0

[10000 rows x 14 columns]

10) Split the data into training and testing

In [27]:

```
from sklearn.model_selection import train_test_split
training_data, testing_data = train_test_split(df, test_size=0.2,
random_state=25)
print(f"No. of training examples: {training_data.shape[0]}")
print(f"No. of testing examples: {testing_data.shape[0]}")

No. of training examples: 8000
No. of testing examples: 2000
```