# **Bank Churn Prediction**

# **Project Introduction**

- Customer relation is key for banking success. XYZ multistate bank is keen on retaining its account holders.
- XYZ bank intend to use Machine Learning techniques to predict, minimize churn based on customer data.

#### **About Dataset**

This dataset is for XYZ multistate bank with following columns:

- Customer ID: A unique identifier for each customer
- Surname: The customer's surname or last name
- Credit Score: A numerical value representing the customer's credit score
- Geography: The country where the customer resides (France, Spain or Germany)
- Gender: The customer's gender (Male or Female)
- Age: The customer's age
- Tenure: The number of years the customer has been with the bank
- Balance: The customer's account balance
- NumOfProducts: The number of bank products the customer uses (e.g., savings account, credit card)
- HasCrCard: Whether the customer has a credit card (1 = yes, 0 = no)
- IsActiveMember: Whether the customer is an active member (1 = yes, 0 = no)
- EstimatedSalary: The estimated salary of the customer
- Exited: Whether the customer has churned (1 = yes, 0 = no)

# **Objectives**

- Data cleansing, data preparation, data preprocessing
- Exploratory Data Analysis(EDA)
- Build and optimize machine learning models using a pipeline and crossvalidation to tune its hyperparameters:

- Random Forest
- Logistic Regression
- Model evaluation using metrix recall, roc auc score, f1-score, confusion matrix
- Extract feature importance from trained pipeline
- Conclusion

```
In [3]: # Import libraries
        import numpy as np
        import pandas as pd
        import matplotlib.pyplot as plt
        %matplotlib inline
        import seaborn as sns
        from sklearn.compose import ColumnTransformer
        from sklearn.impute import SimpleImputer
        from sklearn.preprocessing import StandardScaler, OneHotEncoder
        from sklearn.decomposition import PCA
        from sklearn.model selection import train test split, GridSearchCV, cross val score, StratifiedKFold
        from sklearn.pipeline import Pipeline
        from sklearn.linear model import LogisticRegression
        from sklearn.ensemble import RandomForestClassifier
        from sklearn.metrics import confusion matrix, classification report, ConfusionMatrixDisplay, roc auc score, accuracy
        import warnings
        warnings.filterwarnings('ignore')
In [4]: # Load dataset
        train = pd.read_csv('train.csv')
In [5]: train.shape
Out[5]: (165034, 14)
```

There are total 165034 rows and 14 columns in the dataset.

# **Data Preparation and Explore**

```
In [8]: # replace empty space with Nan value
    train.replace(' ', np.nan, inplace=True)

In [9]: # check missing values for each column
    missing = train.isnull()
    l = missing.columns.values.tolist()
    for col in 1:
        print(missing[col].value_counts(), '\n')
```

False 165034

Name: count, dtype: int64

CustomerId

False 165034

Name: count, dtype: int64

Surname

False 165034

Name: count, dtype: int64

CreditScore

False 165034

Name: count, dtype: int64

Geography

False 165034

Name: count, dtype: int64

Gender

False 165034

Name: count, dtype: int64

Age

False 165034

Name: count, dtype: int64

Tenure

False 165034

Name: count, dtype: int64

Balance

False 165034

Name: count, dtype: int64

NumOfProducts

False 165034

Name: count, dtype: int64

HasCrCard

False 165034

```
Name: count, dtype: int64
        IsActiveMember
        False
                165034
       Name: count, dtype: int64
       EstimatedSalary
        False
                165034
       Name: count, dtype: int64
        Exited
        False
                165034
       Name: count, dtype: int64
In [10]: # check and summarize missing values in each column
         train.isnull().sum()
Out[10]: id
                            0
         CustomerId
                            0
         Surname
         CreditScore
         Geography
         Gender
                            0
         Age
         Tenure
         Balance
         NumOfProducts
                            0
         HasCrCard
         IsActiveMember
         EstimatedSalary
         Exited
         dtype: int64
         No missing data
In [12]: # check duplicates
         len(train.drop_duplicates())
Out[12]: 165034
```

#### No duplicates

```
In [14]: train.count()
Out[14]: id
                             165034
         CustomerId
                             165034
         Surname
                             165034
         CreditScore
                             165034
                             165034
         Geography
         Gender
                             165034
         Age
                             165034
         Tenure
                             165034
          Balance
                             165034
                             165034
         NumOfProducts
         HasCrCard
                             165034
         IsActiveMember
                             165034
          EstimatedSalary
                             165034
          Exited
                             165034
         dtype: int64
In [15]: # drop columns that won't be used for this analysis
         train= train.drop(columns=(['id','CustomerId','Surname']), axis=1)
         train.shape
Out[15]: (165034, 11)
         After data preparation, there are total 165034 rows and 11 columns remain.
In [17]: train.info()
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 165034 entries, 0 to 165033
Data columns (total 11 columns):
```

#	Column	Non-Null Count	Dtype
0	CreditScore	165034 non-null	int64
1	Geography	165034 non-null	object
2	Gender	165034 non-null	object
3	Age	165034 non-null	float64
4	Tenure	165034 non-null	int64
5	Balance	165034 non-null	float64
6	NumOfProducts	165034 non-null	int64
7	HasCrCard	165034 non-null	float64
8	IsActiveMember	165034 non-null	float64
9	EstimatedSalary	165034 non-null	float64
10	Exited	165034 non-null	int64
dtyp	es: float64(5), i	nt64(4), object(2	)

dtypes: float64(5), int64(4), object(2)
memory usage: 13.9+ MB

```
In [18]: # check basic statistical info
        train.describe(include='all')
```

	CreditScore	Geography	Gender	Age	Tenure	Balance	NumOfProducts	HasCrCard	IsAc
count	165034.000000	165034	165034	165034.000000	165034.000000	165034.000000	165034.000000	165034.000000	16
unique	NaN	3	2	NaN	NaN	NaN	NaN	NaN	
top	NaN	France	Male	NaN	NaN	NaN	NaN	NaN	
freq	NaN	94215	93150	NaN	NaN	NaN	NaN	NaN	
mean	656.454373	NaN	NaN	38.125888	5.020353	55478.086689	1.554455	0.753954	
std	80.103340	NaN	NaN	8.867205	2.806159	62817.663278	0.547154	0.430707	
min	350.000000	NaN	NaN	18.000000	0.000000	0.000000	1.000000	0.000000	
25%	597.000000	NaN	NaN	32.000000	3.000000	0.000000	1.000000	1.000000	
50%	659.000000	NaN	NaN	37.000000	5.000000	0.000000	2.000000	1.000000	
75%	710.000000	NaN	NaN	42.000000	7.000000	119939.517500	2.000000	1.000000	
max	850.000000	NaN	NaN	92.000000	10.000000	250898.090000	4.000000	1.000000	

Findings from describe():

Out[18]:

- There are more male customers than female customers.
- The number of French customers is higher than that of German and Spainish.
- Average of credit score is 656, age of customer is 38, tenure is 5 years, bank account balance is \$55478.
- Average of churn percentage is 21%

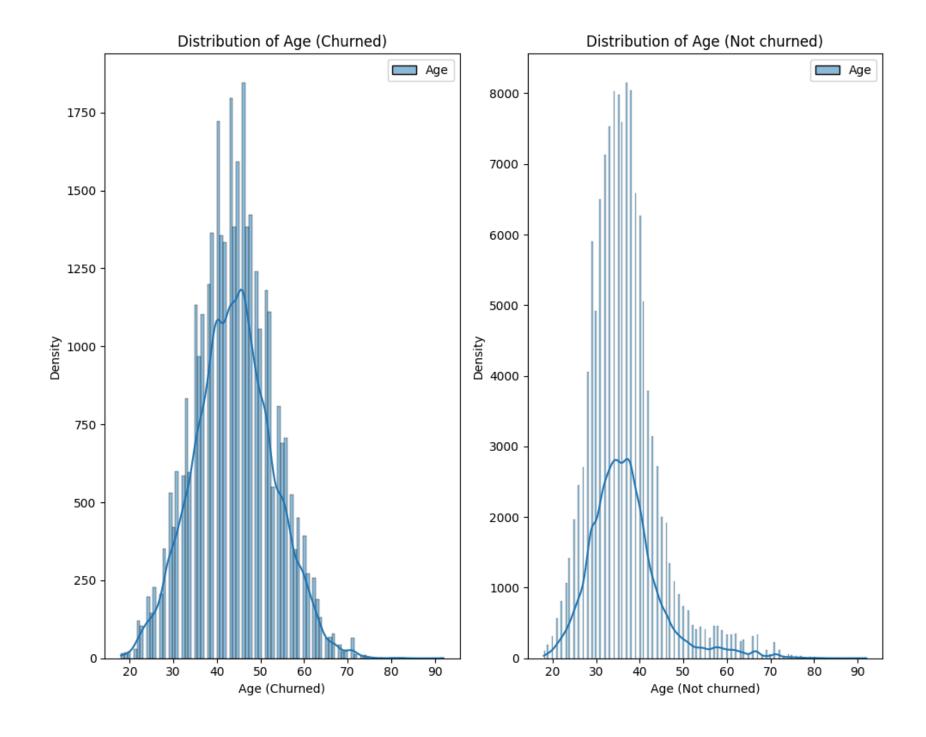
# **Exploratory Data Analysis (EDA)**

## **Age Distribution**

```
In [33]: # 1 --> churned    0 --> not churned
fig, axes = plt.subplots(nrows=1, ncols=2, figsize=(10,8))
sns.histplot(train.loc[train.Exited == 1] [['Age']],kde=True,color='grey', ax=axes[0])
axes[0].set_title('Distribution of Age (Churned)')
```

```
axes[0].set_xlabel('Age (Churned)')
axes[0].set_ylabel('Density')

sns.histplot(train.loc[train.Exited == 0] [['Age']],kde=True,color='grey', ax=axes[1])
axes[1].set_title('Distribution of Age (Not churned)')
axes[1].set_xlabel('Age (Not churned)')
axes[1].set_ylabel('Density')
plt.tight_layout()
plt.savefig('age_dis.png', dpi=300, bbox_inches='tight') # You can change the file format and name
plt.show()
```

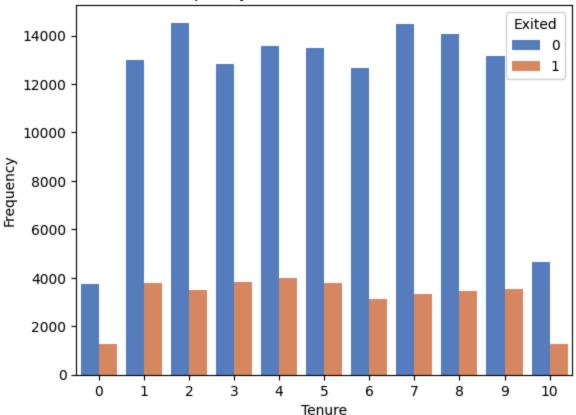


- The histogram on the left indicates that customers aged 40 to 50 has the highest churn rate.
- The histogram on the right indicates that customers aged 30 to 40 are less likely to churn.

### **Tenure Distribution**

```
In [37]: tenure = train.groupby(['Tenure','Exited']).size().reset_index(name='Frequency')
    sns.barplot(x='Tenure', y='Frequency', hue='Exited',data=tenure,palette='muted')
    plt.xlabel('Tenure')
    plt.ylabel('Frequency')
    plt.title('Frequency of churned and not churned')
    plt.savefig('tenure_dis.png', dpi=300, bbox_inches='tight') # You can change the file format and name
    plt.show()
```

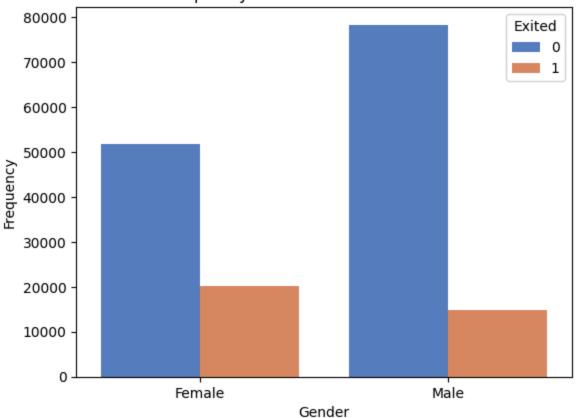
# Frequency of churned and not churned



- Customers who has the least tenure(new customers) and the most tenure (long term/ loyal customers) less likely to churn.
- Customers who has 1 to 9 years tenure have similar churn rate.

## **Gender Distribution**

# Frequency of churned and not churned



- The bar chart shows that the bank has about 32% more male customers than femal customers.
- It also shows that there is no significant difference in churn between genders.

# Target (Churn) Distribution

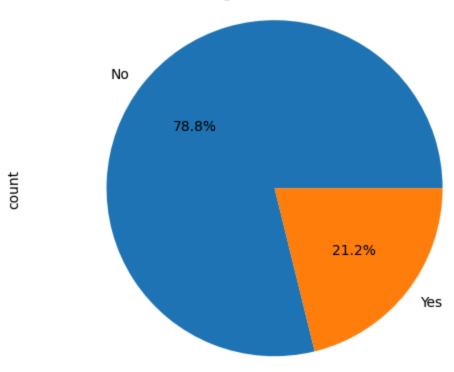
```
In [45]: train.Exited.value_counts()
```

# Out[45]: Exited

0 130113 1 34921

Name: count, dtype: int64

# **Target Distribution**



# **Data Preprocessing**

Split data into training and testing datasets

```
In [53]: X = train.drop('Exited',axis=1)
y = train.Exited
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, stratify=y, random_state=42)
```

Auto detect numerical and categorical columns, assign them to separate features

```
In [56]: # select numerical features
num_features = X_train.select_dtypes('number').columns.tolist()

# Select categorical features
cat_features = X_train.select_dtypes(['object','category']).columns.tolist()
```

Using Pipeline class from Scikit-Learn for streamlining data preprocessing and model training into a single, coherent sequence.

Define separate preprocessing pipelines for both feature types

Combine the transformers into a single column transformer

# Modeling

**Model 1 Random Forest** 

#### Create a model pipeline

Define a parameter grid

Use the grid in a cross validation search to optimize the model

```
In [72]: param_grid_rf = {
          'classifier__n_estimators': [50, 100],
          'classifier__max_depth': [None, 10, 20],
          'classifier__min_samples_split': [2, 5]
}
```

Perform grid search cross-validation and fit the best model to the training data

```
In [75]: cv = StratifiedKFold(n_splits=3, shuffle=True)
```

Train the pipeline model

```
In [78]: rf_model = GridSearchCV(estimator=pipe_rf, param_grid=param_grid_rf, cv=cv, scoring='accuracy', verbose=2)
    rf_model.fit(X_train, y_train)
```

```
Fitting 3 folds for each of 12 candidates, totalling 36 fits
[CV] END classifier max depth=None, classifier min samples split=2, classifier n estimators=50; total time=
                                                                                                              9.0s
[CV] END classifier max depth=None, classifier min samples split=2, classifier n estimators=50; total time=
                                                                                                              7.7s
[CV] END classifier max depth=None, classifier min samples split=2, classifier n estimators=50; total time=
                                                                                                              7.7s
[CV] END classifier max_depth=None, classifier min_samples_split=2, classifier n_estimators=100; total time= 12.1
S
[CV] END classifier max depth=None, classifier min samples split=2, classifier n estimators=100; total time= 12.1
[CV] END classifier max depth=None, classifier min samples split=2, classifier n estimators=100; total time= 14.4
[CV] END classifier max depth=None, classifier min samples split=5, classifier n estimators=50; total time=
                                                                                                              7.9s
[CV] END classifier max depth=None, classifier min samples split=5, classifier n estimators=50; total time=
                                                                                                              8.6s
[CV] END classifier max depth=None, classifier min samples split=5, classifier n estimators=50; total time=
                                                                                                              8.0s
[CV] END classifier max depth=None, classifier min samples split=5, classifier n estimators=100; total time= 16.3
S
[CV] END classifier max depth=None, classifier min samples split=5, classifier n estimators=100; total time= 15.8
S
[CV] END classifier max_depth=None, classifier min_samples_split=5, classifier n_estimators=100; total time= 15.2
[CV] END classifier max depth=10, classifier min samples split=2, classifier n estimators=50; total time=
                                                                                                            4.5s
[CV] END classifier max_depth=10, classifier min_samples_split=2, classifier n estimators=50; total time=
                                                                                                            5.8s
[CV] END classifier max depth=10, classifier min samples split=2, classifier n estimators=50; total time=
                                                                                                            5.3s
[CV] END classifier max depth=10, classifier min samples split=2, classifier n estimators=100; total time=
                                                                                                             9.9s
[CV] END classifier max depth=10, classifier min samples split=2, classifier n estimators=100; total time=
                                                                                                             9.0s
[CV] END classifier max depth=10, classifier min samples split=2, classifier n estimators=100; total time=
                                                                                                             9.0s
[CV] END classifier max depth=10, classifier min samples split=5, classifier n estimators=50; total time=
                                                                                                            3.5s
[CV] END classifier max depth=10, classifier min samples split=5, classifier n estimators=50; total time=
                                                                                                            3.9s
[CV] END classifier max_depth=10, classifier min_samples_split=5, classifier n_estimators=50; total time=
                                                                                                            4.5s
[CV] END classifier max depth=10, classifier min samples split=5, classifier n estimators=100; total time=
                                                                                                             9.5s
[CV] END classifier max depth=10, classifier min samples split=5, classifier n estimators=100; total time=
                                                                                                             9.5s
[CV] END classifier max depth=10, classifier min samples split=5, classifier n estimators=100; total time=
                                                                                                             8.3s
[CV] END classifier max depth=20, classifier min samples split=2, classifier n estimators=50; total time=
                                                                                                            8.1s
[CV] END classifier__max_depth=20, classifier__min_samples_split=2, classifier__n_estimators=50; total time=
                                                                                                            6.4s
[CV] END classifier max depth=20, classifier min samples split=2, classifier n estimators=50; total time=
                                                                                                            6.8s
[CV] END classifier max depth=20, classifier min samples split=2, classifier n estimators=100; total time= 15.3s
[CV] END classifier max depth=20, classifier min samples split=2, classifier n estimators=100; total time= 15.0s
[CV] END classifier max_depth=20, classifier min_samples_split=2, classifier n_estimators=100; total time= 14.0s
[CV] END classifier max depth=20, classifier min samples split=5, classifier n estimators=50; total time=
                                                                                                            6.7s
[CV] END classifier max depth=20, classifier min samples split=5, classifier n estimators=50; total time=
                                                                                                            7.8s
[CV] END classifier max depth=20, classifier min samples split=5, classifier n estimators=50; total time=
                                                                                                            7.9s
[CV] END classifier max depth=20, classifier min samples split=5, classifier n estimators=100; total time= 16.1s
```

Get the model predictions from the grid search estimator on the test data

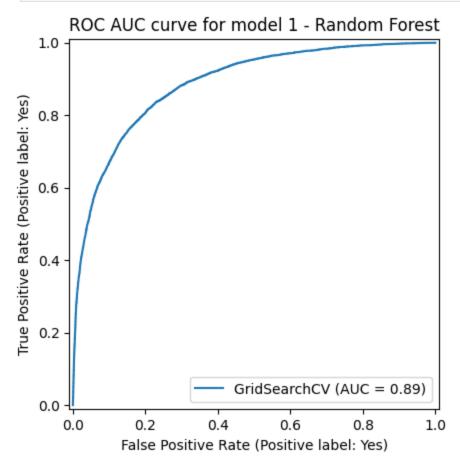
## **Classification Report for Model 1 Random Forest**

```
In [81]: y_pred_rf = rf_model.predict(X_test)
         print(classification_report(y_test, y_pred_rf))
                      precision
                                   recall f1-score
                                                      support
                  No
                           0.88
                                     0.96
                                               0.92
                                                         26023
                 Yes
                           0.77
                                     0.50
                                               0.61
                                                         6984
                                               0.86
                                                         33007
            accuracy
           macro avg
                           0.82
                                     0.73
                                               0.76
                                                         33007
        weighted avg
                           0.85
                                     0.86
                                               0.85
                                                         33007
```

**ROC AUC for Model 1 Random Forest** 

```
In [83]: y_prob_rf = rf_model.predict_proba(X_test)[:,1]
    print('Model 1 (Random Forest) ROC AUC score is:', roc_auc_score(y_test, y_prob_rf))
```

```
In [84]: RocCurveDisplay.from_estimator(rf_model, X_test, y_test)
   plt.title('ROC AUC curve for model 1 - Random Forest')
   plt.savefig('roc_curve_rf_model.png', dpi=300, bbox_inches='tight') # You can change the file format and name
   plt.show()
```

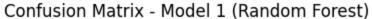


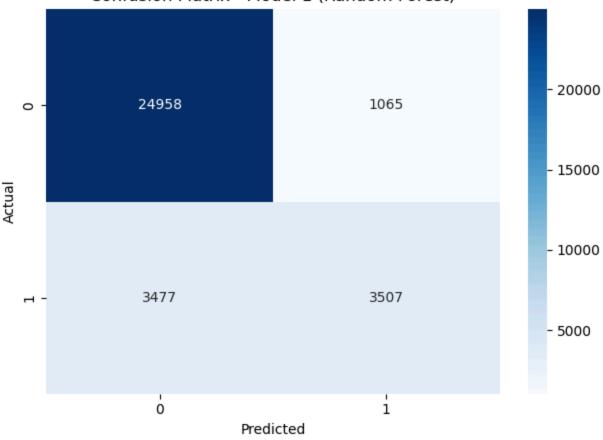
Plot the confusion matrix

```
In [90]: conf_matrix = confusion_matrix(y_test, y_pred_rf)
   plt.figure()
   sns.heatmap(conf_matrix, annot=True, cmap='Blues', fmt='d')

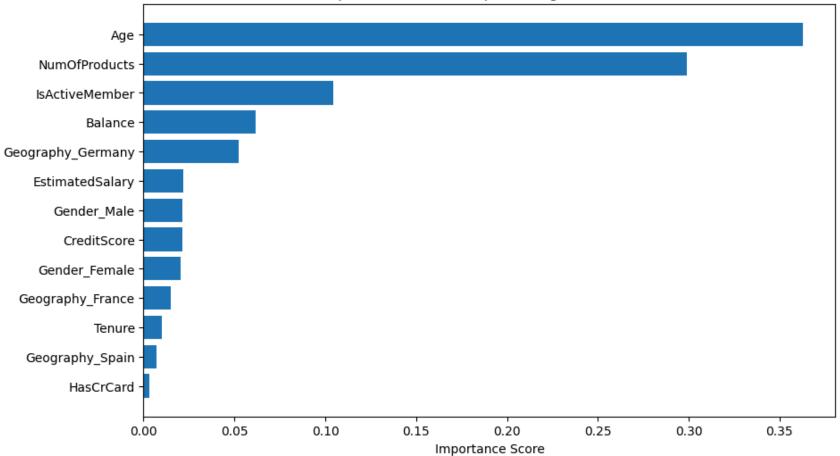
plt.title('Confusion Matrix - Model 1 (Random Forest)')
```

```
plt.xlabel('Predicted')
plt.ylabel('Actual')
plt.tight_layout()
plt.savefig('confusion_matrix_rf_model.png', dpi=300, bbox_inches='tight') # You can change the file format and name
plt.show()
```





# Most Important Features in predicting for Random Forest



Test set accuracy: 86.24%

The Top 5 Importance features in Model 1 Random Forest:

```
importance_df_r = importance_df_r.sort_values(by='Importance', ascending=False)
importance_df_r[:5]
```

$\sim$		1_	г	1	$\overline{}$	4
()	ш.	г				
					U	

	Feature	Importance
1	Age	0.362650
4	NumOfProducts	0.299002
6	IsActiveMember	0.104191
3	Balance	0.061750
9	Geography_Germany	0.052207

The top five important features for model 1 Random Forest are: Age, NumOfProducts, IsActiveMember, Balance, Geography\_Germany.

# **Model 2 Logistic Regression**

Create a model pipeline and define a hyperparameter grid

Use the grid in a cross validation search to optimize the model

Perform grid search cross-validation and fit the best model to the training data

```
In [114... lr_model = GridSearchCV(estimator=pipe_lr, param_grid=param_grid_lr, cv=cv, scoring='accuracy', verbose=2)
# Train the pipeline model
lr_model.fit(X_train, y_train)
```

```
Fitting 3 folds for each of 4 candidates, totalling 12 fits
         [CV] END classifier class weight=None, classifier penalty=11, classifier solver=liblinear; total time=
                                                                                                                   2.5s
        [CV] END classifier__class_weight=None, classifier__penalty=11, classifier__solver=liblinear; total time=
                                                                                                                   2.8s
        [CV] END classifier__class_weight=None, classifier__penalty=11, classifier__solver=liblinear; total time=
                                                                                                                   2.5s
        [CV] END classifier__class_weight=None, classifier__penalty=12, classifier__solver=liblinear; total time=
                                                                                                                   0.2s
         [CV] END classifier class weight=None, classifier penalty=12, classifier solver=liblinear; total time=
                                                                                                                   0.3s
        [CV] END classifier class weight=None, classifier penalty=12, classifier solver=liblinear; total time=
                                                                                                                   0.3s
        [CV] END classifier class weight=balanced, classifier penalty=11, classifier solver=liblinear; total time=
                                                                                                                       0.4s
         [CV] END classifier class weight=balanced, classifier penalty=11, classifier solver=liblinear; total time=
                                                                                                                       9.9s
        [CV] END classifier class weight=balanced, classifier penalty=11, classifier solver=liblinear; total time=
                                                                                                                       0.4s
         [CV] END classifier class weight=balanced, classifier penalty=12, classifier solver=liblinear; total time=
                                                                                                                       0.3s
        [CV] END classifier__class_weight=balanced, classifier__penalty=12, classifier__solver=liblinear; total time=
                                                                                                                       0.4s
        [CV] END classifier class weight=balanced, classifier penalty=12, classifier solver=liblinear; total time=
                                                                                                                       0.4s
Out[114...
                                        GridSearchCV
                                 best estimator : Pipeline
                              preprocessor: ColumnTransformer
                                                            cat
                              num
                    StandardScaler
                                                   ▶ OneHotEncoder
                                 ▶ LogisticRegression
```

### Classification Report for Model 2 Logistic Regression

```
In [116... y_pred_lr = lr_model.predict(X_test)
    print(classification_report(y_test, y_pred_lr))
```

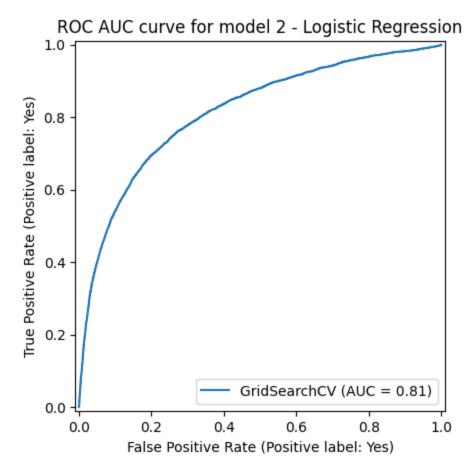
	precision	recall	f1-score	support
No	0.85	0.95	0.90	26023
Yes	0.69	0.38	0.49	6984
accuracy			0.83	33007
macro avg	0.77	0.67	0.70	33007
weighted avg	0.82	0.83	0.81	33007

### ROC AUC for Model 2 Logistic Regression

```
In [119... y_prob_lr = lr_model.predict_proba(X_test)[:,1]
    print('Model 2 (Logistic Regression) ROC AUC score is:', roc_auc_score(y_test, y_prob_lr))

Model 2 (Logistic Regression) ROC AUC score is: 0.8144870655657108

In [121... RocCurveDisplay.from_estimator(lr_model, X_test, y_test)
    plt.title('ROC AUC curve for model 2 - Logistic Regression')
    plt.savefig('roc_curve_lr_model.png', dpi=300, bbox_inches='tight') # You can change the file format and name
    plt.show()
```

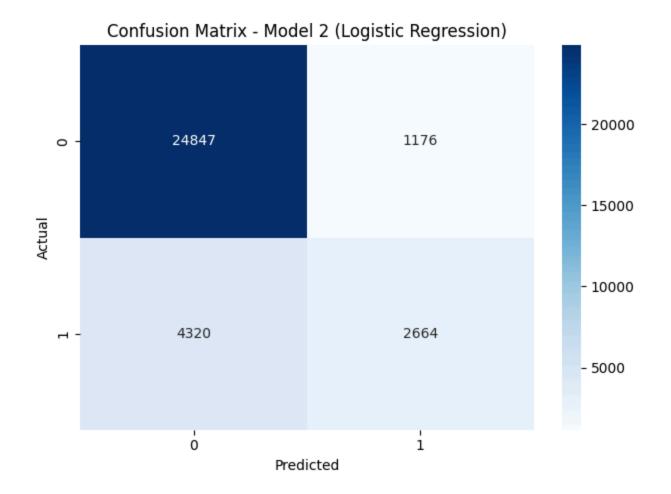


```
In [123... conf_matrix = confusion_matrix(y_test, y_pred_lr)

plt.figure()
    sns.heatmap(conf_matrix, annot=True, cmap='Blues', fmt='d')

plt.title('Confusion Matrix - Model 2 (Logistic Regression)')
    plt.xlabel('Predicted')
    plt.ylabel('Actual')

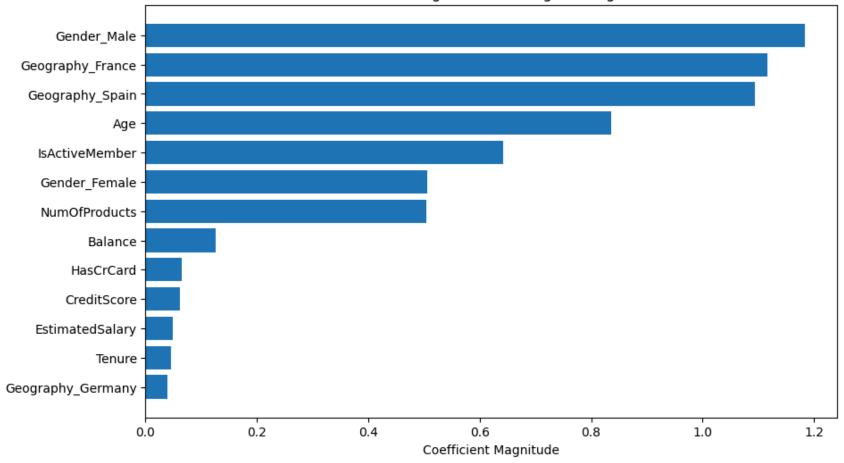
plt.tight_layout()
    plt.savefig('confusion_matrix_lr_model.png', dpi=300, bbox_inches='tight') # You can change the file format and name plt.show()
```



## Feature Importance

```
In [132... # Create a DataFrame for the coefficients
          importance_df_lr = pd.DataFrame({
              'Feature': feature_names,
              'Coefficient': coefficients
          }).sort_values(by='Coefficient', key=abs) # Sort by absolute values
          # Plotting
          plt.figure(figsize=(10, 6))
          plt.barh(importance_df_lr['Feature'], importance_df_lr['Coefficient'].abs())
          #plt.gca().invert_yaxis()
          plt.title('Feature Coefficient magnitudes for Logistic Regression model')
          plt.xlabel('Coefficient Magnitude')
          plt.savefig('feature_importance_lr_model.png', dpi=300, bbox_inches='tight') # You can change the file format and no
          plt.show()
          # Print test score
          test_score_lr = lr_model.best_estimator_.score(X_test, y_test)
          print(f"\nTest set accuracy: {test_score_lr:.2%}")
```

# Feature Coefficient magnitudes for Logistic Regression model



Test set accuracy: 83.35%

The Top 5 Important features in Model 2 Logistic Regression:

```
importance_df_lr['Coefficient']=importance_df_lr['Coefficient'].abs()
importance_df_lr = importance_df_lr.sort_values(by='Coefficient', ascending=False)
importance_df_lr[:5]
```

O.		г	-	$\neg$	-	
Ul	Jτ	П	Т	3	5	

	Feature	Coefficient
12	Gender_Male	1.183180
8	Geography_France	1.116289
10	Geography_Spain	1.093797
1	Age	0.836526
6	IsActiveMember	0.641585

The top five important features for model 2 Logistic Regression are: Gender\_Male, Geography\_France, Geography\_Spain, Age, IsActiveMember.

## Classification Report

#### In [139...

## print(classification\_report(y\_test, y\_pred\_rf))

	precision	recall	f1-score	support
No	0.88	0.96	0.92	26023
Yes	0.77	0.50	0.61	6984
accuracy			0.86	33007
macro avg	0.82	0.73	0.76	33007
weighted avg	0.85	0.86	0.85	33007

#### In [141...

### print(classification\_report(y\_test, y\_pred\_lr))

	precision	recall	f1-score	support
No	0.85	0.95	0.90	26023
Yes	0.69	0.38	0.49	6984
	0.05	0.50	0.15	0301
accuracy			0.83	33007
macro avg	0.77	0.67	0.70	33007
weighted avg	0.82	0.83	0.81	33007

```
In [143... print(f'Model 1 Random Forest accuracy score is: {accuracy_score(y_test, y_pred_rf)*100:.2f}%')
print()
print(f'Model 2 Logistic Regression accuracy score is: {accuracy_score(y_test, y_pred_lr)*100:.2f}%')

Model 1 Random Forest accuracy score is: 86.24%

Model 2 Logistic Regression accuracy score is: 83.35%
```

# Conclusion

• Random Forest model performs better than Logistic Regression model in predicting XYZ bank customer churn with higher roc auc score and accuracy score.

```
In [ ]:
```