

dataset link : <https://www.kaggle.com/code/kmalit/bank-customer-churn-prediction/data>  
(<https://www.kaggle.com/code/kmalit/bank-customer-churn-prediction/data>).

## Importing Modules

In [1]:

```
1 import pandas as pd
2 import numpy as np
3 import matplotlib.pyplot as plt
4 import seaborn as sns
5 import warnings
6 warnings.filterwarnings('ignore')
```

In [2]:

```
1 df=pd.read_csv('Churn_Modelling.csv')
2 df.head()
```

Out[2]:

	RowNumber	CustomerId	Surname	CreditScore	Geography	Gender	Age	Tenure	Bala
0	1	15634602	Hargrave	619	France	Female	42	2	(
1	2	15647311	Hill	608	Spain	Female	41	1	8380i
2	3	15619304	Onio	502	France	Female	42	8	15966(
3	4	15701354	Boni	699	France	Female	39	1	(
4	5	15737888	Mitchell	850	Spain	Female	43	2	12551(

## Variables:

RowNumber — corresponds to the record (row) number and has no effect on the output. This column will be removed.

CustomerId — contains random values and has no effect on customer leaving the bank. This column will be removed.

Surname — the surname of a customer has no impact on their decision to leave the bank. This column will be removed.

CreditScore — can have an effect on customer churn, since a customer with a higher credit score is less likely to leave the bank.

Geography — a customer's location can affect their decision to leave the bank. This column will be removed.

Gender — it's interesting to explore whether gender plays a role in a customer leaving the bank. We'll include this column, too.

Age — this is certainly relevant, since older customers are less likely to leave their bank than younger ones.

Tenure — refers to the number of years that the customer has been a client of the bank. Normally, older clients are more loyal and less likely to leave a bank.

Balance — also a very good indicator of customer churn, as people with a higher balance in their accounts are less likely to leave the bank compared to those with lower balances.

NumOfProducts — refers to the number of products that a customer has purchased through the bank.

HasCrCard — denotes whether or not a customer has a credit card. This column is also relevant, since people with a credit card are less likely to leave the bank. (0=No,1=Yes)

IsActiveMember — active customers are less likely to leave the bank, so we'll keep this. (0=No,1=Yes)

EstimatedSalary — as with balance, people with lower salaries are more likely to leave the bank compared to those with higher salaries.

Exited — whether or not the customer left the bank. This is what we have to predict. (0=No,1=Yes)

Data cleaning

In [3]:

```
1 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
```

It consists of 10000 observations and 12 variables. Independent variables contain information about customers. Dependent variable refers to customer abandonment status.

In [4]:

```
1 df.isnull().sum()
```

Out[4]:

```
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
```

No Null values in Given Dataset.

In [5]:

```
1 df.duplicated().sum()
```

Out[5]:

```
0
```

No duplicate values available

In [6]:

```
1 df.describe()
```

Out[6]:

	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	76485.889288
<b>std</b>	2886.89568	7.193619e+04	96.653299	10.487806	2.892174	62397.405202
<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	97198.540000
<b>75%</b>	7500.25000	1.575323e+07	718.000000	44.000000	7.000000	127644.240000
<b>max</b>	10000.00000	1.581569e+07	850.000000	92.000000	10.000000	250898.090000

We can drop the features which we do not need in model building

In [7]:

```
1 df.drop(['Surname', 'Geography', 'Gender', 'CustomerId', 'RowNumber'], axis=1, inplace=True)
2 df.head()
```

Out[7]:

	CreditScore	Age	Tenure	Balance	NumOfProducts	HasCrCard	IsActiveMember	EstimatedSalary
0	619	42	2	0.00	1	1	1	151603.28
1	608	41	1	83807.86	1	0	1	134216.42
2	502	42	8	159660.80	3	1	0	151603.28
3	699	39	1	0.00	2	0	0	134216.42
4	850	43	2	125510.82	1	1	1	151603.28

## Visualization on Target Column

In [8]:

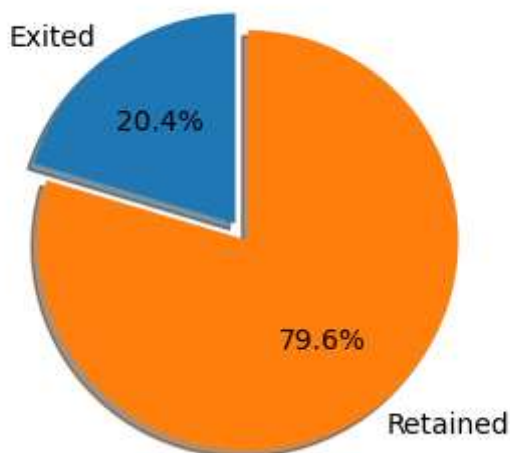
```
1 print(df.Exited.value_counts())
2 labels = 'Exited', 'Retained'
3 sizes = [df.Exited[df['Exited']==1].count(), df.Exited[df['Exited']==0].count()]
4 explode = (0, 0.1)
5 fig1, ax1 = plt.subplots(figsize=(3,3))
6 ax1.pie(sizes, explode=explode, labels=labels, autopct='%1.1f%%',
7         shadow=True, startangle=90)
8 ax1.axis('equal')
9 plt.title("Proportion of customer churned and retained", size = 10)
10 plt.show()
```

0 7963

1 2037

Name: Exited, dtype: int64

Proportion of customer churned and retained



From above observation, approximate 80% (7063 nos.) customers are continue with the bank as a member

## feature scaling

In [9]:

```
1 from sklearn.preprocessing import StandardScaler
2 sc = StandardScaler()
3 df[['CreditScore', 'Age', 'Balance', 'EstimatedSalary']] = sc.fit_transform(df[['CreditScore', 'Age', 'Balance', 'EstimatedSalary']])
4 df.head()
```

Out[9]:

	CreditScore	Age	Tenure	Balance	NumOfProducts	HasCrCard	IsActiveMember	Es
0	-0.326221	0.293517	2	-1.225848	1	1	1	
1	-0.440036	0.198164	1	0.117350	1	0	1	
2	-1.536794	0.293517	8	1.333053	3	1	0	
3	0.501521	0.007457	1	-1.225848	2	0	0	
4	2.063884	0.388871	2	0.785728	1	1	1	

making data balanced by Oversampling and Splitting dataset into x & y

In [10]:

```
1 from imblearn.over_sampling import SMOTE
2 sm = SMOTE(random_state=42)
3 x = df.drop('Exited', axis=1)
4 y = df['Exited']
5 x, y = sm.fit_resample(x, y)
6 x.shape, y.shape
```

Out[10]:

((15926, 8), (15926,))

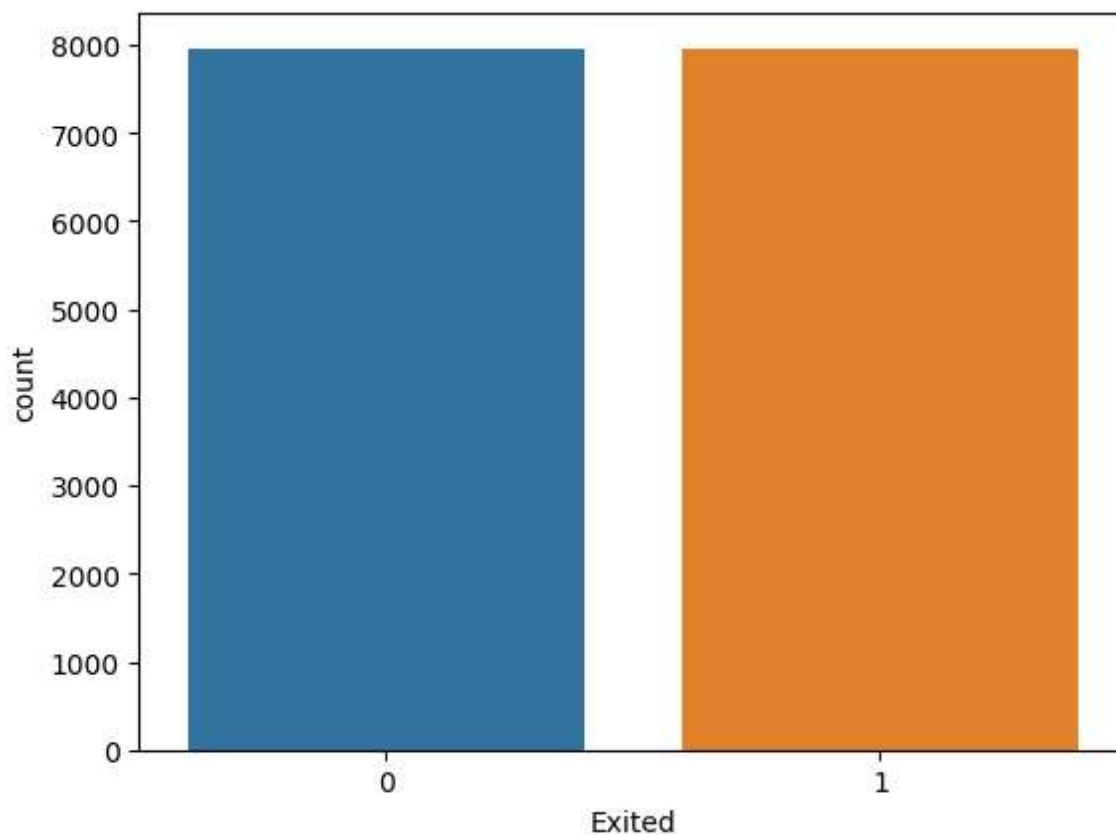
We use Scaling for standardization

In [11]:

```
1 sns.countplot(x=y)
```

Out[11]:

<AxesSubplot: xlabel='Exited', ylabel='count'>



Now my Data is balanced as shown above graph

In [12]:

```
1 x.head()
```

Out[12]:

	CreditScore	Age	Tenure	Balance	NumOfProducts	HasCrCard	IsActiveMember	Es
0	-0.326221	0.293517	2	-1.225848	1	1	1	
1	-0.440036	0.198164	1	0.117350	1	0	1	
2	-1.536794	0.293517	8	1.333053	3	1	0	
3	0.501521	0.007457	1	-1.225848	2	0	0	
4	2.063884	0.388871	2	0.785728	1	1	1	

**Splitting x & y into training and testing**

In [13]:

```
1 from sklearn.model_selection import train_test_split
2 xtrain,xtest,ytrain,ytest =train_test_split(x,y,test_size=0.3,random_state=1)
```

In [14]:

```
1 from sklearn.linear_model import LogisticRegression
2 from sklearn.neighbors import KNeighborsClassifier
3 from sklearn.svm import SVC
4 from sklearn.tree import DecisionTreeClassifier
5 from sklearn import tree
6 from sklearn.metrics import classification_report,accuracy_score
```

In [15]:

```
1 logreg=LogisticRegression()
2 knn=KNeighborsClassifier()
3 svm=SVC()
4 dt=DecisionTreeClassifier()
```

In [16]:

```
1 def mymodel(model):
2     model.fit(xtrain,ytrain) # build model
3     ypred = model.predict(xtest) #predicted value of y
4     train = model.score(xtrain,ytrain)
5     test = model.score(xtest,ytest)
6     print(f'training accuracy {train}')
7     print(f'testing accuracy {test}')
8     print(f'Model Name : {model}')
9     print(classification_report(ytest,ypred))
```

In [17]:

```
1 mymodel(logreg)
```

training accuracy 0.7030857552924291

testing accuracy 0.691084135621599

Model Name : LogisticRegression()

	precision	recall	f1-score	support
0	0.68	0.70	0.69	2322
1	0.70	0.69	0.70	2456
accuracy			0.69	4778
macro avg	0.69	0.69	0.69	4778
weighted avg	0.69	0.69	0.69	4778

In [18]:

```
1 mymodel(svm)
```

training accuracy 0.7704520990312164

testing accuracy 0.756802009208874

Model Name : SVC()

	precision	recall	f1-score	support
0	0.73	0.79	0.76	2322
1	0.78	0.73	0.76	2456
accuracy			0.76	4778
macro avg	0.76	0.76	0.76	4778
weighted avg	0.76	0.76	0.76	4778

In [19]:

```
1 mymodel(knn)
```

training accuracy 0.8878722640832436

testing accuracy 0.8359146086228547

Model Name : KNeighborsClassifier()

	precision	recall	f1-score	support
0	0.89	0.76	0.82	2322
1	0.80	0.91	0.85	2456
accuracy			0.84	4778
macro avg	0.84	0.83	0.83	4778
weighted avg	0.84	0.84	0.83	4778

In [20]:

```
1 mymodel(dt)
```

training accuracy 1.0

testing accuracy 0.8210548346588531

Model Name : DecisionTreeClassifier()

	precision	recall	f1-score	support
0	0.82	0.81	0.81	2322
1	0.82	0.83	0.83	2456
accuracy			0.82	4778
macro avg	0.82	0.82	0.82	4778
weighted avg	0.82	0.82	0.82	4778

## Bagging



In [21]:

```
1 from sklearn.ensemble import BaggingClassifier
2 bg = BaggingClassifier(dt)
3 bg.fit(xtrain,ytrain)
4 ypred = bg.predict(xtest)
5 cr = classification_report(ytest,ypred)
6 train = bg.score(xtrain,ytrain)
7 test = bg.score(xtest,ytest)
8 print(f'training accuracy {train}')
9 print(f'testing accuracy {test}')
10 print(dt)
11 print(cr)
```

training accuracy 0.990043057050592

testing accuracy 0.8539137714524906

DecisionTreeClassifier()

	precision	recall	f1-score	support
0	0.83	0.89	0.86	2322
1	0.89	0.82	0.85	2456
accuracy			0.85	4778
macro avg	0.86	0.85	0.85	4778
weighted avg	0.86	0.85	0.85	4778

## Random Forest

In [22]:

```
1 from sklearn.ensemble import RandomForestClassifier
2 rf = RandomForestClassifier()
3 rf.fit(xtrain,ytrain)
4 ypred = rf.predict(xtest)
5 cr = classification_report(ytest,ypred)
6 train = rf.score(xtrain,ytrain)
7 test = rf.score(xtest,ytest)
8 print(f'training accuracy {train}')
9 print(f'testing accuracy {test}')
10 print(cr)
```

training accuracy 1.0

testing accuracy 0.8737965676015069

	precision	recall	f1-score	support
0	0.87	0.88	0.87	2322
1	0.88	0.87	0.88	2456
accuracy			0.87	4778
macro avg	0.87	0.87	0.87	4778
weighted avg	0.87	0.87	0.87	4778

## Voting Classifier

In [23]:

```
1 models=[]
2 models.append(('logistic regression',logreg))
3 models.append(('KNN',knn))
4 models.append(('Support vector machine',svm))
5 models.append(('Decision Tree',dt))
```

In [24]:

```
1 from sklearn.ensemble import VotingClassifier
2 vc = VotingClassifier(estimators=models)
3 vc.fit(xtrain,ytrain)
4 ypred = vc.predict(xtest)
5 cr = classification_report(ytest,ypred)
6 train = vc.score(xtrain,ytrain)
7 test = vc.score(xtest,ytest)
8 print(f'training accuracy {train}')
9 print(f'testing accuracy {test}')
10 print(cr)
```

training accuracy 0.8486724076067456

testing accuracy 0.7938467978233571

		precision	recall	f1-score	support
	0	0.75	0.85	0.80	2322
	1	0.84	0.74	0.79	2456
	accuracy			0.79	4778
	macro avg	0.80	0.80	0.79	4778
	weighted avg	0.80	0.79	0.79	4778

## Boosting

### a. Adaptive Boosting

In [25]:

```
1 from sklearn.ensemble import AdaBoostClassifier
2 ad = AdaBoostClassifier()
3 ad.fit(xtrain,ytrain)
4 ypred = ad.predict(xtest)
5 cr = classification_report(ytest,ypred)
6 train = ad.score(xtrain,ytrain)
7 test = ad.score(xtest,ytest)
8 print(f'training accuracy {train}')
9 print(f'testing accuracy {test}')
10 print(cr)
```

training accuracy 0.789199856476498

testing accuracy 0.7787777312683131

	precision	recall	f1-score	support
0	0.76	0.79	0.78	2322
1	0.79	0.77	0.78	2456
accuracy			0.78	4778
macro avg	0.78	0.78	0.78	4778
weighted avg	0.78	0.78	0.78	4778

## b. Gradient Boosting

In [26]:

```
1 from sklearn.ensemble import GradientBoostingClassifier
2 gb = GradientBoostingClassifier()
3 gb.fit(xtrain,ytrain)
4 ypred = gb.predict(xtest)
5 cr = classification_report(ytest,ypred)
6 train = gb.score(xtrain,ytrain)
7 test = gb.score(xtest,ytest)
8 print(f'training accuracy {train}')
9 print(f'testing accuracy {test}')
10 print(cr)
```

training accuracy 0.8273232866881952

testing accuracy 0.808915864378401

	precision	recall	f1-score	support
0	0.79	0.83	0.81	2322
1	0.83	0.79	0.81	2456
accuracy			0.81	4778
macro avg	0.81	0.81	0.81	4778
weighted avg	0.81	0.81	0.81	4778

## c. Extreme Gradient Boosting

In [27]:

```
1 from xgboost import XGBClassifier
2 xg = XGBClassifier()
3 xg.fit(xtrain,ytrain)
4 ypred = xg.predict(xtest)
5 cr = classification_report(ytest,ypred)
6 train = xg.score(xtrain,ytrain)
7 test = xg.score(xtest,ytest)
8 print(f'training accuracy {train}')
9 print(f'testing accuracy {test}')
10 print(cr)
```

training accuracy 0.9540724793684966

testing accuracy 0.8890749267475931

	precision	recall	f1-score	support
0	0.86	0.92	0.89	2322
1	0.92	0.86	0.89	2456
accuracy			0.89	4778
macro avg	0.89	0.89	0.89	4778
weighted avg	0.89	0.89	0.89	4778

My model work best on Random Forest i.e. accuracy 84% as compared to others. So we'll build our model on Random Forest.

## Hyper Tunning using GridSearchCV

### HyperTunning on Extreme Gradient Boosting

In [28]:

```
1 from xgboost import XGBClassifier
2 from sklearn.model_selection import GridSearchCV
3 estimator = XGBClassifier(objective= 'binary:logistic',nthread=4,seed=42)
4 parameters = {'max_depth': range (2, 10, 1),'n_estimators': range(60, 220, 40),'lear
```

In [29]:

```
1 gs = GridSearchCV(estimator=estimator,param_grid=parameters,scoring = 'roc_auc',n_jo
```

In [30]:

```
1 gs.fit(xtrain,ytrain)
2 ypred = gs.predict(xtest)
3 cr = classification_report(ytest,ypred)
4 print(cr)
```

Fitting 10 folds for each of 96 candidates, totalling 960 fits

	precision	recall	f1-score	support
--	-----------	--------	----------	---------

	0	0.86	0.91	0.89	2322
	1	0.91	0.86	0.89	2456
accuracy				0.89	4778
macro avg	0.89	0.89	0.89		4778
weighted avg	0.89	0.89	0.89		4778

In [31]:

```
1 gs.best_params_
```

Out[31]:

```
{'learning_rate': 0.1, 'max_depth': 9, 'n_estimators': 180}
```

## Rebuild using xgboost algorithm with best parameters

In [32]:

```
1 from xgboost import XGBClassifier
2 xg = XGBClassifier(learning_rate= 0.1, max_depth= 9, n_estimators= 180)
3 xg.fit(xtrain,ytrain)
4 ypred = xg.predict(xtest)
5 cr = classification_report(ytest,ypred)
6 train = xg.score(xtrain,ytrain)
7 test = xg.score(xtest,ytest)
8 print(f'training accuracy {train}')
9 print(f'testing accuracy {test}')
10 print(cr)
```

training accuracy 0.9777538571941156

testing accuracy 0.8874005860192549

	precision	recall	f1-score	support	
	0	0.86	0.91	0.89	2322
	1	0.91	0.86	0.89	2456
accuracy				0.89	4778
macro avg	0.89	0.89	0.89		4778
weighted avg	0.89	0.89	0.89		4778

So we build model on xgboosting having highest accuracy 89%.

In [35]:

```
1 import pickle
2 pickle.dump(xg , open("predict.pkl" , "wb"))
```

In [36]:

```
1 import pickle
2 pickle.dump(sc , open("stdscl.pkl" , "wb"))
```

## Reporting

The aim of this study was to create classification models for the churn dataset and to predict whether a person abandons the bank by creating models and to obtain maximum accuracy score in the established models. The work done is as follows:

Churn Data Set read.

With Exploratory Data Analysis; The data set's structural data were checked. The types of variables in the dataset were examined. Size information of the dataset was accessed. Descriptive statistics of the data set were examined. It was concluded that there were no missing observations and outliers in the data set.

During Model Building; Logistic Regression, KNN, SVM, DT, Bagging Classifier, Voting Classifier, Random Forests, AdaBoost, GradientBoost, XGBoost like using machine learning models Accuracy Score were calculated. Later Random Forest hyperparameter optimizations (by using GridSearchCV) optimized to increase Accuracy score.

Result; The model created as a result of xgboosting became the model with the maximum Accuracy Score. (0.89)