

Eilidh Pike - Bank Customer Churn Prediction

Task: You are working as Data Scientist in one of the leading banking sector firms. The firm is losing customers at an alarming rate. Your task is to analyse the dataset provided(Exploratory data analysis) and build a machine learning model to predict which customers are likely to churn, so that the company can take proactive measures to retain them.

Aims:

- Identify and visualise which factors will contribute to customer churn
- Build a machine learning model that will predict whether the bank's customers will churn or not.

Dataset Information

- ID — Corresponds to the row number.
- CustomerId — Bank customer Identification Number.
- Surname – Surnames of customers.
- CreditScore - Credit Score of customer.
- Geography - Customer location.
- Gender - Gender of the customer.
- Age - Age of the customer.
- Tenure - How long the customer has been a member of the bank.
- Balance - The balance of a customer.
- NumOfProducts (Number of products) - How many products (accounts and cards) a customer has with the bank.
- HasCrCard (Has credit card) - If the customer has a credit card with the bank.
- IsActiveMember - If the customer is a member of the bank.
- EstimatedSalary - The estimated salary of the customer with the bank.
- Exited - If the customer has left the bank.

In [1]:

```
1 #Libraries
2 import pandas as pd
3 import numpy as np
4 import matplotlib.pyplot as plt
5 import seaborn as sns
6 from sklearn import preprocessing
7 from scipy.stats import chi2_contingency
8 from scipy.stats import pearsonr, spearmanr
9 from sklearn.metrics import confusion_matrix
10 from sklearn.metrics import classification_report
11 from sklearn.metrics import roc_curve, roc_auc_score
12 from sklearn.preprocessing import LabelEncoder
13 from sklearn.preprocessing import StandardScaler
14 from imblearn.over_sampling import SMOTE
15 from sklearn.ensemble import RandomForestClassifier, GradientBoostingClassifier
16 from sklearn.tree import DecisionTreeClassifier
17 import xgboost as xgb
18 import lightgbm as lgb
19 import catboost as cb
20 from sklearn.model_selection import GridSearchCV, train_test_split
21 from sklearn.metrics import accuracy_score, precision_score, recall_score, f1_score
22
23 #Seaborn Parameters
24 sns.set_palette("PRGn")
25 sns.set_style("darkgrid")
```

In [2]:

```
1 # Read the train dataset
```

1) Exploratory Data Analysis (EDA)

In [3]:

```
1 # Displaying the first 5 rows of the df
```

Out[3]:

	id	CustomerId	Surname	CreditScore	Geography	Gender	Age	Tenure	Balance	NumOfProducts	HasCrCard	IsActiveMember	EstimatedSalary
0	0	15674932	Okwudilichukwu	668	France	Male	33.0	3	0.00	2	1.0	0.0	181449.97
1	1	15749177	Okwudiliolisa	627	France	Male	33.0	1	0.00	2	1.0	1.0	49503.50
2	2	15694510	Hsueh	678	France	Male	40.0	10	0.00	2	1.0	0.0	184866.69
3	3	15741417	Kao	581	France	Male	34.0	2	148882.54	1	1.0	1.0	84560.88
4	4	15766172	Chiemenam	716	Spain	Male	33.0	5	0.00	2	1.0	1.0	15068.83

```
In [4]:
Out[4]: (165034, 14)

• Train dataset has 165034 observations and 14 rows

In [5]: 1 # Check missing values
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 165034 entries, 0 to 165033
Data columns (total 14 columns):
#   Column                Non-Null Count  Dtype
---  -
0   id                    165034 non-null  int64
1   CustomerId            165034 non-null  int64
2   Surname               165034 non-null  object
3   CreditScore           165034 non-null  int64
4   Geography             165034 non-null  object
5   Gender               165034 non-null  object
6   Age                  165034 non-null  float64
7   Tenure               165034 non-null  int64
8   Balance              165034 non-null  float64
9   NumOfProducts        165034 non-null  int64
10  HasCrCard             165034 non-null  float64
11  IsActiveMember        165034 non-null  float64
12  EstimatedSalary       165034 non-null  float64
13  Exited                165034 non-null  int64
```

- No missing values in the dataframe.

```
In [6]: 1 # Get unique count for each variable
```

```
Out[6]: id                165034
CustomerId            23221
Surname                2797
CreditScore           457
Geography              3
Gender                 2
Age                   71
Tenure                 11
Balance               30075
NumOfProducts          4
HasCrCard              2
IsActiveMember         2
EstimatedSalary       55298
Exited                 2
dtype: int64
```

```
In [7]: 1 #Checking for duplicate rows
2 duplicate_rows = df[df.duplicated()]
```

Out[7]: 0

```
In [8]:
```

```
Out[8]: Hsia                2456
T'ien                2282
Hs?                 1611
Kao                  1577
Maclean              1577
...
Samaniego            1
Lawley                1
Bonwick              1
Tennant              1
Elkins               1
Name: Surname, Length: 2797, dtype: int64
```

Columns 'id' and 'CustomerId' are merely labels for the observations and have no relevance to customer churn. 'Surname' has been removed as there as multiple instances of the same surname and will lead to either model overfitting or potential bias towards a certain surname.

```
In [9]: 1 # Dropping columns not relevant to predicting customer churn
```

In [10]:

Out[10]:

	count	mean	std	min	25%	50%	75%	max
CreditScore	165034.0	656.454373	80.103340	350.00	597.00	659.0	710.0000	850.00
Age	165034.0	38.125888	8.867205	18.00	32.00	37.0	42.0000	92.00
Tenure	165034.0	5.020353	2.806159	0.00	3.00	5.0	7.0000	10.00
Balance	165034.0	55478.086689	62817.663278	0.00	0.00	0.0	119939.5175	250898.09
NumOfProducts	165034.0	1.554455	0.547154	1.00	1.00	2.0	2.0000	4.00
HasCrCard	165034.0	0.753954	0.430707	0.00	1.00	1.0	1.0000	1.00
IsActiveMember	165034.0	0.497770	0.499997	0.00	0.00	0.0	1.0000	1.00
EstimatedSalary	165034.0	112574.822734	50292.865585	11.58	74637.57	117948.0	155152.4675	199992.48
Exited	165034.0	0.211599	0.408443	0.00	0.00	0.0	0.0000	1.00

Credit Score: The distribution of credit scores appears to be somewhat right-skewed, as the mean is slightly lower than the median, and the maximum value is much higher than the 75th percentile. This skewness suggests that there may be a larger proportion of customers with relatively higher credit scores compared to those with lower credit scores.

Age: The average age of customers in the dataset is around 38 years old, with a standard deviation of approximately 8.87 years. The majority of customers fall within the age range of 32 to 42 years, as the 25th to 75th percentiles (Q1 to Q3) indicate. The youngest customer is 18 years old, while the oldest is 92 years old. This could highlight some potential outliers in the age column. The bank may not be as appealing to older customers- things like the loss of in-person banks and a lack of benefits, or handling of pensions may influence this.

Tenure: On average, customers have been with the bank for around 5 years, with a standard deviation of approximately 2.81 years. The minimum tenure is 0 years, which indicates that the bank is still attracting new customers. The maximum is 10 years - it's not clear if the longest a customer has stayed is 10 years or if this is the data cut-off.

Balance: The average balance across all customers is approximately 55,478, with a large standard deviation of around 62,817. There are customers with a balance of 0 - it might be worth exploring the correlation between balance and active member status. The highest balance observed is 250,898.

Number of Products (NumOfProducts): On average, customers have approximately 1.55 products with the bank, with a standard deviation of around 0.55. The majority of customers have either 1 or 2 products, but some have up to 4 products.

Has Credit Card (HasCrCard): Around 75% of customers in the dataset have a credit card, as indicated by the mean value. This suggests that credit card ownership is quite common among the bank's customers.

Is Active Member (IsActiveMember): Roughly half of the customers are active members, based on the mean value. This indicates that the other half may be inactive or less engaged with the bank's services. We could expect to see a relationship between active members, balance, and number of products.

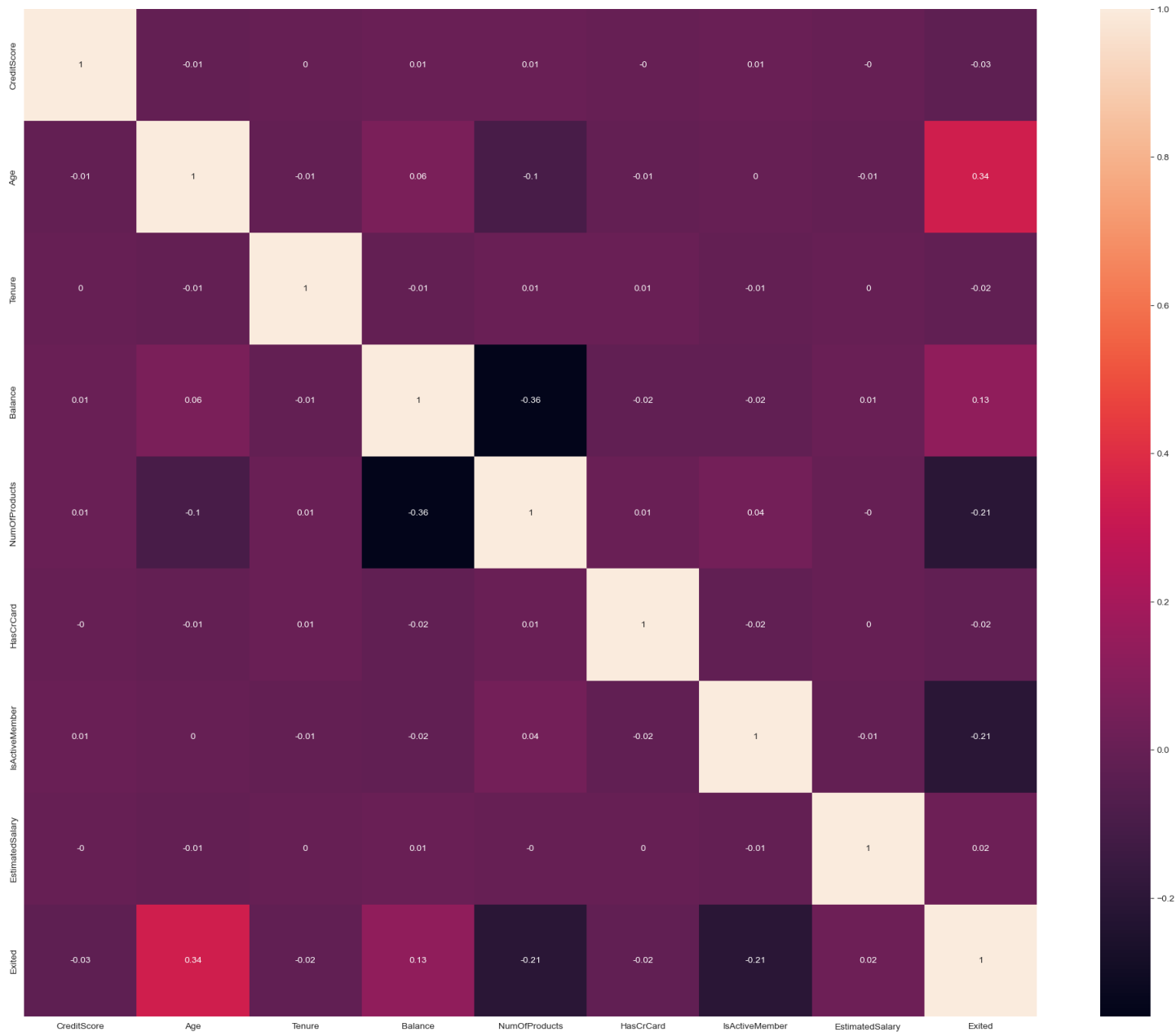
Estimated Salary: The average estimated salary of customers is approximately 112,575, with a standard deviation of around \$50,293. The minimum estimated salary is 11.58, while the maximum is 199,992.

Correlation between variables

```
In [11]: 1 corr = df.corr().round(2)
2         plt.figure(figsize = (25,20))

C:\Users\eilid\AppData\Local\Temp\ipykernel_13236\349184144.py:1: FutureWarning: The default value of numeric_only in DataFrame.corr is deprecated. In a future version, it will default to False. Select only valid columns or specify the value of numeric_only to silence this warning.
  corr = df.corr().round(2)
```

Out[11]: <Axes: >



From a first look, none of our variables display multicollinearity. There are a few relationships of note to the dependent variable:

- Age and Exited have a slightly positive correlation
- Balance and Exited have a slightly positive correlation
- NumofProducts and Exited have a slightly negative correlation
- IsActiveMember and Exited have a slightly negative correlation

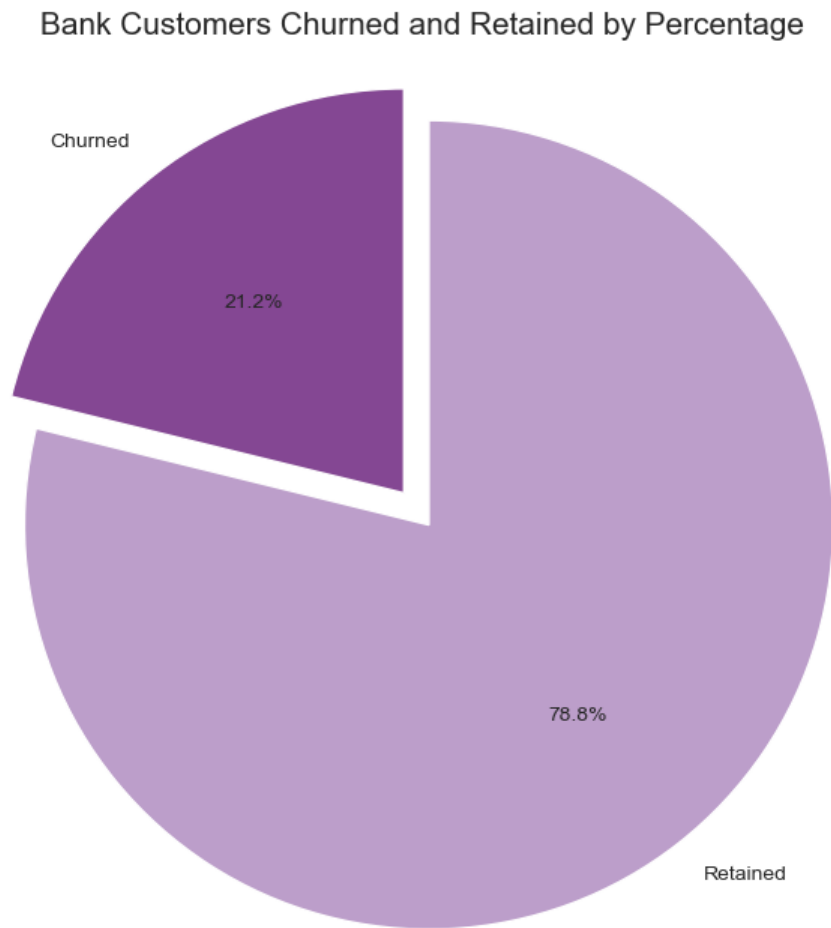
From a first look, we expect these variables to have an influence whether a customer churns.

Target Variable: Exited

The target variable has been previously encoded. The values are:

- 0 = Customer has not churned
- 1 = Customer has churned

```
In [12]: 1 labels = 'Churned', 'Retained'
2 sizes = [df.Exited[df['Exited']==1].count(), df.Exited[df['Exited']==0].count()]
3 explode = (0, 0.1)
4 fig1, ax1 = plt.subplots(figsize=(10, 8))
5 ax1.pie(sizes, explode=explode, labels=labels, autopct='%1.1f%%',
6         shadow=False, startangle=90)
7 ax1.axis('equal')
8 plt.title("Bank Customers Churned and Retained by Percentage", size = 15)
```



- The bank retains 78.8% of its customers and churns 21.2%.
- The target variable is imbalanced with our main concern, churned being very undersampled. In this classification task, this might lead to the models being biased towards predicting the majority class (Retained), and may perform poorly identifying churned customers. Therefore, a class oversampling technique (SMOTE) will be used in the data preprocessing stage.

Continuous Variables: Age, Credit Score, Balance & Estimated Salary

```
In [13]: 1 # Importing the libraries
```

In [14]:

```
1 #Printing the value counts for the continuous columns
2 for column in continuous:
3     print("Value counts for", column, ":")
4     print(df[column].value_counts())
```

Value counts for Age :

```
37.00    9255
38.00    9246
35.00    9118
34.00    8625
36.00    8556
```

...

```
84.00      4
83.00      3
85.00      3
36.44      1
32.34      1
```

Name: Age, Length: 71, dtype: int64

Value counts for CreditScore :

```
850     2532
678     2299
684     1718
667     1658
705     1605
```

...

```
419      1
386      1
358      1
423      1
373      1
```

Name: CreditScore, Length: 457, dtype: int64

Value counts for Balance :

```
0.00      89648
124577.33    88
127864.40     64
122314.50     63
129855.32     59
```

...

```
125824.21      1
158741.56      1
126815.52      1
61172.57       1
110993.29      1
```

Name: Balance, Length: 30075, dtype: int64

Value counts for EstimatedSalary :

```
88890.05     178
140941.47     107
167984.72     100
90876.95      98
129964.94      98
```

...

```
102747.73      1
170593.45      1
109179.48      1
60538.47       1
71173.03       1
```

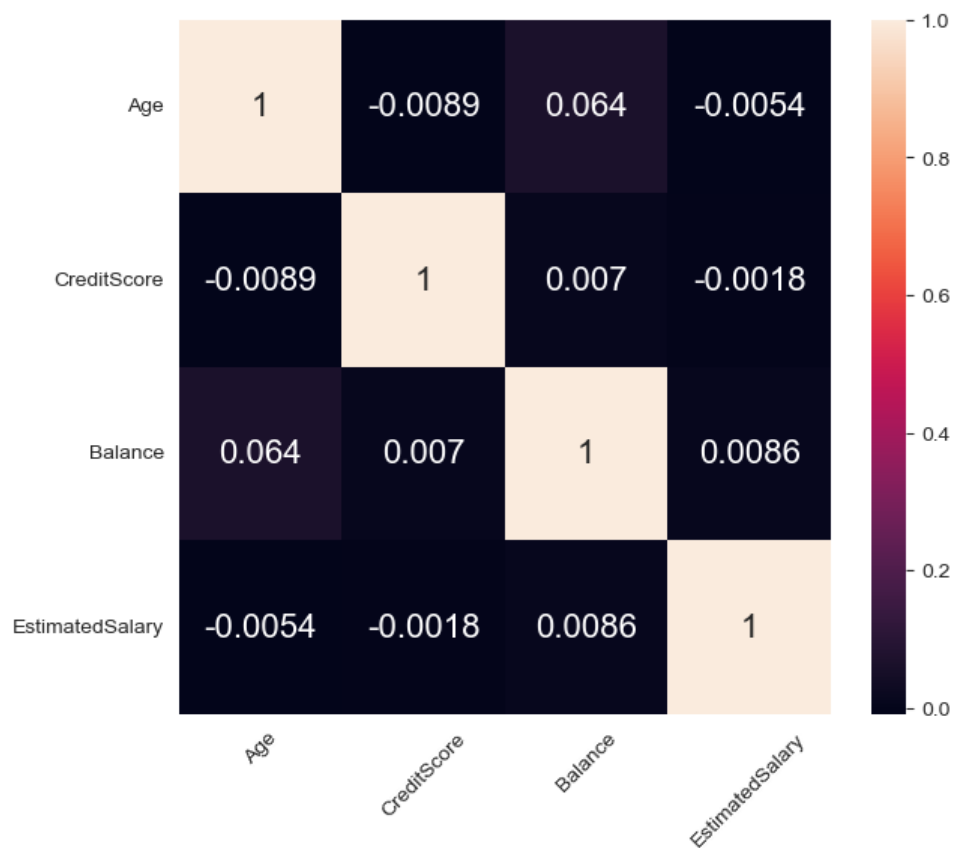
Name: EstimatedSalary, Length: 55298, dtype: int64

- Age: The most common age among customers is 37 years old, with the least common ages being 83, 84, and 85 years. This could indicate that older customers are more likely to churn.
- Credit Score: The most common credit score is 850, but from the descriptive statistics, we know that there are varying frequencies of credit scores, and some only occur once in the dataset.
- Balance: There are 89,684 accounts with 0 balance. This indicates that these are inactive customers, who may be more likely to churn.
- Estimated Salary: There's a wide range of estimated salaries occurring, which is to be expected.

Checking for multicollinearity between continuous variables

In [15]:

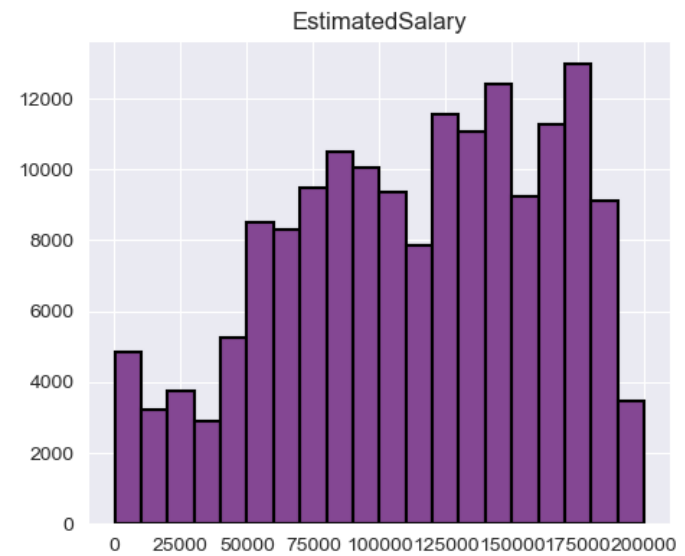
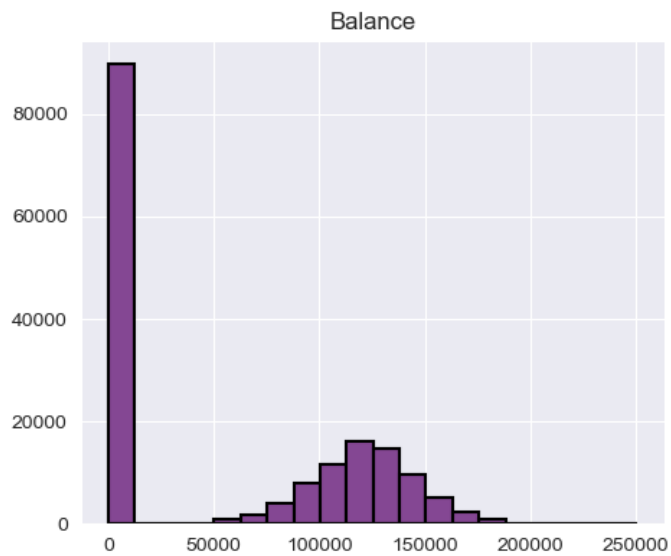
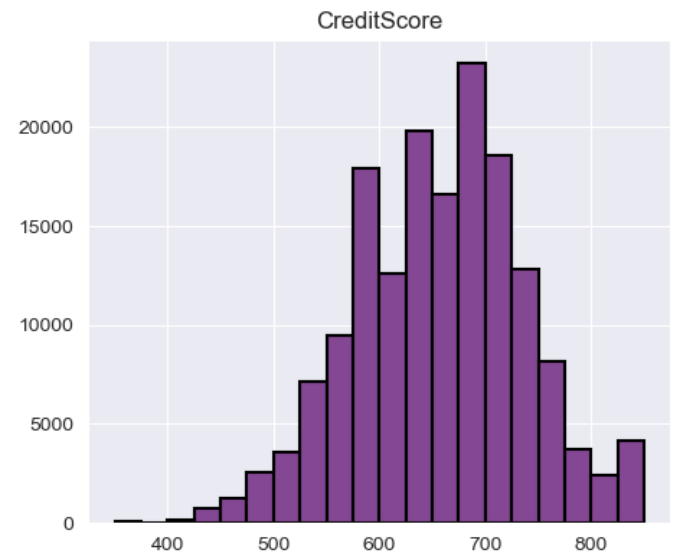
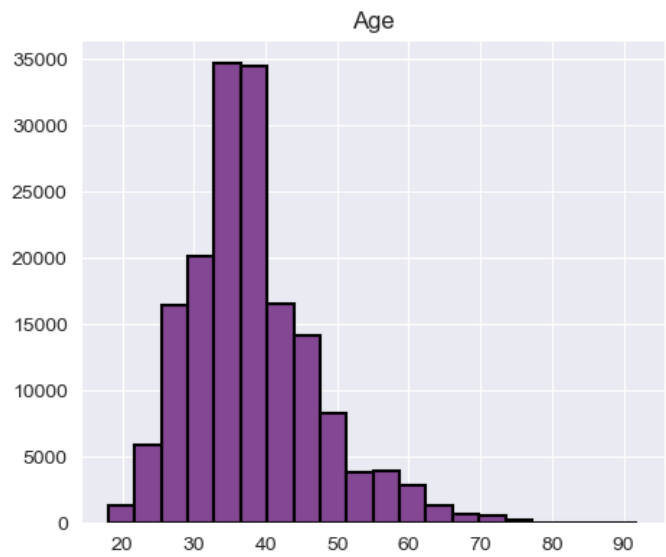
```
1 fig, ax = plt.subplots(figsize=(7, 6))
2
3 sns.heatmap(df[continuous].corr(),
4             annot=True,
5             annot_kws={'fontsize': 16},
6             ax=ax)
7
8 ax.tick_params(axis='x', rotation=45)
```



None of the continuous variables have a multicollinear relationship, so we can keep all these.

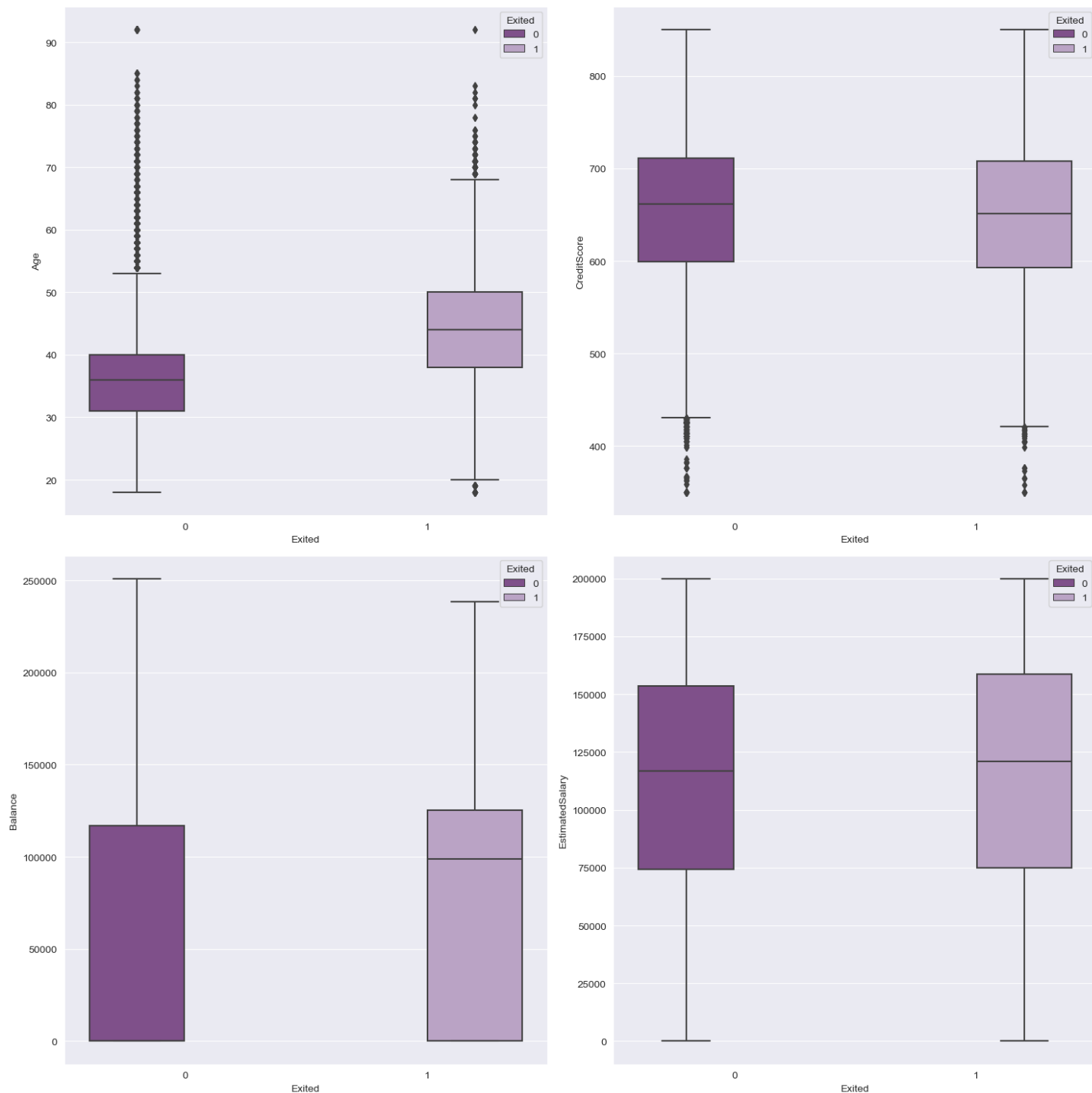
In [16]:

```
1 df[continuous].hist(figsize=(12, 10),  
2     bins=20,  
3     edgecolor='black',  
4     layout=(2, 2),  
5     color='purple')
```



In [17]:

```
1 fig, axes = plt.subplots(2, 2, figsize=(15, 15))
2 axes = axes.flatten()
3 for i, column in enumerate(continuous):
4     sns.boxplot(y=column, x='Exited', hue='Exited', data=df, ax=axes[i])
5 plt.tight_layout()
```



- Age: As inferred from the descriptive statistics, older customers tend to churn more than younger customers. Bank may not be as appealing to older customers- things like the loss of in-person banks and a lack of benefits, or handling of pensions may influence this.
- Credit Score: Although credit score is slightly negatively skewed, there is no significant difference in the credit score column between churned and retained customers.
- Balance: Customers with a higher bank balance are slightly more likely to churn. This suggests no products or policies in place rewarding those with higher bank balances i.e. higher savings interest rate.
- Estimated Salary: Estimated Salary has quite a varied deviation, and the boxplot indicates that those with a higher salary are more likely to churn. This suggests no products or policies in place rewarding those who will bring more earnings to the bank, and will likely have a higher lending potential.

Categorical Variables: Geography, Gender, Tenure, NumOfProducts, HasCrCard & IsActiveMember

In [18]:

```
1 fig, axes = plt.subplots(2, 2, figsize=(15, 15))
```

```
In [19]: 1 #Printing the value counts for the continuous columns
2 for column in categorical:
3     print("Value counts for", column, ":")
4     print(df[column].value_counts())
```

Value counts for Geography :

France	94215
Spain	36213
Germany	34606

Name: Geography, dtype: int64

Value counts for Gender :

Male	93150
Female	71884

Name: Gender, dtype: int64

Value counts for Tenure :

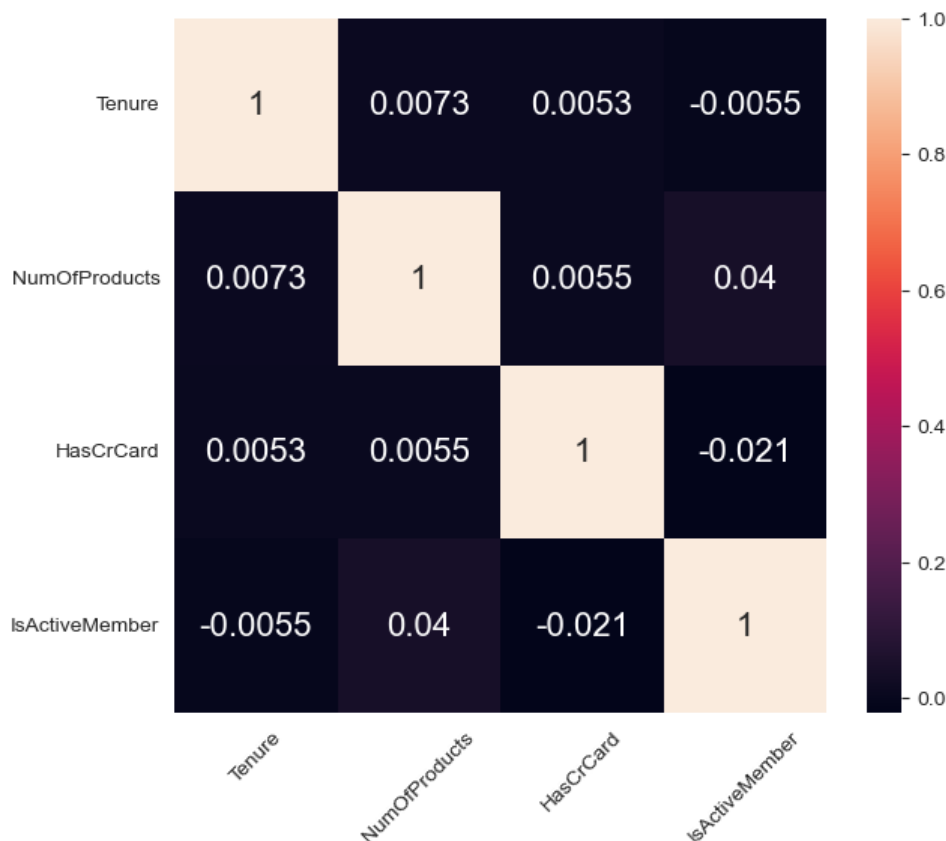
2	18045
7	17810
4	17554
8	17520
5	17268
1	16760
9	16709
3	16630

Checking for multicollinearity between categorical variables

```
In [20]: 1 fig, ax = plt.subplots(figsize=(7, 6))
2
3 sns.heatmap(df[categorical].corr(),
4             annot=True,
5             annot_kws={'fontsize': 16},
6             ax=ax)
7
8 ax.tick_params(axis='x', rotation=45)
```

C:\Users\eilid\AppData\Local\Temp\ipykernel_13236\155750028.py:3: FutureWarning: The default value of numeric_only in DataFrame.corr is deprecated. In a future version, it will default to False. Select only valid columns or specify the value of numeric_only to silence this warning.

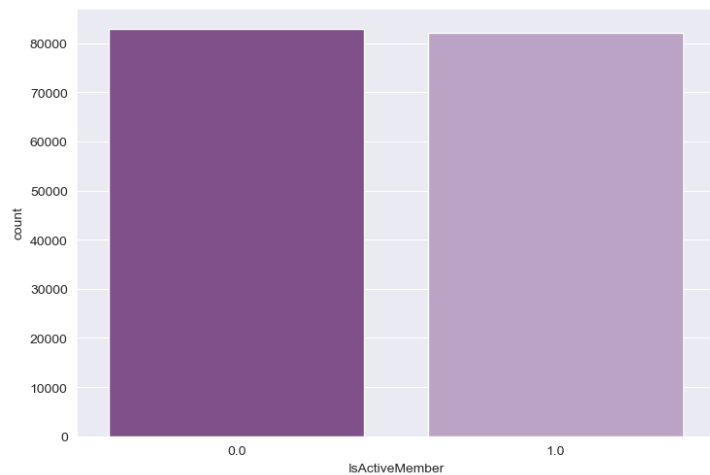
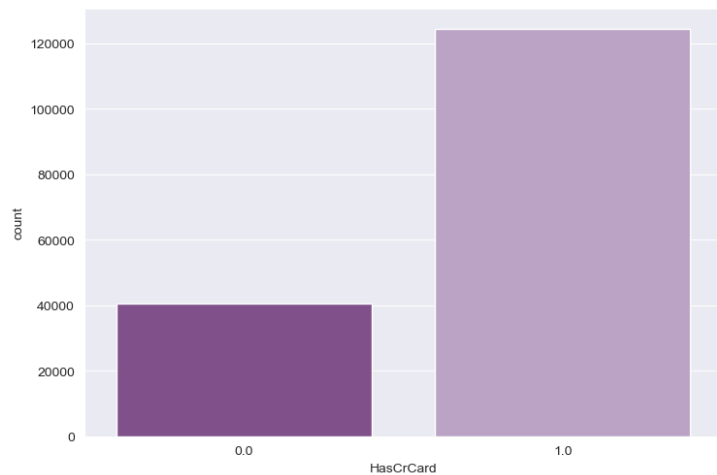
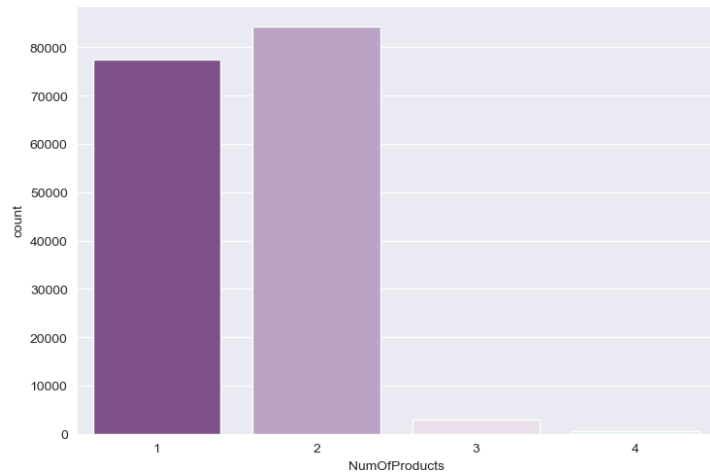
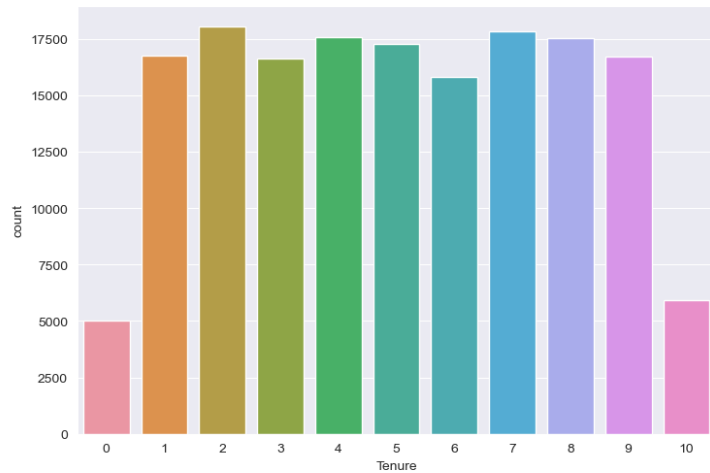
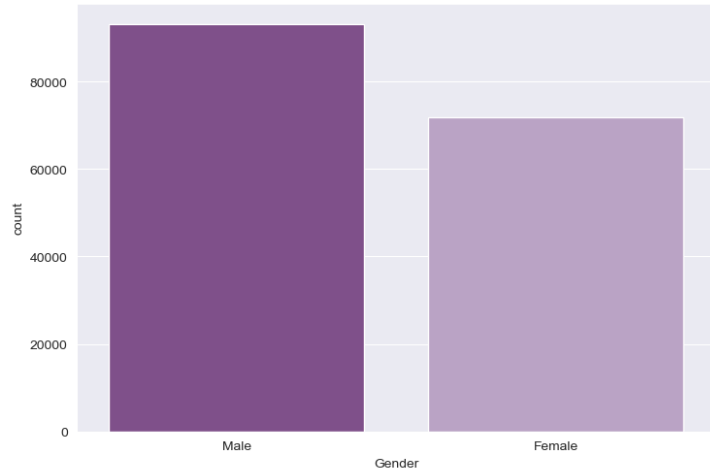
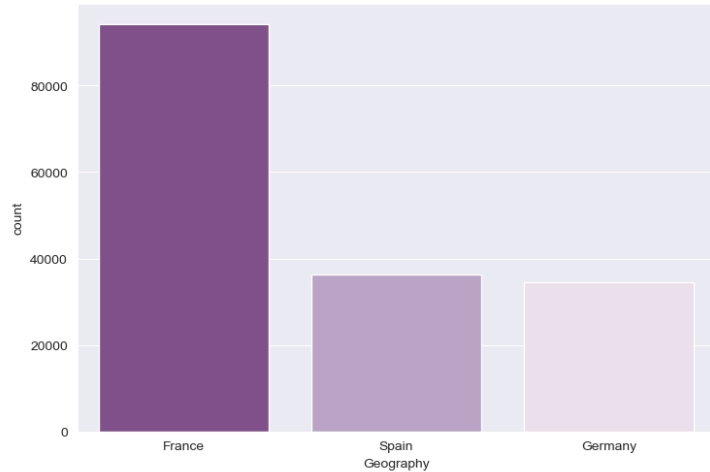
```
sns.heatmap(df[categorical].corr(),
```



- Tenure and IsActiveMember have a slightly negative correlation
- NumOfProducts and IsActiveMember have a slightly positive correlation
- HasCrCard and IsActiveMember have a slightly negative correlation
- Tenure and NumofProducts have a sliglytly positive correlation

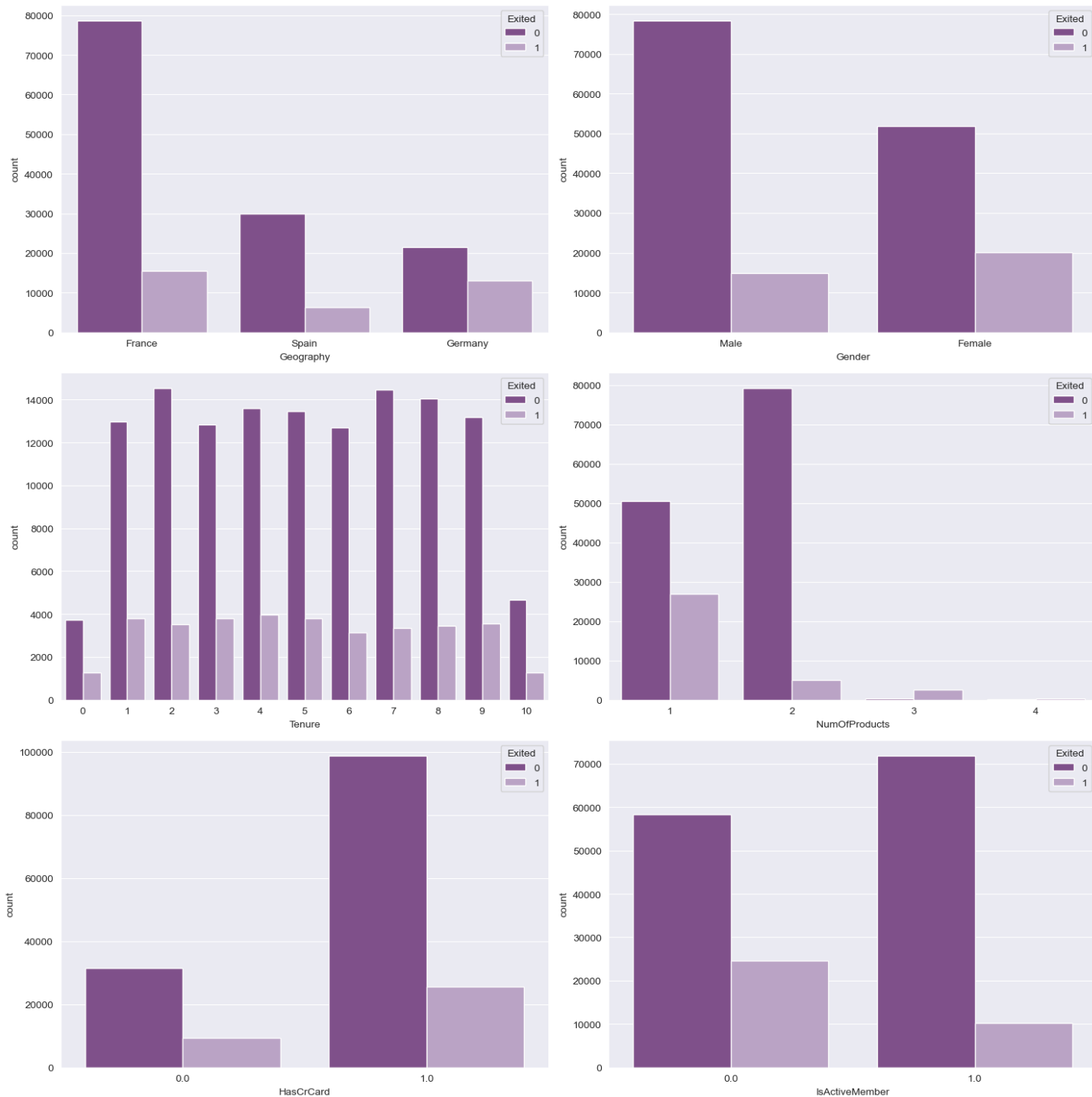
It's worth exploring the feature engineering of some of these variables

```
In [21]: 1 fig, axes = plt.subplots(3, 2, figsize=(15, 15))
2 axes = axes.flatten()
3 for i, column in enumerate(categorical):
4     sns.countplot(x=column, data=df, ax=axes[i])
5
6 plt.tight_layout()
```



In [22]:

```
1 # Set up the layout for subplots
2 fig, axes = plt.subplots(3, 2, figsize=(15, 15))
3 axes = axes.flatten()
4 for i, column in enumerate(categorical):
5     sns.countplot(x=column, hue='Exited', data=df, ax=axes[i])
6 plt.tight_layout()
```



Geography: The majority of customers reside in France, with Spain and Germany representing equal portions of customers. Customers in France are more likely to churn, but Germany has the highest ratio of churned to retained customers. This could indicate a lack of services or higher competition in Germany.

Gender: There are more male than female customers, and female customers are more likely to churn overall.

Tenure: Customers are least likely to churn in the 1st and 10th year, and churn rates are highest during the 2nd and 7th year. This could be due to a bank policy of higher rewards for 1st year holders i.e. cash rewards for switching bank accounts. However, tenure distribution does not show any uniform pattern for churn, and so tenure does not seem to affect the churn rate.

Number of Products (NumOfProducts): The majority of customers have 2 products, with those with 1 product being the most likely to churn. This follows a similar pattern for tenure in that customers may be taking advantage of first year benefits policies offered by the bank and then churning after the first year. Customers with 2 products are the second most likely to churn, but the rate is much lower than that of 1. Interestingly, having 3 products has the highest rate of churning.

Has Credit Card (HasCrCard): The majority of customers have a credit card. Having a credit card does not appear to affect the churn rate.

Is Active Member (IsActiveMember): Roughly half the customers are inactive, and inactive customers are more likely to churn.

Data Preprocessing

Outliers

```
In [23]: 1 #Function to define the outlier thresholds
2 def outlier_thresholds(df, variable, low_quantile=0.05, up_quantile=0.95):
3     quantile_one = df[variable].quantile(low_quantile)
4     quantile_three = df[variable].quantile(up_quantile)
5     interquantile_range = quantile_three - quantile_one
6     up_limit = quantile_three + 1.5 * interquantile_range
7     low_limit = quantile_one - 1.5 * interquantile_range
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65
66
67
68
69
70
71
72
73
74
75
76
77
78
79
80
81
82
83
84
85
86
87
88
89
90
91
92
93
94
95
96
97
98
99
100
```

```
In [24]: 1 # Function to check for outliers based on the above thresholds
2 def has_outliers(df, continuous, plot=False):
3     for col in continuous:
4         low_limit, up_limit = outlier_thresholds(df, col)
5         if df[(df[col] > up_limit) | (df[col] < low_limit)].any(axis=None):
6             number_of_outliers = df[(df[col] > up_limit) | (df[col] < low_limit)].shape[0]
7             print(col, " : ", number_of_outliers, "outliers")
8             #variable_names.append(col)
9             if plot:
10                 sns.boxplot(x=dataframe[col])
11                 plt.show()
12 #Return outliers
13 for var in continuous:
14     print(var, " : ", outlier_thresholds(df, var, 0.05, 0.95))
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65
66
67
68
69
70
71
72
73
74
75
76
77
78
79
80
81
82
83
84
85
86
87
88
89
90
91
92
93
94
95
96
97
98
99
100
```

```
Age has  None Outliers
CreditScore has  None Outliers
Balance has  None Outliers
EstimatedSalary has  None Outliers
```

Feature Selection

The EDA has revealed more features than can be dropped as they do not affect customer churn. As a further measure, all the continous variables undergo both Pearson r and the Spearman p tests, and categorical variables undergo a chi-sqaure test.

Pearson r and the Spearman p for continous variables

```
In [25]: 1 corr_array, p_array = [], []
2 for column in continuous:
3     pearson_corr, pearson_p_value = pearsonr(df[column], df['Exited']) #PearsonR
4     spearman_corr, spearman_p_value = spearmanr(df[column], df['Exited']) #SpearmanP
5     corr_array.append((pearson_corr, spearman_corr))
6     p_array.append((pearson_p_value, spearman_p_value)) # Append results to arrays for a df
7
8 df_corr = pd.DataFrame({
9     'Variable': continuous,
10    'Pearson_corr': [corr[0] for corr in corr_array],
11    'Spearman_corr': [corr[1] for corr in corr_array],
12    'Pearson_p-value': [p_value[0] for p_value in p_array],
13    'Spearman_p-value': [p_value[1] for p_value in p_array]
14 })
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65
66
67
68
69
70
71
72
73
74
75
76
77
78
79
80
81
82
83
84
85
86
87
88
89
90
91
92
93
94
95
96
97
98
99
100
```

```
Out[25]:
```

	Variable	Pearson_corr	Spearman_corr	Pearson_p-value	Spearman_p-value
0	Age	0.340768	0.354665	0.000000e+00	0.000000e+00
1	CreditScore	-0.027383	-0.029120	9.373557e-29	2.663270e-32
2	Balance	0.129743	0.126417	0.000000e+00	0.000000e+00
3	EstimatedSalary	0.018827	0.019997	2.027177e-14	4.494482e-16

As predicted from the initial EDA, 'EstimatedSalary' and 'CreditScore' both scored very close to zero in both Pearson and Spearman correlation efficents inducting they have almost no linear relationship and very little impact on the target variable.

```
In [26]: 1 features_drop = ['EstimatedSalary', 'CreditScore']
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65
66
67
68
69
70
71
72
73
74
75
76
77
78
79
80
81
82
83
84
85
86
87
88
89
90
91
92
93
94
95
96
97
98
99
100
```

Chi-square test for categorical variables

```
In [27]: 1 chi2_array, p_array = [], []
2         for column in categorical:
3
4             crosstab = pd.crosstab(df[column], df['Exited'])
5             chi2, p, dof, expected = chi2_contingency(crosstab)
6             chi2_array.append(chi2) #Create array for df
7             p_array.append(p)
8
9         df_chi = pd.DataFrame({
10             'Variable': categorical,
11             'Chi-square': chi2_array,
12             'p-value': p_array
13         })
```

```
Out[27]:
```

	Variable	Chi-square	p-value
3	NumOfProducts	29130.479017	0.000000e+00
0	Geography	7358.673765	0.000000e+00
5	IsActiveMember	7293.408511	0.000000e+00
1	Gender	3538.452550	0.000000e+00
2	Tenure	265.458532	3.028684e-51
4	HasCrCard	80.780230	2.522693e-19

Similarly as predicted from the initial EDA, 'Tenure' and 'HasCrCard' have a small chi-square and a p-value greater than 0.05, therefore, we can safely drop these features from the dataset.

```
In [28]: 1 features_drop = ['Tenure', 'HasCrCard']
```

```
In [29]: 1 df = df.drop(features_drop, axis=1)
```

```
Out[29]:
```

	Geography	Gender	Age	Balance	NumOfProducts	IsActiveMember	Exited
0	France	Male	33.0	0.00	2	0.0	0
1	France	Male	33.0	0.00	2	1.0	0
2	France	Male	40.0	0.00	2	0.0	0
3	France	Male	34.0	148882.54	1	1.0	0
4	Spain	Male	33.0	0.00	2	1.0	0

Encoding Categorical Features

```
In [30]: 1 # Label encoder selected over one-hot encoding. Note, did try both, and they had no impact on model performance.
2         # Label was kept as no high dimensionality.
3         df['Gender'] = LabelEncoder().fit_transform(df['Gender'])
```

```
In [31]: 1 df = df.drop(features_drop, axis=1)
```

```
Out[31]:
```

	Geography	Gender	Age	Balance	NumOfProducts	IsActiveMember	Exited
0	0	1	33.0	0.00	2	0.0	0
1	0	1	33.0	0.00	2	1.0	0
2	0	1	40.0	0.00	2	0.0	0
3	0	1	34.0	148882.54	1	1.0	0
4	2	1	33.0	0.00	2	1.0	0

Scaling

'Age' ranges from quite small integers (18-92) whereas 'Balance' has much larger integers (0-250,898). The difference between these integers will cause the classification models to misinterpret the data, and so will need to be standard scaled.

```
In [32]: 1 # Standard scaling was used as the 'Age' and 'Balance' variables were on vastly different scales.
2         # It also does not affect the interpretability of the data.
3         scaler = StandardScaler()
4
5         scl_columns = ['Age', 'Balance']
6         df[scl_columns] = scaler.fit_transform(df[scl_columns])
```

Creating 'X_train' & 'Y_train'

```
In [33]: 1 y = df['Exited']
          2 X = df.drop('Exited', 1)

C:\Users\eilid\AppData\Local\Temp\ipykernel_13236\2998808209.py:2: FutureWarning: In a future version of pandas all argument
s of DataFrame.drop except for the argument 'labels' will be keyword-only.
  X = df.drop('Exited', 1)
```

```
In [34]: (165034, 6) (165034,)
```

```
In [35]: 1 X_train, X_test, y_train, y_test = train_test_split(X, y,
          2                                                    test_size=0.20,
```

Using SMOTE to rebalance 'Exited' (Y_train)

```
In [36]: Out[36]: 0    104061
          1    27966
          Name: Exited, dtype: int64
```

```
In [37]: 1 over = SMOTE(sampling_strategy='auto', random_state=42)
          2 X_train, y_train = over.fit_resample(X_train, y_train)
          3
```

```
Out[37]: 0    104061
          1    104061
          Name: Exited, dtype: int64
```

```
In [38]: (208122, 6) (33007, 6) (208122,) (33007,)
```

Modelling

Baseline Models

3 classifiers were selected as the initial baseline; Random Forest, Descion Tree, and Gradient Boosting.

The use of SVM was also explored, but due to time and computational constraints, this could not complete.

In [39]:

```
1 ##### Random Forest Classifier #####
2 random_forest_model = RandomForestClassifier()
3 random_forest_model.fit(X_train, y_train)
4
5 random_forest_predictions = random_forest_model.predict(X_test)
6 random_forest_accuracy = accuracy_score(y_test, random_forest_predictions)
7 random_forest_precision = precision_score(y_test, random_forest_predictions, average='weighted')
8 random_forest_recall = recall_score(y_test, random_forest_predictions, average='weighted')
9 random_forest_f1 = f1_score(y_test, random_forest_predictions, average='weighted')
10
11 random_forest_feature_importance = random_forest_model.feature_importances_
12
13 print("Random Forest Classifier Metrics:")
14 print(f"Accuracy: {random_forest_accuracy:.2f}")
15 print(f"Precision: {random_forest_precision:.2f}")
16 print(f"Recall: {random_forest_recall:.2f}")
17 print(f"F1-score: {random_forest_f1:.2f}")
18 print()
19
20 ##### Decision Tree Classifier #####
21 decision_tree_model = DecisionTreeClassifier()
22 decision_tree_model.fit(X_train, y_train)
23
24 decision_tree_predictions = decision_tree_model.predict(X_test)
25 decision_tree_accuracy = accuracy_score(y_test, decision_tree_predictions)
26 decision_tree_precision = precision_score(y_test, decision_tree_predictions, average='weighted')
27 decision_tree_recall = recall_score(y_test, decision_tree_predictions, average='weighted')
28 decision_tree_f1 = f1_score(y_test, decision_tree_predictions, average='weighted')
29
30 decision_tree_feature_importance = decision_tree_model.feature_importances_
31
32 print("Decision Tree Classifier Metrics:")
33 print(f"Accuracy: {decision_tree_accuracy:.2f}")
34 print(f"Precision: {decision_tree_precision:.2f}")
35 print(f"Recall: {decision_tree_recall:.2f}")
36 print(f"F1-score: {decision_tree_f1:.2f}")
37 print()
38
39 ##### Gradient Boosting #####
40 gradient_boosting_model = GradientBoostingClassifier()
41 gradient_boosting_model.fit(X_train, y_train)
42
43 gradient_boosting_predictions = gradient_boosting_model.predict(X_test)
44 gradient_boosting_accuracy = accuracy_score(y_test, gradient_boosting_predictions)
45 gradient_boosting_precision = precision_score(y_test, gradient_boosting_predictions, average='weighted')
46 gradient_boosting_recall = recall_score(y_test, gradient_boosting_predictions, average='weighted')
47 gradient_boosting_f1 = f1_score(y_test, gradient_boosting_predictions, average='weighted')
48
49 gradient_boosting_feature_importance = gradient_boosting_model.feature_importances_
50
51 print("Gradient Boosting Metrics:")
52 print(f"Accuracy: {gradient_boosting_accuracy:.2f}")
53 print(f"Precision: {gradient_boosting_precision:.2f}")
54 print(f"Recall: {gradient_boosting_recall:.2f}")
55 print(f"F1-score: {gradient_boosting_f1:.2f}")
56
```

Random Forest Classifier Metrics:

Accuracy: 0.80
Precision: 0.83
Recall: 0.80
F1-score: 0.81

Decision Tree Classifier Metrics:

Accuracy: 0.79
Precision: 0.82
Recall: 0.79
F1-score: 0.80

Gradient Boosting Metrics:

Accuracy: 0.81
Precision: 0.85
Recall: 0.81
F1-score: 0.82

Classification Report for Baseline Models


```
In [51]: 1 models = [  
2         ('Random Forest Classifier', random_forest_predictions),  
3         ('Decision Tree Classifier', decision_tree_predictions),  
4         ('Gradient Boosting', gradient_boosting_predictions),  
5     ]  
6  
7  
8 for name, predictions in models:  
9     print(f"Classification Report for {name}:")  
10    print(classification_report(y_test, predictions))
```

Classification Report for Random Forest Classifier:

	precision	recall	f1-score	support
0	0.91	0.83	0.87	26052
1	0.52	0.70	0.60	6955
accuracy			0.80	33007
macro avg	0.72	0.76	0.73	33007
weighted avg	0.83	0.80	0.81	33007

Classification Report for Decision Tree Classifier:

	precision	recall	f1-score	support
0	0.91	0.82	0.86	26052
1	0.50	0.69	0.58	6955
accuracy			0.79	33007
macro avg	0.71	0.75	0.72	33007
weighted avg	0.82	0.79	0.80	33007

Classification Report for Gradient Boosting:

	precision	recall	f1-score	support
0	0.94	0.81	0.87	26052
1	0.53	0.80	0.64	6955
accuracy			0.81	33007
macro avg	0.73	0.80	0.75	33007
weighted avg	0.85	0.81	0.82	33007

Confusion Matrix for Baseline Models

```
In [54]: 1 models = [  
2         ('Random Forest Classifier', random_forest_predictions),  
3         ('Gradient Boosting', gradient_boosting_predictions),  
4         ('Decision Tree Classifier', decision_tree_predictions),  
5     ]  
6  
7 for name, predictions in models:  
8     cm = confusion_matrix(y_test, predictions)  
9     print(f"Confusion Matrix for {name}:")  
10    print(cm)
```

Confusion Matrix for Random Forest Classifier:

[[21550 4502]
[2093 4862]]

Confusion Matrix for Gradient Boosting:

[[21075 4977]
[1393 5562]]

Confusion Matrix for Decision Tree Classifier:

[[21362 4690]
[2173 4782]]

Based on the initial baseline metrics, the Gradient Boosting classifier performs better than the Random Forest and Decision Tree classifiers. In addition, it has the highest precision (53%) and recall (80%), indicating fewer false positives and correctly identifying a larger proportion of actual churned customers.

Because of this, we will now use optimised versions of gradient boosting algorithms (XGBoost, LightGBM, and CatBoost).

XGBoost, LightGBM, and CatBoost Classifiers

In [40]:

```
1 ##### XGBoost Classifier #####
2 xgb_model = xgb.XGBClassifier()
3 xgb_model.fit(X_train, y_train)
4
5 xgb_predictions = xgb_model.predict(X_test)
6 xgb_accuracy = accuracy_score(y_test, xgb_predictions)
7 xgb_precision = precision_score(y_test, xgb_predictions, average='weighted')
8 xgb_recall = recall_score(y_test, xgb_predictions, average='weighted')
9 xgb_f1 = f1_score(y_test, xgb_predictions, average='weighted')
10
11 xgb_feature_importance = xgb_model.feature_importances_
12
13 print("XGBoost Classifier Metrics:")
14 print(f"Accuracy: {xgb_accuracy:.2f}")
15 print(f"Precision: {xgb_precision:.2f}")
16 print(f"Recall: {xgb_recall:.2f}")
17 print(f"F1-score: {xgb_f1:.2f}")
18 print()
19
20 ##### LightGBM Classifier #####
21 lgb_model = lgb.LGBMClassifier(verbose=-1)
22 lgb_model.fit(X_train, y_train)
23
24 lgb_predictions = lgb_model.predict(X_test)
25 lgb_accuracy = accuracy_score(y_test, lgb_predictions)
26 lgb_precision = precision_score(y_test, lgb_predictions, average='weighted')
27 lgb_recall = recall_score(y_test, lgb_predictions, average='weighted')
28 lgb_f1 = f1_score(y_test, lgb_predictions, average='weighted')
29
30 lgb_feature_importance = lgb_model.feature_importances_
31
32 print("LightGBM Classifier Metrics:")
33 print(f"Accuracy: {lgb_accuracy:.2f}")
34 print(f"Precision: {lgb_precision:.2f}")
35 print(f"Recall: {lgb_recall:.2f}")
36 print(f"F1-score: {lgb_f1:.2f}")
37 print()
38
39 ##### CatBoost Classifier #####
40 cb_model = cb.CatBoostClassifier(verbose=False)
41 cb_model.fit(X_train, y_train)
42
43 cb_predictions = cb_model.predict(X_test)
44 cb_accuracy = accuracy_score(y_test, cb_predictions)
45 cb_precision = precision_score(y_test, cb_predictions, average='weighted')
46 cb_recall = recall_score(y_test, cb_predictions, average='weighted')
47 cb_f1 = f1_score(y_test, cb_predictions, average='weighted')
48
49 cb_feature_importance = cb_model.feature_importances_
50
51 print("CatBoost Classifier Metrics:")
52 print(f"Accuracy: {cb_accuracy:.2f}")
53 print(f"Precision: {cb_precision:.2f}")
54 print(f"Recall: {cb_recall:.2f}")
55 print(f"F1-score: {cb_f1:.2f}")
56 print()
57
```

XGBoost Classifier Metrics:

Accuracy: 0.82
Precision: 0.85
Recall: 0.82
F1-score: 0.83

LightGBM Classifier Metrics:

Accuracy: 0.82
Precision: 0.85
Recall: 0.82
F1-score: 0.83

CatBoost Classifier Metrics:

Accuracy: 0.82
Precision: 0.85
Recall: 0.82
F1-score: 0.83

Overall, XGBoost, LightGBM, and CatBoost are only slightly better performing than the baseline models. There is a slight improvement in predicting customer churn.

Confusion Matrix for XGBoost, LightGBM, and CatBoost

```
In [52]: 1 models = [  
2         ('XGBoost Classifier', xgb_predictions),  
3         ('LightGBM Classifier', lgb_predictions),  
4         ('CatBoost Classifier', cb_predictions)  
5     ]  
6  
7     for name, predictions in models:  
8         cm = confusion_matrix(y_test, predictions)  
9         print(f"Confusion Matrix for {name}:")  
10        print(cm)
```

Confusion Matrix for XGBoost Classifier:

```
[[21817  4235]  
 [ 1664  5291]]
```

Confusion Matrix for LightGBM Classifier:

```
[[21645  4407]  
 [ 1558  5397]]
```

Confusion Matrix for CatBoost Classifier:

```
[[21967  4085]  
 [ 1704  5251]]
```

Classification Report for XGBoost, LightGBM, and CatBoost

```
In [53]: 1 models = [  
2         ('XGBoost Classifier', xgb_predictions),  
3         ('LightGBM Classifier', lgb_predictions),  
4         ('CatBoost Classifier', cb_predictions)  
5     ]  
6  
7     for name, predictions in models:  
8         print(f"Classification Report for {name}:")  
9         print(classification_report(y_test, predictions))
```

Classification Report for XGBoost Classifier:

	precision	recall	f1-score	support
0	0.93	0.84	0.88	26052
1	0.56	0.76	0.64	6955
accuracy			0.82	33007
macro avg	0.74	0.80	0.76	33007
weighted avg	0.85	0.82	0.83	33007

Classification Report for LightGBM Classifier:

	precision	recall	f1-score	support
0	0.93	0.83	0.88	26052
1	0.55	0.78	0.64	6955
accuracy			0.82	33007
macro avg	0.74	0.80	0.76	33007
weighted avg	0.85	0.82	0.83	33007

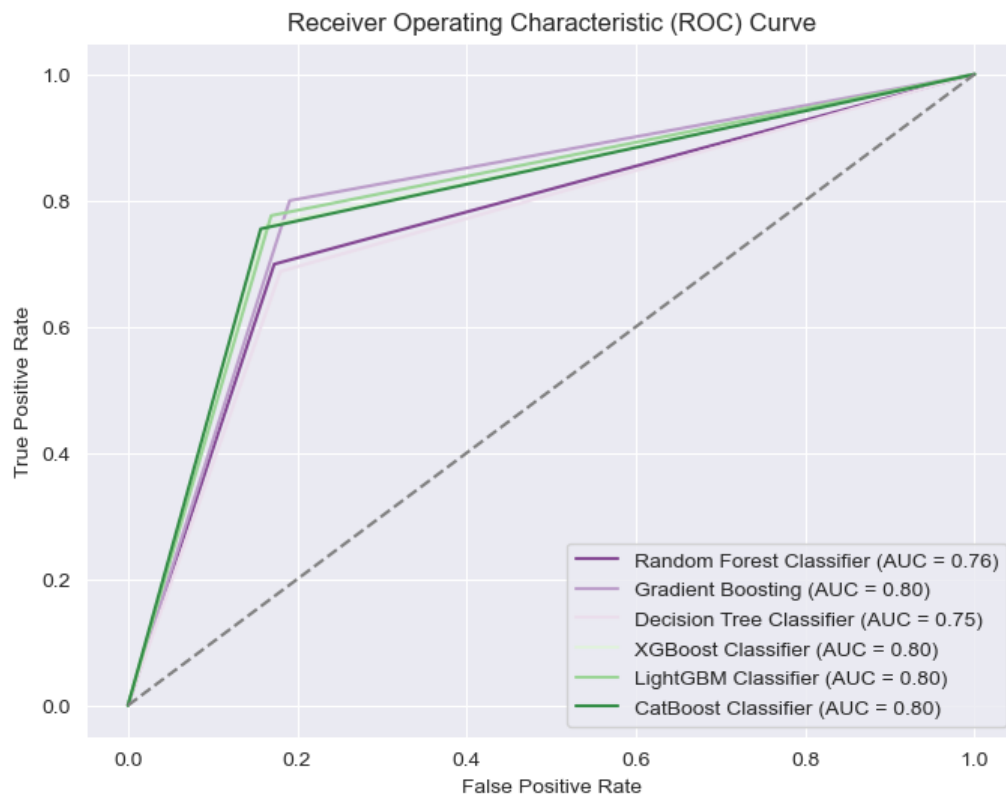
Classification Report for CatBoost Classifier:

	precision	recall	f1-score	support
0	0.93	0.84	0.88	26052
1	0.56	0.75	0.64	6955
accuracy			0.82	33007
macro avg	0.75	0.80	0.76	33007
weighted avg	0.85	0.82	0.83	33007

```

In [43]: 1 # Calculate ROC curve and AUC score for each model
2 models = [
3     ('Random Forest Classifier', random_forest_predictions),
4     ('Gradient Boosting', gradient_boosting_predictions),
5     ('Decision Tree Classifier', decision_tree_predictions),
6     ('XGBoost Classifier', xgb_predictions),
7     ('LightGBM Classifier', lgb_predictions),
8     ('CatBoost Classifier', cb_predictions)
9 ]
10
11 plt.figure(figsize=(8, 6))
12
13 for name, predictions in models:
14     fpr, tpr, _ = roc_curve(y_test, predictions)
15     auc_score = roc_auc_score(y_test, predictions)
16     plt.plot(fpr, tpr, label=f'{name} (AUC = {auc_score:.2f})')
17
18 plt.plot([0, 1], [0, 1], linestyle='--', color='gray')
19 plt.xlabel('False Positive Rate')
20 plt.ylabel('True Positive Rate')
21 plt.title('Receiver Operating Characteristic (ROC) Curve')
22 plt.legend(loc='lower right')
23 plt.grid(True)
24 plt.show()

```



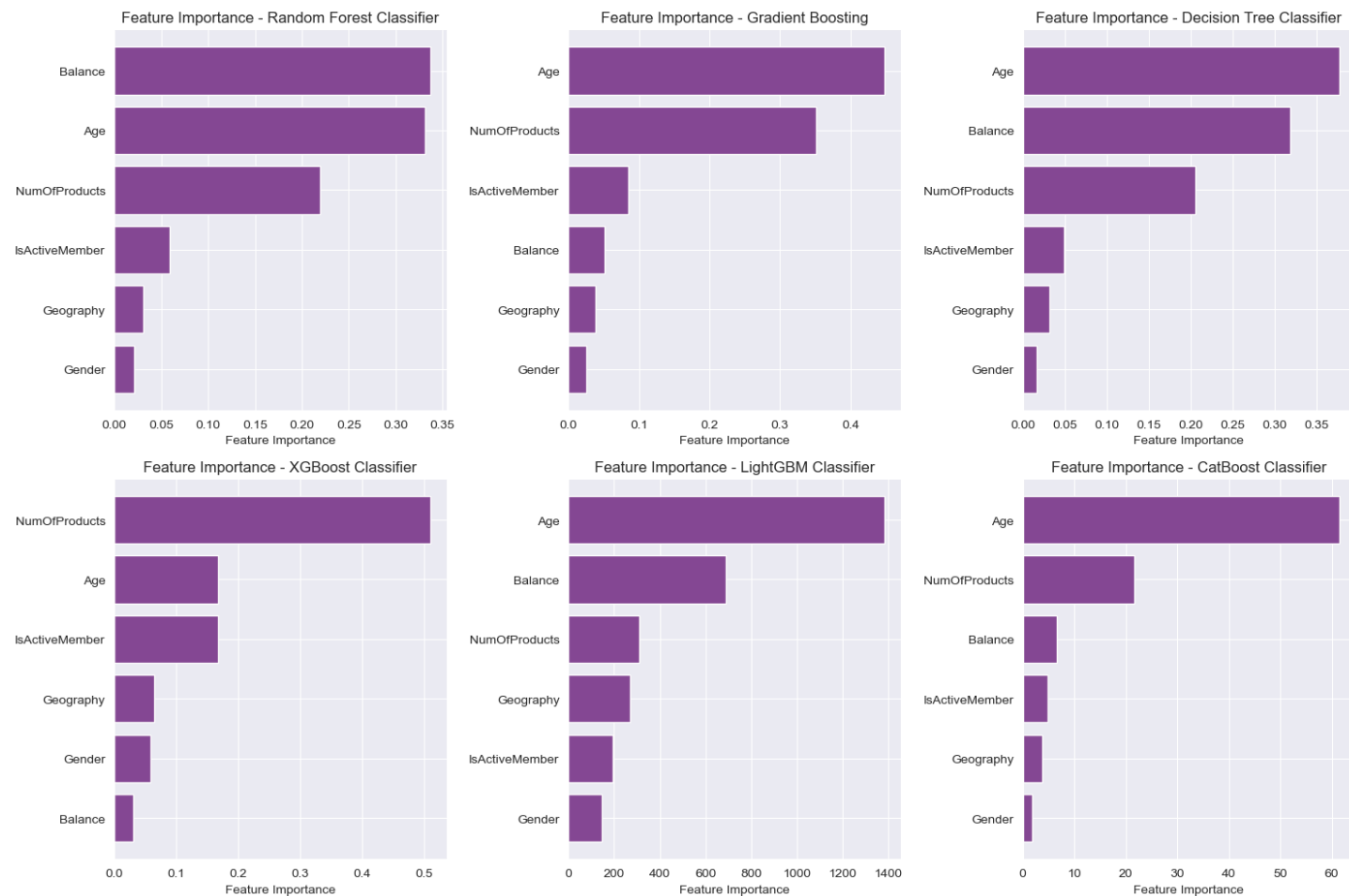
Feature Importance

In [44]:

```

1 feature_names = X_train.columns.tolist()
2
3 # Define a function to plot feature importance
4 def plot_feature_importance(feature_importance, feature_names, model_name, ax):
5     sorted_idx = np.argsort(feature_importance)
6     ax.barh(range(len(feature_importance)), feature_importance[sorted_idx], align='center')
7     ax.set_yticks(range(len(feature_importance)))
8     ax.set_yticklabels([feature_names[i] for i in sorted_idx])
9     ax.set_xlabel('Feature Importance')
10    ax.set_title(f'Feature Importance - {model_name}')
11
12 # Assuming you have trained models and extracted feature importance already
13 # and feature_names extracted from X_train
14
15 # Example feature importance for each model
16 feature_importance_data = [
17     ('Random Forest Classifier', random_forest_feature_importance),
18     ('Gradient Boosting', gradient_boosting_feature_importance),
19     ('Decision Tree Classifier', decision_tree_feature_importance),
20     ('XGBoost Classifier', xgb_feature_importance),
21     ('LightGBM Classifier', lgb_feature_importance),
22     ('CatBoost Classifier', cb_feature_importance)
23 ]
24
25 fig, axs = plt.subplots(2, 3, figsize=(15, 10))
26
27 for (model_name, feature_importance), ax in zip(feature_importance_data, axs.flatten()):
28     plot_feature_importance(feature_importance, feature_names, model_name, ax)
29
30 plt.tight_layout()
31 plt.show()

```



Model Tuning using GridSearchCV

Note: This was not used in the final tuning model, but will be used in future versions.

```
In [46]: 1 # Define parameter grids for each model
2 xgb_param_grid = {
3     'n_estimators': [100, 200],
4     'learning_rate': [0.1, 0.05],
5     'max_depth': [3, 4],
6 }
7
8 lgb_param_grid = {
9     'n_estimators': [100, 200],
10    'learning_rate': [0.1, 0.05],
11    'max_depth': [3, 4],
12 }
13
14 cb_param_grid = {
15     'iterations': [100, 200],
16     'learning_rate': [0.1, 0.05],
17     'depth': [3, 4],
18 }
19
20 models = [
21     ('XGBoost Classifier', xgb.XGBClassifier(), xgb_param_grid),
22     ('LightGBM Classifier', lgb.LGBMClassifier(), lgb_param_grid),
23     ('CatBoost Classifier', cb.CatBoostClassifier(verbose=False), cb_param_grid)
24 ]
25
26 for name, model, param_grid in models:
27     grid_search = GridSearchCV(model, param_grid, cv=3, scoring='accuracy', n_jobs=-1)
28     grid_search.fit(X_train[:len(X_train)//10], y_train[:len(y_train)//10]) # Using 1/10th of the data
29
30     best_estimator = grid_search.best_estimator_
31     best_params = grid_search.best_params_
32
33     refined_param_grid = {key: [best_params[key]] for key in best_params}
34     refined_grid_search = GridSearchCV(best_estimator, refined_param_grid, cv=3, scoring='accuracy', n_jobs=-1)
35     refined_grid_search.fit(X_train, y_train)
36
37     print(f"Best parameters for {name}: {refined_grid_search.best_params_}")
38     print(f"Best score for {name}: {refined_grid_search.best_score_:.2f}")
39     print()
40
41     val_score = refined_grid_search.best_estimator_.score(X_val, y_val)
42     print(f"Validation score for {name}: {val_score:.2f}")
43     print()
44 ..
```

```
Best parameters for XGBoost Classifier: {'learning_rate': 0.1, 'max_depth': 4, 'n_estimators': 200}
Best score for XGBoost Classifier: 0.83
```

```
Validation score for XGBoost Classifier: 0.82
```

```
[LightGBM] [Info] Number of positive: 8333, number of negative: 8316
[LightGBM] [Info] Auto-choosing row-wise multi-threading, the overhead of testing was 0.000205 seconds.
You can set `force_row_wise=true` to remove the overhead.
And if memory is not enough, you can set `force_col_wise=true`.
[LightGBM] [Info] Total Bins 518
[LightGBM] [Info] Number of data points in the train set: 16649, number of used features: 6
[LightGBM] [Info] [binary:BoostFromScore]: pavg=0.500511 -> initscore=0.002042
[LightGBM] [Info] Start training from score 0.002042
[LightGBM] [Warning] No further splits with positive gain, best gain: -inf
[LightGBM] [Warning] No further splits with positive gain, best gain: -inf
[LightGBM] [Warning] No further splits with positive gain, best gain: -inf
[LightGBM] [Warning] No further splits with positive gain, best gain: -inf
[LightGBM] [Warning] No further splits with positive gain, best gain: -inf
[LightGBM] [Warning] No further splits with positive gain, best gain: -inf
[LightGBM] [Warning] No further splits with positive gain, best gain: -inf
[LightGBM] [Warning] No further splits with positive gain, best gain: -inf
[LightGBM] [Warning] No further splits with positive gain, best gain: -inf
[LightGBM] [Warning] No further splits with positive gain, best gain: -inf
```

Conclusions

- All three classifiers have similar performance in predicting churned customers, with precision, recall, F1-score, and accuracy hovering around 55-56%, 75-78%, 64%, and 82% respectively.
- There is no significant difference in performance between XGBoost, LightGBM, and CatBoost based on these metrics.

Future work

- Feature Engineering - exploring the relationship of age and tenure, combining credit score, salary, and balance, or just salary and balance
- Hyperparameter Tuning - The GridSearchCv is very rudimentary. Would like to expand the search for parameters and use better fine-tuned hyper parameters.
- Also realised that there was a mistake made with the Pearson and Spearman correlation tests, as these measure the relationship to other variables, not the targeted variable. Future versions will fix this and revise the model engineering if needed.

In []: