

## Importing the libraries

```
import numpy as np
import pandas as pd

import matplotlib.pyplot as plt

from google.colab import drive
drive.mount('/content/drive')
```

## Importing the Dataset

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

data.head()

Run this cell to mount your Google Drive.  
Learn more

Dismiss



				CreditScore	Geography	Gender	Age	Tenure	Balance
0	1	15634602	Hargrave	619	France	Female	42	2	0.0
1	2	15647311	Hill	608	Spain	Female	41	1	83807.0
2	3	15619304	Onio	502	France	Female	42	8	159660.0
3	4	15701354	Boni	699	France	Female	39	1	0.0

```
data.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
```

```
data.describe()
```

	RowNumber	CustomerId	CreditScore	Age	Tenure	Balance
<b>count</b>	10000.00000	1.000000e+04	10000.000000	10000.000000	10000.000000	10000.000000
<b>mean</b>	5000.50000	1.569094e+07	650.528800	38.921800	5.012800	76485.800000
<b>std</b>	2886.89568	7.193619e+04	96.653299	10.487806	2.892174	62397.400000
<b>min</b>	1.00000	1.556570e+07	350.000000	18.000000	0.000000	0.000000
<b>25%</b>	2500.75000	1.562853e+07	584.000000	32.000000	3.000000	0.000000
<b>50%</b>	5000.50000	1.569074e+07	652.000000	37.000000	5.000000	97198.500000
<b>75%</b>	7500.25000	1.575323e+07	718.000000	44.000000	7.000000	127644.200000
<b>max</b>	10000.00000	1.591500e+07	850.000000	92.000000	10.000000	250898.000000

Run this cell to mount your Google Drive.  
Learn more

Dismiss

```
data.tail()
```

	RowNumber	CustomerId	Surname	CreditScore	Geography	Gender	Age	Tenure	Balance
<b>9995</b>	9996	15606229	Obijiaku	771	France	Male	39	10	70000
<b>9996</b>	9997	15569892	Johnstone	516	France	Male	35	10	83000
<b>9997</b>	9998	15584532	Liu	709	France	Female	36	10	81000
<b>9998</b>	9999	15682355	Sabbatini	772	Germany	Male	42	10	93000
<b>9999</b>	10000	15628319	Walker	792	France	Female	28	10	79000



```
# Checking if our dataset contains any NULL values
```

```
data.isnull().sum()
```

```
RowNumber      0
CustomerId      0
Surname         0
CreditScore     0
Geography       0
Gender          0
Age             0
Tenure          0
Balance         0
```

```
NumOfProducts      0
HasCrCard           0
IsActiveMember      0
EstimatedSalary     0
Exited              0
dtype: int64
```

## Data Analysis

```
data['Gender'].value_counts()
```

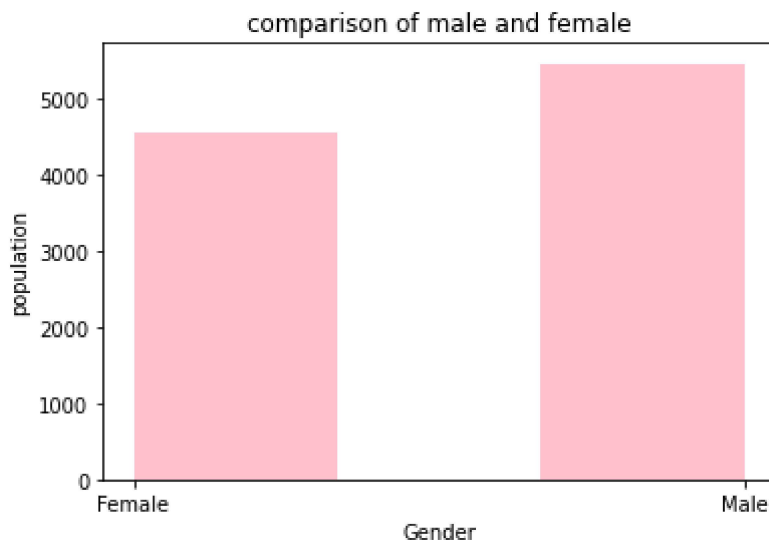
```
Male      5457
Female    4543
Name: Gender, dtype: int64
```

```
# Plotting the features of the dataset to see the correlation between them
```

```
plt.hist(data['Gender'], color = 'pink')
plt.title('Gender Distribution')
plt.xlabel('Gender')
plt.ylabel('Population')
plt.show()
```

Run this cell to mount your Google Drive.  
Learn more

Dismiss



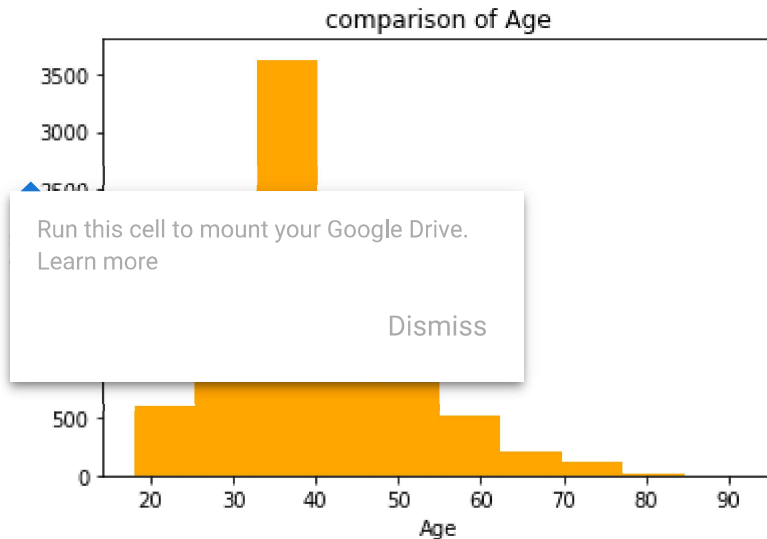
```
data['Age'].value_counts()
```

```
37    478
38    477
35    474
36    456
34    447
...
92     2
82     1
88     1
```

```
85      1
83      1
Name: Age, Length: 70, dtype: int64
```

```
# comparison of age in the dataset
```

```
plt.hist(x = data.Age, bins = 10, color = 'orange')
plt.title('comparison of Age')
plt.xlabel('Age')
plt.ylabel('population')
plt.show()
```

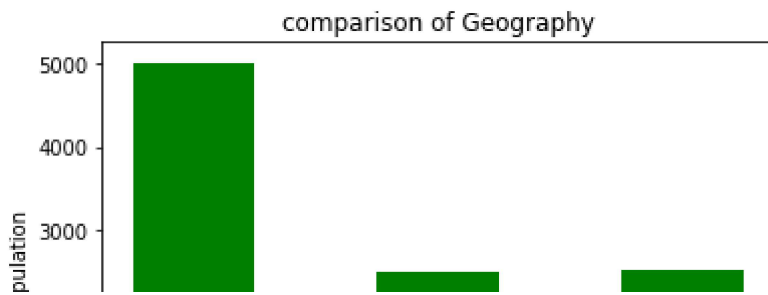


```
data['Geography'].value_counts()
```

```
France      5014
Germany     2509
Spain       2477
Name: Geography, dtype: int64
```

```
# comparison of geography
```

```
plt.hist(x = data.Geography, bins = 5, color = 'green')
plt.title('comparison of Geography')
plt.xlabel('Geography')
plt.ylabel('population')
plt.show()
```



```
data['HasCrCard'].value_counts()
```

```
1    7055
0    2945
Name: HasCrCard, dtype: int64
```

```
data['IsActiveMember'].value_counts()
```

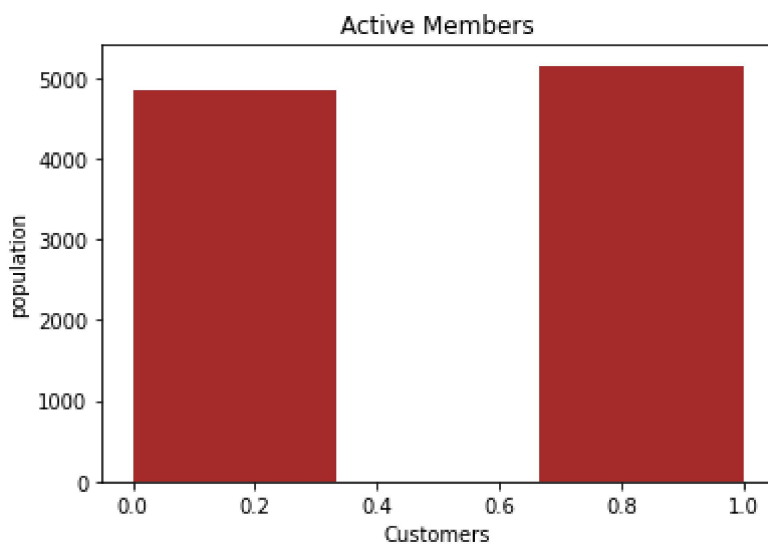
```
1    5151
```

Run this cell to mount your Google Drive.  
Learn more

Dismiss

```
# How many customers have ?
```

```
plt.hist(x = data.IsActiveMember, bins = 3, color = 'brown')
plt.title('Active Members')
plt.xlabel('Customers')
plt.ylabel('population')
plt.show()
```



```
# plotting a pie chart
```

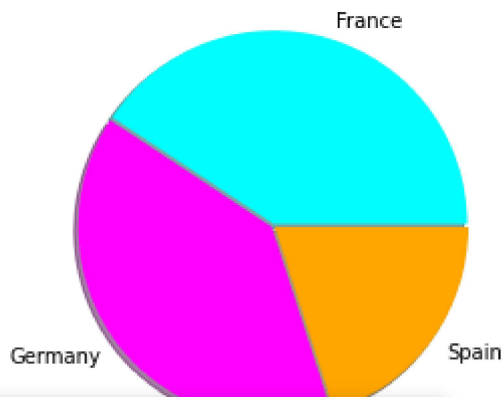
```
labels = 'France', 'Germany', 'Spain'
colors = ['cyan', 'magenta', 'orange']
sizes = [311, 300, 153]
```

```
explode = [ 0.01, 0.01, 0.01]
```

```
plt.pie(sizes, colors = colors, labels = labels, explode = explode, shadow = True)
```

```
plt.axis('equal')
```

```
plt.show()
```



Run this cell to mount your Google Drive.  
Learn more

**Data F**

Dismiss

```
# Removing the unnecessary features from the dataset
```

```
data = data.drop(['CustomerId', 'Surname', 'RowNumber'], axis = 1)
```

```
print(data.columns)
```

```
Index(['CreditScore', 'Geography', 'Gender', 'Age', 'Tenure', 'Balance',
       'NumOfProducts', 'HasCrCard', 'IsActiveMember', 'EstimatedSalary',
       'Exited'],
      dtype='object')
```

```
data.shape
```

```
# splitting the dataset into x(independent variables) and y(dependent variables)
```

```
x = data.iloc[:,0:10]
```

```
y = data.iloc[:,10]
```

```
print(x.shape)
```

```
print(y.shape)
```

```
print(x.columns)
```

```
#print(y)
```

```
(10000, 10)
```

```
(10000,)
Index(['CreditScore', 'Geography', 'Gender', 'Age', 'Tenure', 'Balance',
      'NumOfProducts', 'HasCrCard', 'IsActiveMember', 'EstimatedSalary'],
      dtype='object')
```

```
# Encoding Categorical variables into numerical variables
# One Hot Encoding
```

```
x = pd.get_dummies(x)
```

```
x.head()
```

	CreditScore	Age	Tenure	Balance	NumOfProducts	HasCrCard	IsActiveMember
0	619	42	2	0.00	1	1	1
1	608	41	1	83807.86	1	0	1
				50.80	3	1	0
				0.00	2	0	0
				10.82	1	1	1

Run this cell to mount your Google Drive.  
Learn more

Dismiss



```
x.shape
```

```
(10000, 13)
```

```
# splitting the data into training and testing set
```

```
from sklearn.model_selection import train_test_split
x_train, x_test, y_train, y_test = train_test_split(x, y, test_size = 0.25, random_state = 0)
```

```
print(x_train.shape)
print(y_train.shape)
print(x_test.shape)
print(y_test.shape)
```

```
(7500, 13)
(7500,)
(2500, 13)
(2500,)
```

[Colab paid products](#) - [Cancel contracts here](#)

Run this cell to mount your Google Drive.  
[Learn more](#)

Dismiss

Completed at 7:21 AM

