

Bank Customer Churn Prediction

The aim of this project is to analyze the bank customers' demographics and financial information which includes customer's age, gender, credit score, balance and many others to predict whether the customer will leave the bank or not.

About the dataset The dataset is taken from Kaggle. It contains 10000 rows and 14 columns. The objective of the dataset is to predict whether the customer will leave the bank or not, based on the customer's demographics and financial information included in the dataset.

The dataset has several factors that can influence the customer to leave the bank, which are termed as independent variables. The target variable is the customer's decision to leave the bank, which is termed as dependent variable.

```
In [1]: #Importing the libraries
import numpy as np
import pandas as pd
import matplotlib.pyplot as plt
import seaborn as sns

In [3]: #Loading the dataset
df = pd.read_csv('churn.csv')
df.head()

Out[3]:
   RowNumber  CustomerId  Surname  CreditScore  Geography  Gender  Age  Tenure  Balance  NumOfProducts  HasCreditCard  IsActiveMember  EstimatedSalary  Exited
0          0           1  15634602  Hargrave         619      France  Female  42      2      0.00              1              1              1      101348.88      1
1          1           2      1159264      Hill         608      Spain  Female  41      1      83807.86              1              0              1      112642.36      0
2          2           3  146126294      Owa         602      France  Female  42      8  159602.80              3              1              0      112932.37      1
3          3          4  15701354      Boni         699      France  Female  39      1      0.00              2              0              0      93626.63      0
4          4          5  1577888      Mitchell      650      Spain  Female  43      2  125510.82              1              1              1      79084.10      0

In [5]: #Checking the shape of the dataset
df.shape

Out[5]:
(10088, 14)

Dropping the unnecessary columns - RowNumber, CustomerId, Surname

In [6]: #Drop columns
df = df.drop(['RowNumber', 'CustomerId', 'Surname'], axis=1)

Checking for Null/Missing values

In [7]: #Is there any null value?
df.isnull().sum()

Out[7]:
CreditScore      0
Geography        0
Gender            0
Age              0
Tenure           0
Balance          0
NumOfProducts   0
HasCreditCard    0
IsActiveMember   0
EstimatedSalary  0
Exited           0
dtype: int64

Checking the data types of the columns

In [8]: #Column data types
df.dtypes

Out[8]:
CreditScore      int64
Geography        object
Age              int64
Tenure           int64
Balance          float64
NumOfProducts    int64
HasCreditCard    int64
IsActiveMember   float64
EstimatedSalary  float64
Exited           int64
dtype: object

Checking for duplicate values

In [9]: #Duplicate values
df.duplicated().sum()

Out[9]:
0

Renaming the column 'Exited' to 'Churn'

In [10]: #Rename column
df.rename(columns={'Exited':'Churn'}, inplace=True)

Descriptive Statistics

In [11]: #Descriptive statistics
df.describe()

Out[11]:
   CreditScore      Age      Tenure      Balance  NumOfProducts  HasCreditCard  IsActiveMember  EstimatedSalary      Churn
count  10000.000000  10000.000000  10000.000000  10000.000000  10000.000000  10000.000000  10000.000000  10000.000000  10000.000000
mean     65.528000   38.821800     5.012800   76485.889208      1.520200      0.705500      0.515100  100009.239811      0.203700
std      96.653299   10.487806     2.892174   62297.452002      0.581654      0.459044      0.499797   57510.492818      0.402769
min      50.000000   18.000000      0.000000      0.000000      0.000000      0.000000      0.000000     11.580000      0.000000
25%      58.400000   32.000000      3.000000      0.000000      1.000000      0.000000      0.000000   51002.110000      0.000000
50%      62.000000   37.000000      5.000000   97395.540000      1.000000      1.000000      1.000000  100193.915000      0.000000
75%      73.000000   41.000000      7.000000  127044.240000      2.000000      1.000000      1.000000  146886.475000      0.000000
max     150.000000   52.000000     10.000000  259989.990000      4.000000      1.000000      1.000000  199992.480000      1.000000
```

```
In [12]: df.head()

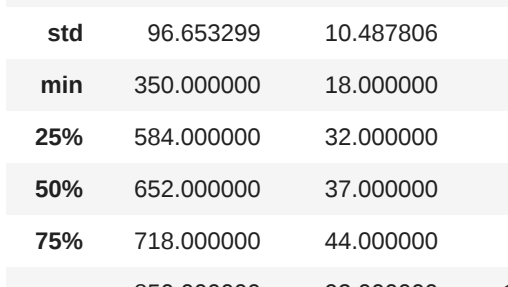
Out[12]:
   CreditScore  Geography  Gender  Age  Tenure  Balance  NumOfProducts  HasCreditCard  IsActiveMember  EstimatedSalary  Churn
0           619      France  Female  42      2      0.00              1              1              1      101348.88      1
1           608      Spain  Female  41      1   83807.86              1              0              1      112642.36      0
2           602      France  Female  42      8  159602.80              3              1              0      112932.37      1
3           699      France  Female  39      1      0.00              2              0              0      93626.63      0
4           650      Spain  Female  43      2  125510.82              1              1              1      79084.10      0
```

Explorative Data Analysis

Explorative data analysis. We will be looking at the distribution of the data, the correlation between features and the target variable and the relationship between the features and the target variable. I will start by looking at the distribution of the data, followed by the relationship between the features and the target variable.

Pie Chart for Customer Churn

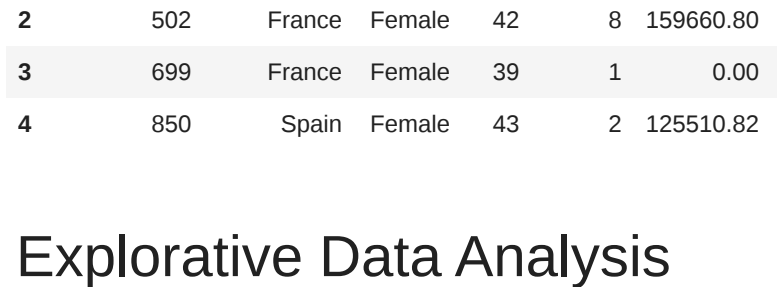
```
In [51]: #Pie chart
plt.figure(figsize=(10,6))
sns.countplot(x='Churn', data=df, hue='Churn')
plt.title('Gender Distribution')
plt.xlabel('Gender')
plt.ylabel('Churn Percentage')
```



The pie chart clearly visualizes the customer churn in the dataset. The majority of the customers in the dataset continue to use the services of the bank with only 20.4% of the customers churning.

Gender

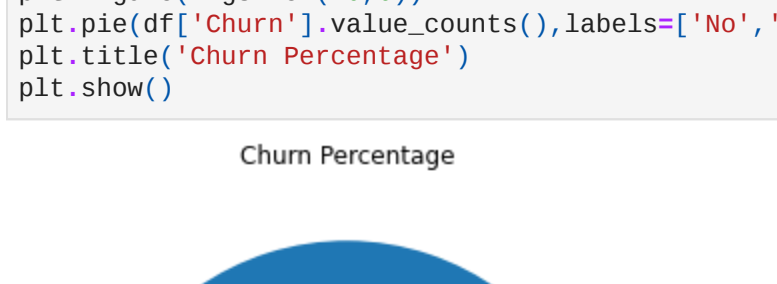
```
In [52]: #gender and customer churn
sns.countplot(x='Gender', data=df, hue='Churn')
plt.title('Gender Distribution')
plt.xlabel('Gender')
plt.ylabel('Count')
plt.show()
```



As shown in the graph, majority of the customers are male. But upon looking at the customer churn, we can see that females have more tendency to churn as compared to males. However there is not much difference between the churn count of the two genders so we cannot have a hypothesis regarding the customer churn based on the gender of the customer.

Age Distribution

```
In [16]: #Histogram for age distribution
sns.histplot(data=df, x='Age', hue='Churn', multiple='stack', kde=True)
<AxesSubplot: xlabel='Age', ylabel='Count'>
```



This histogram visualizes the age distribution and the churn count of the customers. The majority of the customers are from age group 30-40 years old. However the customer churn count is highest for the customers' age 40 and 50. In addition to that customers from age group 20-25 years old count for the lowest churn count. Therefore, age plays a significant role in customer churn, where late adults are more likely to churn as compared to young adults with minimal churn count.

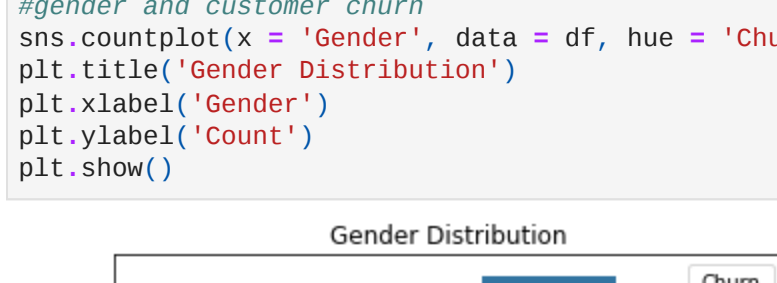
Credit Score

```
In [17]: fig, ax = plt.subplots(1,2,figsize=(15, 5))
sns.boxplot(x='Churn', y='CreditScore', data=df, ax=ax[0])
sns.violinplot(x='Churn', y='CreditScore', data=df, ax=ax[1])
<AxesSubplot: xlabel='Churn', ylabel='CreditScore'>
```



Customer location

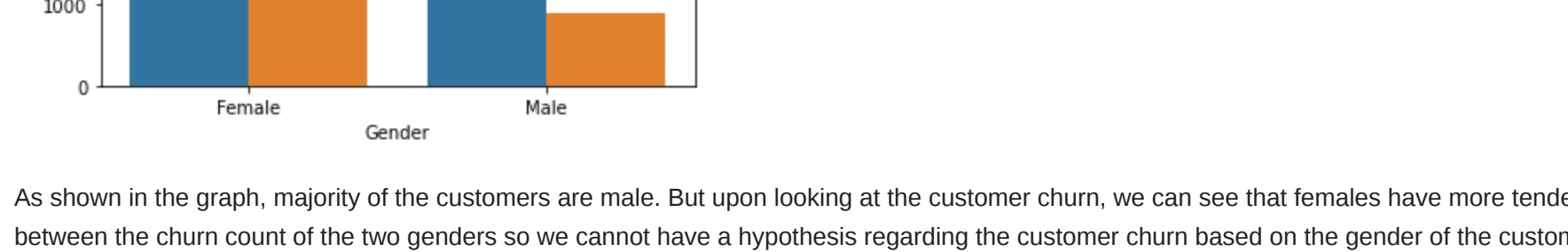
```
In [18]: #sns.countplot(x='Geography', hue='Churn', data=df)
plt.title('Geography and Churn')
plt.xlabel('Geography')
plt.ylabel('Count')
plt.show()
```



This graphs shows the number of customers from their respective countries along with their churn count. Majority of the customers are from France, followed by Spain and Germany. However in contrast to that Germany has the highest number of customer count followed by France and Spain. From this we can infer that German customers are more likely to churn than the customers from other countries.

Tenure

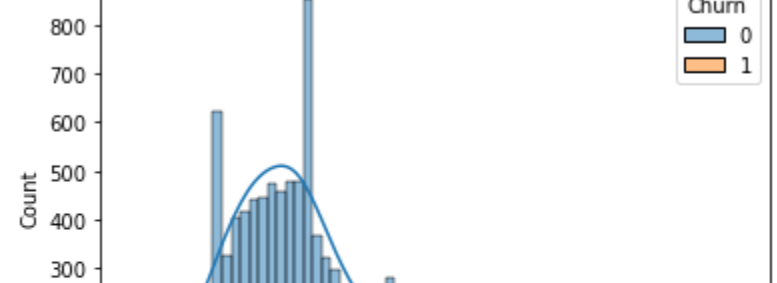
```
In [19]: fig, ax = plt.subplots(1,2,figsize=(15,5))
sns.countplot(x='Tenure', data=df, ax=ax[0])
sns.countplot(x='Tenure', hue='Churn', data=df, ax=ax[1])
<AxesSubplot: xlabel='Tenure', ylabel='count'>
```



Tenure refers to the time (in years) that a customer has been a client of the bank. Majority of the customers in the dataset have a tenure between 1-9 years, having equal distribution among them. There are very few customers with a tenure of less than 1 years or more than 9 years. Looking at the churn of these customers based on their tenure, it can be observed that customers with tenure 1-9 years have higher churn count with 3 years tenure followed those with 3 years tenure. However customers more than 9 years or tenure counts for the least churn. This is because the customers with higher tenure are more loyal to the bank and less likely to churn.

Bank Balance

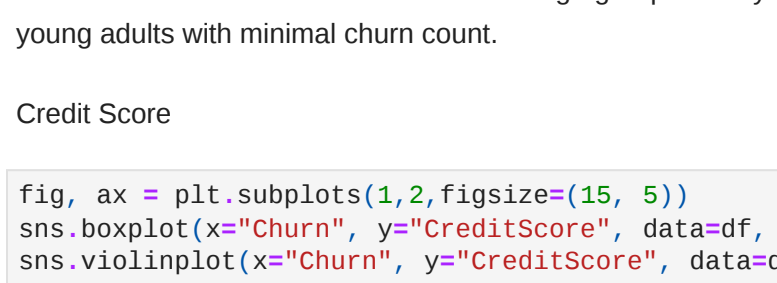
```
In [20]: sns.histplot(data=df, x='Balance', hue='Churn', multiple='stack', kde=True)
<AxesSubplot: xlabel='Balance', ylabel='count'>
```



A huge number of customers have zero bank balance which also resulted in them leaving the bank. However, customer having bank balance between 10000 to 150000 are more likely to leave the bank after the customers with zero bank balance.

Number of products purchased

```
In [21]: sns.countplot(x='NumOfProducts', hue='Churn', data=df)
<AxesSubplot: xlabel='NumOfProducts', ylabel='count'>
```



In the dataset, we have customers in four categories according to the number of products purchased. The customers with purchase of 1 or 2 products are highest in number and have low churn count in comparison to the non churn customers in the category. However, in the category where customers have purchased 3 or 4 products the number of leaving customers is much higher than the non leaving customers. Therefore, the number of product purchased is a good indicator of customer churn

Customers with/without credit card

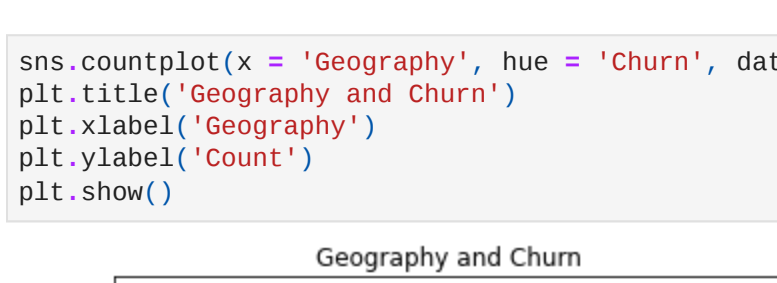
```
In [22]: sns.countplot(x=df['HasCreditCard'], hue=df['Churn'])
<AxesSubplot: xlabel='HasCreditCard', ylabel='count'>
```



Majority of the customers have credit cards i.e. nearly 70% of the customers have credit cards leaving 30% of the customers who do not have credit cards. Moreover, the number of customers leaving the bank are more whom have a credit card.

Active Members

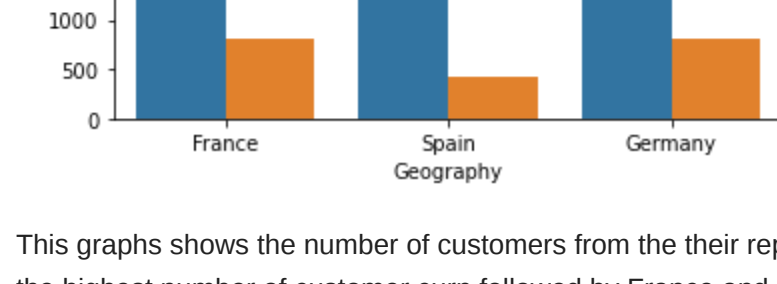
```
In [23]: sns.countplot(x='IsActiveMember', hue='Churn', data=df)
<AxesSubplot: xlabel='IsActiveMember', ylabel='count'>
```



As expected, the churn count is higher for non active members as compared to the active members of the bank. This is because the active members are more satisfied with the services of the bank and hence they are less likely to leave the bank. Therefore, the bank should focus on the non active members and try to improve their services to retain them.

Estimated Salary

```
In [25]: sns.histplot(data=df, x='EstimatedSalary', hue='Churn', multiple='stack', palette='Set2')
<AxesSubplot: xlabel='EstimatedSalary', ylabel='count'>
```



Data Preprocessing-2

Label encoding the variables

```
In [27]: #label encoding
variables = ['Geography', 'Gender']
from sklearn.preprocessing import LabelEncoder
le=LabelEncoder()
for i in variables:
    le.fit(df[i].unique())
    df[i]=le.transform(df[i])
print(df[i].unique())

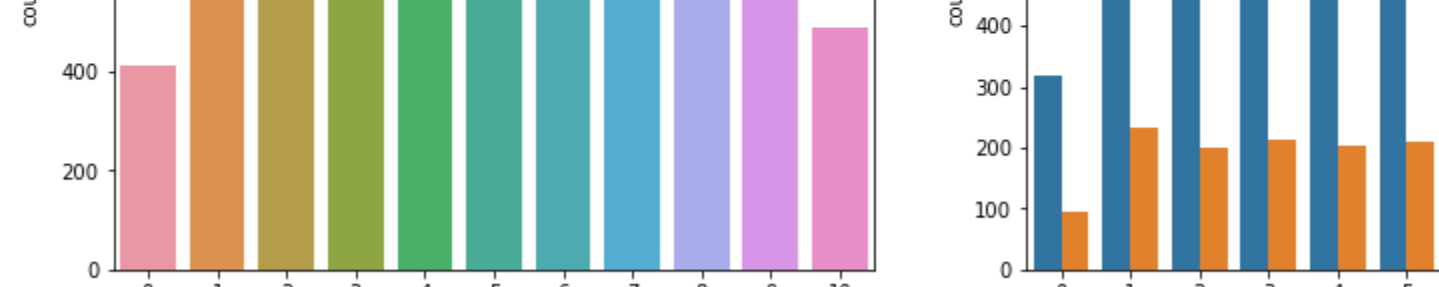
Geography [0 1]
Gender [0 1]
```

Normalization

```
In [28]: #normalize the continuous variables
from sklearn.preprocessing import StandardScaler
scaler = StandardScaler()
df[['CreditScore', 'Balance', 'EstimatedSalary']] = scaler.fit_transform(df[['CreditScore', 'Balance', 'EstimatedSalary']])
```

Confusion Matrix Heatmap

```
In [29]: plt.figure(figsize=(12,12))
sns.heatmap(df.confusion_matrix(), annot=True, cmap='coolwarm')
plt.title('Confusion Matrix')
```



There is no significant correlation among the variables. So, lets proceed to model building

Train Test Split

```
In [31]: #train test split
from sklearn.model_selection import train_test_split
X_train,X_test,y_train,y_test=train_test_split(df.drop('Churn',axis=1),df['Churn'],test_size=0.3,random_state=42)
```

Churn Prediction For predicting the churn of customers, depending on the data of the customers, we will use the following models:

Decision Tree Classifier Random Forest Classifier

Decision Tree Classifier

Using GridSearchCV to find the best parameters for the model.

```
In [32]: from sklearn.tree import DecisionTreeClassifier
from sklearn.model_selection import GridSearchCV
#creating decision tree classifier object
dtree = DecisionTreeClassifier()
```

```
#defining parameter range
param_grid = {
    'max_depth': [2,4,6,8,10,12,14,16,18,20],
    'min_samples_leaf': [1,2,3,4,5,6,7,8,9,10],
    'criterion': ['gini', 'entropy'],
    'random_state': [0,42]}

#creating grid search object
grid_dtree = GridSearchCV(dtree, param_grid, cv = 5, scoring = 'roc_auc', n_jobs = -1, verbose = 1)
```

```
#fitting the grid search object to the training data
grid_dtree.fit(X_train, y_train)

#printing the best parameters
print('Best parameters found: ', grid_dtree.best_params_)
```

Fitting 5 folds for each of 480 candidates, totalling 2880 fits
Best parameters found: {'criterion': 'gini', 'max_depth': 6, 'min_samples_leaf': 10, 'random_state': 42}

Adding the parameters to the model

```
In [33]: dtree = DecisionTreeClassifier(criterion='gini', max_depth=6, random_state=42, min_samples_leaf=10)
dtree
```

```
Out[33]: DecisionTreeClassifier(max_depth=6, min_samples_leaf=10, random_state=42)
```

Random Forest Classifier

```
In [36]: #training the model
dtree.fit(X_train, y_train)

#training accuracy
dtree.score(X_train, y_train)
```

```
Out[36]: 0.8787428571428571
```

Predicting the customer churn from Test set

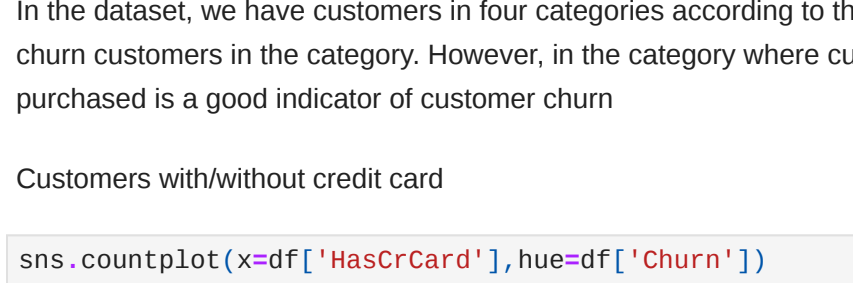
```
In [39]: rfc_pred = rfc.predict(X_test)
```

Model Evaluation

Decision Tree Classifier

Confusion Matrix Heatmap

```
In [40]: #confusion matrix heatmap
from sklearn.metrics import confusion_matrix
plt.figure(figsize=(8,6))
sns.heatmap(confusion_matrix(y_test,dtree_pred),annot=True,fmt='d',cmap='Blues')
plt.xlabel('Predicted')
plt.ylabel('Actual')
plt.title('Confusion Matrix for Decision Tree')
plt.show()
```



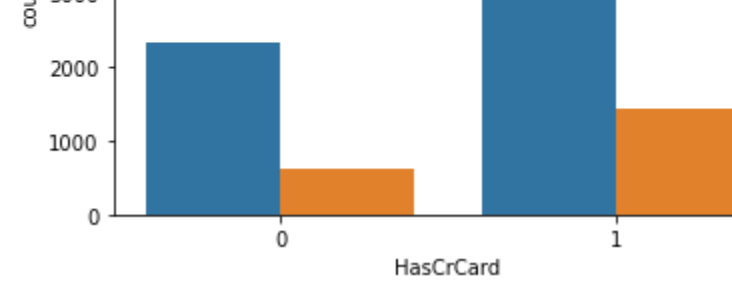
The True Positive shows the count of correctly classified data points whereas the False Positive elements are those that are misclassified by the model. The higher the True Positive values of the confusion matrix the better, indicating many correct predictions

Distribution Plot

```
In [41]: ax = sns.displot(y_test, hist=False, color='r', label='Actual Value')
sns.displot(rfc_pred, hist=False, color='b', label='Fitted Values', ax=ax)

C:\Users\Wp\anaconda3\lib\site-packages\seaborn\distributions.py:2619: FutureWarning: 'displot' 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 'kdeplot' (an axes-level function for kernel density plots).
C:\Users\Wp\anaconda3\lib\site-packages\seaborn\distributions.py:2619: FutureWarning: 'displot' 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 'kdeplot' (an axes-level function for kernel density plots).
warnings.warn(msg, FutureWarning)

<AxesSubplot: xlabel='Churn', ylabel='Density'>
```



The more overlapping of two colors, the more accurate the model is.

Classification Report

```
In [43]: from sklearn.metrics import classification_report
print(classification_report(y_test, dtree_pred))
```

	precision	recall	f1-score	support
0	0.87	0.98	0.92	2416
1	0.86	0.39	0.52	584
accuracy	0.83	0.68	0.76	3000
macro avg	0.83	0.68	0.72	3000
weighted avg	0.85	0.66	0.74	3000

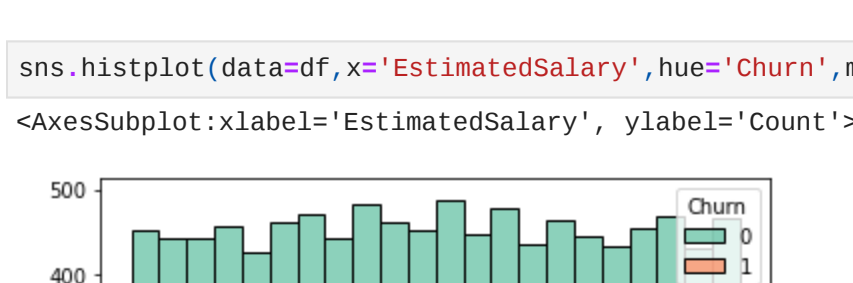
```
In [44]: from sklearn.metrics import accuracy_score, mean_absolute_error, r2_score
print("Accuracy Score: ", accuracy_score(y_test, dtree_pred))
print("Mean Absolute Error: ", mean_absolute_error(y_test, dtree_pred))
print("R2 Score: ", r2_score(y_test, dtree_pred))
```

Accuracy Score: 0.868
Mean Absolute Error: 0.13066666666666666
R2 Score: 0.15480823333333333

Random Forest Classifier

Confusion Matrix Heatmap

```
In [46]: plt.figure(figsize=(8,6))
sns.heatmap(confusion_matrix(y_test,rfc_pred),annot=True,fmt='d',cmap='Blues')
plt.xlabel('Predicted')
plt.ylabel('Actual')
plt.title('Confusion Matrix for Random Forest Classifier')
plt.show()
```



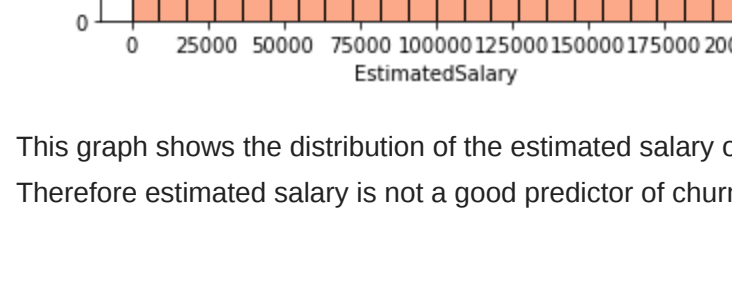
The True Positive shows the count of correctly classified data points whereas the False Positive elements are those that are misclassified by the model. The higher the True Positive values of the confusion matrix the better, indicating many correct predictions

Distribution Plot

```
In [48]: ax = sns.displot(y_test, hist=False, color='r', label='Actual Value')
sns.displot(rfc_pred, hist=False, color='b', label='Fitted Values', ax=ax)

C:\Users\Wp\anaconda3\lib\site-packages\seaborn\distributions.py:2619: FutureWarning: 'displot' 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 'kdeplot' (an axes-level function for kernel density plots).
C:\Users\Wp\anaconda3\lib\site-packages\seaborn\distributions.py:2619: FutureWarning: 'displot' 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 'kdeplot' (an axes-level function for kernel density plots).
warnings.warn(msg, FutureWarning)

<AxesSubplot: xlabel='Churn', ylabel='Density'>
```



Classification Report

```
In [49]: from sklearn.metrics import classification_report
print(classification_report(y_test, rfc_pred))
```

	precision	recall	f1-score	support
0	0.87	0.98	0.92	2416
1	0.82	0.41	0.55	584
accuracy	0.85	0.70	0.77	3000
macro avg	0.85	0.70	0.74	3000
weighted avg	0.86	0.67	0.76	3000

```
In [50]: print("Accuracy Score: ", accuracy_score(y_test, rfc_pred))
print("Mean Absolute Error: ", mean_absolute_error(y_test, rfc_pred))
print("R2 Score: ", r2_score(y_test, rfc_pred))
```

Accuracy Score: 0.886
Mean Absolute Error: 0.102
R2 Score: 0.15480823333333333

Conclusion From the exploratory data analysis, I have concluded that the churn count of the customers depends upon the following factors:

1.Age 2.Geography 3.Tenure 4.Balance 5.Number of Products 6.Has Credit Card 7.IsActive Member Coming to the classification models. I have used the following models:

Decision Tree Classifier Random Forest Classifier Both the models were hyperparameter tuned using GridSearchCV. Both the models have nearly equal accuracy score. But, the Random Forest Classifier has a better accuracy score than the Decision Tree Classifier.

In []:

In []: