

GOAL: A bank is investigating a very high rate of customer leaving the bank. Our goal is to investigate and predict which of the customers are more likely to leave the bank soon based on the given dataset.

```
In [116]: import numpy as np
import pandas as pd
import matplotlib.pyplot as plt
import seaborn as sns
import tensorflow as tf
from scipy.stats import skew
from sklearn.model_selection import train_test_split
from sklearn.preprocessing import StandardScaler
from sklearn.metrics import classification_report
from tensorflow.keras.models import Sequential
from tensorflow.keras.layers import Dense
from sklearn.metrics import confusion_matrix
import warnings
warnings.filterwarnings('ignore')
```

```
In [17]: df = pd.read_csv("/content/drive/MyDrive/Colab Notebooks/Churn_Modelling.csv")
```

EDA and PREPROCESSING :

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

```
Out[18]:
```

	RowNumber	CustomerId	Surname	CreditScore	Geography	Gender	Age	Tenure	Balance	I
0	1	15634602	Hargrave	619	France	Female	42	2	0.00	
1	2	15647311	Hill	608	Spain	Female	41	1	83807.86	
2	3	15619304	Onio	502	France	Female	42	8	159660.80	
3	4	15701354	Boni	699	France	Female	39	1	0.00	
4	5	15737888	Mitchell	850	Spain	Female	43	2	125510.82	

```
In [19]: df.drop(["RowNumber"], axis=1, inplace=True)
```

```
In [20]: df.drop(["CustomerId"], axis=1, inplace=True)
```

Irrelevant columns RowNumber and CustomerID are dropped.

```
In [21]: df.shape
```

```
Out[21]: (10000, 12)
```

```
In [22]: df.info()
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 10000 entries, 0 to 9999
Data columns (total 12 columns):
 #   Column           Non-Null Count  Dtype  
 --- 
 0   Surname          10000 non-null   object  
 1   CreditScore      10000 non-null   int64  
 2   Geography         10000 non-null   object  
 3   Gender            10000 non-null   object  
 4   Age               10000 non-null   int64  
 5   Tenure            10000 non-null   int64  
 6   Balance           10000 non-null   float64
 7   IsFraud          10000 non-null   int64  
 8   ProductCategory  10000 non-null   object  
 9   Month             10000 non-null   int64  
 10  Day              10000 non-null   int64  
 11  WeekofYear       10000 non-null   int64
```

```

5   Tenure          10000 non-null  int64
6   Balance         10000 non-null  float64
7   NumOfProducts   10000 non-null  int64
8   HasCrCard       10000 non-null  int64
9   IsActiveMember  10000 non-null  int64
10  EstimatedSalary 10000 non-null  float64
11  Exited          10000 non-null  int64
dtypes: float64(2), int64(7), object(3)
memory usage: 937.6+ KB

```

In [23]: `df.isnull().sum()`

```

Out[23]: 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

```

In [24]: *#Finding the special characters in the data frame*
`df.isin(['?']).sum(axis=0)`

```

Out[24]: 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

```

We can see that there are no null values in our dataset.

In [25]: *#Separating numeric and categoric columns:*
`df_num = df.select_dtypes(exclude=['object'])`
`df_cat = df.select_dtypes(include=['object'])`

In [26]: `df_num.head()`

```

Out[26]:   CreditScore  Age  Tenure  Balance  NumOfProducts  HasCrCard  IsActiveMember  EstimatedSal
          0            619    42       2      0.00             1           1              1            101348
          1            608    41       1    83807.86             1           0              1            112542
          2            502    42       8   159660.80             3           1              0            113931
          3            699    39       1      0.00             2           0              0            93826
          4            850    43       2   125510.82             1           1              1            79084

```

In [27]: `df_cat.head()`

Out[27]:

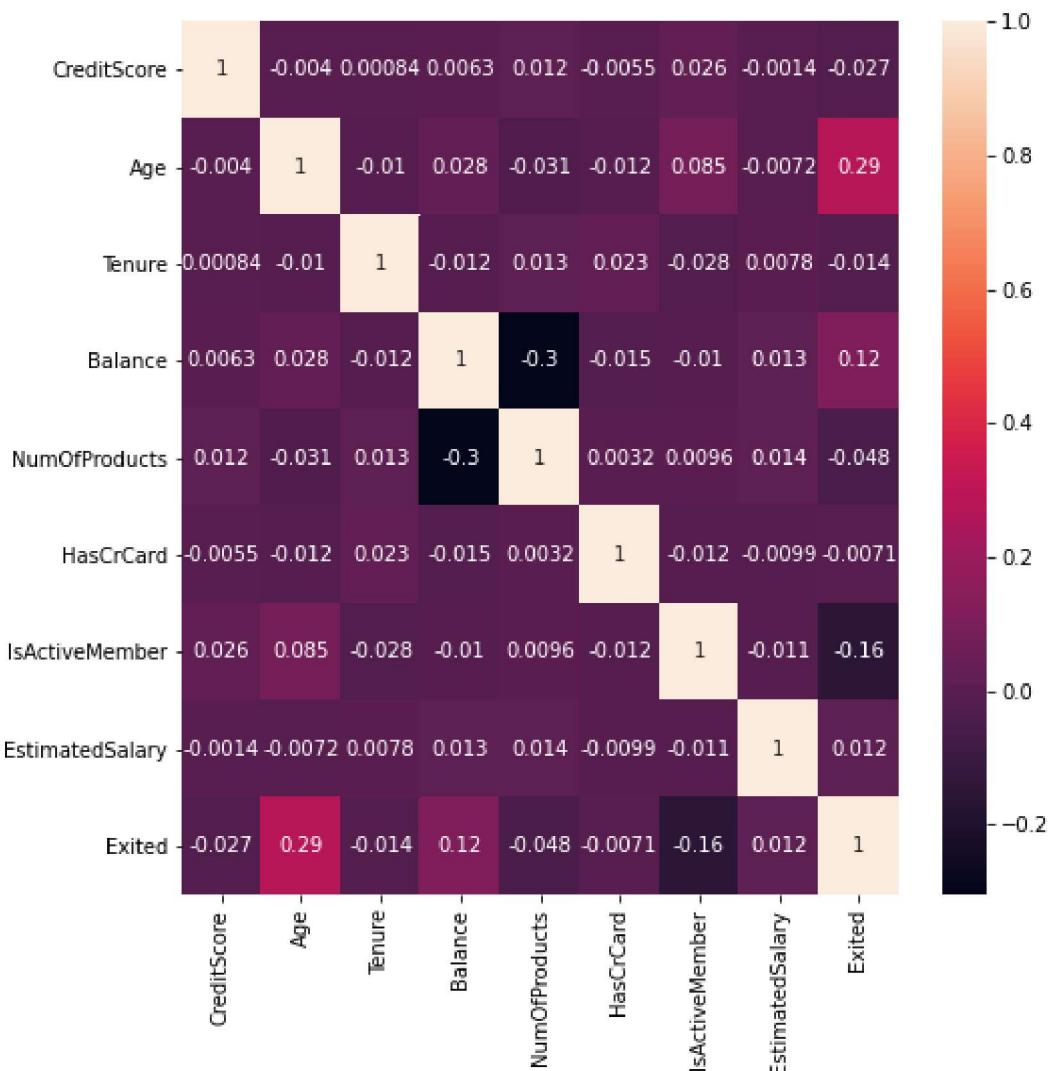
Surname Geography Gender

0	Hargrave	France	Female
1	Hill	Spain	Female
2	Onio	France	Female
3	Boni	France	Female
4	Mitchell	Spain	Female

CHECKING MULTICOLLINEARITY :

In [28]:

```
plt.figure(figsize=(8,8))
sns.heatmap(df.corr(), annot=True)
plt.show()
```

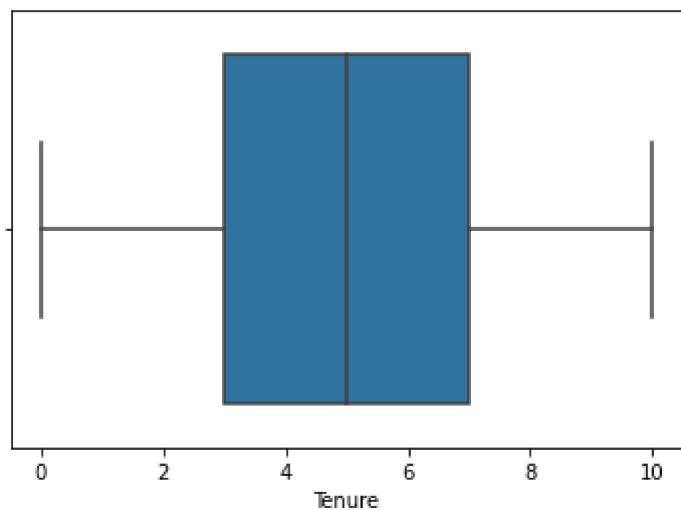
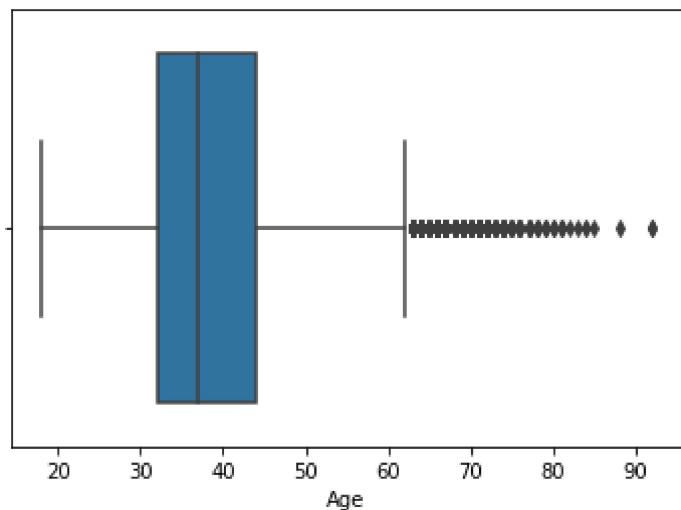
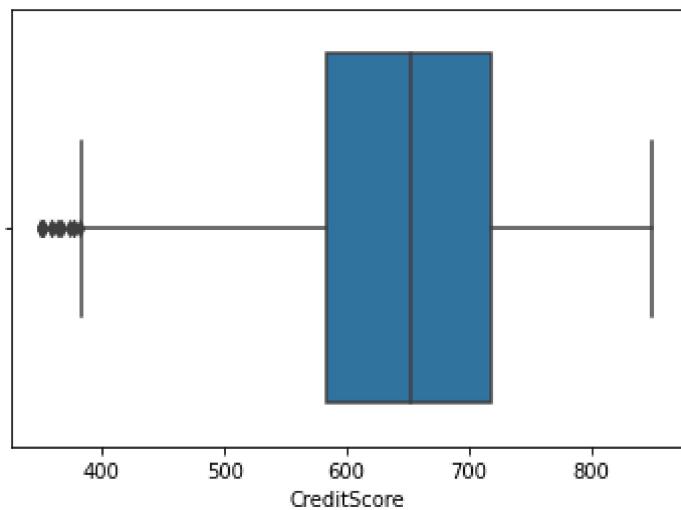


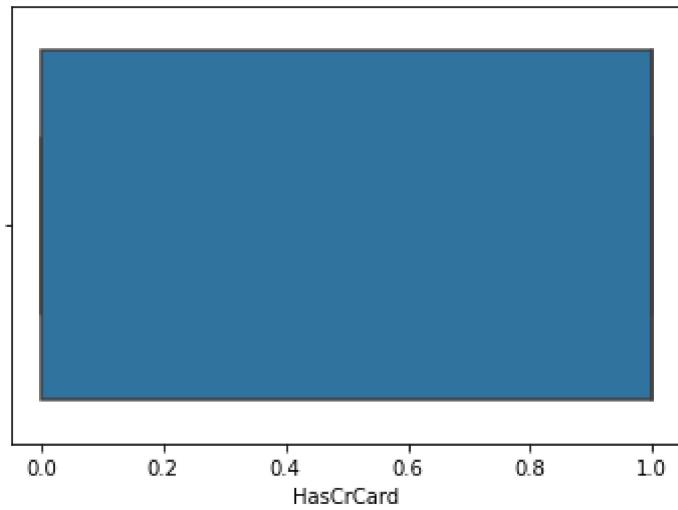
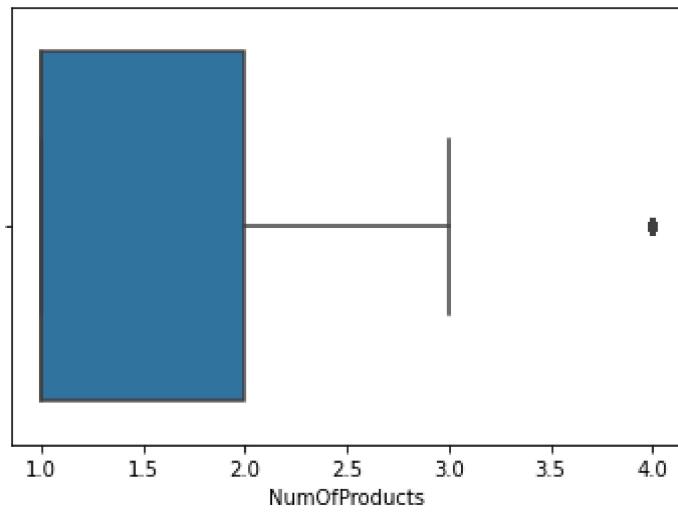
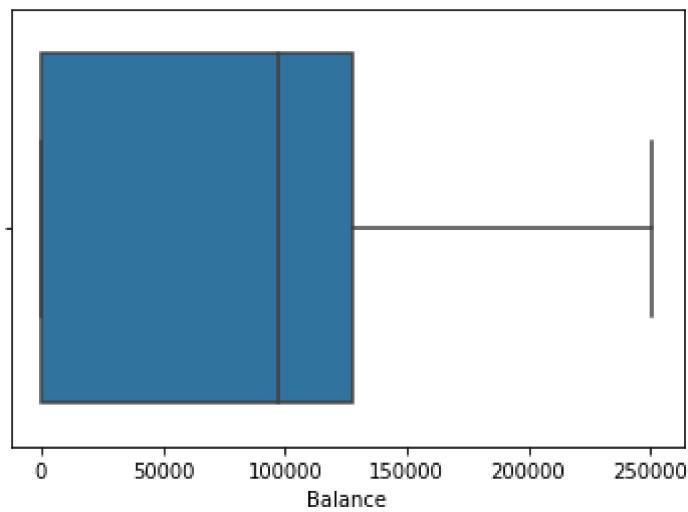
We can see that there is no significant multicollinearity between the coulmns.

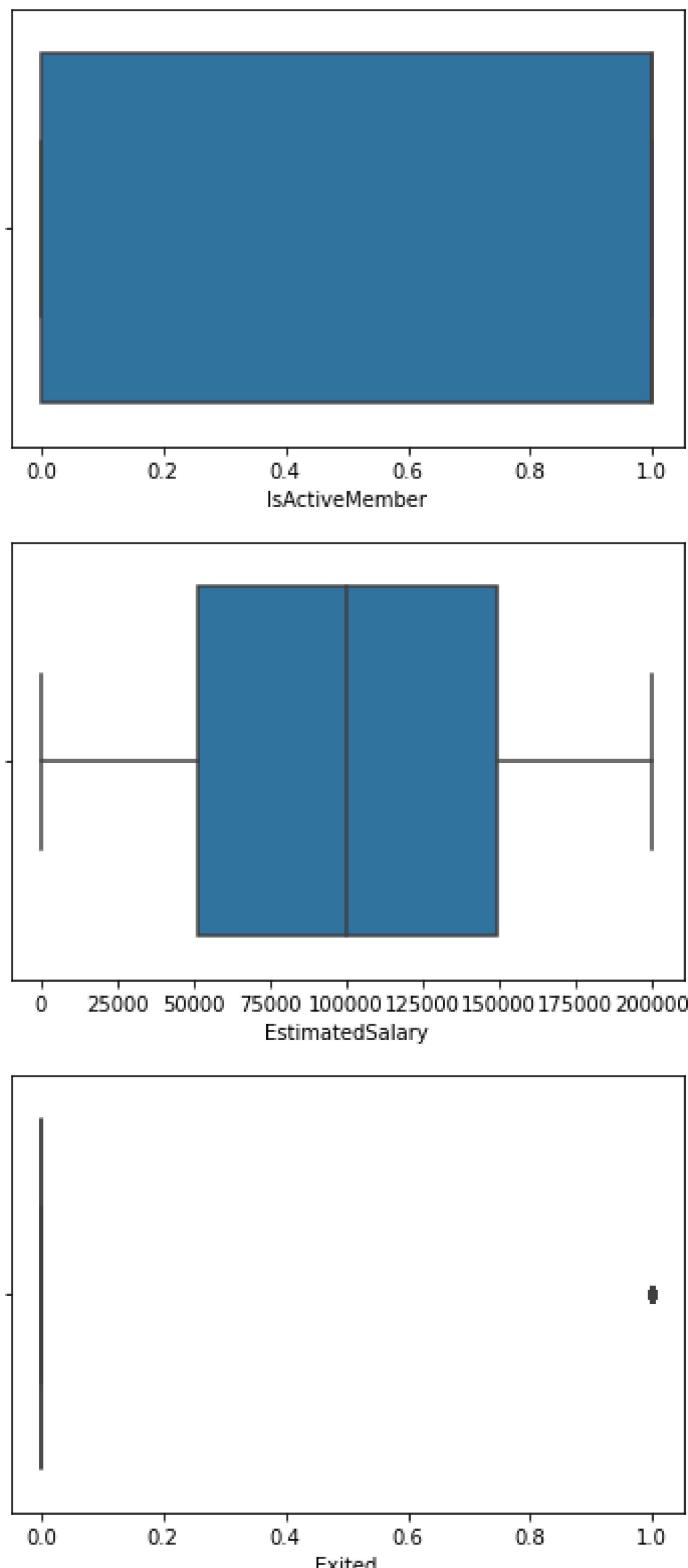
UNIVARIATE ANALYSIS

In [29]:

```
for col in df_num:
    plt.figure()
    sns.boxplot(df_num[col])
    plt.show()
```



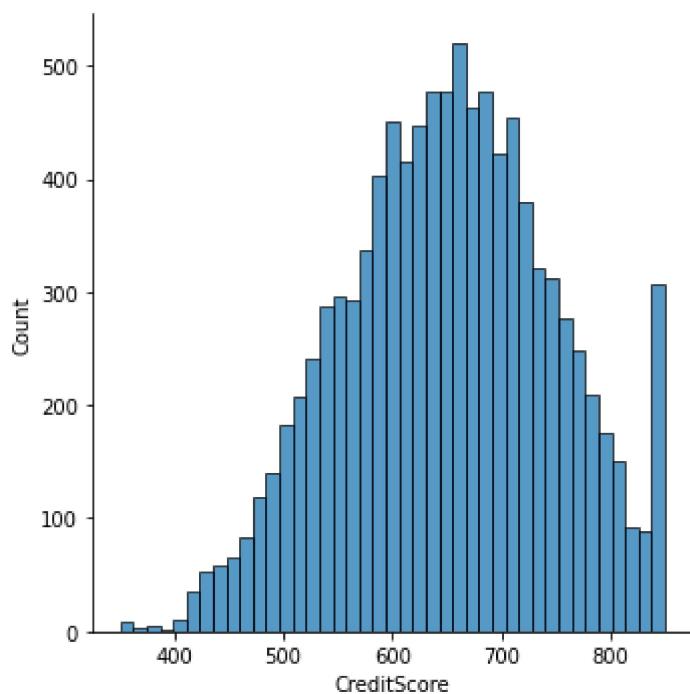




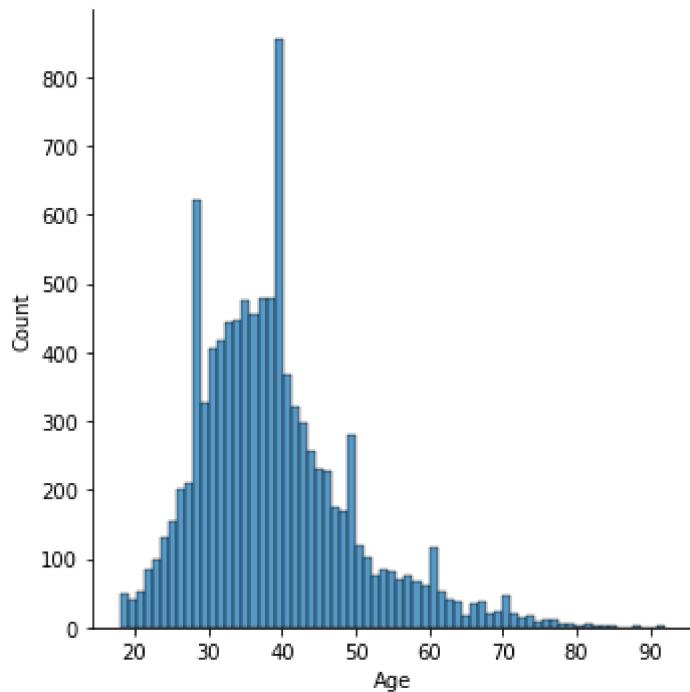
There are outliers in CreditScore, Age and NumofProduct colums.

```
In [30]: for col in df_num:  
    print(col, "-", skew(df_num[col]))  
    plt.figure()  
    sns.displot(df_num[col])  
    plt.show()
```

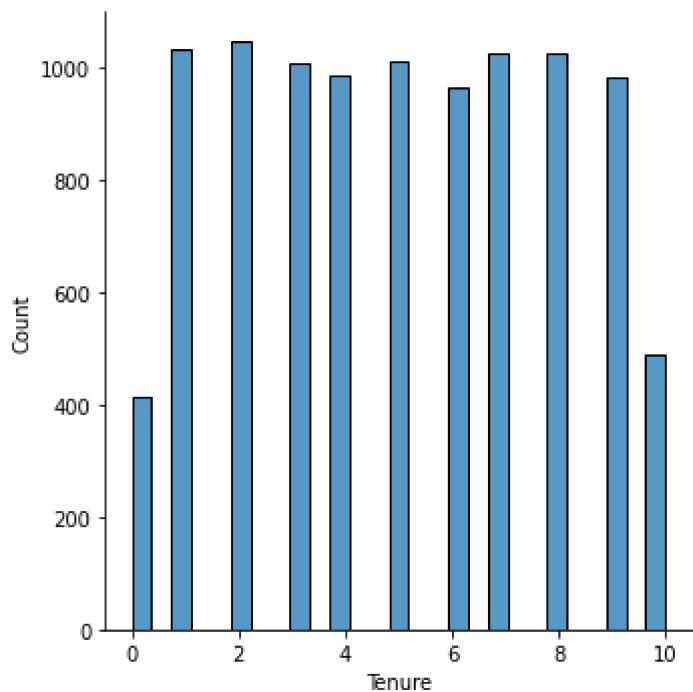
```
CreditScore - -0.07159586676212397  
<Figure size 432x288 with 0 Axes>
```



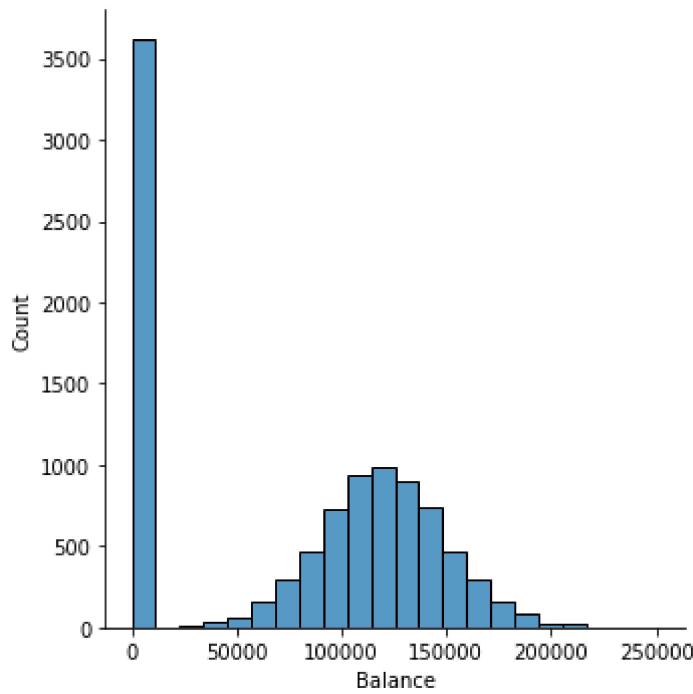
Age - 1.0111685586628079
<Figure size 432x288 with 0 Axes>



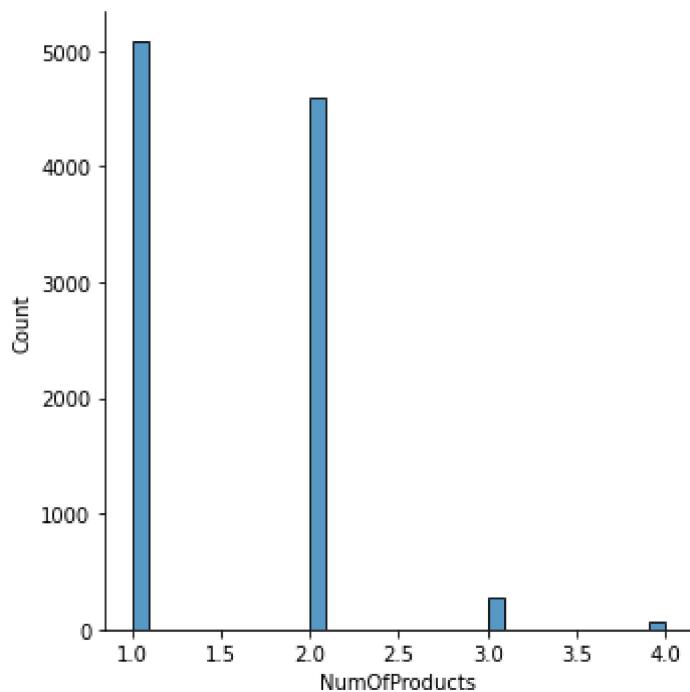
Tenure - 0.010989809189781041
<Figure size 432x288 with 0 Axes>



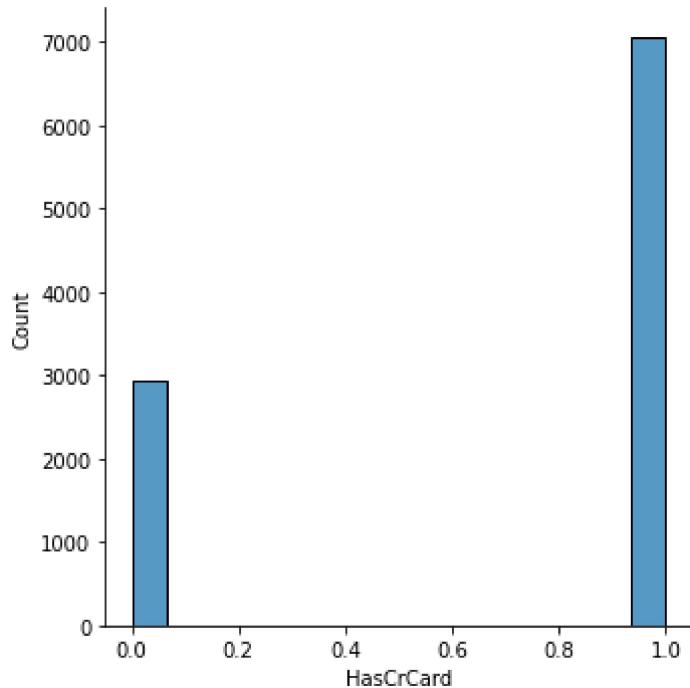
Balance - -0.14108754375291138
<Figure size 432x288 with 0 Axes>



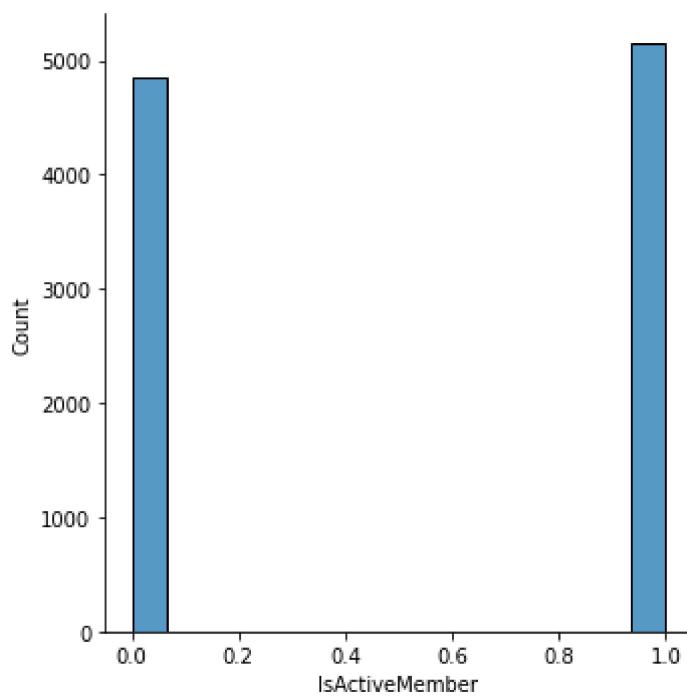
NumOfProducts - 0.745456048438949
<Figure size 432x288 with 0 Axes>



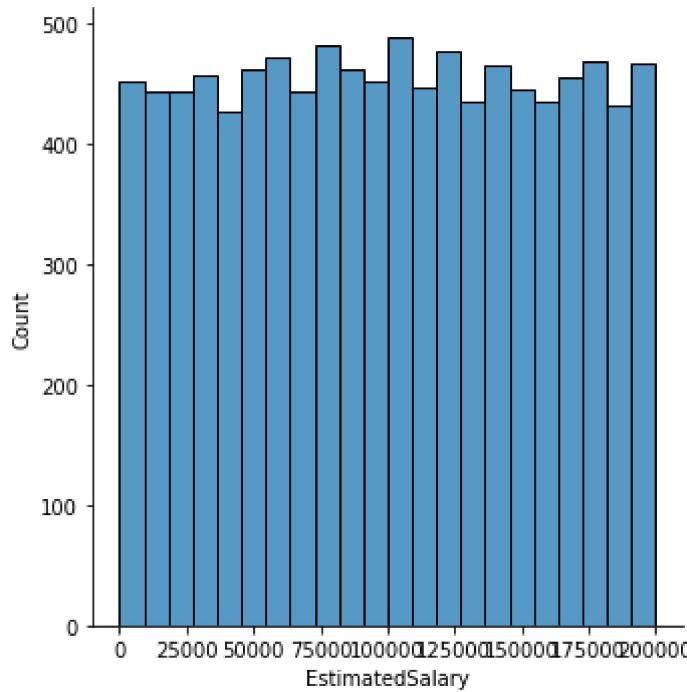
HasCrCard - -0.9016763178640548
<Figure size 432x288 with 0 Axes>



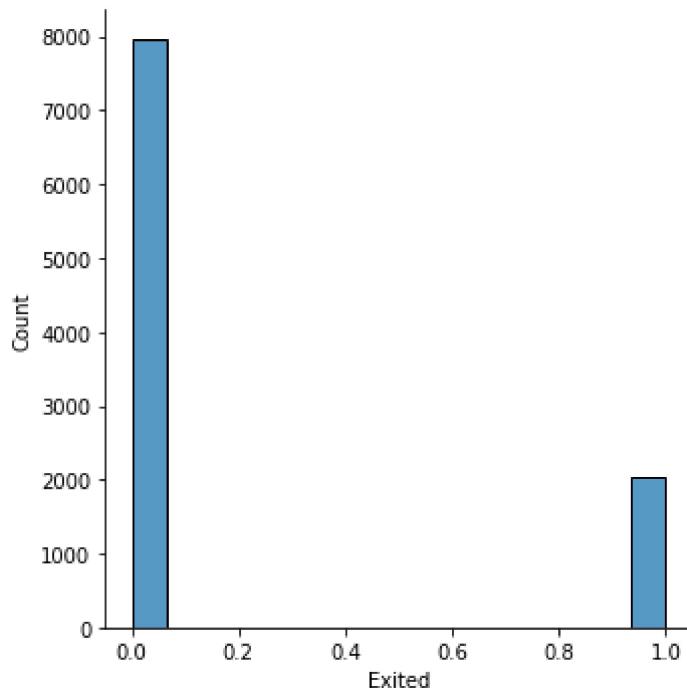
IsActiveMember - -0.06042756246298516
<Figure size 432x288 with 0 Axes>



EstimatedSalary - 0.0020850448448748848
<Figure size 432x288 with 0 Axes>

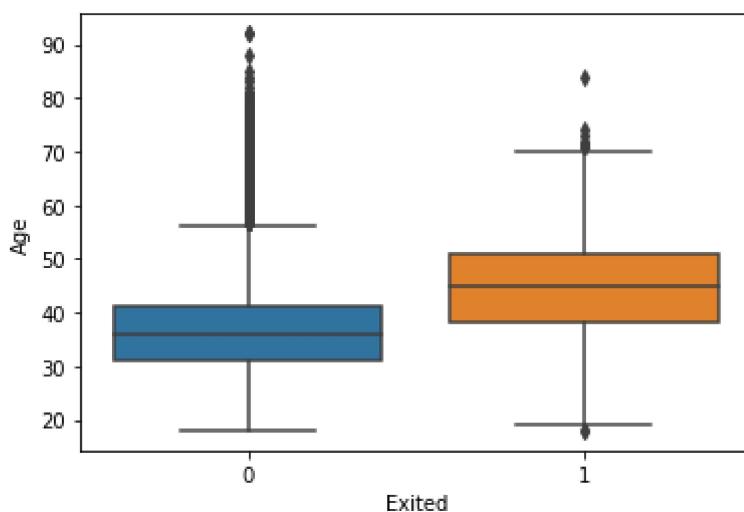
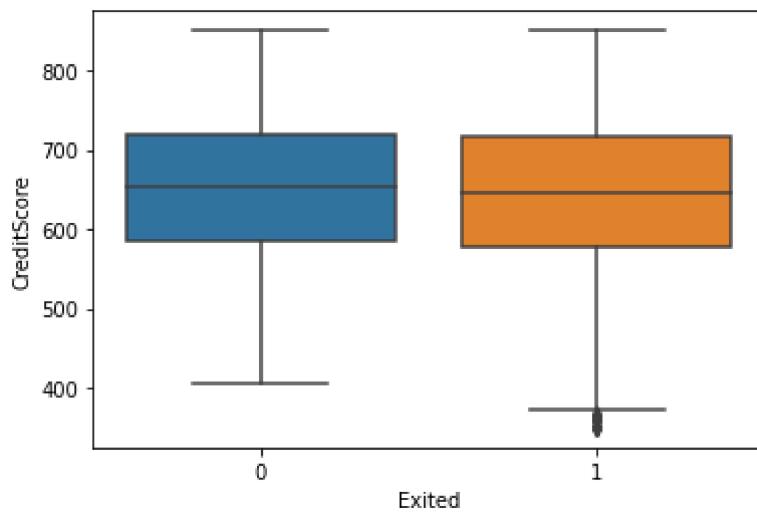


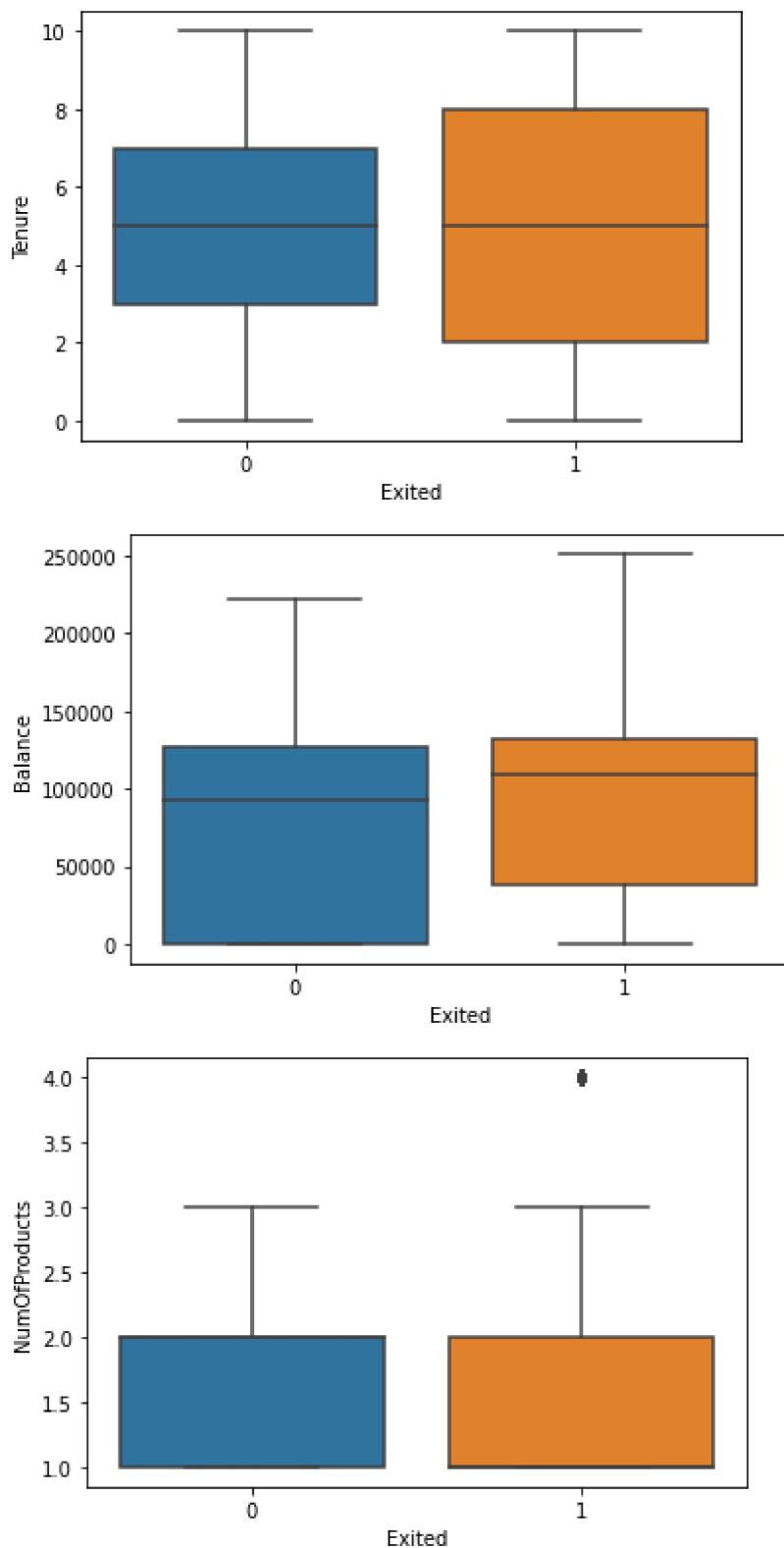
Exited - 1.4713899141398699
<Figure size 432x288 with 0 Axes>

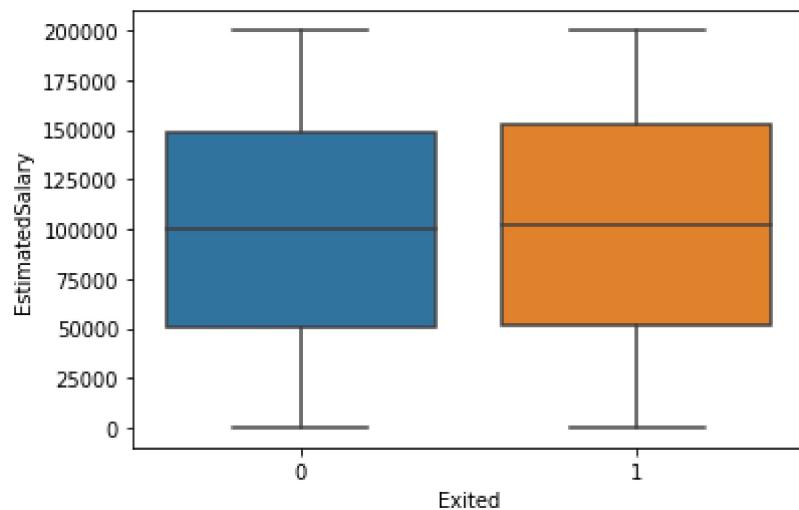
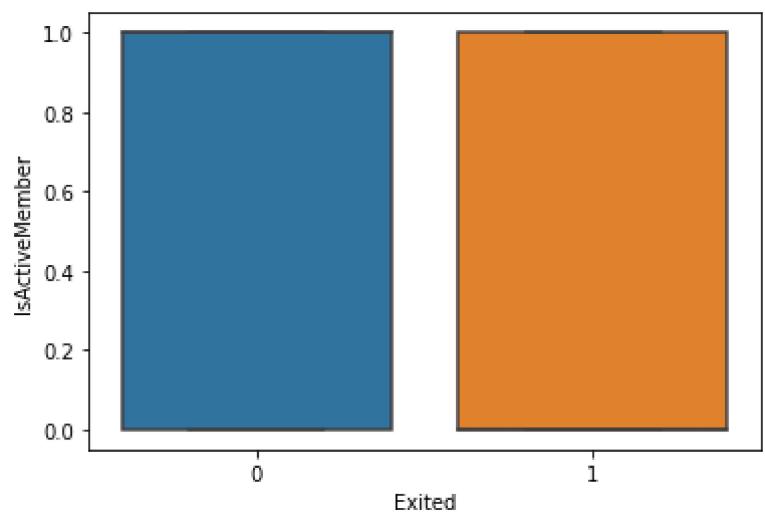
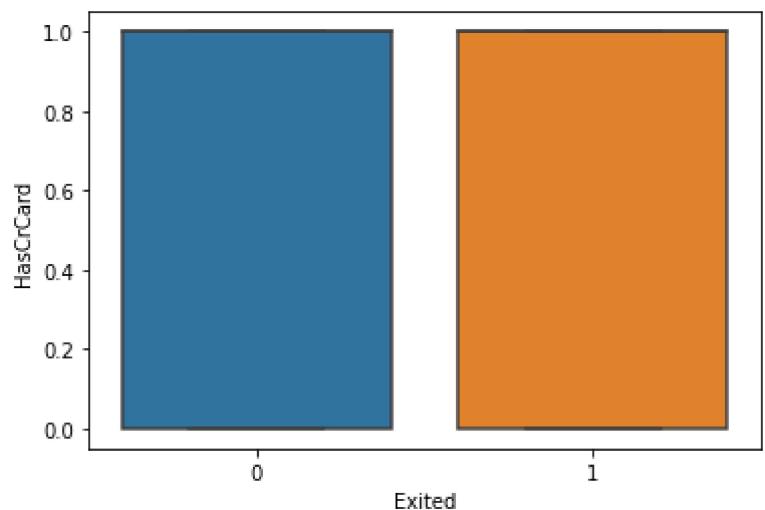


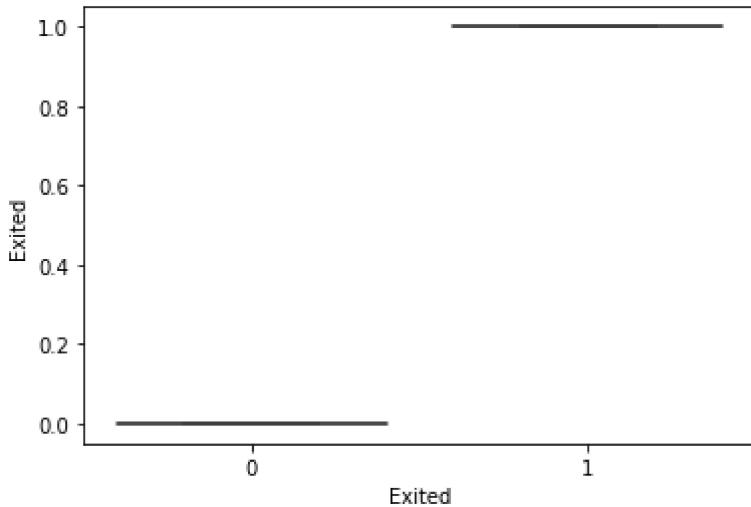
BIVARIATE ANALYSIS :

```
In [31]: for col in df_num:  
    plt.figure()  
    sns.boxplot(x='Exited',y = df_num[col],data=df_num)  
    plt.show()
```





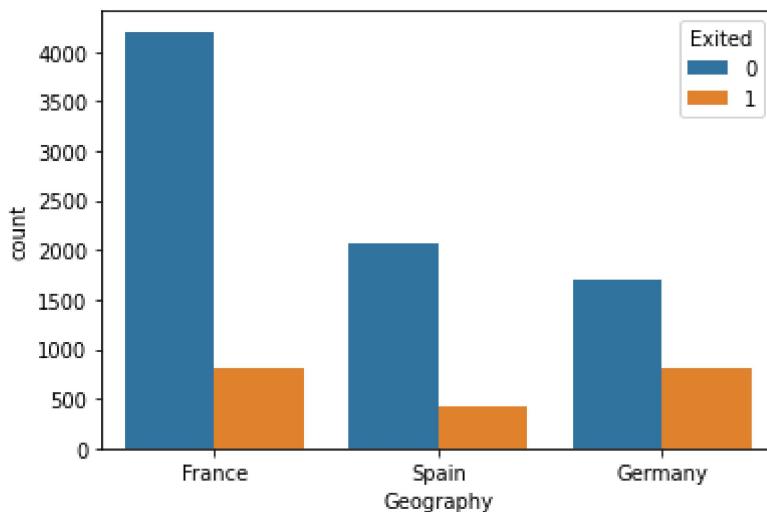




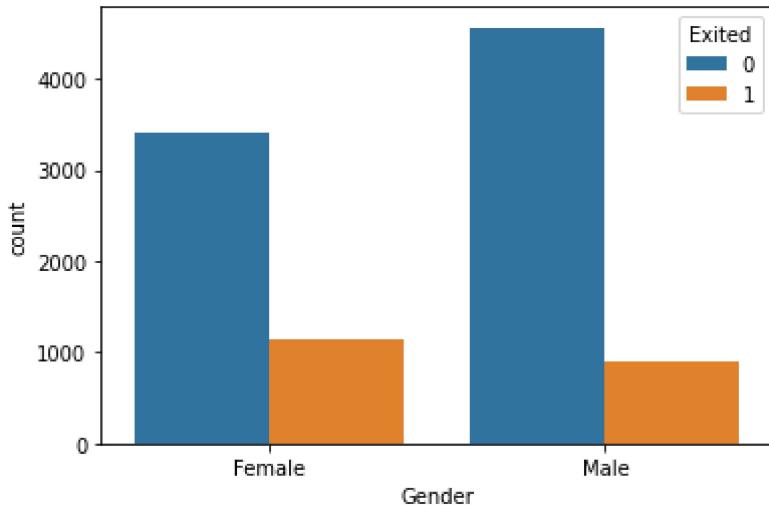
From the boxplots we can conclude that:

- 1> There is no significant difference in the credit score distribution between retained and churned customers.
- 2> Customers mostly between the age group of 35 to 50 are churning and those between 30 to 40 age group are not. The bank may need to review their target market or review the strategy for retention between the different age groups.
- 3> Customers having shorter tenure period have stayed with the bank while those having a larger tenure period have Exited.
- 4> The bank is losing customers with significant bank balances.
- 5> NumOfProducts, HasCrCard, IsActiveMember and EstimatedSalary do not have any significant relation with Customer churn.

```
In [32]: plt.figure()
sns.countplot('Geography', hue = 'Exited', data = df)
plt.show()
```



```
In [33]: plt.figure()
sns.countplot('Gender', hue = 'Exited', data = df)
plt.show()
```

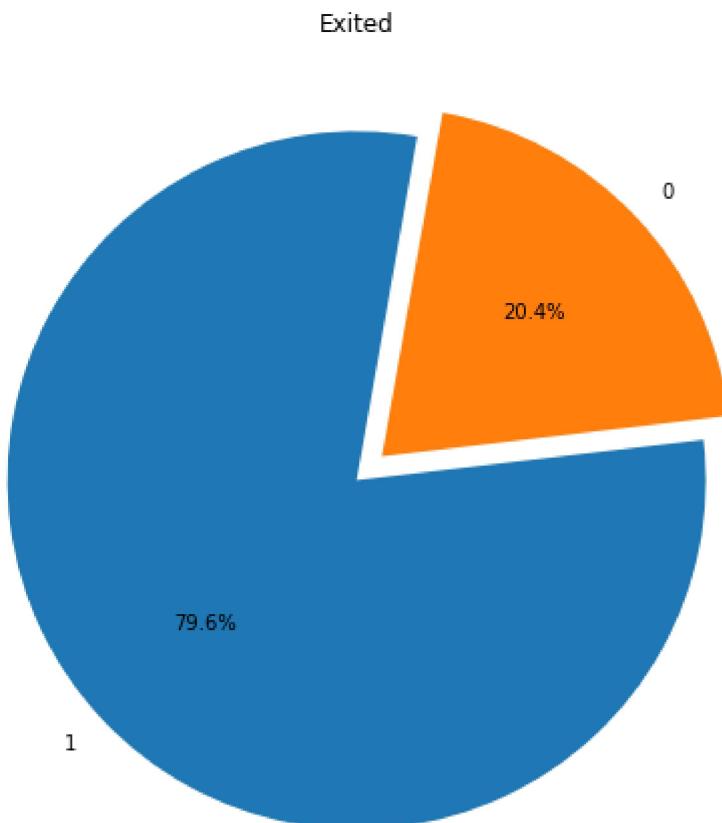


From the count plots we can conclude that :

- 1> Of the three countries, Germany has the highest Churning rate where as France has the highest number of customer who stayed.
- 2>The number of female customers churning is greater than that of male.

CHECKING COLUMNS FOR IMBALANCED DATA :

```
In [34]: explode = (0, 0.1)
plt.figure(figsize=(8,8))
plt.pie(df_num['Exited'].value_counts(), labels=df_num['Exited'].unique(), explode =
plt.title('Exited')
plt.show()
```



We can see that there is imbalanced data in Exited column with almost 80 percent of customers churning and only 20 percent retained.

Label Encoding :

```
In [35]: from sklearn.preprocessing import LabelEncoder
for col in df_cat:
    le = LabelEncoder()
    df_cat[col] = le.fit_transform(df_cat[col])
```

```
In [36]: #Concatenating the numeric and categoric columns
df_new = pd.concat([df_num, df_cat], axis=1)
```

```
In [37]: df_new.head()
```

```
Out[37]:
```

	CreditScore	Age	Tenure	Balance	NumOfProducts	HasCrCard	IsActiveMember	EstimatedSal
0	619	42	2	0.00		1	1	101348
1	608	41	1	83807.86		1	0	112542
2	502	42	8	159660.80		3	1	113931
3	699	39	1	0.00		2	0	93826
4	850	43	2	125510.82		1	1	79084

```
In [38]: df_new.shape
```

```
Out[38]: (10000, 12)
```

```
In [39]: y = df_new[["Exited"]]
X = df_new.drop(["Exited"],axis=1)
```

```
In [40]: for col in X:
    if skew(X[col] >0.5) or skew(X[col] <0.5):
        X[col]= np.sqrt(X[col])
```

```
In [129...]: #Splitting training and testint data.
X_train,X_test,y_train,y_test = train_test_split(X,y,test_size=0.3,random_state=1)
```

```
In [130...]: from sklearn.preprocessing import StandardScaler
ss = StandardScaler()
X_train_ss = ss.fit_transform(X_train)
X_test_ss = ss.transform(X_test)
```

BASELINE MODEL :

```
In [131...]: base_model = Sequential()
base_model.add(Dense(16,input_shape=(X.shape[1],),activation = 'relu'))
base_model.add(Dense(8,activation='relu'))
base_model.add(Dense(4,activation='relu'))
base_model.add(Dense(1,activation='sigmoid'))
```

```
In [132...]: base_model.summary()
```

```
Model: "sequential_8"
```

Layer (type)	Output Shape	Param #
<hr/>		

dense_32 (Dense)	(None, 16)	192
dense_33 (Dense)	(None, 8)	136
dense_34 (Dense)	(None, 4)	36
dense_35 (Dense)	(None, 1)	5
<hr/>		
Total params: 369		
Trainable params: 369		
Non-trainable params: 0		

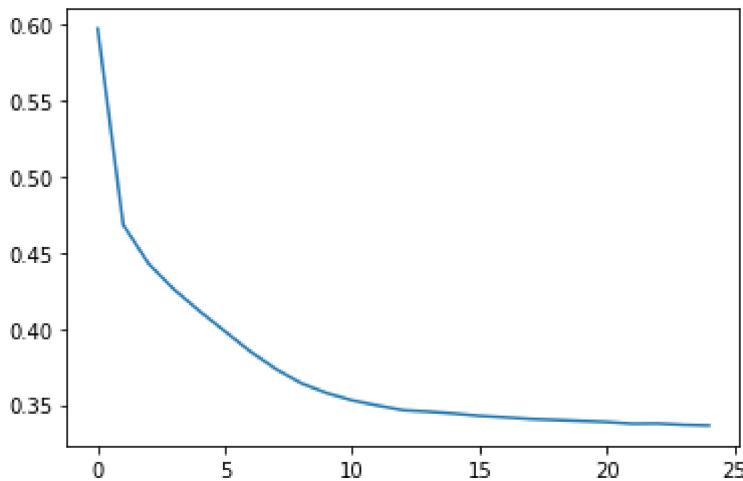
```
In [133...]: base_model.compile(optimizer = 'adam', loss = 'binary_crossentropy')
```

```
In [134...]: base_model_history = base_model.fit(X_train_ss, y_train, epochs = 25, batch_size = 5)
```

Epoch 1/25
140/140 [=====] - 0s 754us/step - loss: 0.6723
Epoch 2/25
140/140 [=====] - 0s 889us/step - loss: 0.4766
Epoch 3/25
140/140 [=====] - 0s 721us/step - loss: 0.4434
Epoch 4/25
140/140 [=====] - 0s 687us/step - loss: 0.4348
Epoch 5/25
140/140 [=====] - 0s 687us/step - loss: 0.4161
Epoch 6/25
140/140 [=====] - 0s 720us/step - loss: 0.3967
Epoch 7/25
140/140 [=====] - 0s 682us/step - loss: 0.3936
Epoch 8/25
140/140 [=====] - 0s 693us/step - loss: 0.3819
Epoch 9/25
140/140 [=====] - 0s 714us/step - loss: 0.3593
Epoch 10/25
140/140 [=====] - 0s 663us/step - loss: 0.3430
Epoch 11/25
140/140 [=====] - 0s 678us/step - loss: 0.3630
Epoch 12/25
140/140 [=====] - 0s 900us/step - loss: 0.3529
Epoch 13/25
140/140 [=====] - 0s 718us/step - loss: 0.3629
Epoch 14/25
140/140 [=====] - 0s 805us/step - loss: 0.3451
Epoch 15/25
140/140 [=====] - 0s 704us/step - loss: 0.3407
Epoch 16/25
140/140 [=====] - 0s 785us/step - loss: 0.3465
Epoch 17/25
140/140 [=====] - 0s 720us/step - loss: 0.3430
Epoch 18/25
140/140 [=====] - 0s 829us/step - loss: 0.3498
Epoch 19/25
140/140 [=====] - 0s 847us/step - loss: 0.3423
Epoch 20/25
140/140 [=====] - 0s 772us/step - loss: 0.3438
Epoch 21/25
140/140 [=====] - 0s 806us/step - loss: 0.3354
Epoch 22/25
140/140 [=====] - 0s 681us/step - loss: 0.3325
Epoch 23/25
140/140 [=====] - 0s 794us/step - loss: 0.3366
Epoch 24/25
140/140 [=====] - 0s 711us/step - loss: 0.3235
Epoch 25/25
140/140 [=====] - 0s 874us/step - loss: 0.3417

```
In [135...]: plt.plot(base_model_history.history["loss"])
```

Out[135]: [`<matplotlib.lines.Line2D at 0x7f9193b849d0>`]



In [136]: `y_pred_base = base_model.predict(X_test_ss)`

In [137]: `y_pred_base = np.where(y_pred_base>0.5,1,0)`

In [138]: `print(classification_report(y_test,y_pred_base))`

	precision	recall	f1-score	support
0	0.86	0.97	0.91	2373
1	0.77	0.43	0.55	627
accuracy			0.85	3000
macro avg	0.82	0.70	0.73	3000
weighted avg	0.84	0.85	0.84	3000

Our main aim is to predict the customers that will possibly churn so that the bank can retain them. So if a customer who is going to churn is wrongly predicted by the model as not going to churn, the bank won't make strategy to retain that customer and will lose him. Hence the False negative i.e. recall measures on the 1's is of more importance to us than the overall accuracy score of the model. For the base line model as we can see the value of recall is very low.

Handeling imbalanced data

Oversampling

In [53]: `from imblearn.over_sampling import RandomOverSampler`

In [54]: `ros = RandomOverSampler(random_state=1)`

In [55]: `X_sample_1 ,y_sample_1 = ros.fit_resample(X_train_ss,y_train)`

In [56]: `pd.Series(y_sample_1).value_counts()`

Out[56]:

1	5590
0	5590
	dtype: int64

Under Sampling

```
In [57]: from imblearn.under_sampling import RandomUnderSampler
In [58]: rus = RandomUnderSampler(random_state=1)
In [59]: X_sample_2, y_sample_2 = rus.fit_resample(X_train_ss, y_train)
In [ ]: pd.Series(y_sample_2).value_counts()
```

ANN WITH OVERSAMPLING:

```
In [139... model = Sequential()
model.add(Dense(16, input_shape=(X.shape[1],), activation="relu"))
model.add(Dense(8, activation="relu"))
model.add(Dense(4, activation="relu"))
model.add(Dense(1, activation="sigmoid"))

In [140... model.compile(loss="binary_crossentropy", optimizer="adam")

In [141... model.summary()

Model: "sequential_9"
-----
```

Layer (type)	Output Shape	Param #
dense_36 (Dense)	(None, 16)	192
dense_37 (Dense)	(None, 8)	136
dense_38 (Dense)	(None, 4)	36
dense_39 (Dense)	(None, 1)	5

```
Total params: 369
Trainable params: 369
Non-trainable params: 0
```

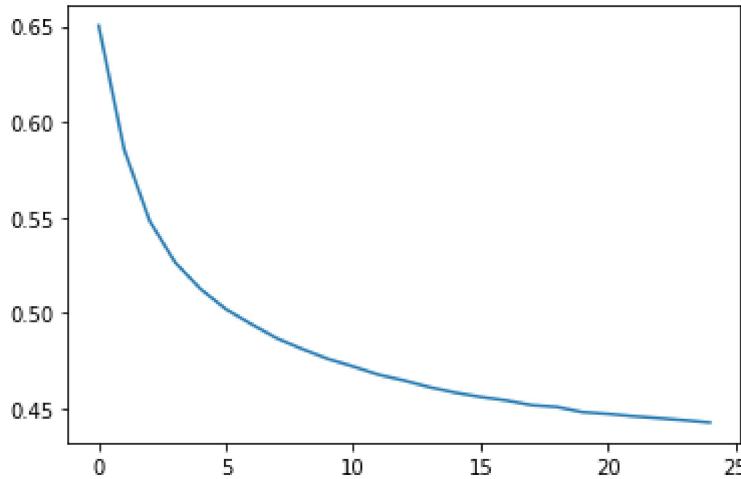
```
In [104... history = model.fit(X_sample_1,y_sample_1,epochs=25,batch_size=50)

Epoch 1/25
224/224 [=====] - 1s 778us/step - loss: 0.6690
Epoch 2/25
224/224 [=====] - 0s 746us/step - loss: 0.6000
Epoch 3/25
224/224 [=====] - 0s 752us/step - loss: 0.5580
Epoch 4/25
224/224 [=====] - 0s 789us/step - loss: 0.5303
Epoch 5/25
224/224 [=====] - 0s 783us/step - loss: 0.5106
Epoch 6/25
224/224 [=====] - 0s 781us/step - loss: 0.5035
Epoch 7/25
224/224 [=====] - 0s 822us/step - loss: 0.4999
Epoch 8/25
224/224 [=====] - 0s 804us/step - loss: 0.4977
Epoch 9/25
224/224 [=====] - 0s 795us/step - loss: 0.4797
Epoch 10/25
224/224 [=====] - 0s 772us/step - loss: 0.4795
Epoch 11/25
```

```
224/224 [=====] - 0s 735us/step - loss: 0.4658
Epoch 12/25
224/224 [=====] - 0s 739us/step - loss: 0.4682
Epoch 13/25
224/224 [=====] - 0s 760us/step - loss: 0.4615
Epoch 14/25
224/224 [=====] - 0s 764us/step - loss: 0.4612
Epoch 15/25
224/224 [=====] - 0s 819us/step - loss: 0.4551
Epoch 16/25
224/224 [=====] - 0s 787us/step - loss: 0.4558
Epoch 17/25
224/224 [=====] - 0s 805us/step - loss: 0.4597
Epoch 18/25
224/224 [=====] - 0s 833us/step - loss: 0.4429
Epoch 19/25
224/224 [=====] - 0s 775us/step - loss: 0.4492
Epoch 20/25
224/224 [=====] - 0s 757us/step - loss: 0.4511
Epoch 21/25
224/224 [=====] - 0s 814us/step - loss: 0.4489
Epoch 22/25
224/224 [=====] - 0s 772us/step - loss: 0.4538
Epoch 23/25
224/224 [=====] - 0s 791us/step - loss: 0.4408
Epoch 24/25
224/224 [=====] - 0s 791us/step - loss: 0.4407
Epoch 25/25
224/224 [=====] - 0s 765us/step - loss: 0.4477
```

In [105...]: `plt.plot(history.history["loss"])`

Out[105...]: [`<matplotlib.lines.Line2D at 0x7f919b6e5390>`]



In [106...]: `y_pred = model.predict(X_sample_1)`

In [107...]: `y_pred = np.where(y_pred >= 0.5, 1, 0)`

In [108...]: `print(classification_report(y_sample_1,y_pred))`

	precision	recall	f1-score	support
0	0.78	0.84	0.81	5590
1	0.82	0.76	0.79	5590
accuracy			0.80	11180
macro avg	0.80	0.80	0.80	11180
weighted avg	0.80	0.80	0.80	11180

As we can see that the values of Precision, Recall and F1score have highly increased from

that of the baseline model. Also Modelling with Oversampling has given us a descent recall of 76% which is the paramount for our Goal.

ANN WITH UNDERSAMPLING :

In [109...]

```
model = Sequential()
model.add(Dense(16, input_shape=(X.shape[1],), activation="relu"))
model.add(Dense(8, activation="relu"))
model.add(Dense(4, activation="relu"))
model.add(Dense(1, activation="sigmoid"))
```

In [110...]

```
model.compile(loss="binary_crossentropy", optimizer="adam")
```

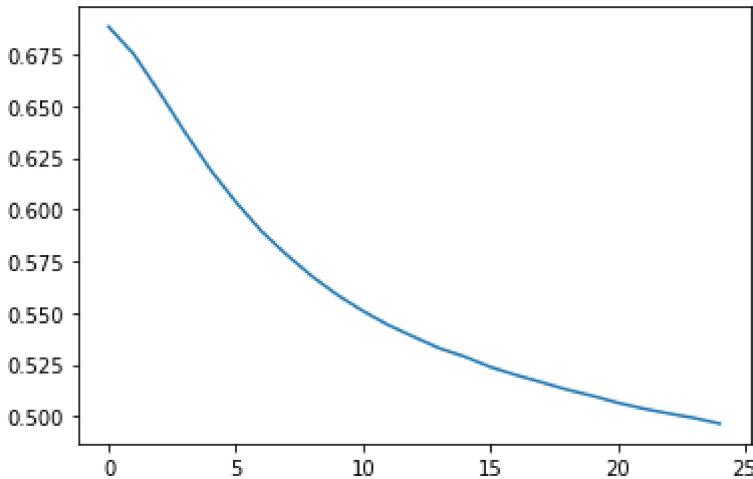
In [111...]

```
history = model.fit(X_sample_2, y_sample_2, epochs=25, batch_size=50)
```

```
Epoch 1/25
57/57 [=====] - 0s 792us/step - loss: 0.6922
Epoch 2/25
57/57 [=====] - 0s 737us/step - loss: 0.6792
Epoch 3/25
57/57 [=====] - 0s 728us/step - loss: 0.6625
Epoch 4/25
57/57 [=====] - 0s 756us/step - loss: 0.6466
Epoch 5/25
57/57 [=====] - 0s 728us/step - loss: 0.6135
Epoch 6/25
57/57 [=====] - 0s 752us/step - loss: 0.5957
Epoch 7/25
57/57 [=====] - 0s 738us/step - loss: 0.5924
Epoch 8/25
57/57 [=====] - 0s 789us/step - loss: 0.5811
Epoch 9/25
57/57 [=====] - 0s 693us/step - loss: 0.5760
Epoch 10/25
57/57 [=====] - 0s 804us/step - loss: 0.5527
Epoch 11/25
57/57 [=====] - 0s 796us/step - loss: 0.5575
Epoch 12/25
57/57 [=====] - 0s 748us/step - loss: 0.5426
Epoch 13/25
57/57 [=====] - 0s 773us/step - loss: 0.5415
Epoch 14/25
57/57 [=====] - 0s 767us/step - loss: 0.5368
Epoch 15/25
57/57 [=====] - 0s 886us/step - loss: 0.5259
Epoch 16/25
57/57 [=====] - 0s 715us/step - loss: 0.5214
Epoch 17/25
57/57 [=====] - 0s 742us/step - loss: 0.5263
Epoch 18/25
57/57 [=====] - 0s 731us/step - loss: 0.5141
Epoch 19/25
57/57 [=====] - 0s 845us/step - loss: 0.5046
Epoch 20/25
57/57 [=====] - 0s 711us/step - loss: 0.5080
Epoch 21/25
57/57 [=====] - 0s 754us/step - loss: 0.5068
Epoch 22/25
57/57 [=====] - 0s 763us/step - loss: 0.4970
Epoch 23/25
57/57 [=====] - 0s 721us/step - loss: 0.5008
Epoch 24/25
57/57 [=====] - 0s 797us/step - loss: 0.4902
Epoch 25/25
57/57 [=====] - 0s 1ms/step - loss: 0.4974
```

```
In [112]: plt.plot(history.history["loss"])
```

```
Out[112]: []
```



```
In [113]: y_pred = model.predict(X_sample_2)
```

```
In [114]: y_pred = np.where(y_pred >= 0.5, 1, 0)
```

```
In [115]: print(classification_report(y_sample_2,y_pred))
```

	precision	recall	f1-score	support
0	0.77	0.77	0.77	1410
1	0.77	0.77	0.77	1410
accuracy			0.77	2820
macro avg	0.77	0.77	0.77	2820
weighted avg	0.77	0.77	0.77	2820

As we can see Modelling with Undersampling is giving us better precision, recall and f1score than the Base model but, not as good as Modelling with Oversampling.

CONCLUSION:

From the review of the fitted models above, the best model that gives a decent balance of the recall and precision is ANN WITH OVERSAMPLING, where we are getting a precision on 1's of 82 percent, recall of 76 percent, f1-score of 79 percent and accuracy of 80 percent. Out of all customers that the model thinks will churn, 82 percent do actually churn.