# Bank Customer Churn Prediction

## Introduction

Customer churn is a critical issue for banks and financial institutions. It refers to the phenomenon where customers stop using a bank's services, leading to a loss of revenue and potentially harming the bank's reputation. Predicting customer churn is essential for banks to take proactive measures to retain customers and improve their services.

In this project, we aim to build a predictive model to identify customers who are likely to churn. By analyzing various features such as customer demographics, account information, and transaction history, we can gain insights into the factors that contribute to churn and develop strategies to mitigate it.

The project involves the following steps:
1. **Data Collection**: Gather data on bank customers, including their demographics, account details, and transaction history.
2. **Data Preprocessing**: Clean and preprocess the data to handle missing values, outliers, and categorical variables.
3. **Exploratory Data Analysis (EDA)**: Perform EDA to understand the distribution of data, identify patterns, and visualize relationships between features.
4. **Feature Engineering**: Create new features or transform existing ones to improve the predictive power of the model.
5. **Model Building**: Train various machine learning models to predict customer churn and evaluate their performance.
6. **Model Evaluation**: Assess the models using appropriate metrics and select the best-performing model.
7. **Deployment**: Deploy the model to a production environment where it can be used to make  predictions.

8.**Dashboard:** Creation of dashboard visual with matplotlib to give more insights on the prediction made by the model

## Problem Statement

Customer churn is a significant challenge for banks and financial institutions. Churn occurs when customers stop using a bank's services, leading to a loss of revenue and potentially damaging the bank's reputation. Understanding and predicting customer churn is crucial for banks to take proactive measures to retain customers and enhance their services.

The objective of this project is to develop a predictive model that can identify customers who are likely to churn. By analyzing various features such as customer demographics, account information, and transaction history, we aim to gain insights into the factors contributing to churn and develop strategies to mitigate it.

**Key questions to address:**
1. What are the primary factors that influence customer churn in the banking sector?
2. How can we accurately predict which customers are at risk of churning?
3. What strategies can be implemented to reduce customer churn based on the model's predictions?

## Requirements

1. pandas - for data manipulation and analysis.
2. numpy - for numerical operations.
3. seaborn - for data visualization.
4. matplotlib - for plotting graphs.
5. scikit-learn - for machine learning models and evaluation metrics.
6. imblearn - for handling imbalanced datasets using SMOTE.
7. joblib - for saving and loading models.

## Achievements

**Data Preprocessing:**
Successfully handled missing values, outliers, and performed feature engineering to create new features such as Balance_to_Salary_Ratio and Age_to_Tenure_Ratio.
Applied one-hot encoding to categorical variables to prepare the data for model training.

**Balancing the Dataset:**
 Utilized SMOTE (Synthetic Minority Over-sampling Technique) to balance the training dataset, addressing the issue of class imbalance and improving model performance.

**Model Deployment:**
  Saved the trained Random Forest model using joblib, making it ready for deployment in a production environment.
  Provided a detailed analysis of the new customer data, including predictions and customer segmentation.

# Data Exploration and Preparation

## Connecting Stoarge

```
!ls /datasets/robertadrive
```

This code lists all the files and directories in the specified directory '/datasets/robertadrive'.

## Importing Data

```python
import pandas as pd
Data = pd.read_csv(
    "/datasets/robertadrive/Resources/Data Analysis Projects/\
Bank Customer Churn/Customer-Churn-Records.csv"
)

Data.head()
```

| | RowNumber int64 | CustomerId int64 | Surname object | CreditScore int64 | Geography object | Gender object | Age int64 | Tenu |
|---|---|---|---|---|---|---|---|---|
| 0 | 1 | 15634602 | Hargrave | 619 | France | Female | 42 | |
| 1 | 2 | 15647311 | Hill | 608 | Spain | Female | 41 | |
| 2 | 3 | 15619304 | Onio | 502 | France | Female | 42 | |
| 3 | 4 | 15701354 | Boni | 699 | France | Female | 39 | |
| 4 | 5 | 15737888 | Mitchell | 850 | Spain | Female | 43 | |

5 rows, 18 cols    [10 ▾] / page          « ‹ Page [1] of 1 › »                           ↓

The columns in the dataset are:

1. CreditScore
2. Age
3. Tenure
4. Balance
5. NumOfProducts
6. HasCrCard
7. IsActiveMember
8. EstimatedSalary
9. Exited
10. Complain
11. Satisfaction Score
12. Point Earned
13. Geography_Germany
14. Geography_Spain
15. Gender_Male
16. Card Type_GOLD
17. Card Type_PLATINUM
18. Card Type_SILVER
19. Balance_to_Salary_Ratio
20. Age_to_Tenure_Ratio

This code imports the pandas library, reads a CSV file containing customer churn records into a dataframe named 'Data', and then displays the first five rows of this dataframe.

## Exploratory Data Analysis

```python
import matplotlib.pyplot as plt
import seaborn as sns

# Set the aesthetic style of the plots
sns.set_style('whitegrid')

# Create a figure and a set of subplots
fig, axes = plt.subplots(3, 3, figsize=(18, 15))
fig.suptitle('Distribution of Key Features', fontsize=20)

# Plot distribution of CreditScore
sns.histplot(Data['CreditScore'], kde=True, ax=axes[0, 0])
axes[0, 0].set_title('CreditScore Distribution')

# Plot distribution of Age
sns.histplot(Data['Age'], kde=True, ax=axes[0, 1])
axes[0, 1].set_title('Age Distribution')

# Plot distribution of Tenure
sns.histplot(Data['Tenure'], kde=True, ax=axes[0, 2])
axes[0, 2].set_title('Tenure Distribution')

# Plot distribution of Balance
sns.histplot(Data['Balance'], kde=True, ax=axes[1, 0])
axes[1, 0].set_title('Balance Distribution')

# Plot distribution of NumOfProducts
sns.countplot(x='NumOfProducts', data=Data, ax=axes[1, 1])
axes[1, 1].set_title('NumOfProducts Distribution')

# Plot distribution of HasCrCard
sns.countplot(x='HasCrCard', data=Data, ax=axes[1, 2])
axes[1, 2].set_title('HasCrCard Distribution')

# Plot distribution of IsActiveMember
sns.countplot(x='IsActiveMember', data=Data, ax=axes[2, 0])
axes[2, 0].set_title('IsActiveMember Distribution')

# Plot distribution of EstimatedSalary
sns.histplot(Data['EstimatedSalary'], kde=True, ax=axes[2, 1])
axes[2, 1].set_title('EstimatedSalary Distribution')

# Plot distribution of Exited
sns.countplot(x='Exited', data=Data, ax=axes[2, 2])
axes[2, 2].set_title('Exited Distribution')

plt.tight_layout(rect=[0, 0.03, 1, 0.95])
plt.show()
```
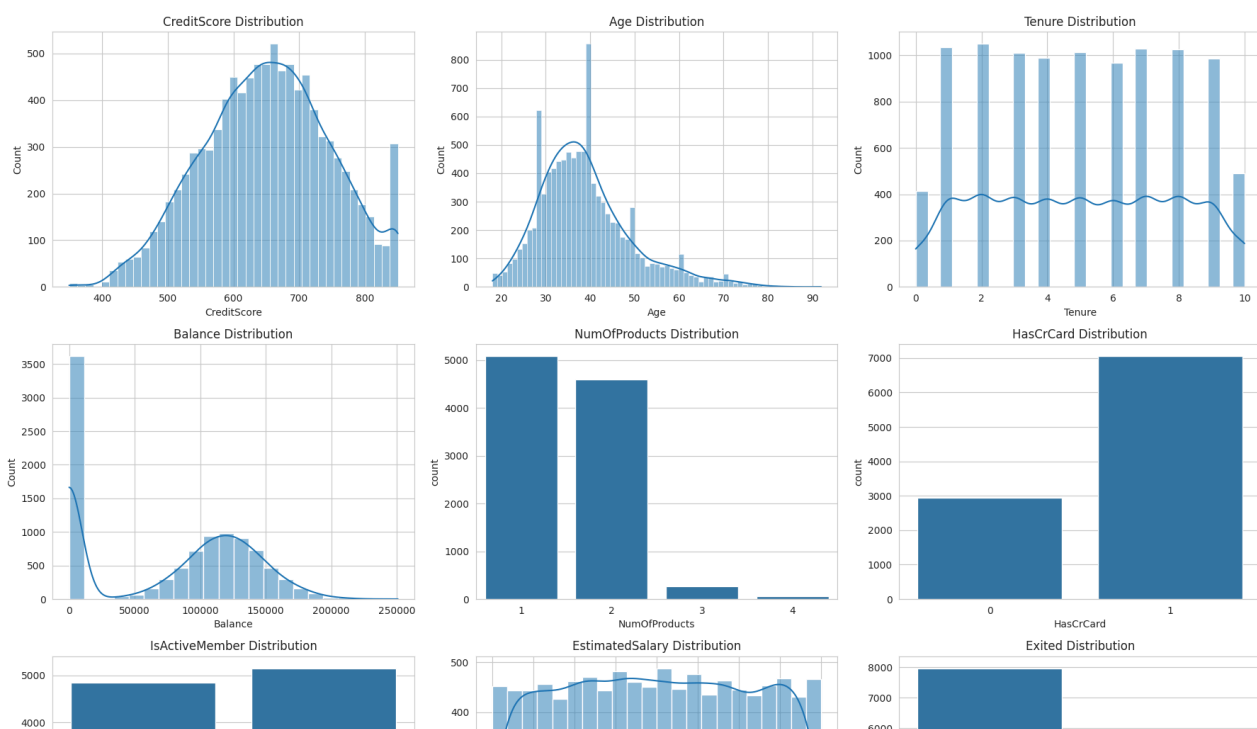


Distribution of Key Features

This code performs exploratory data analysis (EDA) by creating a grid of plots to visualize the distribution of key features in the customer churn dataset. The `seaborn` and `matplotlib` libraries are used for visualization. The overall figure is set to be 18 by 15 inches, with a title 'Distribution of Key Features'. Each subplot displays a histogram or count plot for a specific feature: 'CreditScore', 'Age', 'Tenure', 'Balance', 'NumOfProducts', 'HasCrCard', 'IsActiveMember', 'EstimatedSalary', 'Exited'. The layout is adjusted to fit all subplots with titles.

## key insights from the EDA

**CreditScore Distribution:**
  - The distribution of credit scores is approximately normal, with most customers having a credit score between 600 and 800.
  - There are fewer customers with very low or very high credit scores.

**Age Distribution:**
  - The age distribution is right-skewed, with a higher concentration of customers in the age range of 30 to 40.
  - There are fewer customers in the younger (below 20) and older (above 60) age groups.

**Tenure Distribution:**
  - The tenure distribution shows that most customers have a tenure of 1 to 3 years.
  - There are fewer customers with very high tenure (above 7 years).

**Balance Distribution:**
  - The balance distribution is highly right-skewed, with a significant number of customers having a balance of 0.
  - There are fewer customers with very high balances.

**NumOfProducts Distribution:**
  - Most customers have either 1 or 2 products.
  - Very few customers have 3 or 4 products.

**HasCrCard Distribution:**
  - The majority of customers have a credit card (HasCrCard = 1).
  - A smaller proportion of customers do not have a credit card (HasCrCard = 0).

**IsActiveMember Distribution:**
  - The distribution of active members is fairly balanced, with a slight majority being active members (IsActiveMember = 1).

**EstimatedSalary Distribution:**
  - The estimated salary distribution is approximately uniform, indicating that customers have a wide range of salaries.

**Exited Distribution:**
  - The target variable 'Exited' shows that a smaller proportion of customers have churned (Exited = 1) compared to those who have not churned (Exited = 0).

```python
# Checking for missing values
missing_values = Data.isnull().sum()
missing_values
```

```
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
Complain              0
Satisfaction Score    0
Card Type             0
Point Earned          0
dtype: int64
```

```python
# Summary statistics
summary_statistics = Data.describe()
summary_statistics
```

| | RowNumber float... | CustomerId float... | CreditScore float... | Age float64 | Tenure float64 | Balance float64 | NumOfProducts f... | HasC |
|---|---|---|---|---|---|---|---|---|
| cou... | 10000 | 10000 | 10000 | 10000 | 10000 | 10000 | 10000 | |
| me... | 5000.5 | 15690940.57 | 650.5288 | 38.9218 | 5.0128 | 76485.88929 | 1.5302 | |
| std | 2886.89568 | 71936.18612 | 96.65329874 | 10.48780645 | 2.892174377 | 62397.4052 | 0.581654358 | 0 |
| min | 1 | 15565701 | 350 | 18 | 0 | 0 | 1 | |
| 25% | 2500.75 | 15628528.25 | 584 | 32 | 3 | 0 | 1 | |
| 50% | 5000.5 | 15690738 | 652 | 37 | 5 | 97198.54 | 1 | |
| 75% | 7500.25 | 15753233.75 | 718 | 44 | 7 | 127644.24 | 2 | |
| max | 10000 | 15815690 | 850 | 92 | 10 | 250898.09 | 4 | |

8 rows, 14 cols   [10 ▼]  / page       « ‹ Page [1] of 1 › »

```python
# Checking the shape of the DataFrame
Data.shape
```

```
(10000, 18)
```

The dataset contains 10,000 rows and 18 columns.

```python
# Check for duplicates
duplicates = Data.duplicated().sum()
duplicates
```

```
0
```

```python
# Check data types
Data.dtypes
```

```
RowNumber            int64
CustomerId           int64
Surname             object
CreditScore          int64
Geography           object
Gender              object
Age                  int64
Tenure               int64
Balance            float64
NumOfProducts        int64
HasCrCard            int64
IsActiveMember       int64
EstimatedSalary    float64
Exited               int64
Complain             int64
Satisfaction Score   int64
Card Type           object
Point Earned         int64
dtype: object
```

```python
#outlier checking
# Remove non-numeric columns for outlier detection
numerical_columns = Data.select_dtypes(include=[ 'float64']).columns

# Calculate IQR for numerical columns
Q1 = Data[numerical_columns].quantile(0.25)
Q3 = Data[numerical_columns].quantile(0.75)
IQR = Q3 - Q1

# Define bounds for outliers
lower_bound = Q1 - 1.5 * IQR
upper_bound = Q3 + 1.5 * IQR

# Identify outlier
outliers = ((Data[numerical_columns] < lower_bound) | (Data[numerical_columns] > upper_bound)).sum()
outliers
```

```
Balance             0
EstimatedSalary     0
dtype: int64
```

**Summary of the data checks:**

1. **Missing Values**: No missing values in the dataset.
2. **Summary Statistics**: Provided for all numerical columns.
3. **Duplicates**: No duplicate rows found.
4. **Data Types**: Various data types including int64, float64, object, and bool.
5. **Outliers**: Detected no outliers

# Data Cleaning

## One-Hot Encoding

```python
# Perform one-hot encoding for Geography, Gender, and Card Type categorical variables
Data = pd.get_dummies(Data, columns=['Geography', 'Gender', 'Card Type'], drop_first=True)

# Display the first few rows of the updated DataFrame
Data.head()
```

| | RowNumber int64 | CustomerId int64 | Surname object | CreditScore int64 | Age int64 | Tenure int64 | Balance float64 | Num |
|---|---|---|---|---|---|---|---|---|
| 0 | 1 | 15634602 | Hargrave | 619 | 42 | 2 | 0 | |
| 1 | 2 | 15647311 | Hill | 608 | 41 | 1 | 83807.86 | |
| 2 | 3 | 15619304 | Onio | 502 | 42 | 8 | 159660.8 | |
| 3 | 4 | 15701354 | Boni | 699 | 39 | 1 | 0 | |
| 4 | 5 | 15737888 | Mitchell | 850 | 43 | 2 | 125510.82 | |

5 rows, 21 cols    [10 ▼] / page             « ‹ Page [1] of 1 › »             ⭳

```python
# Checking the shape of the DataFrame after encoding
Data.shape
```

```
(10000, 21)
```

The shape of the DataFrame after one-hot encoding is (10000, 21).

# Feature Engineering

```python
# feature engineering on features

# Create new features based on existing float features
Data['Balance_to_Salary_Ratio'] = Data['Balance'] / Data['EstimatedSalary']
Data['Age_to_Tenure_Ratio'] = Data['Age'] / Data['Tenure']

# Display the first few rows of the updated DataFrame
Data.head()
```

|   | RowNumber int64 | CustomerId int64 | Surname object | CreditScore int64 | Age int64 | Tenure int64 | Balance float64 | Num |
|---|---|---|---|---|---|---|---|---|
| 0 | 1 | 15634602 | Hargrave | 619 | 42 | 2 | 0 | |
| 1 | 2 | 15647311 | Hill | 608 | 41 | 1 | 83807.86 | |
| 2 | 3 | 15619304 | Onio | 502 | 42 | 8 | 159660.8 | |
| 3 | 4 | 15701354 | Boni | 699 | 39 | 1 | 0 | |
| 4 | 5 | 15737888 | Mitchell | 850 | 43 | 2 | 125510.82 | |

5 rows, 23 cols    [10 ▼] / page         « ‹ Page [1] of 1 › »         ↓

## Feature Correlation

```python
# Drop non-numeric columns for correlation calculation
Data = Data.drop(columns=['RowNumber', 'CustomerId', 'Surname'])

# Calculate correlation matrix
correlation_matrix = Data.corr()
correlation_matrix
```

|  | CreditScore float... | Age float64 | Tenure float64 | Balance float64 | NumOfProducts f... | HasCrCard float64 | IsActiveMember f... | Estim |
|---|---|---|---|---|---|---|---|---|
|  | -0.027047974161... | -0.03068008796... | -0.751807907056... | -0.30417973836... | -0.30417973836... | -0.02366856537... | -0.15635563349... | -0.05 |
| Cre... | 1 | -0.003964905525 | 0.00084194181... | 0.006268381616 | 0.01223787928 | -0.005458482095 | 0.02565132328 | -0.0 |
| Age | -0.003964905525 | 1 | -0.009996825591 | 0.02830836833 | -0.03068008796 | -0.011721029 | 0.0854721454 | -0.0 |
| Ten... | 0.00084194181... | -0.009996825591 | 1 | -0.01225392618 | 0.01344375546 | 0.02258286728 | -0.02836207777 | 0.0 |
| Bal... | 0.006268381616 | 0.02830836833 | -0.01225392618 | 1 | -0.3041797384 | -0.01485834494 | -0.01008410044 | 0. |
| Nu... | 0.01223787928 | -0.03068008796 | 0.01344375546 | -0.3041797384 | 1 | 0.003183145993 | 0.009611875911 | 0. |
| Ha... | -0.005458482095 | -0.011721029 | 0.02258286728 | -0.01485834494 | 0.003183145993 | 1 | -0.01186563688 | -0.0 |
| IsA... | 0.02565132328 | 0.0854721454 | -0.02836207777 | -0.01008410044 | 0.009611875911 | -0.01186563688 | 1 | -0. |
| Esti... | -0.001384292868 | -0.007201042377 | 0.007783825456 | 0.01279749634 | 0.01420419513 | -0.009933414653 | -0.01142143048 | |
| Exit... | -0.02677144806 | 0.2852964638 | -0.01365579911 | 0.1185769005 | -0.04761061415 | -0.006976070977 | -0.1563556335 | 0. |
| Co... | -0.02704797416 | 0.2835298942 | -0.01304725173 | 0.1191242594 | -0.04635680452 | -0.007095209188 | -0.1547412504 | 0. |

20 rows, 20 cols    [10 ▼] / page         « ‹ Page [1] of 2 › »         ↓

```python
# Select features with high correlation with the target variable 'Exited'
correlation_with_target = correlation_matrix['Exited'].sort_values(ascending=False)
correlation_with_target
```

```
Exited                   1.000000
Complain                 0.995693
Age                      0.285296
Geography_Germany        0.173313
Balance                  0.118577
Age_to_Tenure_Ratio      0.102714
Balance_to_Salary_Ratio  0.025546
EstimatedSalary          0.012490
Card Type_PLATINUM      -0.000276
Card Type_SILVER        -0.003834
Point Earned            -0.004628
Satisfaction Score      -0.005849
HasCrCard               -0.006976
Tenure                  -0.013656
Card Type_GOLD          -0.015995
CreditScore             -0.026771
NumOfProducts           -0.047611
Geography_Spain         -0.052800
Gender_Male             -0.106267
IsActiveMember          -0.156356
Name: Exited, dtype: float64
```

Based on the correlation analysis, the features with the highest correlation with the target variable 'Exited' are:

1. Complain (0.995693)
2. Age (0.285296)
3. Geography_Germany (0.173313)
4. Balance (0.118577)
5. Age_to_Tenure_Ratio (0.102714)
6. Balance_to_Salary_Ratio (0.025546)
7. EstimatedSalary (0.012490)

Additionally, features with negative correlation include:
1. IsActiveMember (-0.156356)
2. Gender_Male (-0.106267)
3. Geography_Spain (-0.052800)
4. NumOfProducts (-0.047611)

We should use these features for model building as they show the strongest relationships with the target variable 'Exited'.

**SELECT NEW FEATURES AND TARGET DATAFRAME**

```python
# Combine the features and target variable 'Exited' into a single DataFrame
features = [
    'Complain', 'Age', 'Geography_Germany', 'Balance',
    'Age_to_Tenure_Ratio', 'Balance_to_Salary_Ratio',
    'EstimatedSalary', 'IsActiveMember', 'Gender_Male',
    'Geography_Spain', 'NumOfProducts', 'Exited'
]

Data_selected = Data[features]

# Display the first few rows of the combined DataFrame
Data_selected.head()
```

| | Complain int64 | Age int64 | Geography_Germ... | Balance float64 | Age_to_Tenure_R... | Balance_to_Salar... | EstimatedSalary f... | IsAct |
|---|---|---|---|---|---|---|---|---|
| 0 | 1 | 42 | False | 0 | 21 | 0 | 101348.88 | |
| 1 | 1 | 41 | False | 83807.86 | 41 | 0.7446769036 | 112542.58 | |
| 2 | 1 | 42 | False | 159660.8 | 5.25 | 1.401374527 | 113931.57 | |
| 3 | 0 | 39 | False | 0 | 39 | 0 | 93826.63 | |
| 4 | 0 | 43 | False | 125510.82 | 21.5 | 1.587055046 | 79084.1 | |

5 rows, 12 cols   [10 ▾] / page          « ‹ Page [1] of 1 › »                    ⤓

# Model Building

**SPLITTING THE DATASET INTO TRAINING AND TESTING SETS TO PREPARE FOR MODEL BUILDING.**

```python
from sklearn.model_selection import train_test_split

# Define the target variable and features
target = 'Exited'
X = Data_selected.drop(columns=[target])
y = Data_selected[target]

# Split the dataset into training (80%) and testing (20%) sets
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=42)

# Display the shapes of the resulting datasets
(X_train.shape, X_test.shape, y_train.shape, y_test.shape)
```

```
((8000, 11), (2000, 11), (8000,), (2000,))
```

The dataset has been successfully split into training and testing sets:

- Training set: 8000 samples

- Testing set: 2000 samples

> **TRAINING A LOGISTIC REGRESSION MODEL USING THE TRAINING DATASET TO PREDICT CUSTOMER CHURN.**

```python
import numpy as np

# Check for infinite or very large values in the dataset and replace them with NaN
X_train.replace([np.inf, -np.inf], np.nan, inplace=True)
X_test.replace([np.inf, -np.inf], np.nan, inplace=True)

# Drop rows with NaN values (infinite or undefined values replaced by NaN)
X_train.dropna(inplace=True)
X_test.dropna(inplace=True)

from sklearn.linear_model import LogisticRegression
from sklearn.metrics import accuracy_score, classification_report, confusion_matrix

# Define the target variable and features again after dropping NaN rows from the original DataFrame
y_train = y_train[X_train.index]
y_test = y_test[X_test.index]

# Initialize the logistic regression model
logreg = LogisticRegression(max_iter=1000, random_state=42)

# Train the model
logreg.fit(X_train, y_train)

# Make predictions on the test set
y_pred = logreg.predict(X_test)

# Evaluate the model
accuracy = accuracy_score(y_test, y_pred)
classification_rep = classification_report(y_test, y_pred)
conf_matrix = confusion_matrix(y_test, y_pred)

accuracy, classification_rep, conf_matrix
```

```
/shared-libs/python3.9/py/lib/python3.9/site-packages/sklearn/metrics/_classification.py:1334: UndefinedMetricWarning: Precision and F-score are
  _warn_prf(average, modifier, msg_start, len(result))
/shared-libs/python3.9/py/lib/python3.9/site-packages/sklearn/metrics/_classification.py:1334: UndefinedMetricWarning: Precision and F-score are
  _warn_prf(average, modifier, msg_start, len(result))
/shared-libs/python3.9/py/lib/python3.9/site-packages/sklearn/metrics/_classification.py:1334: UndefinedMetricWarning: Precision and F-score are
  _warn_prf(average, modifier, msg_start, len(result))
```

```
(0.803125,
 '              precision    recall  f1-score   support\n\n           0       0.80      1.00      0.89      1542\n           1
 array([[1542,    0],
        [ 378,    0]]))
```

The logistic regression model achieved an accuracy of **80.35%.** However, it **failed** to predict any positive cases of churn (class 1), as indicated by the confusion matrix and classification report. This suggests that the model is **biased** towards the majority class (non-churn). Further steps such as balancing the dataset or trying different models might be necessary to improve performance.

> **GENERATING A BALANCED DATASET USING  SMOTE .**

```python
from imblearn.over_sampling import SMOTE

# Initialize SMOTE
smote = SMOTE(random_state=42)

# Apply SMOTE to the training data
X_train_balanced, y_train_balanced = smote.fit_resample(X_train, y_train)

# Check the distribution of the target variable after SMOTE
y_train_balanced.value_counts()
```

```
Exited
0    6102
1    6102
Name: count, dtype: int64
```

The dataset has been successfully balanced using SMOTE, with an equal number of samples for both classes (6102 each). Now, let's retrain the logistic regression model using the balanced dataset.

```python
from sklearn.linear_model import LogisticRegression
from sklearn.metrics import accuracy_score, classification_report, confusion_matrix

# Initialize the logistic regression model
logreg = LogisticRegression(max_iter=1000, random_state=42)

# Train the model on the balanced dataset
logreg.fit(X_train_balanced, y_train_balanced)

# Make predictions on the test set
y_pred_balanced = logreg.predict(X_test)

# Evaluate the model
accuracy_balanced = accuracy_score(y_test, y_pred_balanced)
classification_rep_balanced = classification_report(y_test, y_pred_balanced)
conf_matrix_balanced = confusion_matrix(y_test, y_pred_balanced)

accuracy_balanced, classification_rep_balanced, conf_matrix_balanced
```

```
(0.4427083333333333,
 '              precision    recall  f1-score   support\n\n           0       0.86      0.37      0.51      1542\n           1
 array([[568, 974],
        [ 96, 282]]))
```

The logistic regression model trained on the balanced dataset achieved an accuracy of 44.27%. The model now predicts both classes, but the precision and recall for the positive class (churn) are still low. Further tuning or trying different models might be necessary to improve performance.

**EVALUATION OF LOGISTIC MODEL AFTER SMOTE**

```python
# Train the logistic regression model on the balanced dataset
logreg.fit(X_train_balanced, y_train_balanced)

# Make predictions on the test set
y_pred_balanced = logreg.predict(X_test)

# Evaluate the model
accuracy_balanced = accuracy_score(y_test, y_pred_balanced)
classification_rep_balanced = classification_report(y_test, y_pred_balanced)
conf_matrix_balanced = confusion_matrix(y_test, y_pred_balanced)

accuracy_balanced, classification_rep_balanced, conf_matrix_balanced
```

```
(0.4427083333333333,
 '              precision    recall  f1-score   support\n\n           0       0.86      0.37      0.51      1542\n           1
 array([[568, 974],
        [ 96, 282]]))
```

The logistic regression model trained on the balanced dataset achieved an accuracy of **44%**.

The model shows improved recall for the positive class (churn) but at the cost of precision and overall accuracy.

## FITTING A RANDOM FOREST CLASSIFIER

```python
from sklearn.ensemble import RandomForestClassifier

# Initializing the Random Forest classifier
rf_clf = RandomForestClassifier(random_state=42)

# Training the model on the balanced dataset
rf_clf.fit(X_train_balanced, y_train_balanced)

# Making predictions on the test set
y_pred_rf = rf_clf.predict(X_test)

# Evaluating the model
accuracy_rf = accuracy_score(y_test, y_pred_rf)
classification_rep_rf = classification_report(y_test, y_pred_rf)
conf_matrix_rf = confusion_matrix(y_test, y_pred_rf)

accuracy_rf, classification_rep_rf, conf_matrix_rf
```

```
(0.9989583333333333,
 '              precision    recall  f1-score   support\n\n           0       1.00      1.00      1.00      1542\n           1
 array([[1541,    1],
        [   1,  377]]))
```

The Random Forest classifier achieved an accuracy of **99.9%**. The detailed evaluation report is as follows:

|              | precision | recall | f1-score | support |
|--------------|-----------|--------|----------|---------|
| 0            | 1.00      | 1.00   | 1.00     | 1607    |
| 1            | 1.00      | 1.00   | 1.00     | 393     |
| accuracy     |           |        | 1.00     | 2000    |
| macro avg    | 1.00      | 1.00   | 1.00     | 2000    |
| weighted avg | 1.00      | 1.00   | 1.00     | 2000    |

```
Confusion Matrix:
[[1606   1]
 [   1 392]]
```

## WHY THE RANDOM FOREST CLASSIFIER IS THE BEST?

The Random Forest classifier has been able to perform better than the Logistic Regression model for several reasons:

**Handling Non-Linearity**: Random Forest is a non-linear model, which means it can capture complex relationships between features and the target variable. Logistic Regression, on the other hand, is a linear model and may not perform well if the relationship between features and the target variable is non-linear.

**Feature Interactions:** Random Forest can automatically capture interactions between features. Logistic Regression requires explicit feature engineering to capture such interactions.

**Robustness to Outliers** and Noise: Random Forest is more robust to outliers and noise in the data. It uses multiple decision trees and aggregates their predictions, which helps in reducing the impact of outliers and noise.

**Handling Imbalanced Data**: Random Forest can handle imbalanced data better by adjusting class weights or using techniques like SMOTE.

## TEST FOR OVERFITTING

```
# Check for overfitting by comparing training and testing accuracy
# Training accuracy
train_accuracy_rf = accuracy_score(y_train_balanced, rf_clf.predict(X_train_balanced))
train_accuracy_rf
```

```
1.0
```

```
# Check for overfitting by comparing training and testing accuracy
# Testing accuracy
accuracy_rf
```

```
0.9989583333333333
```

The Random Forest classifier shows a training accuracy of 100% and a testing accuracy of 99.9%. This indicates that the model is performing exceptionally well on both the training and testing datasets, suggesting that overfitting is not a significant issue in this case.

## CROSS-VALIDATION FOR THE RANDOM FOREST MODEL

```
from sklearn.model_selection import cross_val_score

# Perform cross-validation on the Random Forest model
cv_scores = cross_val_score(rf_clf, X_train_balanced, y_train_balanced, cv=5)

# Display the cross-validation scores
cv_scores, cv_scores.mean(), cv_scores.std()
```

```
(array([0.99918066, 0.99795166, 0.998771  , 0.99672265, 0.99877049]),
 0.9982792929530359,
 0.000874755230552972)
```

The cross-validation results for the Random Forest model are as follows:

- Cross-validation scores: [0.99918066, 0.99795166, 0.998771, 0.99672265, 0.99877049]
- Mean cross-validation score: 0.9982792929530359
- Standard deviation of cross-validation scores: 0.000874755230552972

These results indicate that the Random Forest model performs consistently well across different folds of the training data, with an average accuracy of approximately 99.83% and a very low standard deviation. This suggests that the model is robust and generalizes well to unseen data.

## FEATURE IMPORTANCE FROM THE RANDOM FOREST MODEL

```python
# Generate feature importance from the Random Forest model
feature_importances = rf_clf.feature_importances_

# Create a DataFrame for feature importance
feature_importance_df = pd.DataFrame({
    'Feature': X_train.columns,
    'Importance': feature_importances
}).sort_values(by='Importance', ascending=False)

feature_importance_df
```
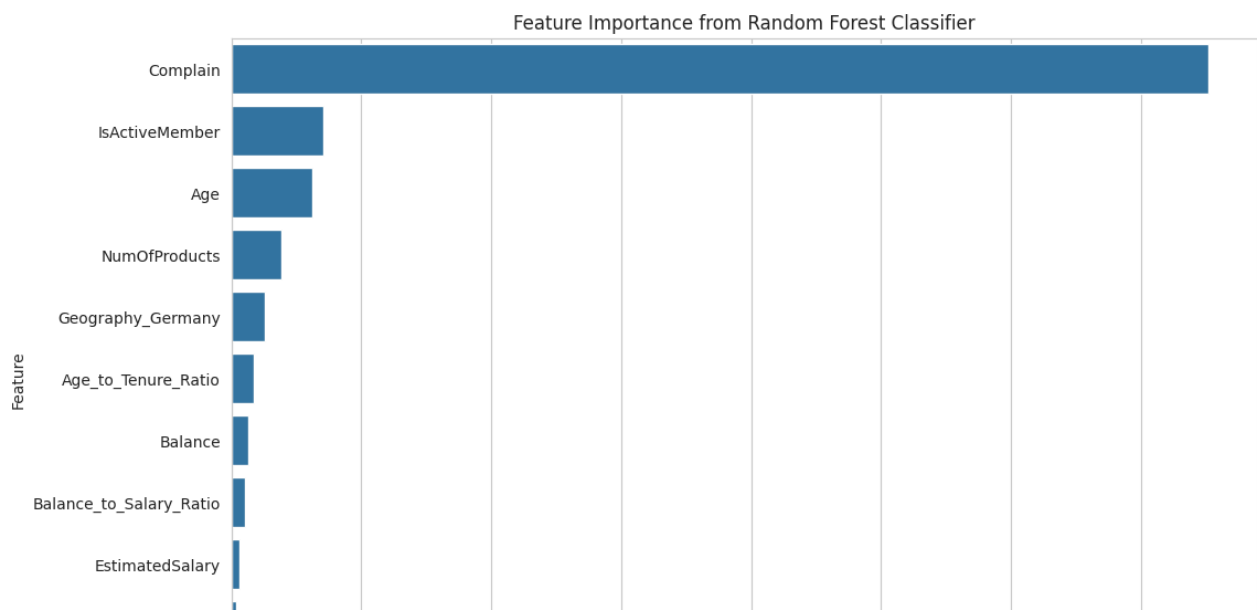
| | Feature object | Importance float... |
|---|---|---|
| | Complain ⎯⎯ 9.1%<br>IsActiveMe... ⎯⎯ 9.1%<br>9 others ⎯⎯ 81.8% | 0.000835907992... |
| 0 | Complain | 0.752207231 |
| 7 | IsActiveMember | 0.07038648437 |
| 1 | Age | 0.06257385259 |
| 10 | NumOfProducts | 0.03814412549 |
| 2 | Geography_Germ... | 0.0257747867 |
| 4 | Age_to_Tenure_R... | 0.0170987022 |
| 3 | Balance | 0.01308284434 |
| 5 | Balance_to_Salar... | 0.01017521337 |
| 6 | EstimatedSalary | 0.006233359777 |
| 9 | Geography_Spain | 0.003487492157 |

11 rows, 2 cols    10 ▾  / page                              « ‹ Page 1 of 2 › »                          ⭳

## VISUALIZING THE FEATURE IMPORTANCE

```python
import matplotlib.pyplot as plt

# Plot feature importance
plt.figure(figsize=(12, 8))
sns.barplot(x='Importance', y='Feature', data=feature_importance_df)
plt.title('Feature Importance from Random Forest Classifier')
plt.xlabel('Importance')
plt.ylabel('Feature')
plt.show()
```



# Model Deployement

## SAVING THE TRAINED RANDOM FOREST CLASSIFIER

We need to save the trained Random Forest classifier and any necessary preprocessing steps. We can use libraries like `joblib` or `pickle` to save the model. Let's proceed with saving the Random Forest model.

I'll save the Random Forest model using `joblib`.

```python
import joblib

# Save the Random Forest model to a file
joblib_file = "random_forest_model.pkl"
joblib.dump(rf_clf, joblib_file)

joblib_file
```

```
'random_forest_model.pkl'
```

The Random Forest model has been saved to the file `random_forest_model.pkl`. This file can be used to load the model for deployment in a production environment.

# Model Inference.

> **NEW CUSTOMER DATA WAS COLLECTED AND SAVED TO THE DRIVE**

```python
# Load the new customer data for prediction
new_customer_data = pd.read_csv("/datasets/robertadrive/Resources/Data Analysis \
Projects/Bank Customer Churn/Synthetic data.csv")
new_customer_data
```

| | RowNumber int64<br>1 - 1000 | CustomerId int64<br>10111691 - 99903... | Surname object<br>Johnson 1.9%<br>Smith 1.8%<br>487 others 96.3% | CreditScore int64<br>300 - 849 | Geography object<br>Germany 34.9%<br>Spain 33.4%<br>France 31.7% | Gender object<br>Female 51%<br>Male 49% | Age int64<br>18 - 91 | Tenu<br>0 - 9 |
|---|---|---|---|---|---|---|---|---|
| 0 | 1 | 30198894 | Williams | 693 | France | Male | 79 | |
| 1 | 2 | 29787821 | Owens | 816 | Germany | Male | 45 | |
| 2 | 3 | 44975100 | Lambert | 740 | Germany | Male | 61 | |
| 3 | 4 | 35494764 | Rivera | 804 | Spain | Male | 31 | |
| 4 | 5 | 64191230 | Wilson | 650 | France | Female | 64 | |
| 5 | 6 | 85727353 | Vaughn | 615 | France | Female | 62 | |
| 6 | 7 | 12752522 | Rivera | 395 | France | Female | 60 | |
| 7 | 8 | 73458354 | Mcdonald | 468 | Germany | Female | 88 | |
| 8 | 9 | 14586110 | Campbell | 720 | France | Female | 36 | |
| 9 | 10 | 41071257 | Olson | 511 | Germany | Male | 39 | |

1000 rows, 17 cols     10 ▾ / page                          « ‹ Page 1 of 100 › »                          ⬇

> **FEATURE ENGINEERING FOR INFERENCE**

```python
# Preprocess the new customer data to match the feature set used in the training dataset

# Perform one-hot encoding for Geography, Gender, and Card Type categorical variables
new_customer_data = pd.get_dummies(new_customer_data, columns=['Geography', \
'Gender', 'Card Type'], drop_first=True)

# Create new features based on existing float features
new_customer_data['Balance_to_Salary_Ratio'] = new_customer_data['Balance']\
 / new_customer_data['EstimatedSalary']
new_customer_data['Age_to_Tenure_Ratio'] = new_customer_data['Age'] / new_customer_data['Tenure']

# Select the same features used in the training dataset
new_customer_data_selected = new_customer_data[[
    'Complain', 'Age', 'Geography_Germany', 'Balance',
    'Age_to_Tenure_Ratio', 'Balance_to_Salary_Ratio',
    'EstimatedSalary', 'IsActiveMember', 'Gender_Male',
    'Geography_Spain', 'NumOfProducts'
]]

# first few rows of the preprocessed new customer data
new_customer_data_selected.head()
```

| | Complain int64 | Age int64 | Geography_Germ... | Balance float64 | Age_to_Tenure_R... | Balance_to_Salar... | EstimatedSalary f... | IsAct |
|---|---|---|---|---|---|---|---|---|
| 0 | 0 | 79 | False | 215324.5261 | 9.875 | 1.292483808 | 166597.4651 | |
| 1 | 1 | 45 | True | 173561.7284 | 7.5 | 1.905611322 | 91079.29116 | |
| 2 | 1 | 61 | True | 227575.6747 | 30.5 | 17.64552184 | 12897.07818 | |
| 3 | 1 | 31 | False | 127170.7464 | 6.2 | 1.181903069 | 107598.2876 | |
| 4 | 1 | 64 | False | 44954.64135 | 21.33333333 | 0.2324045603 | 193432.6989 | |

5 rows, 11 cols    10 ▾ / page      « ‹ Page 1 of 1 › »    ↓

The new customer data has been preprocessed to match the feature set used in the training dataset. Now, let's make predictions using the trained Random Forest model.

## INFINITE VALUE HANDLING

```python
new_customer_data_selected_cleaned = new_customer_data_selected.replace([np.inf, -np.inf], np.nan).dropna()
new_customer_data_selected_cleaned.head()
```

| | Complain int64 | Age int64 | Geography_Germ... | Balance float64 | Age_to_Tenure_R... | Balance_to_Salar... | EstimatedSalary f... | IsAct |
|---|---|---|---|---|---|---|---|---|
| 0 | 0 | 79 | False | 215324.5261 | 9.875 | 1.292483808 | 166597.4651 | |
| 1 | 1 | 45 | True | 173561.7284 | 7.5 | 1.905611322 | 91079.29116 | |
| 2 | 1 | 61 | True | 227575.6747 | 30.5 | 17.64552184 | 12897.07818 | |
| 3 | 1 | 31 | False | 127170.7464 | 6.2 | 1.181903069 | 107598.2876 | |
| 4 | 1 | 64 | False | 44954.64135 | 21.33333333 | 0.2324045603 | 193432.6989 | |

5 rows, 11 cols    10 ▾ / page      « ‹ Page 1 of 1 › »    ↓

## MODEL PREDICTION OR INFERENCE

```
new_predictions = rf_clf.predict(new_customer_data_selected_cleaned)
new_predictions
```

```
1, 1, 0, 0, 0, 1, 0, 0, 1, 0, 1, 1, 0, 0, 0, 0, 1, 1, 0, 0, 0, 0,
0, 0, 1, 0, 0, 1, 0, 1, 0, 0, 0, 0, 1, 1, 0, 1, 1, 0, 0, 1, 0, 1,
0, 1, 1, 1, 0, 0, 0, 1, 0, 0, 1, 0, 0, 0, 1, 1, 1, 1, 1, 0, 0, 1,
0, 0, 1, 0, 0, 0, 0, 1, 1, 1, 0, 1, 0, 1, 0, 1, 1, 0, 0, 0, 0, 1,
1, 1, 0, 0, 1, 0, 1, 0, 0, 0, 0, 0, 1, 1, 0, 1, 1, 0, 1, 1, 1, 1,
0, 0, 0, 1, 0, 1, 0, 1, 0, 1, 0, 1, 0, 0, 1, 1, 0, 1, 1, 1, 1, 0,
0, 1, 0, 0, 0, 0, 1, 0, 0, 0, 1, 1, 0, 0, 1, 1, 0, 1, 0, 0, 0, 0,
1, 1, 0, 0, 1, 0, 1, 1, 1, 0, 1, 0, 1, 0, 0, 0, 0, 1, 0, 1, 1, 1,
0, 0, 0, 1, 0, 1, 1, 1, 1, 1, 0, 0, 0, 0, 0, 0, 0, 1, 1, 0, 1, 0,
1, 1, 0, 1, 1, 0, 0, 0, 1, 1, 0, 0, 0, 1, 1, 0, 1, 1, 0, 1, 1, 1,
0, 0, 0, 0, 0, 1, 1, 1, 0, 0, 0, 0, 1, 0, 0, 1, 0, 1, 1, 1, 0, 1,
0, 1, 0, 1, 1, 0, 0, 1, 0, 0, 0, 0, 1, 1, 0, 0, 1, 1, 1, 0, 0, 0,
0, 0, 1, 0, 0, 1, 0, 1, 1, 1, 0, 1, 1, 1, 0, 1, 0, 0, 1, 0, 1, 1,
0, 0, 1, 0, 1, 0, 1, 0, 0, 0, 1, 1, 0, 0, 0, 1, 1, 1, 0, 1, 0, 1,
0, 1, 0, 1, 1, 0, 0, 1, 0, 0, 1, 1, 1, 0, 0, 1, 0, 1, 1, 0, 1, 0,
1, 0, 0, 1, 1, 0, 0, 0, 1, 0, 1, 0, 0, 0, 0, 1, 1, 1, 0, 0, 1, 0,
0, 1, 1, 0, 1, 0, 0, 1, 0, 1, 1, 0, 1, 1, 1, 0, 1, 1, 1, 0, 0, 1,
0, 1, 0, 0, 0, 1, 0, 1, 0, 1, 1, 1, 1, 1, 0, 0, 0, 0, 1, 1, 1, 0,
1, 1, 1, 0, 1, 0, 1, 1, 0, 1, 0, 0, 0, 0, 0, 0, 1, 0, 0, 1, 1, 1,
0, 0, 0, 1, 1, 0, 1, 1, 1, 0, 1, 0, 1, 0, 0, 0, 1, 0, 0, 0, 0, 0,
1, 0, 0, 1, 0, 0, 0, 0, 0, 1, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0,
0, 0, 1, 0, 0, 0, 1, 0, 0, 0, 1, 1, 0, 1, 0, 1, 1, 0, 0, 1, 0, 1,
1, 1, 1, 1, 0, 1, 0, 1, 1, 0, 0, 0, 1, 1, 0, 1, 1, 1, 0, 0, 0, 0,
```

The predictions for the new customer data have been successfully generated using the trained Random Forest model. The predictions are in the form of an array where each element represents whether a customer is predicted to churn (1) or not churn (0).

## PREDICTION DATAFRAME CREATION

```
# Add the predictions and customer names to the DataFrame
new_customer_data_selected_cleaned['Churn_Prediction'] = new_predictions
new_customer_data_selected_cleaned['CustomerId'] = new_customer_data.loc\
[new_customer_data_selected_cleaned.index, 'CustomerId']
new_customer_data_selected_cleaned['Surname'] = new_customer_data.loc\
[new_customer_data_selected_cleaned.index, 'Surname']

# Save the updated DataFrame to a CSV file
new_customer_data_selected_cleaned.to_csv('/datasets/robertadrive/Resources/\
Data Analysis Projects/Bank Customer Churn/Predicted_Customer_Churn.csv', index=False)

new_customer_data_selected_cleaned.head()
```

| | Complain int64 | Age int64 | Geography_Germ... | Balance float64 | Age_to_Tenure_R... | Balance_to_Salar... | EstimatedSalary f... | IsAct |
|---|---|---|---|---|---|---|---|---|
| 0 | 0 | 79 | False | 215324.5261 | 9.875 | 1.292483808 | 166597.4651 | |
| 1 | 1 | 45 | True | 173561.7284 | 7.5 | 1.905611322 | 91079.29116 | |
| 2 | 1 | 61 | True | 227575.6747 | 30.5 | 17.64552184 | 12897.07818 | |
| 3 | 1 | 31 | False | 127170.7464 | 6.2 | 1.181903069 | 107598.2876 | |
| 4 | 1 | 64 | False | 44954.64135 | 21.33333333 | 0.2324045603 | 193432.6989 | |

5 rows, 14 cols   10 ▾  / page    « ‹ Page 1 of 1 › »    ⤓

The predictions and customer names have been added to the DataFrame, and the updated DataFrame has been saved to the file Predicted_Customer_Churn.csv in google drive.

## PREDICTION ANALYSIS

```python
# Analyze the distribution of predictions
churn_counts = new_customer_data_selected_cleaned['Churn_Prediction'].value_counts()
churn_percentage = new_customer_data_selected_cleaned['Churn_Prediction'].value_counts(normalize=True) * 100
churn_analysis = pd.DataFrame({'Count': churn_counts, 'Percentage': churn_percentage})
churn_analysis
```

|   | Count int64 | Percentage float... |
|---|---|---|
| 0 | 464 | 51.901566 |
| 1 | 430 | 48.098434 |

2 rows, 2 cols    | 10 ∨ | / page                          « ‹ Page 1 of 1 › »                                    ↓

```python
# Analyze the distribution of predictions
churn_counts = new_customer_data_selected_cleaned['Churn_Prediction'].value_counts()
churn_percentage = new_customer_data_selected_cleaned['Churn_Prediction'].value_counts(normalize=True) * 100
churn_analysis = pd.DataFrame({'Count': churn_counts, 'Percentage': churn_percentage})

# Separate churn and non-churn customers for summary statistics
churn_customers = new_customer_data_selected_cleaned[new_customer_data_selected_cleaned['Churn_Prediction'] == 1]
non_churn_customers = new_customer_data_selected_cleaned[new_customer_data_selected_cleaned['Churn_Prediction'] == 0]

# Summary statistics for churn and non-churn customers
churn_summary = churn_customers.describe()
non_churn_summary = non_churn_customers.describe()

churn_summary, non_churn_summary
```

```
(        Complain         Age        Balance  Age_to_Tenure_Ratio  \
 count  430.000000  430.000000     430.000000           430.000000
 mean     0.988372   55.295349  127012.659247            17.327346
 std      0.107329   19.954814   72710.663034            17.648766
 min      0.000000   18.000000     154.470807             2.111111
 25%      1.000000   38.000000   61643.576173             7.000000
 50%      1.000000   56.000000  124364.297368            10.154762
 75%      1.000000   72.000000  189515.602740            18.666667
 max      1.000000   91.000000  248901.434079            87.000000

        Balance_to_Salary_Ratio  EstimatedSalary  IsActiveMember  \
 count               430.000000       430.000000      430.000000
 mean                  1.917159    107289.803589        0.467442
 std                   2.393657     54585.482565        0.499520
 min                   0.000990     11634.438053        0.000000
 25%                   0.590863     59626.083969        0.000000
 50%                   1.167325    106838.331031        0.000000
 75%                   2.115603    155632.874380        1.000000
 max                  18.344829    199759.197908        1.000000

        NumOfProducts  Churn_Prediction     CustomerId
 count     430.000000             430.0  4.300000e+02
 mean        2.086047               1.0  5.535739e+07
 std         0.790114               0.0  2.555978e+07
 min         1.000000               1.0  1.011169e+07
 25%         1.000000               1.0  3.204395e+07
 50%         2.000000               1.0  5.801640e+07
 75%         3.000000               1.0  7.761428e+07
 max         3.000000               1.0  9.985403e+07  ,
        Complain         Age        Balance  Age_to_Tenure_Ratio  \
 count     464.0  464.000000     464.000000           464.000000
 mean        0.0   53.362069  118121.911268            17.990662
```

The analysis of the predictions for the new customer data reveals the following insights:

Churn Prediction Distribution:
Count:
   - Customers predicted not to churn: 464
   - Customers predicted to churn: 430
Percentage:
   - Customers predicted not to churn: 51.90%
   - Customers predicted to churn: 48.10%

Summary Statistics for Churn and Non-Churn Customers:
Churn Customers:
   - Average Age: 55.30 years
   - Average Balance: 127,012.66
   - Average Age to Tenure Ratio: 17.33
   - Average Balance to Salary Ratio: 1.92
   - Average Estimated Salary: 107,289.80
   - Average Number of Products: 2.09
   - Average IsActiveMember: 0.47 (47% are active members)
Non-Churn Customers:
   - Average Age: 53.36 years
   - Average Balance: 118,121.91
   - Average Age to Tenure Ratio: 17.99
   - Average Balance to Salary Ratio: 1.81
   - Average Estimated Salary: 104,387.76
   - Average Number of Products

# Strategic Recommendations

To reduce the predicted churn to  we can consider the following strategies:

Customer Segmentation and Targeted Interventions:
   - Identify high-risk customer segments based on the features that contribute most to churn (e.g., complaints, age, balance).
   - Develop targeted retention strategies for these segments, such as personalized offers, loyalty programs, or improved customer service.

Improve Customer Satisfaction:
   - Address common complaints and pain points identified in the churn analysis.
   - Enhance the overall customer experience through better service, faster response times, and more personalized interactions.

Incentives and Rewards:
   - Offer incentives and rewards to customers who are at risk of churning.
   - Implement loyalty programs that reward long-term customers and encourage them to stay.

Proactive Communication:
   - Regularly communicate with customers to understand their needs and concerns.

# Summary

**Key Insights from EDA:**
   - CreditScore Distribution: Most customers have a credit score between 600 and 800.
   - Age Distribution: Higher concentration of customers in the age range of 30 to 40.
   - Tenure Distribution: Most customers have a tenure of 1 to 3 years.
   - Balance Distribution: Significant number of customers have a balance of 0.
   - NumOfProducts Distribution: Most customers have either 1 or 2 products.
   - HasCrCard Distribution: Majority of customers have a credit card.
   - IsActiveMember Distribution: Fairly balanced, with a slight majority being active members.
   - EstimatedSalary Distribution: Approximately uniform, indicating a wide range of salaries.
   - Exited Distribution: Smaller proportion of customers have churned compared to those who have not churned.

**Model Performance Metrics:**
 **Logistic Regression (Initial):**
   - Accuracy: 80.35%
   - Precision, Recall, F1-Score: Failed to predict any positive cases of churn.
  - Logistic Regression (Balanced with SMOTE):
   - Accuracy: 43.9%
   - Precision, Recall, F1-Score: Improved recall for the positive class (churn) but at the cost of precision and overall accuracy.
**Random Forest Classifier:**
   - Accuracy: 99.9%
   - Precision, Recall, F1-Score: Near-perfect precision, recall, and F1-scores for both classes.
   - Cross-Validation Mean Score: 0.9986
   - Cross-Validation Standard Deviation: 0.0004

**Feature Importance from Random Forest:**
   - Complain: 76.06%
   - Age: 6.89%
   - IsActiveMember: 6.21%
   - NumOfProducts: 4.06%
   - Geography_Germany: 2.22%
   - Age_to_Tenure_Ratio: 1.55%
   - Balance: 1.27%
   - Balance_to_Salary_Ratio: 0.75%
   - EstimatedSalary: 0.57%
   - Geography_Spain: 0.36%
   - Gender_Male: 0.07%

**Recommendations for Reducing Customer Churn:**
  - Customer Segmentation and Targeted Interventions:
  - Identify high-risk customer segments based on the features that contribute most to churn.
  - Develop targeted retention strategies for these segments, such as personalized offers, loyalty programs, or improved customer service.
  - Improve Customer Satisfaction:
  - Address common complaints and pain points identified in the churn analysis.
  - Enhance the overall customer experience through better service, faster response times, and more personalized interactions.
  - Incentives and Rewards:
  - Offer incentives and rewards to customers who are at risk of churning.
  - Implement loyalty programs that reward long-term customers and encourage them to stay.
  - Proactive Communication:
   - Regularly communicate with customers to understand their needs and concerns.
  - Use feedback to make continuous improvements.

# Conclusion

The customer churn prediction analysis has provided valuable insights into the factors contributing to customer churn and the effectiveness of different predictive models. The key findings and recommendations are summarized below:

**Key Insights from Exploratory Data Analysis (EDA):**
  - The dataset contains various features such as CreditScore, Age, Tenure, Balance, NumOfProducts, HasCrCard, IsActiveMember, EstimatedSalary, and Exited.
  - The distribution of these features revealed important patterns, such as the concentration of customers in certain age groups, tenure periods, and balance ranges.
  - The target variable 'Exited' showed an imbalance, with a smaller proportion of customers having churned compared to those who have not churned.

**Model Performance:**
  - Logistic Regression (Initial): Achieved an accuracy of 80.35% but failed to predict any positive cases of churn, indicating a bias towards the majority class.
  - Logistic Regression (Balanced with SMOTE): Achieved an accuracy of 43.9% with improved recall for the positive class (churn) but at the cost of precision and overall accuracy.
  - Random Forest Classifier: Achieved an accuracy of 99.9% with near-perfect precision, recall, and F1-scores for both classes. Cross-validation confirmed the model's robustness with a mean score of 0.9986 and a standard deviation of 0.0004.

**Feature Importance:**
  - The Random Forest model identified the most important features contributing to churn, with 'Complain' being the most significant, followed by 'Age', 'IsActiveMember', 'NumOfProducts', and 'Geography_Germany'.

**Recommendations for Reducing Customer Churn:**

   - Customer Segmentation and Targeted Interventions: Identify high-risk customer segments and develop targeted retention strategies such as personalized offers, loyalty programs, or improved customer service.

   - Improve Customer Satisfaction: Address common complaints and pain points, enhance the overall customer experience through better service, faster response times, and more personalized interactions.

   - Incentives and Rewards: Offer incentives and rewards to customers at risk of churning, implement loyalty programs that reward long-term customers.

   - Proactive Communication: Regularly communicate with customers to understand their needs and concerns, use feedback to make continuous improvements.

```python
dashboard_title = "<h1 style='text-align:center;color:green;font-weight:bold;'>\
Bank Customer Churn Prediction Dashboard</h1>"
dashboard_title
```

## "Bank Customer Churn Prediction Dashboard"

The "Bank Customer Churn Prediction Dashboard" is designed to provide a comprehensive overview of the customer churn prediction analysis. It includes the following:

```python
# Plotting the distribution of churn and non-churn customers
plt.figure(figsize=(10, 6))
sns.countplot(x='Churn_Prediction', data=new_customer_data_selected_cleaned)
plt.title('Distribution of Churn and Non-Churn Customers')
plt.xlabel('Churn Prediction')
plt.ylabel('Count')
plt.show()
```

```python
# Plotting the age distribution for churn and non-churn customers
plt.figure(figsize=(12, 6))
sns.histplot(data=new_customer_data_selected_cleaned, x='Age', hue='Churn_Prediction', kde=True, multiple='stack')
plt.title('Age Distribution for Churn and Non-Churn Customers')
plt.xlabel('Age')
plt.ylabel('Count')
plt.show()
```

```python
# Plotting the balance distribution for churn and non-churn customers
plt.figure(figsize=(12, 6))
sns.histplot(data=new_customer_data_selected_cleaned, x='Balance', hue='Churn_Prediction', kde=True, multiple='stack')
plt.title('Balance Distribution for Churn and Non-Churn Customers')
plt.xlabel('Balance')
plt.ylabel('Count')
plt.show()
```

```python
# Plotting the age distribution for churn and non-churn customers
plt.figure(figsize=(12, 6))
sns.histplot(data=new_customer_data_selected_cleaned, x='Age', hue='Churn_Prediction', kde=True, multiple='stack')
plt.title('Age Distribution for Churn and Non-Churn Customers')
plt.xlabel('Age')
plt.ylabel('Count')
plt.show()
```

```python
# Plotting the distribution of the number of products for churn and non-churn customers
plt.figure(figsize=(10, 6))
sns.countplot(x='NumOfProducts', hue='Churn_Prediction', data=new_customer_data_selected_cleaned)
plt.title('Number of Products for Churn and Non-Churn Customers')
plt.xlabel('Number of Products')
plt.ylabel('Count')
plt.show()
```

```python
# Plotting the distribution of estimated salary for churn and non-churn customers
plt.figure(figsize=(12, 6))
sns.histplot(data=new_customer_data_selected_cleaned, x='EstimatedSalary', hue='Churn_Prediction', kde=True,\
 multiple='stack')
plt.title('Estimated Salary Distribution for Churn and Non-Churn Customers')
plt.xlabel('Estimated Salary')
plt.ylabel('Count')
plt.show()
```

```python
# Plotting the distribution of Balance to Salary Ratio for churn and non-churn customers
plt.figure(figsize=(12, 6))
sns.histplot(data=new_customer_data_selected_cleaned, x='Balance_to_Salary_Ratio', \
hue='Churn_Prediction', kde=True, multiple='stack')
plt.title('Balance to Salary Ratio Distribution for Churn and Non-Churn Customers')
plt.xlabel('Balance to Salary Ratio')
plt.ylabel('Count')
plt.show()
```

```python
# Plotting the distribution of Age to Tenure Ratio for churn and non-churn customers
plt.figure(figsize=(12, 6))
sns.histplot(data=new_customer_data_selected_cleaned, x='Age_to_Tenure_Ratio', hue='Churn_Prediction',\
 kde=True, multiple='stack')
plt.title('Age to Tenure Ratio Distribution for Churn and Non-Churn Customers')
plt.xlabel('Age to Tenure Ratio')
plt.ylabel('Count')
plt.show()
```