

MSIN0097 Individual Coursework - Credit Card Churn Prediction-2

April 20, 2024

Word Count: 1392

```
[303]: import pandas as pd
import numpy as np
import matplotlib.pyplot as plt
import seaborn as sns
import tensorflow as tf
import warnings

from pandas.plotting import scatter_matrix
from sklearn.preprocessing import StandardScaler, OneHotEncoder
from sklearn.model_selection import train_test_split
from sklearn.compose import ColumnTransformer
from sklearn.linear_model import LogisticRegression
from sklearn.metrics import accuracy_score, confusion_matrix, \
    classification_report, roc_auc_score
from sklearn.metrics import roc_curve, precision_recall_curve, roc_auc_score, \
    PrecisionRecallDisplay, precision_recall_fscore_support
from sklearn.svm import SVC
from sklearn.ensemble import RandomForestClassifier
from xgboost import XGBClassifier
from sklearn.neural_network import MLPClassifier
from numpy import mean
from sklearn.datasets import make_classification
from sklearn.model_selection import GridSearchCV, cross_val_score
from sklearn.model_selection import RepeatedStratifiedKFold
from xgboost import XGBClassifier
from tensorflow.keras import layers, models, losses, metrics
from tensorflow.keras.datasets import mnist
from tensorflow import keras
from tensorflow.keras.optimizers import SGD
from keras.models import Sequential
from sklearn.model_selection import cross_val_score
from sklearn.model_selection import KFold
warnings.filterwarnings("ignore")
```

Business Problem: A consumer credit card bank is facing the problem of customer attrition. The dataset contains a series of customer data including demographic characteristics and and historical

activities, along with a variable 'Attrition_Flag' to label existing customer or attrited customer. The dataset is downloaded from Kaggle: <https://www.kaggle.com/datasets/anwarsan/credit-card-bank-churn/data>

Overall Strategy: The object is to build a predictive model(supervised learning) to predict whether a customer will become attrited or not based on the information in the dataset. For those customers who are labeled as potentially attrited customer by the model, the bank could implement a series of customer retention strategies to prevent customer loss. For example, the bank could offer them discount on annual fee/interest rate or limited time bonus etc.

1 Get the Data

1.1 Download the Data

```
[255]: data=pd.read_csv('credit_card_churn.csv')
data = data.iloc[:, 1:-2] #Delete last two column advised by the dataset
      ↪publication
```

1.2 Take a quick look at the data structure

```
[256]: data.head()
```

```
[256]:
```

	Attrition_Flag	Customer_Age	Gender	Dependent_count	Education_Level	\
0	Existing Customer	45	M	3	High School	
1	Existing Customer	49	F	5	Graduate	
2	Existing Customer	51	M	3	Graduate	
3	Existing Customer	40	F	4	High School	
4	Existing Customer	40	M	3	Uneducated	

	Marital_Status	Income_Category	Card_Category	Months_on_book	\
0	Married	\$60K - \$80K	Blue	39	
1	Single	Less than \$40K	Blue	44	
2	Married	\$80K - \$120K	Blue	36	
3	Unknown	Less than \$40K	Blue	34	
4	Married	\$60K - \$80K	Blue	21	

	Total_Relationship_Count	Months_Inactive_12_mon	Contacts_Count_12_mon	\
0	5	1	3	
1	6	1	2	
2	4	1	0	
3	3	4	1	
4	5	1	0	

	Credit_Limit	Total_Revolving_Bal	Avg_Open_To_Buy	Total_Amt_Chng_Q4_Q1	\
0	12691.0	777	11914.0	1.335	
1	8256.0	864	7392.0	1.541	
2	3418.0	0	3418.0	2.594	
3	3313.0	2517	796.0	1.405	

4	4716.0	0	4716.0	2.175
---	--------	---	--------	-------

	Total_Trans_Amt	Total_Trans_Ct	Total_Ct_Chng_Q4_Q1	Avg_Utilization_Ratio
0	1144	42	1.625	0.061
1	1291	33	3.714	0.105
2	1887	20	2.333	0.000
3	1171	20	2.333	0.760
4	816	28	2.500	0.000

[257]: data.info()

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 10127 entries, 0 to 10126
Data columns (total 20 columns):
#   Column                                Non-Null Count  Dtype
---  -
0   Attrition_Flag                       10127 non-null  object
1   Customer_Age                         10127 non-null  int64
2   Gender                               10127 non-null  object
3   Dependent_count                      10127 non-null  int64
4   Education_Level                      10127 non-null  object
5   Marital_Status                       10127 non-null  object
6   Income_Category                     10127 non-null  object
7   Card_Category                       10127 non-null  object
8   Months_on_book                      10127 non-null  int64
9   Total_Relationship_Count            10127 non-null  int64
10  Months_Inactive_12_mon              10127 non-null  int64
11  Contacts_Count_12_mon              10127 non-null  int64
12  Credit_Limit                        10127 non-null  float64
13  Total_Revolving_Bal                 10127 non-null  int64
14  Avg_Open_To_Buy                     10127 non-null  float64
15  Total_Amt_Chng_Q4_Q1                10127 non-null  float64
16  Total_Trans_Amt                     10127 non-null  int64
17  Total_Trans_Ct                      10127 non-null  int64
18  Total_Ct_Chng_Q4_Q1                 10127 non-null  float64
19  Avg_Utilization_Ratio                10127 non-null  float64
dtypes: float64(5), int64(9), object(6)
memory usage: 1.5+ MB

No Null value found in the dataset
```

[258]: data.describe()

```
[258]:
```

	Customer_Age	Dependent_count	Months_on_book	\
count	10127.000000	10127.000000	10127.000000	
mean	46.325960	2.346203	35.928409	
std	8.016814	1.298908	7.986416	
min	26.000000	0.000000	13.000000	

25%	41.000000	1.000000	31.000000
50%	46.000000	2.000000	36.000000
75%	52.000000	3.000000	40.000000
max	73.000000	5.000000	56.000000

	Total_Relationship_Count	Months_Inactive_12_mon	\
count	10127.000000	10127.000000	
mean	3.812580	2.341167	
std	1.554408	1.010622	
min	1.000000	0.000000	
25%	3.000000	2.000000	
50%	4.000000	2.000000	
75%	5.000000	3.000000	
max	6.000000	6.000000	

	Contacts_Count_12_mon	Credit_Limit	Total_Revolving_Bal	\
count	10127.000000	10127.000000	10127.000000	
mean	2.455317	8631.953698	1162.814061	
std	1.106225	9088.776650	814.987335	
min	0.000000	1438.300000	0.000000	
25%	2.000000	2555.000000	359.000000	
50%	2.000000	4549.000000	1276.000000	
75%	3.000000	11067.500000	1784.000000	
max	6.000000	34516.000000	2517.000000	

	Avg_Open_To_Buy	Total_Amt_Chng_Q4_Q1	Total_Trans_Amt	Total_Trans_Ct	\
count	10127.000000	10127.000000	10127.000000	10127.000000	
mean	7469.139637	0.759941	4404.086304	64.858695	
std	9090.685324	0.219207	3397.129254	23.472570	
min	3.000000	0.000000	510.000000	10.000000	
25%	1324.500000	0.631000	2155.500000	45.000000	
50%	3474.000000	0.736000	3899.000000	67.000000	
75%	9859.000000	0.859000	4741.000000	81.000000	
max	34516.000000	3.397000	18484.000000	139.000000	

	Total_Ct_Chng_Q4_Q1	Avg_Utilization_Ratio
count	10127.000000	10127.000000
mean	0.712222	0.274894
std	0.238086	0.275691
min	0.000000	0.000000
25%	0.582000	0.023000
50%	0.702000	0.176000
75%	0.818000	0.503000
max	3.714000	0.999000

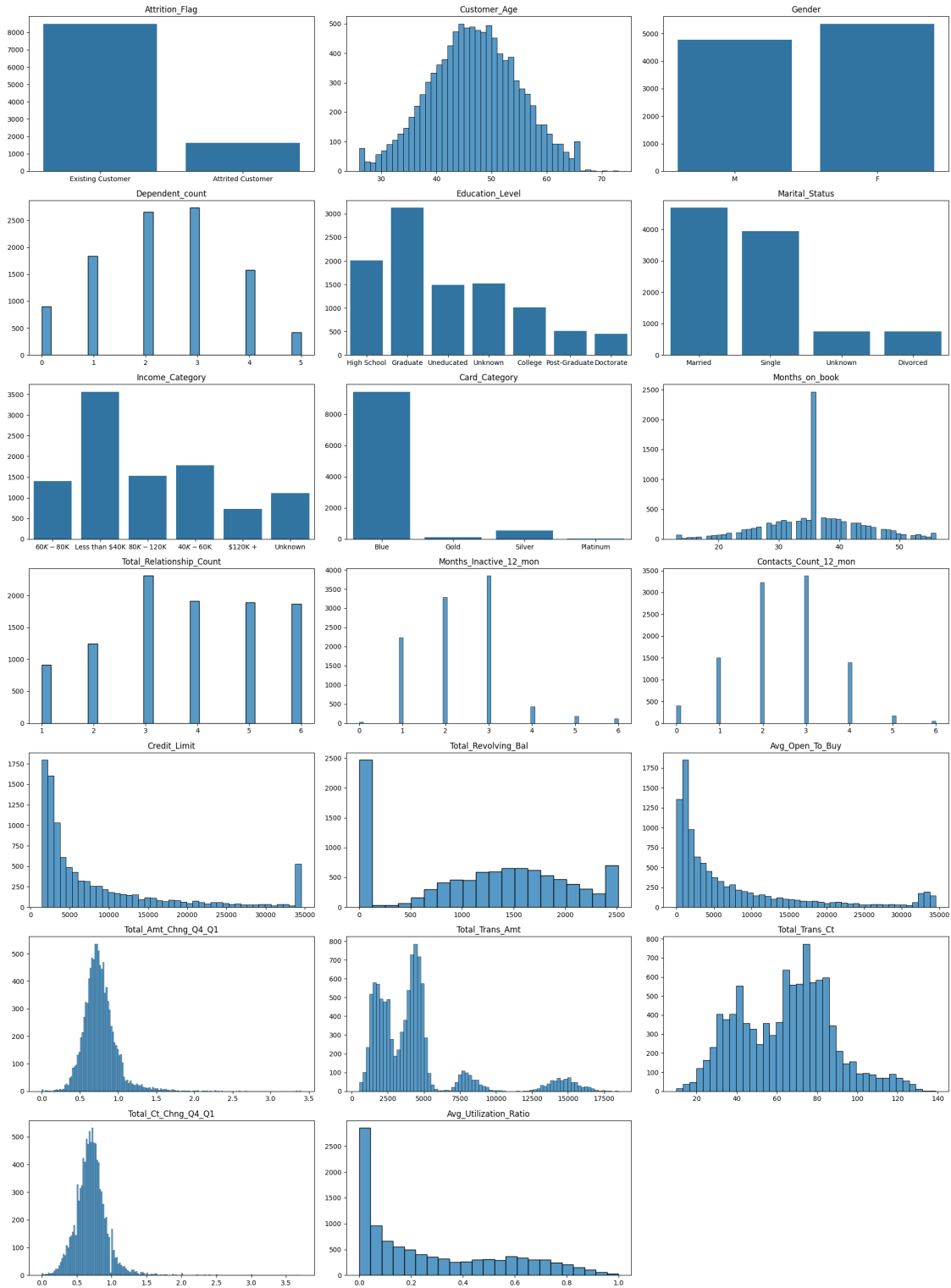
```
[259]: # Plot distributions for all variables
plt.figure(figsize=(20, 27))
```

```

for i, col in enumerate(data.columns):
    plt.subplot(7, 3, i+1)
    if data[col].dtype == 'object':
        sns.countplot(x=data[col], data=data) # countplot for categorical
        ↪ variables
    else:
        sns.histplot(data[col]) # histplot for numerical variables
    plt.title(col)
    plt.xlabel('')
    plt.ylabel('')

plt.tight_layout()
plt.show()

```



```
[260]: # Count the number of existing customers and calculate the percentage of
        ↪existing customers in the dataset
existing_count=data[(data['Attrition_Flag']=='Existing Customer')]
existing_per = existing_count.shape[0]/data.shape[0]
attrition_per = 1-existing_per
print(existing_per)
print(attrition_per)
```

```
0.8393403772094401
```

```
0.1606596227905599
```

From the first plot, I noticed that the dataset may have imbalanced classes since 84% of the sample are existing customers and only 16% are attrited customers. With the imbalanced classes, a model could get a pretty high accuracy just by predicting the majority class, but it may fail to capture the minority class. Below are some techniques to solve the imbalanced classes:

1. Resampling techniques like undersampling or oversampling (undersampling discards potentially useful information and oversampling increases the likelihood of overfitting since it replicates the minority class event)
2. Use various performance metrics: accuracy can be misleading for imbalanced datasets. Confusion matrix, precision, recall, F1 score, and area under ROC curve could be used for model evaluation as well
3. Penalize algorithms (cost-sensitive training) that increase the cost of classification mistakes in the minority class like Penalized-SVM and Weighted XGBoost
4. Use ensemble algorithms which performs well on imbalanced data like random forests and gradient boosted trees

Technique 2,3,4 will be used in this project to solve imbalanced classes after analyzing the pros and cons of the 4 options.

1.3 Create a test set

```
[261]: # Split the data into train(80%) and test(20%) set
        # The train set will be used in model training, and the test set will be used
        ↪to generate predicted value and assess model performance
train_set, test_set = train_test_split(data, test_size=0.2, random_state=42) #
        ↪set the random seed to ensure the same output every time
print(len(train_set))
print(len(test_set))
```

```
8101
```

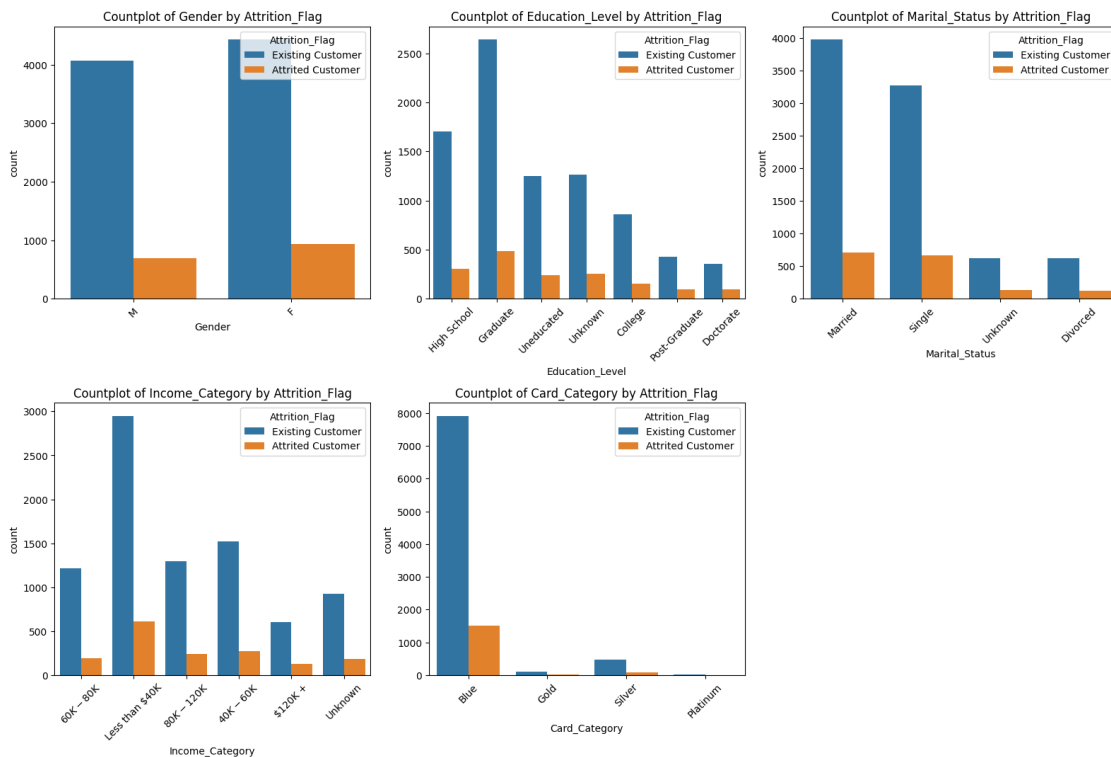
```
2026
```

2 Discover and Visualize the Data to Gain Insights

2.1 Countplot plots of categorical variables by attrition flags

I want to see whether categorical variables have the same distribution for existing and attrited customers

```
[262]: df_categorical = data.select_dtypes(include='object').
        drop(columns=['Attrition_Flag']) # Generate a set of categorical variables
plt.figure(figsize=(16, 15))
for i, column in enumerate(df_categorical.columns, 1):
    plt.subplot(3, 3, i)
    sns.countplot(x=column, hue='Attrition_Flag', data=data)
    plt.title(f'Countplot of {column} by Attrition_Flag')
    plt.xticks(rotation=45)
plt.tight_layout()
plt.show()
```



Some preliminary insights from EDA:

1. Married people are more likely to hold the card, followed by singles
2. Customers with graduate degree are more likely to hold the card, compared with other degrees
3. People whose income less than \$40k are most likely to subscribe to the credit card

2.2 Box plots of numerical variables by attrition flags

I want to see whether numerical variables have the same distribution for existing and attrited customers

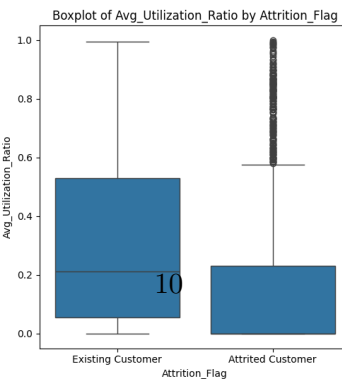
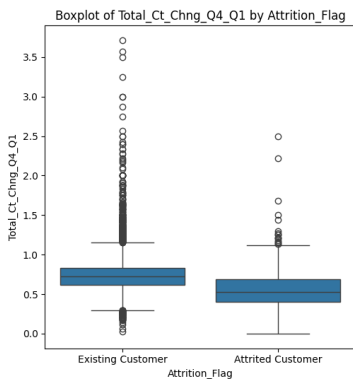
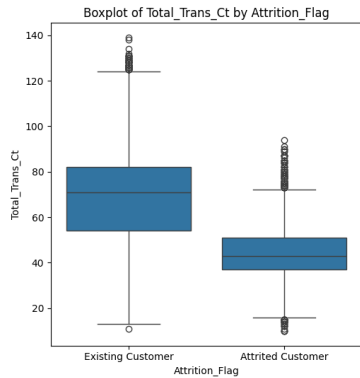
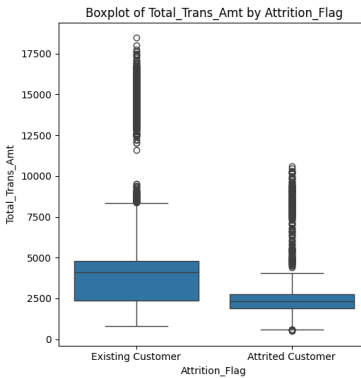
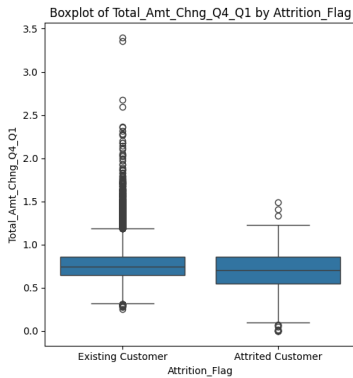
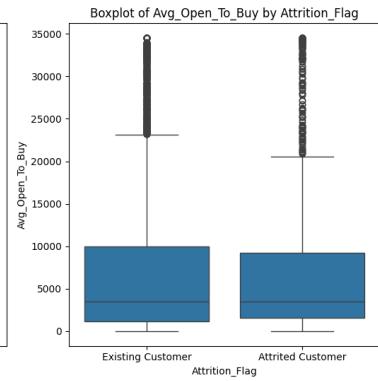
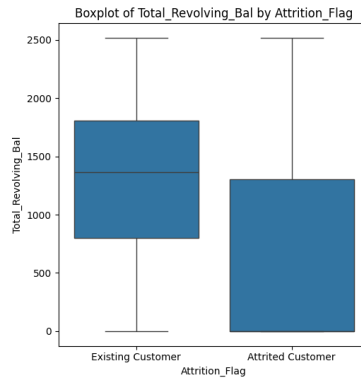
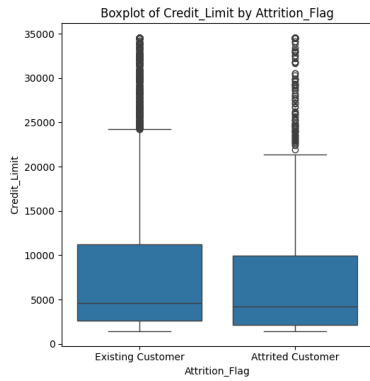
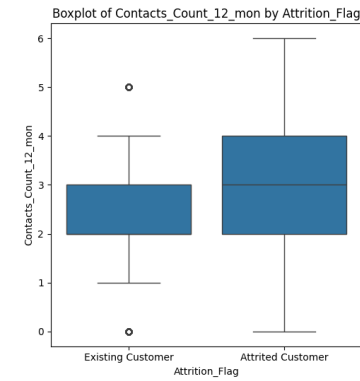
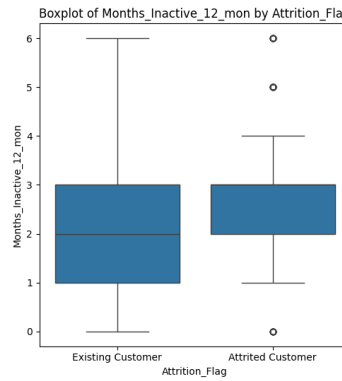
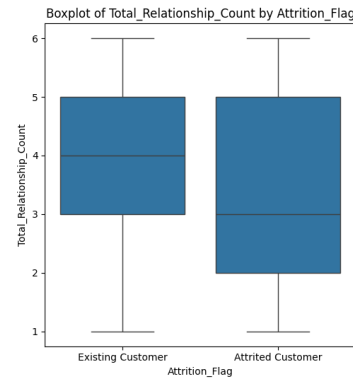
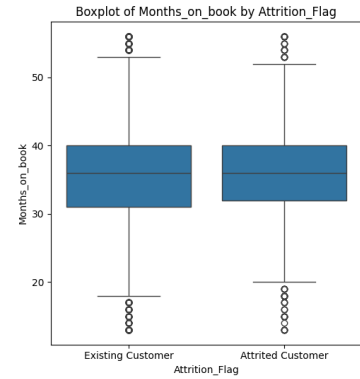
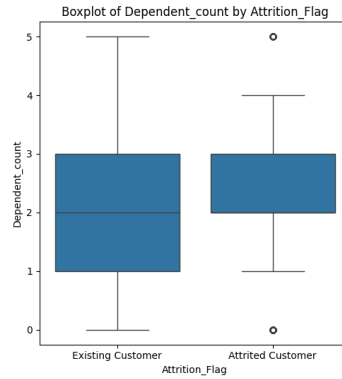
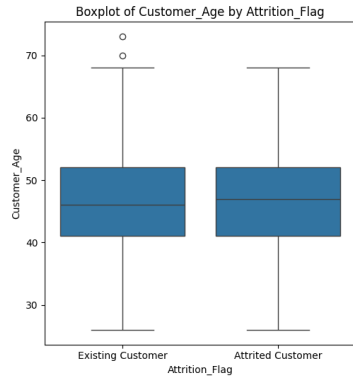
[263]:


```

df_numerical = data.select_dtypes(include=['int64', 'float64']) # Generate a
↳ set of numerical variables
plt.figure(figsize=(16, 27))
num_cols = len(df_numerical.columns)
rows = 5
cols = num_cols // rows + (num_cols % rows > 0)

for i, column in enumerate(df_numerical.columns, 1):
    plt.subplot(rows, cols, i)
    sns.boxplot(x='Attrition_Flag', y=column, data=data)
    plt.title(f'Boxplot of {column} by Attrition_Flag')
plt.tight_layout()
plt.show()

```

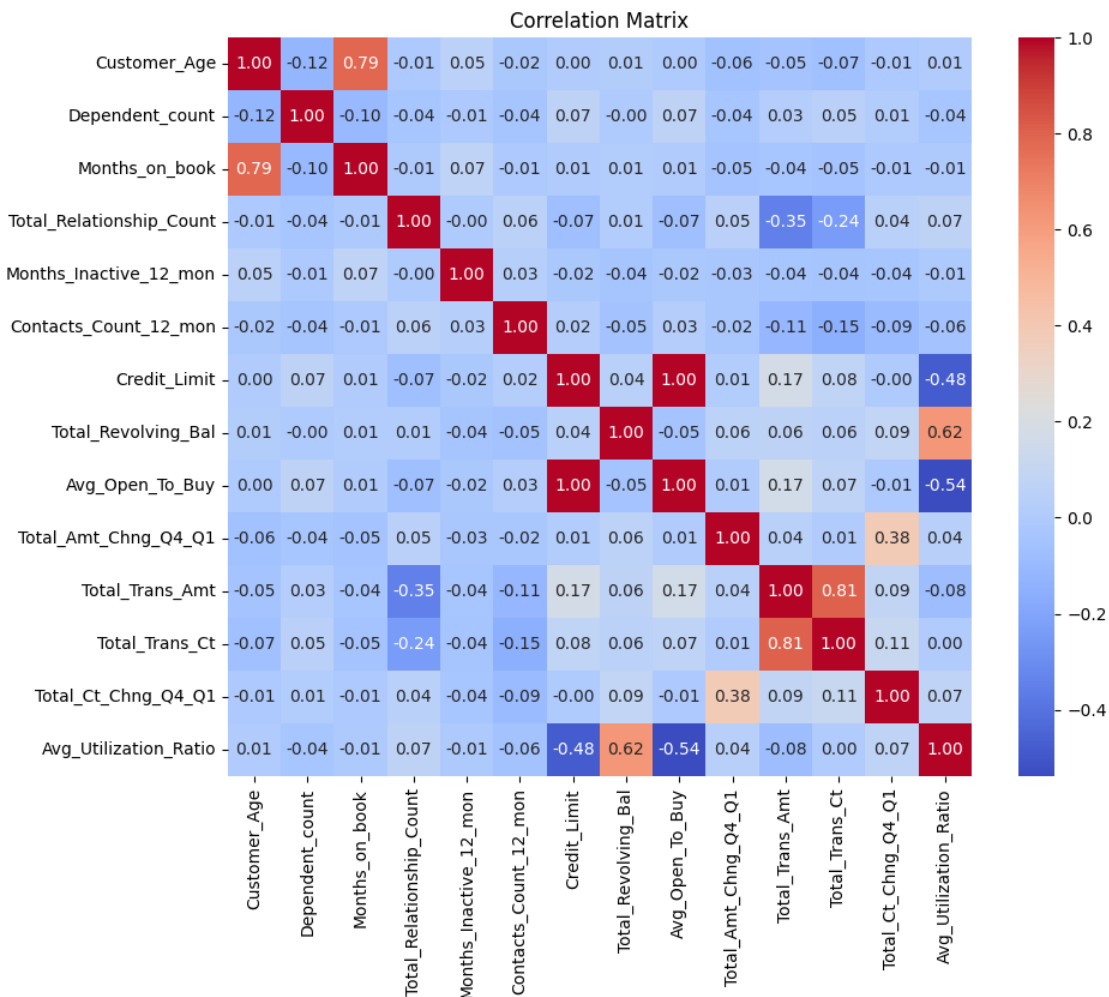


Some preliminary insights from EDA:

1. Attrited customers tend to have less product with the bank
2. Attrited customers tend to have more contacts with the bank in last 12 month
3. Attrited customers tend to have lower revolving balance and less transactions

2.3 Looking for Correlations

```
[264]: # Plot correation matrix for all numerical variables
corr_matrix = df_numerical.corr()
plt.figure(figsize=(10, 8))
sns.heatmap(corr_matrix, annot=True, cmap='coolwarm', fmt=".2f",
            annot_kws={"size": 10})
plt.title('Correlation Matrix')
plt.show()
```

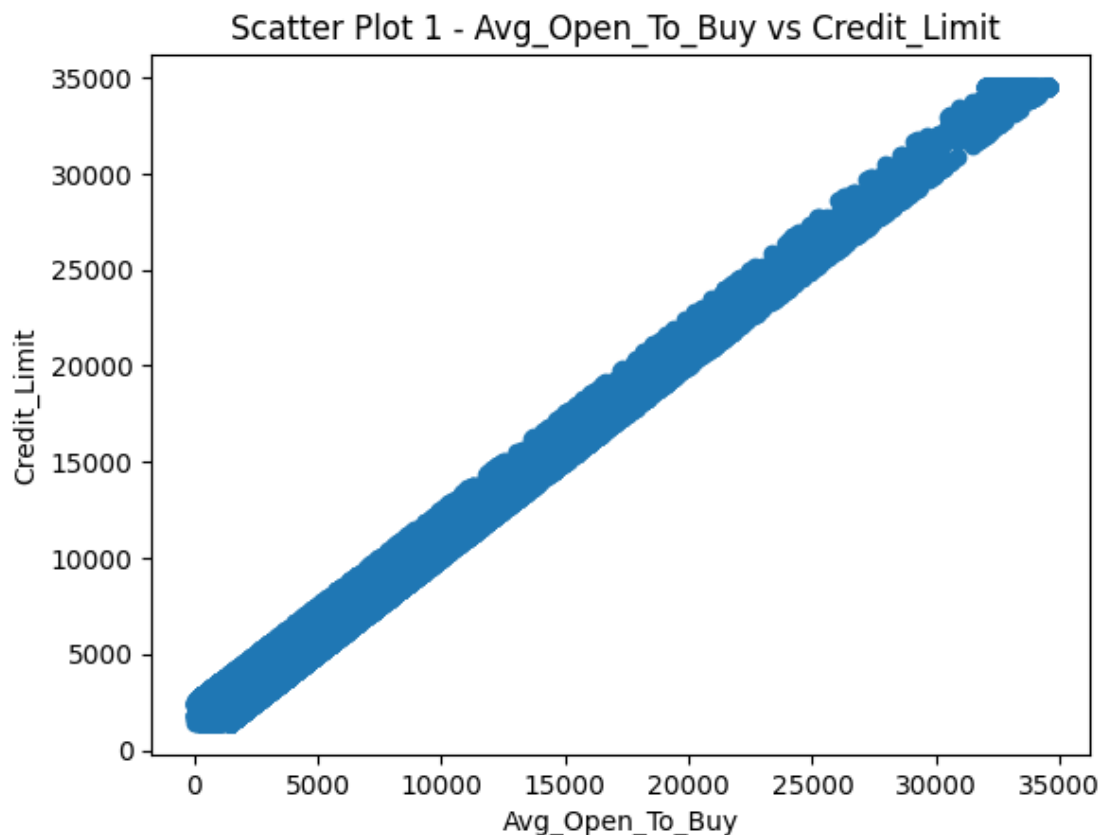


Three pairs of variables have high correlation(≥ 0.7) : Avg_Open_To_Buy & Credit_Limit, Months_on_book & Customer_Age, Total_Trans_Ct & Total_Trans_Amt. Scatter plot will be plotted to further investigate the correlation.

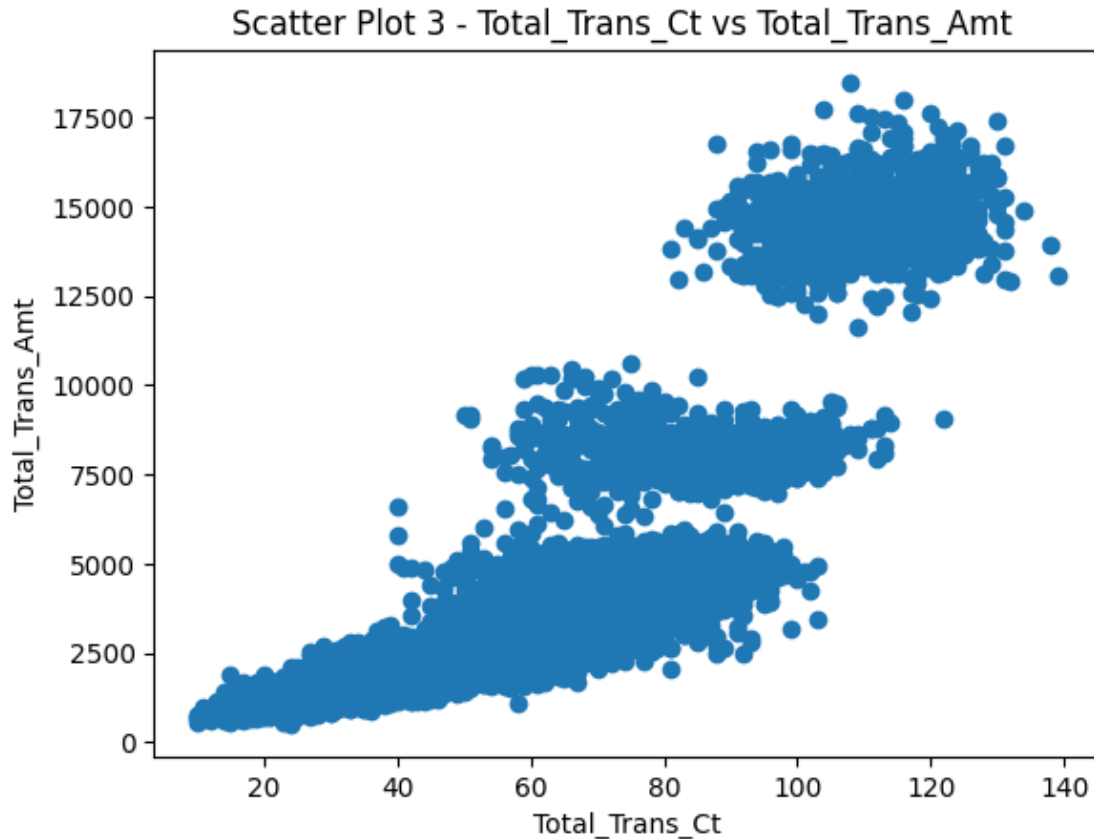
```
[265]: plt.scatter(df_numerical['Avg_Open_To_Buy'], df_numerical['Credit_Limit'])
plt.xlabel('Avg_Open_To_Buy')
plt.ylabel('Credit_Limit')
plt.title('Scatter Plot 1 - Avg_Open_To_Buy vs Credit_Limit ')
plt.show()

plt.scatter(df_numerical['Months_on_book'], df_numerical['Customer_Age'])
plt.xlabel('Months_on_book')
plt.ylabel('Customer_Age')
plt.title('Scatter Plot 2 - Months_on_book vs Customer_Age ')
plt.show()

plt.scatter(df_numerical['Total_Trans_Ct'], df_numerical['Total_Trans_Amt'])
plt.xlabel('Total_Trans_Ct')
plt.ylabel('Total_Trans_Amt')
plt.title('Scatter Plot 3 - Total_Trans_Ct vs Total_Trans_Amt ')
plt.show()
```







Due to the high colinearity indicated in the scatter plots above, I decide to delete Avg_Open_To_Buy, Customer_Age, and Total_Trans_Amt to reduce redundancy and improve stability and generalization of the model

```
[266]: columns_to_exclude = ['Customer_Age', 'Avg_Open_To_Buy', 'Total_Trans_Amt']
data_new = data.drop(columns=columns_to_exclude)
```

3 Prepare the Data for Machine Learning Algorithms

3.1 Deal with missing values

```
[267]: #Count number of rows that have at least one 'Unkown' value
filtered_df=data_new[(data_new['Income_Category']=='Unknown')|
                     (data_new['Education_Level']=='Unknown')|
                     (data_new['Marital_Status']=='Unknown')]
num_rows = filtered_df.shape[0]
print(num_rows)
```

3046

Education_Level, Marital_Status, and Income_Category have 'Unknown' value, and there are

3046(~30% of the data) rows with at least one 'Unknown' value. There are several strategies dealing with 'Unknown' values: 1. Delete rows with 'Unknown' value 2. Replace with the mean, median, or mode of the respective feature. This is typically used for numerical features 3. Use machine learning algorithms to predict missing values based on other features. i.e. XGBoost

Given large size of the dataset, I will choose option 1 to delete rows with missing values. If good results couldn't be obtained with the new dataset, I will seek other techniques for these 'Unknown' values.

```
[268]: # Option 1 : Delete rows with 'Unknown' value
data_opt1 = data_new[(data_new['Income_Category'] != 'Unknown') &
                     (data_new['Education_Level'] != 'Unknown') &
                     (data_new['Marital_Status'] != 'Unknown')].copy()
```

```
[269]: # ## Option 2: Replace 'Unknown' value with mode
# data_opt2=data_new.copy()
# mode_education= data_opt2['Education_Level'].mode()[0]
# mode_marital= data_opt2['Marital_Status'].mode()[0]
# mode_income= data_opt2['Income_Category'].mode()[0]
# data_opt2['Education_Level']=data_opt2['Education_Level'].
#   ↳replace('Unknown',mode_education)
# data_opt2['Marital_Status']=data_opt2['Marital_Status'].
#   ↳replace('Unknown',mode_marital)
# data_opt2['Income_Category']=data_opt2['Income_Category'].
#   ↳replace('Unknown',mode_income)
# data_opt2.info()
```

3.2 Deal with categorical variables

```
[271]: # Create rpreprocessor for x variables
categorical_features = ['Gender', 'Education_Level', 'Marital_Status',
#   ↳'Income_Category', 'Card_Category']
numeric_features = ['Dependent_count', 'Months_on_book',
#   ↳'Total_Relationship_Count',
#   ↳'Months_Inactive_12_mon', 'Contacts_Count_12_mon', 'Credit_Limit',
#   ↳'Total_Revolving_Bal', 'Total_Amt_Chng_Q4_Q1', 'Total_Trans_Ct',
#   ↳'Total_Ct_Chng_Q4_Q1', 'Avg_Utilization_Ratio']
preprocessor = ColumnTransformer(
    transformers=
    [
        ('num', StandardScaler(), numeric_features), # standardize numerical
#   ↳variables
        ('cat', OneHotEncoder(), categorical_features) # encode categorical
#   ↳variables
    ])

```

```
# Since LabelEncoder() will encode on alphabet order, attrited customers will
↳ be encoded 0.
# I want to label attrited customers as positive(1), so I manually encode the y
↳ variable.
data_opt1['Attrition_Flag']=data_opt1['Attrition_Flag'].replace({'Existing_
↳ Customer':0,'Attrited Customer':1}).astype(int)
```

[272]: data_opt1.info()

```
<class 'pandas.core.frame.DataFrame'>
Index: 7081 entries, 0 to 10126
Data columns (total 17 columns):
#   Column                                Non-Null Count  Dtype
---  -
0   Attrition_Flag                        7081 non-null   int64
1   Gender                               7081 non-null   object
2   Dependent_count                       7081 non-null   int64
3   Education_Level                       7081 non-null   object
4   Marital_Status                       7081 non-null   object
5   Income_Category                       7081 non-null   object
6   Card_Category                         7081 non-null   object
7   Months_on_book                       7081 non-null   int64
8   Total_Relationship_Count              7081 non-null   int64
9   Months_Inactive_12_mon                7081 non-null   int64
10  Contacts_Count_12_mon                 7081 non-null   int64
11  Credit_Limit                          7081 non-null   float64
12  Total_Revolving_Bal                   7081 non-null   int64
13  Total_Amt_Chng_Q4_Q1                  7081 non-null   float64
14  Total_Trans_Ct                        7081 non-null   int64
15  Total_Ct_Chng_Q4_Q1                   7081 non-null   float64
16  Avg_Utilization_Ratio                 7081 non-null   float64
dtypes: float64(4), int64(8), object(5)
memory usage: 995.8+ KB
```

4 Select and Train a Model

[273]: # Define the explanatory and dependent variables and the train and test set
columns_to_exclude = ['Attrition_Flag']
X = data_opt1.drop(columns=columns_to_exclude)
y=data_opt1['Attrition_Flag']
#Split the data set into train(80%) and test(20%) dataset
X_train,X_test,y_train,y_test=train_test_split(X,y,test_size=0.
↳ 2,random_state=42)
print(len(X_train),len(X_test)) # Print the size of train and test dataset

5664 1417


```
[274]: #Preprocess the X_train and X_test datasets
X_train_preproc = preprocessor.fit_transform(X_train)
X_test_preproc = preprocessor.transform(X_test)
```

4.1 Logistic Regression

```
[275]: # Fit the logistic regression model using the train set
log_reg=LogisticRegression()
log_reg.fit(X_train_preproc, y_train)
```

```
[275]: LogisticRegression()
```

```
[276]: #Use the trained model to predict y value
y_pred_lr=log_reg.predict(X_test_preproc)

#Calculate y_score for AUC
y_score_lr = log_reg.predict_proba(X_test_preproc)[: , 1]
```

```
[277]: #Generate performance metrics
accuracy_lr = accuracy_score(y_test, y_pred_lr)
conf_matrix_lr = confusion_matrix(y_test, y_pred_lr)
classification_rep_lr = classification_report(y_test, y_pred_lr)
auc_lr=roc_auc_score(y_test, y_score_lr)

print(f"Test accuracy: \n {accuracy_lr:.2%}")
print ("Confusion Matrix : \n", conf_matrix_lr)
print("\nClassification Report: : \n",classification_rep_lr)
print(f"AUC: \n {auc_lr:.2%}")
```

Test accuracy:

90.26%

Confusion Matrix :

```
[[1169  36]
 [ 102 110]]
```

Classification Report: :

	precision	recall	f1-score	support
0	0.92	0.97	0.94	1205
1	0.75	0.52	0.61	212
accuracy			0.90	1417
macro avg	0.84	0.74	0.78	1417
weighted avg	0.89	0.90	0.89	1417

AUC:

90.92%

The logistic regression model has a 90.26% accuracy which is pretty high. However, I noticed that the model has a very high f1-score for the majority class(0 or existing customers) but a low f1-score for the minority class(1 or attrited customers). The object of the project is to model to identify the potential attrited customers, so I want to reduce the false negative cases(attrited but predicted as existing). As a result, a model with a high recall rate ($TP/(TP+FN)$) is desired. Clearly, the logistic model with 52% recall rate is not satisfying.

4.2 Penalized-SVM

Increasing the cost of classification mistakes in the minority class is a technique to deal with dataset with imbalanced classes. A popular algorithm is Penalized-SVM. During training, the argument `class_weight='balanced'` could be used to penalize mistakes on the minority class by an amount proportional to how under-represented it is.

```
[278]: # Fit the model
svc = SVC(class_weight='balanced', probability=True)
#class_weight='balanced' to penalize mistakes on the minority class by an
    ↳ amount proportional to how under-represented it is.
svc.fit(X_train_preproc, y_train)
```

```
[278]: SVC(class_weight='balanced', probability=True)
```

```
[279]: # Generate the predicted y value and calculate the y_score
y_pred_svc=svc.predict(X_test_preproc)
y_score_svc= svc.predict_proba(X_test_preproc)[:, 1]
```

```
[280]: #Generate performance metrics
accuracy_svc = accuracy_score(y_test, y_pred_svc)
conf_matrix_svc = confusion_matrix(y_test, y_pred_svc)
classification_rep_svc = classification_report(y_test, y_pred_svc)
auc_svc=roc_auc_score(y_test, y_score_svc)

print(f"Test accuracy: \n {accuracy_svc:.2%}")
print ("Confusion Matrix : \n", conf_matrix_svc)
print("\nClassification Report: : \n",classification_rep_svc)
print(f"AUC: \n {auc_svc:.2%}")
```

Test accuracy:

89.98%

Confusion Matrix :

```
[[1094  111]
```

```
[  31  181]]
```

Classification Report: :

	precision	recall	f1-score	support
0	0.97	0.91	0.94	1205
1	0.62	0.85	0.72	212

accuracy			0.90	1417
macro avg	0.80	0.88	0.83	1417
weighted avg	0.92	0.90	0.91	1417

AUC:
95.11%

Even though the test accuracy of the Penalized-SVM model is a little bit lower than the Logistic Regression model, the AUC slightly increases and the recall rate increases significantly by 33%. However, the precision($TP/(TP+FP)$) of the model is only 62%, which means ~40% of the attributed customers predicted by the model are existing customers. Offering discount/bonus to these customers may lead to profit reduction.

4.3 Random Forest

```
[281]: # Fit the model
random_forest = RandomForestClassifier(random_state=42)
random_forest.fit(X_train_preproc, y_train)
```

```
[281]: RandomForestClassifier(random_state=42)
```

```
[282]: # Generate the predicted y value and calculate the y_score
y_pred_rf=random_forest.predict(X_test_preproc)
y_score_rf = random_forest.predict_proba(X_test_preproc)[:, 1]
```

```
[283]: #Generate performance metrics
accuracy_rf = accuracy_score(y_test, y_pred_rf)
conf_matrix_rf = confusion_matrix(y_test, y_pred_rf)
classification_rep_rf = classification_report(y_test, y_pred_rf)
auc_rf=roc_auc_score(y_test, y_score_rf)

print(f"Test accuracy: \n {accuracy_rf:.2%}")
print ("Confusion Matrix : \n", conf_matrix_rf)
print("\nClassification Report: : \n",classification_rep_rf)
print(f"AUC: \n {auc_rf:.2%}")
```

Test accuracy:
92.94%

Confusion Matrix :
[[1181 24]
[76 136]]

Classification Report: :

	precision	recall	f1-score	support
0	0.94	0.98	0.96	1205
1	0.85	0.64	0.73	212
accuracy			0.93	1417

macro avg	0.89	0.81	0.85	1417
weighted avg	0.93	0.93	0.93	1417

AUC:

95.79%

Even though the Random Forest model has a higher accuracy and AUC compared to the Penalized-SVM model, the recall rate (64%) is pretty low.

4.4 XGBoost

Extreme Gradient Boosting(XGBoost) is a machine-learning algorithm based on the gradient boosting(GBM) algorithm. However, a few differences of XGBoost make it better than GBM in terms of performance and speed.

1. Regularization: XGBoost implements regularization in its algorithm to avoid overfitting, whereas GBM doesn't.
2. Parallelization: GBM tends to have a slower training time than the XGBoost because the latter algorithm implements parallelization during the training process.
3. Missing Data Handling: XGBoost has its own in-built missing data handler, whereas GBM doesn't.
4. In-Built Cross-Validation: XGBoost has an in-built Cross-Validation that could improve the model generalization and robustness.

```
[285]: # Fit the model
xgboost = XGBClassifier()
xgboost.fit(X_train_preproc,y_train)
```

```
[285]: XGBClassifier(base_score=None, booster=None, callbacks=None,
                    colsample_bylevel=None, colsample_bynode=None,
                    colsample_bytree=None, device=None, early_stopping_rounds=None,
                    enable_categorical=False, eval_metric=None, feature_types=None,
                    gamma=None, grow_policy=None, importance_type=None,
                    interaction_constraints=None, learning_rate=None, max_bin=None,
                    max_cat_threshold=None, max_cat_to_onehot=None,
                    max_delta_step=None, max_depth=None, max_leaves=None,
                    min_child_weight=None, missing=nan, monotone_constraints=None,
                    multi_strategy=None, n_estimators=None, n_jobs=None,
                    num_parallel_tree=None, random_state=None, ...)
```

```
[286]: # Generate the predicted y value and calculate the y_score
y_pred_xgb = xgboost.predict(X_test_preproc)
y_score_xgb = xgboost.predict_proba(X_test_preproc)[: , 1]
```

```
[287]: #Generate performance metrics
accuracy_xgb = accuracy_score(y_test, y_pred_xgb)
conf_matrix_xgb = confusion_matrix(y_test, y_pred_xgb)
classification_rep_xgb = classification_report(y_test, y_pred_xgb)
auc_xgb=roc_auc_score(y_test, y_score_xgb)
```

```

print(f"Test accuracy: \n {accuracy_xgb:.2%}")
print ("Confusion Matrix : \n", conf_matrix_xgb)
print("\nClassification Report: : \n",classification_rep_xgb)
print(f"AUC: \n {auc_xgb:.2%}")

```

Test accuracy:
93.72%

Confusion Matrix :
[[1167 38]
[51 161]]

Classification Report: :

	precision	recall	f1-score	support
0	0.96	0.97	0.96	1205
1	0.81	0.76	0.78	212
accuracy			0.94	1417
macro avg	0.88	0.86	0.87	1417
weighted avg	0.94	0.94	0.94	1417

AUC:
96.89%

XGBoost has a highest test accuracy, f1 score, and AUC among all models.

4.5 Multilayer Perceptrons (MLP)

```

[288]: # Define MLP model
tf.random.set_seed(42)
mlp = models.Sequential([
    layers.Dense(64, activation='relu', input_shape=(X_train_preproc.
↪shape[1],)),
    layers.Dense(32, activation='relu'),
    layers.Dense(1, activation='sigmoid')
])

# Compile the model
mlp.compile(optimizer='sgd',
            loss=keras.losses.BinaryCrossentropy(),
            metrics=[keras.metrics.AUC()]) # Use AUC as the metrics

# Train the model
mlp.fit(X_train_preproc, y_train, epochs=250, batch_size=128)

```

Epoch 1/250

```

45/45          0s 449us/step -
auc_332: 0.4361 - loss: 0.6412
Epoch 2/250
45/45          0s 392us/step -
auc_332: 0.5488 - loss: 0.4784
Epoch 3/250
45/45          0s 313us/step -
auc_332: 0.6676 - loss: 0.4205
Epoch 4/250
45/45          0s 324us/step -
auc_332: 0.7454 - loss: 0.3904
Epoch 5/250
45/45          0s 331us/step -
auc_332: 0.7866 - loss: 0.3695
Epoch 6/250
45/45          0s 317us/step -
auc_332: 0.8103 - loss: 0.3529
Epoch 7/250
45/45          0s 359us/step -
auc_332: 0.8259 - loss: 0.3391
Epoch 8/250
45/45          0s 314us/step -
auc_332: 0.8365 - loss: 0.3274
Epoch 9/250
45/45          0s 318us/step -
auc_332: 0.8447 - loss: 0.3174
Epoch 10/250
45/45          0s 315us/step -
auc_332: 0.8514 - loss: 0.3090
Epoch 11/250
45/45          0s 319us/step -
auc_332: 0.8569 - loss: 0.3018
Epoch 12/250
45/45          0s 327us/step -
auc_332: 0.8618 - loss: 0.2956
Epoch 13/250
45/45          0s 329us/step -
auc_332: 0.8654 - loss: 0.2903
Epoch 14/250
45/45          0s 341us/step -
auc_332: 0.8693 - loss: 0.2858
Epoch 15/250
45/45          0s 314us/step -
auc_332: 0.8722 - loss: 0.2818
Epoch 16/250
45/45          0s 332us/step -
auc_332: 0.8744 - loss: 0.2783
Epoch 17/250

```

```

45/45          0s 330us/step -
auc_332: 0.8772 - loss: 0.2752
Epoch 18/250
45/45          0s 332us/step -
auc_332: 0.8798 - loss: 0.2725
Epoch 19/250
45/45          0s 326us/step -
auc_332: 0.8815 - loss: 0.2700
Epoch 20/250
45/45          0s 317us/step -
auc_332: 0.8836 - loss: 0.2678
Epoch 21/250
45/45          0s 322us/step -
auc_332: 0.8851 - loss: 0.2658
Epoch 22/250
45/45          0s 320us/step -
auc_332: 0.8863 - loss: 0.2640
Epoch 23/250
45/45          0s 322us/step -
auc_332: 0.8878 - loss: 0.2623
Epoch 24/250
45/45          0s 331us/step -
auc_332: 0.8893 - loss: 0.2608
Epoch 25/250
45/45          0s 319us/step -
auc_332: 0.8905 - loss: 0.2594
Epoch 26/250
45/45          0s 309us/step -
auc_332: 0.8917 - loss: 0.2580
Epoch 27/250
45/45          0s 319us/step -
auc_332: 0.8929 - loss: 0.2568
Epoch 28/250
45/45          0s 312us/step -
auc_332: 0.8942 - loss: 0.2556
Epoch 29/250
45/45          0s 313us/step -
auc_332: 0.8956 - loss: 0.2545
Epoch 30/250
45/45          0s 307us/step -
auc_332: 0.8966 - loss: 0.2534
Epoch 31/250
45/45          0s 317us/step -
auc_332: 0.8976 - loss: 0.2525
Epoch 32/250
45/45          0s 316us/step -
auc_332: 0.8977 - loss: 0.2515
Epoch 33/250

```

```

45/45          0s 319us/step -
auc_332: 0.8987 - loss: 0.2506
Epoch 34/250
45/45          0s 310us/step -
auc_332: 0.8993 - loss: 0.2498
Epoch 35/250
45/45          0s 339us/step -
auc_332: 0.9002 - loss: 0.2489
Epoch 36/250
45/45          0s 325us/step -
auc_332: 0.9011 - loss: 0.2481
Epoch 37/250
45/45          0s 328us/step -
auc_332: 0.9018 - loss: 0.2474
Epoch 38/250
45/45          0s 328us/step -
auc_332: 0.9025 - loss: 0.2466
Epoch 39/250
45/45          0s 320us/step -
auc_332: 0.9031 - loss: 0.2459
Epoch 40/250
45/45          0s 330us/step -
auc_332: 0.9036 - loss: 0.2452
Epoch 41/250
45/45          0s 335us/step -
auc_332: 0.9043 - loss: 0.2445
Epoch 42/250
45/45          0s 324us/step -
auc_332: 0.9052 - loss: 0.2439
Epoch 43/250
45/45          0s 310us/step -
auc_332: 0.9058 - loss: 0.2432
Epoch 44/250
45/45          0s 301us/step -
auc_332: 0.9063 - loss: 0.2426
Epoch 45/250
45/45          0s 327us/step -
auc_332: 0.9069 - loss: 0.2420
Epoch 46/250
45/45          0s 316us/step -
auc_332: 0.9074 - loss: 0.2414
Epoch 47/250
45/45          0s 307us/step -
auc_332: 0.9079 - loss: 0.2407
Epoch 48/250
45/45          0s 305us/step -
auc_332: 0.9084 - loss: 0.2402
Epoch 49/250

```


45/45 0s 320us/step -
auc_332: 0.9091 - loss: 0.2396
Epoch 50/250
45/45 0s 304us/step -
auc_332: 0.9096 - loss: 0.2390
Epoch 51/250
45/45 0s 330us/step -
auc_332: 0.9102 - loss: 0.2384
Epoch 52/250
45/45 0s 305us/step -
auc_332: 0.9107 - loss: 0.2378
Epoch 53/250
45/45 0s 320us/step -
auc_332: 0.9113 - loss: 0.2373
Epoch 54/250
45/45 0s 302us/step -
auc_332: 0.9117 - loss: 0.2367
Epoch 55/250
45/45 0s 323us/step -
auc_332: 0.9124 - loss: 0.2361
Epoch 56/250
45/45 0s 302us/step -
auc_332: 0.9129 - loss: 0.2356
Epoch 57/250
45/45 0s 330us/step -
auc_332: 0.9134 - loss: 0.2350
Epoch 58/250
45/45 0s 298us/step -
auc_332: 0.9137 - loss: 0.2345
Epoch 59/250
45/45 0s 340us/step -
auc_332: 0.9143 - loss: 0.2339
Epoch 60/250
45/45 0s 307us/step -
auc_332: 0.9148 - loss: 0.2334
Epoch 61/250
45/45 0s 308us/step -
auc_332: 0.9152 - loss: 0.2328
Epoch 62/250
45/45 0s 294us/step -
auc_332: 0.9158 - loss: 0.2323
Epoch 63/250
45/45 0s 303us/step -
auc_332: 0.9163 - loss: 0.2318
Epoch 64/250
45/45 0s 308us/step -
auc_332: 0.9167 - loss: 0.2312
Epoch 65/250

```

45/45          0s 299us/step -
auc_332: 0.9174 - loss: 0.2307
Epoch 66/250
45/45          0s 384us/step -
auc_332: 0.9177 - loss: 0.2302
Epoch 67/250
45/45          0s 392us/step -
auc_332: 0.9181 - loss: 0.2297
Epoch 68/250
45/45          0s 323us/step -
auc_332: 0.9185 - loss: 0.2292
Epoch 69/250
45/45          0s 299us/step -
auc_332: 0.9190 - loss: 0.2287
Epoch 70/250
45/45          0s 299us/step -
auc_332: 0.9194 - loss: 0.2282
Epoch 71/250
45/45          0s 305us/step -
auc_332: 0.9199 - loss: 0.2277
Epoch 72/250
45/45          0s 325us/step -
auc_332: 0.9204 - loss: 0.2271
Epoch 73/250
45/45          0s 295us/step -
auc_332: 0.9209 - loss: 0.2266
Epoch 74/250
45/45          0s 296us/step -
auc_332: 0.9212 - loss: 0.2261
Epoch 75/250
45/45          0s 295us/step -
auc_332: 0.9218 - loss: 0.2256
Epoch 76/250
45/45          0s 320us/step -
auc_332: 0.9221 - loss: 0.2251
Epoch 77/250
45/45          0s 292us/step -
auc_332: 0.9227 - loss: 0.2246
Epoch 78/250
45/45          0s 290us/step -
auc_332: 0.9231 - loss: 0.2241
Epoch 79/250
45/45          0s 294us/step -
auc_332: 0.9234 - loss: 0.2236
Epoch 80/250
45/45          0s 294us/step -
auc_332: 0.9237 - loss: 0.2231
Epoch 81/250

```

```

45/45          0s 297us/step -
auc_332: 0.9241 - loss: 0.2226
Epoch 82/250
45/45          0s 299us/step -
auc_332: 0.9245 - loss: 0.2221
Epoch 83/250
45/45          0s 296us/step -
auc_332: 0.9250 - loss: 0.2216
Epoch 84/250
45/45          0s 300us/step -
auc_332: 0.9255 - loss: 0.2212
Epoch 85/250
45/45          0s 297us/step -
auc_332: 0.9258 - loss: 0.2207
Epoch 86/250
45/45          0s 293us/step -
auc_332: 0.9262 - loss: 0.2202
Epoch 87/250
45/45          0s 357us/step -
auc_332: 0.9266 - loss: 0.2197
Epoch 88/250
45/45          0s 308us/step -
auc_332: 0.9271 - loss: 0.2192
Epoch 89/250
45/45          0s 300us/step -
auc_332: 0.9273 - loss: 0.2187
Epoch 90/250
45/45          0s 300us/step -
auc_332: 0.9284 - loss: 0.2183
Epoch 91/250
45/45          0s 302us/step -
auc_332: 0.9287 - loss: 0.2178
Epoch 92/250
45/45          0s 298us/step -
auc_332: 0.9292 - loss: 0.2173
Epoch 93/250
45/45          0s 296us/step -
auc_332: 0.9296 - loss: 0.2168
Epoch 94/250
45/45          0s 297us/step -
auc_332: 0.9299 - loss: 0.2163
Epoch 95/250
45/45          0s 323us/step -
auc_332: 0.9303 - loss: 0.2159
Epoch 96/250
45/45          0s 303us/step -
auc_332: 0.9307 - loss: 0.2154
Epoch 97/250

```

```

45/45          0s 293us/step -
auc_332: 0.9310 - loss: 0.2149
Epoch 98/250
45/45          0s 297us/step -
auc_332: 0.9313 - loss: 0.2144
Epoch 99/250
45/45          0s 302us/step -
auc_332: 0.9318 - loss: 0.2139
Epoch 100/250
45/45          0s 299us/step -
auc_332: 0.9323 - loss: 0.2135
Epoch 101/250
45/45          0s 298us/step -
auc_332: 0.9327 - loss: 0.2130
Epoch 102/250
45/45          0s 295us/step -
auc_332: 0.9330 - loss: 0.2125
Epoch 103/250
45/45          0s 410us/step -
auc_332: 0.9332 - loss: 0.2121
Epoch 104/250
45/45          0s 304us/step -
auc_332: 0.9336 - loss: 0.2116
Epoch 105/250
45/45          0s 302us/step -
auc_332: 0.9340 - loss: 0.2111
Epoch 106/250
45/45          0s 295us/step -
auc_332: 0.9343 - loss: 0.2107
Epoch 107/250
45/45          0s 295us/step -
auc_332: 0.9345 - loss: 0.2102
Epoch 108/250
45/45          0s 293us/step -
auc_332: 0.9349 - loss: 0.2098
Epoch 109/250
45/45          0s 294us/step -
auc_332: 0.9353 - loss: 0.2093
Epoch 110/250
45/45          0s 299us/step -
auc_332: 0.9355 - loss: 0.2089
Epoch 111/250
45/45          0s 318us/step -
auc_332: 0.9360 - loss: 0.2084
Epoch 112/250
45/45          0s 296us/step -
auc_332: 0.9363 - loss: 0.2080
Epoch 113/250

```

```

45/45          0s 294us/step -
auc_332: 0.9365 - loss: 0.2075
Epoch 114/250
45/45          0s 300us/step -
auc_332: 0.9368 - loss: 0.2071
Epoch 115/250
45/45          0s 297us/step -
auc_332: 0.9370 - loss: 0.2066
Epoch 116/250
45/45          0s 296us/step -
auc_332: 0.9372 - loss: 0.2062
Epoch 117/250
45/45          0s 297us/step -
auc_332: 0.9375 - loss: 0.2058
Epoch 118/250
45/45          0s 300us/step -
auc_332: 0.9378 - loss: 0.2053
Epoch 119/250
45/45          0s 306us/step -
auc_332: 0.9381 - loss: 0.2049
Epoch 120/250
45/45          0s 298us/step -
auc_332: 0.9385 - loss: 0.2045
Epoch 121/250
45/45          0s 302us/step -
auc_332: 0.9388 - loss: 0.2040
Epoch 122/250
45/45          0s 297us/step -
auc_332: 0.9390 - loss: 0.2036
Epoch 123/250
45/45          0s 300us/step -
auc_332: 0.9395 - loss: 0.2032
Epoch 124/250
45/45          0s 299us/step -
auc_332: 0.9397 - loss: 0.2028
Epoch 125/250
45/45          0s 297us/step -
auc_332: 0.9399 - loss: 0.2023
Epoch 126/250
45/45          0s 291us/step -
auc_332: 0.9402 - loss: 0.2019
Epoch 127/250
45/45          0s 293us/step -
auc_332: 0.9404 - loss: 0.2015
Epoch 128/250
45/45          0s 331us/step -
auc_332: 0.9407 - loss: 0.2011
Epoch 129/250

```

```

45/45          0s 301us/step -
auc_332: 0.9409 - loss: 0.2007
Epoch 130/250
45/45          0s 302us/step -
auc_332: 0.9412 - loss: 0.2003
Epoch 131/250
45/45          0s 301us/step -
auc_332: 0.9415 - loss: 0.1999
Epoch 132/250
45/45          0s 296us/step -
auc_332: 0.9418 - loss: 0.1995
Epoch 133/250
45/45          0s 297us/step -
auc_332: 0.9421 - loss: 0.1991
Epoch 134/250
45/45          0s 295us/step -
auc_332: 0.9426 - loss: 0.1987
Epoch 135/250
45/45          0s 296us/step -
auc_332: 0.9429 - loss: 0.1984
Epoch 136/250
45/45          0s 302us/step -
auc_332: 0.9431 - loss: 0.1980
Epoch 137/250
45/45          0s 296us/step -
auc_332: 0.9433 - loss: 0.1976
Epoch 138/250
45/45          0s 305us/step -
auc_332: 0.9435 - loss: 0.1972
Epoch 139/250
45/45          0s 312us/step -
auc_332: 0.9438 - loss: 0.1968
Epoch 140/250
45/45          0s 301us/step -
auc_332: 0.9440 - loss: 0.1965
Epoch 141/250
45/45          0s 299us/step -
auc_332: 0.9444 - loss: 0.1961
Epoch 142/250
45/45          0s 297us/step -
auc_332: 0.9445 - loss: 0.1957
Epoch 143/250
45/45          0s 298us/step -
auc_332: 0.9448 - loss: 0.1954
Epoch 144/250
45/45          0s 300us/step -
auc_332: 0.9451 - loss: 0.1950
Epoch 145/250

```

```

45/45          0s 296us/step -
auc_332: 0.9453 - loss: 0.1946
Epoch 146/250
45/45          0s 293us/step -
auc_332: 0.9457 - loss: 0.1943
Epoch 147/250
45/45          0s 291us/step -
auc_332: 0.9458 - loss: 0.1939
Epoch 148/250
45/45          0s 299us/step -
auc_332: 0.9460 - loss: 0.1936
Epoch 149/250
45/45          0s 301us/step -
auc_332: 0.9462 - loss: 0.1932
Epoch 150/250
45/45          0s 297us/step -
auc_332: 0.9464 - loss: 0.1929
Epoch 151/250
45/45          0s 301us/step -
auc_332: 0.9468 - loss: 0.1926
Epoch 152/250
45/45          0s 297us/step -
auc_332: 0.9470 - loss: 0.1922
Epoch 153/250
45/45          0s 299us/step -
auc_332: 0.9474 - loss: 0.1919
Epoch 154/250
45/45          0s 350us/step -
auc_332: 0.9475 - loss: 0.1916
Epoch 155/250
45/45          0s 305us/step -
auc_332: 0.9478 - loss: 0.1912
Epoch 156/250
45/45          0s 297us/step -
auc_332: 0.9480 - loss: 0.1909
Epoch 157/250
45/45          0s 301us/step -
auc_332: 0.9482 - loss: 0.1906
Epoch 158/250
45/45          0s 299us/step -
auc_332: 0.9484 - loss: 0.1903
Epoch 159/250
45/45          0s 296us/step -
auc_332: 0.9485 - loss: 0.1899
Epoch 160/250
45/45          0s 298us/step -
auc_332: 0.9488 - loss: 0.1896
Epoch 161/250

```

```

45/45          0s 295us/step -
auc_332: 0.9487 - loss: 0.1893
Epoch 162/250
45/45          0s 303us/step -
auc_332: 0.9489 - loss: 0.1890
Epoch 163/250
45/45          0s 297us/step -
auc_332: 0.9491 - loss: 0.1887
Epoch 164/250
45/45          0s 295us/step -
auc_332: 0.9493 - loss: 0.1884
Epoch 165/250
45/45          0s 298us/step -
auc_332: 0.9495 - loss: 0.1881
Epoch 166/250
45/45          0s 303us/step -
auc_332: 0.9497 - loss: 0.1878
Epoch 167/250
45/45          0s 305us/step -
auc_332: 0.9498 - loss: 0.1875
Epoch 168/250
45/45          0s 302us/step -
auc_332: 0.9499 - loss: 0.1872
Epoch 169/250
45/45          0s 298us/step -
auc_332: 0.9501 - loss: 0.1868
Epoch 170/250
45/45          0s 299us/step -
auc_332: 0.9503 - loss: 0.1865
Epoch 171/250
45/45          0s 295us/step -
auc_332: 0.9504 - loss: 0.1862
Epoch 172/250
45/45          0s 292us/step -
auc_332: 0.9507 - loss: 0.1859
Epoch 173/250
45/45          0s 339us/step -
auc_332: 0.9508 - loss: 0.1856
Epoch 174/250
45/45          0s 292us/step -
auc_332: 0.9510 - loss: 0.1853
Epoch 175/250
45/45          0s 288us/step -
auc_332: 0.9512 - loss: 0.1850
Epoch 176/250
45/45          0s 298us/step -
auc_332: 0.9513 - loss: 0.1847
Epoch 177/250

```



```

45/45          0s 309us/step -
auc_332: 0.9514 - loss: 0.1844
Epoch 178/250
45/45          0s 294us/step -
auc_332: 0.9516 - loss: 0.1842
Epoch 179/250
45/45          0s 298us/step -
auc_332: 0.9517 - loss: 0.1839
Epoch 180/250
45/45          0s 300us/step -
auc_332: 0.9518 - loss: 0.1836
Epoch 181/250
45/45          0s 319us/step -
auc_332: 0.9520 - loss: 0.1833
Epoch 182/250
45/45          0s 299us/step -
auc_332: 0.9522 - loss: 0.1830
Epoch 183/250
45/45          0s 297us/step -
auc_332: 0.9523 - loss: 0.1827
Epoch 184/250
45/45          0s 302us/step -
auc_332: 0.9526 - loss: 0.1824
Epoch 185/250
45/45          0s 297us/step -
auc_332: 0.9528 - loss: 0.1822
Epoch 186/250
45/45          0s 293us/step -
auc_332: 0.9529 - loss: 0.1819
Epoch 187/250
45/45          0s 296us/step -
auc_332: 0.9531 - loss: 0.1816
Epoch 188/250
45/45          0s 292us/step -
auc_332: 0.9532 - loss: 0.1813
Epoch 189/250
45/45          0s 299us/step -
auc_332: 0.9534 - loss: 0.1811
Epoch 190/250
45/45          0s 296us/step -
auc_332: 0.9536 - loss: 0.1808
Epoch 191/250
45/45          0s 296us/step -
auc_332: 0.9538 - loss: 0.1805
Epoch 192/250
45/45          0s 297us/step -
auc_332: 0.9538 - loss: 0.1803
Epoch 193/250

```

```

45/45          0s 293us/step -
auc_332: 0.9539 - loss: 0.1800
Epoch 194/250
45/45          0s 290us/step -
auc_332: 0.9541 - loss: 0.1797
Epoch 195/250
45/45          0s 297us/step -
auc_332: 0.9546 - loss: 0.1795
Epoch 196/250
45/45          0s 295us/step -
auc_332: 0.9548 - loss: 0.1792
Epoch 197/250
45/45          0s 294us/step -
auc_332: 0.9548 - loss: 0.1790
Epoch 198/250
45/45          0s 300us/step -
auc_332: 0.9551 - loss: 0.1787
Epoch 199/250
45/45          0s 293us/step -
auc_332: 0.9549 - loss: 0.1785
Epoch 200/250
45/45          0s 293us/step -
auc_332: 0.9550 - loss: 0.1782
Epoch 201/250
45/45          0s 330us/step -
auc_332: 0.9552 - loss: 0.1780
Epoch 202/250
45/45          0s 312us/step -
auc_332: 0.9553 - loss: 0.1778
Epoch 203/250
45/45          0s 297us/step -
auc_332: 0.9554 - loss: 0.1775
Epoch 204/250
45/45          0s 289us/step -
auc_332: 0.9555 - loss: 0.1773
Epoch 205/250
45/45          0s 297us/step -
auc_332: 0.9557 - loss: 0.1770
Epoch 206/250
45/45          0s 293us/step -
auc_332: 0.9558 - loss: 0.1768
Epoch 207/250
45/45          0s 290us/step -
auc_332: 0.9559 - loss: 0.1766
Epoch 208/250
45/45          0s 296us/step -
auc_332: 0.9560 - loss: 0.1763
Epoch 209/250

```

```

45/45          0s 293us/step -
auc_332: 0.9561 - loss: 0.1761
Epoch 210/250
45/45          0s 302us/step -
auc_332: 0.9562 - loss: 0.1759
Epoch 211/250
45/45          0s 298us/step -
auc_332: 0.9564 - loss: 0.1756
Epoch 212/250
45/45          0s 306us/step -
auc_332: 0.9565 - loss: 0.1754
Epoch 213/250
45/45          0s 283us/step -
auc_332: 0.9566 - loss: 0.1752
Epoch 214/250
45/45          0s 297us/step -
auc_332: 0.9567 - loss: 0.1750
Epoch 215/250
45/45          0s 302us/step -
auc_332: 0.9569 - loss: 0.1747
Epoch 216/250
45/45          0s 300us/step -
auc_332: 0.9569 - loss: 0.1745
Epoch 217/250
45/45          0s 298us/step -
auc_332: 0.9568 - loss: 0.1743
Epoch 218/250
45/45          0s 311us/step -
auc_332: 0.9569 - loss: 0.1741
Epoch 219/250
45/45          0s 301us/step -
auc_332: 0.9570 - loss: 0.1739
Epoch 220/250
45/45          0s 301us/step -
auc_332: 0.9572 - loss: 0.1736
Epoch 221/250
45/45          0s 301us/step -
auc_332: 0.9573 - loss: 0.1734
Epoch 222/250
45/45          0s 302us/step -
auc_332: 0.9574 - loss: 0.1732
Epoch 223/250
45/45          0s 301us/step -
auc_332: 0.9575 - loss: 0.1730
Epoch 224/250
45/45          0s 296us/step -
auc_332: 0.9576 - loss: 0.1728
Epoch 225/250

```

```

45/45          0s 296us/step -
auc_332: 0.9578 - loss: 0.1725
Epoch 226/250
45/45          0s 300us/step -
auc_332: 0.9579 - loss: 0.1723
Epoch 227/250
45/45          0s 310us/step -
auc_332: 0.9580 - loss: 0.1721
Epoch 228/250
45/45          0s 299us/step -
auc_332: 0.9581 - loss: 0.1719
Epoch 229/250
45/45          0s 301us/step -
auc_332: 0.9582 - loss: 0.1717
Epoch 230/250
45/45          0s 302us/step -
auc_332: 0.9585 - loss: 0.1715
Epoch 231/250
45/45          0s 298us/step -
auc_332: 0.9585 - loss: 0.1713
Epoch 232/250
45/45          0s 300us/step -
auc_332: 0.9586 - loss: 0.1711
Epoch 233/250
45/45          0s 321us/step -
auc_332: 0.9588 - loss: 0.1709
Epoch 234/250
45/45          0s 313us/step -
auc_332: 0.9589 - loss: 0.1707
Epoch 235/250
45/45          0s 296us/step -
auc_332: 0.9590 - loss: 0.1705
Epoch 236/250
45/45          0s 296us/step -
auc_332: 0.9590 - loss: 0.1703
Epoch 237/250
45/45          0s 301us/step -
auc_332: 0.9592 - loss: 0.1701
Epoch 238/250
45/45          0s 299us/step -
auc_332: 0.9592 - loss: 0.1699
Epoch 239/250
45/45          0s 310us/step -
auc_332: 0.9594 - loss: 0.1697
Epoch 240/250
45/45          0s 296us/step -
auc_332: 0.9595 - loss: 0.1695
Epoch 241/250

```

```

45/45          0s 296us/step -
auc_332: 0.9596 - loss: 0.1693
Epoch 242/250
45/45          0s 297us/step -
auc_332: 0.9597 - loss: 0.1691
Epoch 243/250
45/45          0s 298us/step -
auc_332: 0.9598 - loss: 0.1689
Epoch 244/250
45/45          0s 304us/step -
auc_332: 0.9599 - loss: 0.1687
Epoch 245/250
45/45          0s 317us/step -
auc_332: 0.9600 - loss: 0.1685
Epoch 246/250
45/45          0s 316us/step -
auc_332: 0.9601 - loss: 0.1683
Epoch 247/250
45/45          0s 299us/step -
auc_332: 0.9602 - loss: 0.1681
Epoch 248/250
45/45          0s 297us/step -
auc_332: 0.9603 - loss: 0.1679
Epoch 249/250
45/45          0s 302us/step -
auc_332: 0.9604 - loss: 0.1677
Epoch 250/250
45/45          0s 297us/step -
auc_332: 0.9605 - loss: 0.1676

```

[288]: <keras.src.callbacks.history.History at 0x2a871ac50>

```

[289]: # Generate the predicted probability of y
y_score_mlp = mlp.predict(X_test_preproc)

# Convert probability to y value
y_pred_mlp = (y_score_mlp > 0.5).astype(int)

# Generate the test loss and AUC
test_loss, auc_mlp=mlp.evaluate(X_test_preproc, y_test)

```

```

45/45          0s 483us/step
45/45          0s 315us/step -
auc_332: 0.9499 - loss: 0.1866

```

```

[290]: #Generate performance metrics
accuracy_mlp = accuracy_score(y_test, y_pred_mlp)
conf_matrix_mlp = confusion_matrix(y_test, y_pred_mlp)

```

```

classification_rep_mlp = classification_report(y_test, y_pred_mlp)
# auc_mlp=roc_auc_score(y_test, y_score_mlp)

print(f"Test accuracy: \n {accuracy_mlp:.2%}")
print ("Confusion Matrix : \n", conf_matrix_mlp)
print("\nClassification Report: : \n",classification_rep_mlp)
print(f"AUC: \n {auc_mlp:.2%}")

```

Test accuracy:

92.45%

Confusion Matrix :

[[1171 34]

[73 139]]

Classification Report: :

	precision	recall	f1-score	support
0	0.94	0.97	0.96	1205
1	0.80	0.66	0.72	212
accuracy			0.92	1417
macro avg	0.87	0.81	0.84	1417
weighted avg	0.92	0.92	0.92	1417

AUC:

94.01%

4.6 Better Model Evaluation Using Cross-Validation

```

[291]: kf = KFold(n_splits=5, shuffle=True, random_state=42)

models = [log_reg, svc, random_forest, xgboost]
scores = ['accuracy', 'recall', 'precision', 'f1', 'roc_auc']
score_dict = {}
score_dict['Model']=['Logistic Regression', 'Penalized-SVM', 'Random Forest',
↳'XGBoost', 'MLP']

for score in scores:
    score_dict[score] = []
    for model in models:
        mean_score = cross_val_score(model, X_train_preproc, y_train, cv=kf,
↳scoring=score).mean()
        score_dict[score].append(mean_score)

```

```

[292]: # Create a function that returns the Keras model
def create_model():
    model = Sequential([

```

```

        Dense(64, activation='relu', input_shape=(X_train_preproc.shape[1],)),
        Dense(32, activation='relu'),
        Dense(1, activation='sigmoid')
    ])
    model.compile(optimizer='sgd',
                  loss=keras.losses.BinaryCrossentropy(),
                  metrics=[keras.metrics.AUC()])
    return model

# Create a KerasClassifier
model = KerasClassifier(model=create_model,
                        epochs=250,
                        batch_size=128,
                        verbose=0)

# Define the cross-validation process to be used inside cross_val_score
cv_mlp = StratifiedKFold(n_splits=5, shuffle=True, random_state=42)

# Evaluate model using cross validation
for score in scores:
    mean_score = cross_val_score(model, X_train_preproc, y_train, cv=cv_mlp,
    ↪scoring=score).mean()
    score_dict[score].append(mean_score)

```

```

[293]: result_df = pd.DataFrame(data = score_dict)
      result_df.head()

```

```

[293]:

```

	Model	accuracy	recall	precision	f1	roc_auc
0	Logistic Regression	0.894244	0.504234	0.749398	0.602205	0.894754
1	Penalized-SVM	0.887359	0.833813	0.607558	0.702056	0.940904
2	Random Forest	0.924081	0.608033	0.878381	0.718005	0.951609
3	XGBoost	0.933616	0.719342	0.840516	0.774982	0.962683
4	MLP	0.917370	0.645955	0.790351	0.704159	0.935716

Based on the table above, I want to fine-tune the Random Forest, XGBoost, and MLP models. Even though all these models have high AUCs, these three models have highest f1 score. In this project, the bank wants to identify potential attrited customers and take actions to avoid customer loss, but the bank doesn't want to give 'win back offer' to too many existing customers since the offers reduces their profit. As a result, I want to find a balance between recall and precision in the model, so f1 score is used as the metric for the model evaluation.

5 Fine-Tune the Model

I'll use grid search technique to find the hyperparameters of the model that maximize f1 score.

5.1 Grid Search

5.1.1 Random Forest

```
[213]: # Define model
model = RandomForestClassifier(random_state=42, class_weight='balanced')

# Define grid
param_grid = {
    'n_estimators': [150, 200, 250],
    'max_depth': [80, 90, 100],
    'min_samples_split': [8, 10, 12],
    'min_samples_leaf': [3, 4, 5]
}

# Define grid search based on Recall and AUC
grid_search = GridSearchCV(estimator=model, param_grid=param_grid, cv=5,
    ↪scoring='f1')

grid_search=grid_search.fit(X_train_preproc, y_train)

[214]: print("Best f1 score: %f using %s" % (grid_search.best_score_, grid_search.
    ↪best_params_))
```

Best f1 score: 0.759088 using {'max_depth': 80, 'min_samples_leaf': 5, 'min_samples_split': 12, 'n_estimators': 150}

5.1.2 XGBoost

In this section, I want to search the optimized `scale_pos_weight` which is a hyperparameter in XGBoost with the effect of weighing the balance of positive examples (minority class), relative to negative examples (majority class) when boosting decision trees. It has by default value 1. A sensible default value to set for the `scale_pos_weight` hyperparameter is the inverse of the class distribution. For example, for a dataset with a 1 to 100 ratio for examples in the minority to majority classes, the `scale_pos_weight` can be set to 100. This will give classification errors made by the model on the minority class 100 times more impact, and in turn, 100 times more correction than errors made on the majority class.

```
[180]: # Define model
model = XGBClassifier()
# Define grid - all integers between 1 and 40
weights = list(range(1, 41))
param_grid = dict(scale_pos_weight=weights)
# Define evaluation procedure
cv = RepeatedStratifiedKFold(n_splits=5, n_repeats=3, random_state=42)
# Define grid search based on different scoring method(recall, precision, f1,
    ↪score, accuracy, AUC)
grid_recall = GridSearchCV(estimator=model, param_grid=param_grid, cv=cv,
    ↪scoring='recall')
```



```

grid_precision = GridSearchCV(estimator=model, param_grid=param_grid, cv=cv,
    ↳scoring='precision')
grid_f1 = GridSearchCV(estimator=model, param_grid=param_grid, cv=cv,
    ↳scoring='f1')

# Execute the grid search
grid_result_recall = grid_recall.fit(X_train_preproc,y_train)
grid_result_precision = grid_precision.fit(X_train_preproc,y_train)
grid_result_f1 = grid_f1.fit(X_train_preproc,y_train)

# Save the mean core by different weight
means_recall = grid_result_recall.cv_results_['mean_test_score']
means_precision = grid_result_precision.cv_results_['mean_test_score']
means_f1 = grid_result_f1.cv_results_['mean_test_score']

```

```

[181]: # Plot the mean score by different weight
plt.figure()
plt.plot(weights, means_recall, color='gray', lw=2, label='Recall')
plt.plot(weights, means_precision, color='darkorange', lw=2, label='Precision' )
plt.plot(weights, means_f1, color='green', lw=2, label='F1 Score')

plt.xlabel('Scale_pos_weight')
plt.ylabel('Mean Score')
plt.title('Grid Search - Mean Score')
plt.legend(loc='center left', bbox_to_anchor=(1, 0.5))
plt.show()

```



From the plot above, I observed that averaged recall($TP/(TP+FN)$) increases as the `scale_pos_weight` increases, which means the model becomes more powerful to identify the attrited customers as `scale_pos_weight` increases. However, the precision decreases as the `scale_pos_weight` increases because the increase in penalty on minority samples may lead the model to classify more samples as positive, thereby increasing the false positive rate and decreasing precision, and also accuracy and AUC.

In this project, the bank wants to identify potential attrited customers and take actions to avoid customer loss, but the bank doesn't want to give 'win back offer' to too many existing customers since the offers reduces their profit. As a result, I want to find a balance between recall and precision in the model, and the f1 score could be a good metric as it is the weighted average of recall and precision.

```
[182]: # Report the best configuration
print("Best F1 score: %f using %s" % (grid_result_f1.best_score_,
    ↪grid_result_f1.best_params_))
```

Best F1 score: 0.789974 using {'scale_pos_weight': 3}

5.1.3 MLP

```
[294]: def create_model(num_neurons=32, learning_rate=0.01, num_hidden_layers=1):
    model = Sequential()
    model.add(Dense(num_neurons, activation='relu',
    ↪input_shape=(X_train_preproc.shape[1],)))
    for _ in range(num_hidden_layers):
        model.add(Dense(num_neurons, activation='relu'))
    model.add(Dense(1, activation='sigmoid'))

    optimizer = keras.optimizers.SGD(learning_rate=learning_rate)
    model.compile(optimizer=optimizer,
                  loss=keras.losses.BinaryCrossentropy(),
                  metrics=[keras.metrics.AUC()])

    return model

# Wrap the Keras model in a scikit-learn estimator
model = KerasClassifier(model=create_model,num_hidden_layers=1, num_neurons=32,
    ↪learning_rate=0.01,
                        epochs=250, batch_size=128, verbose=0)

param_grid = {
    'num_neurons': [32, 64, 128],
    'learning_rate': [0.001, 0.01, 0.1],
    'num_hidden_layers': [1, 2]
}
```

```
grid_f1 = GridSearchCV(estimator=model, param_grid=param_grid, cv=5,
    ↪scoring='f1')
```

```
[295]: grid_result_f1 = grid_f1.fit(X_train_preproc,y_train)
print("Best F1 score: %f using %s" % (grid_result_f1.best_score_,
    ↪grid_result_f1.best_params_))
```

Best F1 score: 0.725710 using {'learning_rate': 0.01, 'num_hidden_layers': 2, 'num_neurons': 64}

5.2 Evaluate the Model on the Test Set

5.2.1 Random Forest

```
[218]: random_forest =
    ↪RandomForestClassifier(random_state=42,max_depth=80,min_samples_leaf=5,
    ↪min_samples_split= 12, n_estimators=150,class_weight='balanced')
random_forest.fit(X_train_preproc, y_train)
```

```
[218]: RandomForestClassifier(class_weight='balanced', max_depth=80,
    min_samples_leaf=5, min_samples_split=12,
    n_estimators=150, random_state=42)
```

```
[219]: # Generate the predicted y value and calculate the y_score
y_pred_rf=random_forest.predict(X_test_preproc)
y_score_rf = random_forest.predict_proba(X_test_preproc)[:, 1]
```

```
[220]: #Generate performance metrics
accuracy_rf = accuracy_score(y_test, y_pred_rf)
conf_matrix_rf = confusion_matrix(y_test, y_pred_rf)
classification_rep_rf = classification_report(y_test, y_pred_rf)
auc_rf=roc_auc_score(y_test, y_score_rf)

print(f"Test accuracy: \n {accuracy_rf:.2%}")
print ("Confusion Matrix : \n", conf_matrix_rf)
print("\nClassification Report: : \n",classification_rep_rf)
print(f"AUC: \n {auc_rf:.2%}")
```

Test accuracy:

92.73%

Confusion Matrix :

```
[[1154  51]
 [ 52 160]]
```

Classification Report: :

	precision	recall	f1-score	support
0	0.96	0.96	0.96	1205
1	0.76	0.75	0.76	212

accuracy			0.93	1417
macro avg	0.86	0.86	0.86	1417
weighted avg	0.93	0.93	0.93	1417

AUC:
95.71%

5.2.2 XGBoost

```
[233]: xgboost = XGBClassifier(scale_pos_weight=3)
xgboost.fit(X_train_preproc,y_train)
```

```
[233]: XGBClassifier(base_score=None, booster=None, callbacks=None,
    colsample_bylevel=None, colsample_bynode=None,
    colsample_bytree=None, device=None, early_stopping_rounds=None,
    enable_categorical=False, eval_metric=None, feature_types=None,
    gamma=None, grow_policy=None, importance_type=None,
    interaction_constraints=None, learning_rate=None, max_bin=None,
    max_cat_threshold=None, max_cat_to_onehot=None,
    max_delta_step=None, max_depth=None, max_leaves=None,
    min_child_weight=None, missing=nan, monotone_constraints=None,
    multi_strategy=None, n_estimators=None, n_jobs=None,
    num_parallel_tree=None, random_state=None, ...)
```

```
[234]: # Generate the predicted y value and calculate the y_score
y_pred_xgb = xgboost.predict(X_test_preproc)
y_score_xgb = xgboost.predict_proba(X_test_preproc)[:, 1]
```

```
[235]: #Generate performance metrics
accuracy_xgb = accuracy_score(y_test, y_pred_xgb)
conf_matrix_xgb = confusion_matrix(y_test, y_pred_xgb)
classification_rep_xgb = classification_report(y_test, y_pred_xgb)
auc_xgb=roc_auc_score(y_test, y_score_xgb)

print(f"Test accuracy: \n {accuracy_xgb:.2%}")
print ("Confusion Matrix : \n", conf_matrix_xgb)
print("\nClassification Report: : \n",classification_rep_xgb)
print(f"AUC: \n {auc_xgb:.2%}")
```

Test accuracy:
94.00%

Confusion Matrix :

```
[[1159  46]
 [ 39 173]]
```

Classification Report: :

	precision	recall	f1-score	support
--	-----------	--------	----------	---------

0	0.97	0.96	0.96	1205
1	0.79	0.82	0.80	212
accuracy			0.94	1417
macro avg	0.88	0.89	0.88	1417
weighted avg	0.94	0.94	0.94	1417

AUC:
97.12%

5.2.3 MLP

```
[299]: mlp = models.Sequential([
        layers.Dense(64, activation='relu', input_shape=(X_train_preproc.
        ↪shape[1],)),
        layers.Dense(64, activation='relu'),
        layers.Dense(64, activation='relu'),
        layers.Dense(1, activation='sigmoid')
    ])
    # Compile the model
    mlp.compile(optimizer=SGD(learning_rate=0.01),
                loss=keras.losses.BinaryCrossentropy(),
                metrics=[keras.metrics.AUC()]) # Use AUC as the metrics
    # Train the model
    mlp.fit(X_train_preproc, y_train, epochs=250, batch_size=128)
```

```
Epoch 1/250
45/45          0s 439us/step -
auc_449: 0.4091 - loss: 0.7034
Epoch 2/250
45/45          0s 449us/step -
auc_449: 0.4992 - loss: 0.5195
Epoch 3/250
45/45          0s 390us/step -
auc_449: 0.6120 - loss: 0.4464
Epoch 4/250
45/45          0s 403us/step -
auc_449: 0.7012 - loss: 0.4118
Epoch 5/250
45/45          0s 410us/step -
auc_449: 0.7594 - loss: 0.3919
Epoch 6/250
45/45          0s 418us/step -
auc_449: 0.7946 - loss: 0.3774
Epoch 7/250
45/45          0s 402us/step -
auc_449: 0.8159 - loss: 0.3649
```

Epoch 8/250
45/45 0s 404us/step -
auc_449: 0.8300 - loss: 0.3537
Epoch 9/250
45/45 0s 421us/step -
auc_449: 0.8403 - loss: 0.3433
Epoch 10/250
45/45 0s 412us/step -
auc_449: 0.8485 - loss: 0.3336
Epoch 11/250
45/45 0s 404us/step -
auc_449: 0.8552 - loss: 0.3246
Epoch 12/250
45/45 0s 421us/step -
auc_449: 0.8604 - loss: 0.3164
Epoch 13/250
45/45 0s 401us/step -
auc_449: 0.8646 - loss: 0.3089
Epoch 14/250
45/45 0s 399us/step -
auc_449: 0.8687 - loss: 0.3020
Epoch 15/250
45/45 0s 413us/step -
auc_449: 0.8722 - loss: 0.2958
Epoch 16/250
45/45 0s 404us/step -
auc_449: 0.8758 - loss: 0.2902
Epoch 17/250
45/45 0s 413us/step -
auc_449: 0.8790 - loss: 0.2852
Epoch 18/250
45/45 0s 406us/step -
auc_449: 0.8817 - loss: 0.2806
Epoch 19/250
45/45 0s 404us/step -
auc_449: 0.8844 - loss: 0.2766
Epoch 20/250
45/45 0s 423us/step -
auc_449: 0.8867 - loss: 0.2729
Epoch 21/250
45/45 0s 414us/step -
auc_449: 0.8886 - loss: 0.2697
Epoch 22/250
45/45 0s 484us/step -
auc_449: 0.8909 - loss: 0.2667
Epoch 23/250
45/45 0s 400us/step -
auc_449: 0.8926 - loss: 0.2640

Epoch 24/250
45/45 0s 424us/step -
auc_449: 0.8946 - loss: 0.2615
Epoch 25/250
45/45 0s 404us/step -
auc_449: 0.8960 - loss: 0.2592
Epoch 26/250
45/45 0s 409us/step -
auc_449: 0.8979 - loss: 0.2571
Epoch 27/250
45/45 0s 401us/step -
auc_449: 0.8994 - loss: 0.2552
Epoch 28/250
45/45 0s 386us/step -
auc_449: 0.9008 - loss: 0.2534
Epoch 29/250
45/45 0s 396us/step -
auc_449: 0.9019 - loss: 0.2517
Epoch 30/250
45/45 0s 395us/step -
auc_449: 0.9033 - loss: 0.2502
Epoch 31/250
45/45 0s 404us/step -
auc_449: 0.9045 - loss: 0.2487
Epoch 32/250
45/45 0s 395us/step -
auc_449: 0.9056 - loss: 0.2473
Epoch 33/250
45/45 0s 396us/step -
auc_449: 0.9065 - loss: 0.2460
Epoch 34/250
45/45 0s 387us/step -
auc_449: 0.9079 - loss: 0.2447
Epoch 35/250
45/45 0s 390us/step -
auc_449: 0.9090 - loss: 0.2435
Epoch 36/250
45/45 0s 463us/step -
auc_449: 0.9099 - loss: 0.2424
Epoch 37/250
45/45 0s 424us/step -
auc_449: 0.9108 - loss: 0.2413
Epoch 38/250
45/45 0s 409us/step -
auc_449: 0.9117 - loss: 0.2402
Epoch 39/250
45/45 0s 403us/step -
auc_449: 0.9123 - loss: 0.2392

```

Epoch 40/250
45/45          0s 399us/step -
auc_449: 0.9132 - loss: 0.2382
Epoch 41/250
45/45          0s 391us/step -
auc_449: 0.9141 - loss: 0.2372
Epoch 42/250
45/45          0s 400us/step -
auc_449: 0.9147 - loss: 0.2363
Epoch 43/250
45/45          0s 401us/step -
auc_449: 0.9155 - loss: 0.2354
Epoch 44/250
45/45          0s 393us/step -
auc_449: 0.9162 - loss: 0.2345
Epoch 45/250
45/45          0s 401us/step -
auc_449: 0.9170 - loss: 0.2336
Epoch 46/250
45/45          0s 404us/step -
auc_449: 0.9176 - loss: 0.2327
Epoch 47/250
45/45          0s 391us/step -
auc_449: 0.9182 - loss: 0.2319
Epoch 48/250
45/45          0s 405us/step -
auc_449: 0.9189 - loss: 0.2310
Epoch 49/250
45/45          0s 393us/step -
auc_449: 0.9195 - loss: 0.2302
Epoch 50/250
45/45          0s 402us/step -
auc_449: 0.9201 - loss: 0.2294
Epoch 51/250
45/45          0s 406us/step -
auc_449: 0.9210 - loss: 0.2286
Epoch 52/250
45/45          0s 391us/step -
auc_449: 0.9218 - loss: 0.2278
Epoch 53/250
45/45          0s 402us/step -
auc_449: 0.9223 - loss: 0.2271
Epoch 54/250
45/45          0s 390us/step -
auc_449: 0.9230 - loss: 0.2263
Epoch 55/250
45/45          0s 387us/step -
auc_449: 0.9237 - loss: 0.2255

```



```

Epoch 56/250
45/45          0s 382us/step -
auc_449: 0.9242 - loss: 0.2248
Epoch 57/250
45/45          0s 392us/step -
auc_449: 0.9248 - loss: 0.2240
Epoch 58/250
45/45          0s 402us/step -
auc_449: 0.9253 - loss: 0.2233
Epoch 59/250
45/45          0s 404us/step -
auc_449: 0.9259 - loss: 0.2226
Epoch 60/250
45/45          0s 390us/step -
auc_449: 0.9264 - loss: 0.2219
Epoch 61/250
45/45          0s 392us/step -
auc_449: 0.9271 - loss: 0.2212
Epoch 62/250
45/45          0s 390us/step -
auc_449: 0.9277 - loss: 0.2205
Epoch 63/250
45/45          0s 389us/step -
auc_449: 0.9281 - loss: 0.2198
Epoch 64/250
45/45          0s 385us/step -
auc_449: 0.9285 - loss: 0.2192
Epoch 65/250
45/45          0s 395us/step -
auc_449: 0.9289 - loss: 0.2185
Epoch 66/250
45/45          0s 379us/step -
auc_449: 0.9294 - loss: 0.2179
Epoch 67/250
45/45          0s 391us/step -
auc_449: 0.9300 - loss: 0.2172
Epoch 68/250
45/45          0s 374us/step -
auc_449: 0.9303 - loss: 0.2166
Epoch 69/250
45/45          0s 391us/step -
auc_449: 0.9307 - loss: 0.2160
Epoch 70/250
45/45          0s 380us/step -
auc_449: 0.9312 - loss: 0.2153
Epoch 71/250
45/45          0s 390us/step -
auc_449: 0.9319 - loss: 0.2147

```

Epoch 72/250
45/45 0s 382us/step -
auc_449: 0.9323 - loss: 0.2141
Epoch 73/250
45/45 0s 381us/step -
auc_449: 0.9326 - loss: 0.2135
Epoch 74/250
45/45 0s 388us/step -
auc_449: 0.9330 - loss: 0.2129
Epoch 75/250
45/45 0s 386us/step -
auc_449: 0.9334 - loss: 0.2123
Epoch 76/250
45/45 0s 383us/step -
auc_449: 0.9337 - loss: 0.2117
Epoch 77/250
45/45 0s 433us/step -
auc_449: 0.9341 - loss: 0.2111
Epoch 78/250
45/45 0s 401us/step -
auc_449: 0.9347 - loss: 0.2105
Epoch 79/250
45/45 0s 374us/step -
auc_449: 0.9358 - loss: 0.2099
Epoch 80/250
45/45 0s 381us/step -
auc_449: 0.9362 - loss: 0.2093
Epoch 81/250
45/45 0s 391us/step -
auc_449: 0.9366 - loss: 0.2087
Epoch 82/250
45/45 0s 389us/step -
auc_449: 0.9369 - loss: 0.2081
Epoch 83/250
45/45 0s 373us/step -
auc_449: 0.9373 - loss: 0.2075
Epoch 84/250
45/45 0s 377us/step -
auc_449: 0.9374 - loss: 0.2069
Epoch 85/250
45/45 0s 384us/step -
auc_449: 0.9376 - loss: 0.2063
Epoch 86/250
45/45 0s 379us/step -
auc_449: 0.9381 - loss: 0.2058
Epoch 87/250
45/45 0s 392us/step -
auc_449: 0.9384 - loss: 0.2052

Epoch 88/250
45/45 0s 390us/step -
auc_449: 0.9389 - loss: 0.2046
Epoch 89/250
45/45 0s 384us/step -
auc_449: 0.9393 - loss: 0.2040
Epoch 90/250
45/45 0s 388us/step -
auc_449: 0.9396 - loss: 0.2034
Epoch 91/250
45/45 0s 380us/step -
auc_449: 0.9400 - loss: 0.2029
Epoch 92/250
45/45 0s 391us/step -
auc_449: 0.9405 - loss: 0.2023
Epoch 93/250
45/45 0s 390us/step -
auc_449: 0.9409 - loss: 0.2018
Epoch 94/250
45/45 0s 396us/step -
auc_449: 0.9413 - loss: 0.2012
Epoch 95/250
45/45 0s 397us/step -
auc_449: 0.9415 - loss: 0.2007
Epoch 96/250
45/45 0s 386us/step -
auc_449: 0.9420 - loss: 0.2001
Epoch 97/250
45/45 0s 396us/step -
auc_449: 0.9423 - loss: 0.1996
Epoch 98/250
45/45 0s 387us/step -
auc_449: 0.9427 - loss: 0.1991
Epoch 99/250
45/45 0s 409us/step -
auc_449: 0.9432 - loss: 0.1986
Epoch 100/250
45/45 0s 390us/step -
auc_449: 0.9435 - loss: 0.1980
Epoch 101/250
45/45 0s 381us/step -
auc_449: 0.9440 - loss: 0.1975
Epoch 102/250
45/45 0s 397us/step -
auc_449: 0.9445 - loss: 0.1970
Epoch 103/250
45/45 0s 395us/step -
auc_449: 0.9448 - loss: 0.1965

Epoch 104/250
45/45 0s 448us/step -
auc_449: 0.9451 - loss: 0.1959
Epoch 105/250
45/45 0s 393us/step -
auc_449: 0.9454 - loss: 0.1954
Epoch 106/250
45/45 0s 389us/step -
auc_449: 0.9457 - loss: 0.1949
Epoch 107/250
45/45 0s 389us/step -
auc_449: 0.9459 - loss: 0.1944
Epoch 108/250
45/45 0s 379us/step -
auc_449: 0.9463 - loss: 0.1938
Epoch 109/250
45/45 0s 383us/step -
auc_449: 0.9465 - loss: 0.1933
Epoch 110/250
45/45 0s 409us/step -
auc_449: 0.9469 - loss: 0.1928
Epoch 111/250
45/45 0s 385us/step -
auc_449: 0.9472 - loss: 0.1923
Epoch 112/250
45/45 0s 415us/step -
auc_449: 0.9476 - loss: 0.1917
Epoch 113/250
45/45 0s 387us/step -
auc_449: 0.9480 - loss: 0.1912
Epoch 114/250
45/45 0s 389us/step -
auc_449: 0.9482 - loss: 0.1907
Epoch 115/250
45/45 0s 394us/step -
auc_449: 0.9486 - loss: 0.1902
Epoch 116/250
45/45 0s 386us/step -
auc_449: 0.9488 - loss: 0.1897
Epoch 117/250
45/45 0s 398us/step -
auc_449: 0.9489 - loss: 0.1892
Epoch 118/250
45/45 0s 389us/step -
auc_449: 0.9492 - loss: 0.1887
Epoch 119/250
45/45 0s 375us/step -
auc_449: 0.9495 - loss: 0.1882

Epoch 120/250
45/45 0s 390us/step -
auc_449: 0.9497 - loss: 0.1878
Epoch 121/250
45/45 0s 377us/step -
auc_449: 0.9500 - loss: 0.1873
Epoch 122/250
45/45 0s 376us/step -
auc_449: 0.9504 - loss: 0.1868
Epoch 123/250
45/45 0s 383us/step -
auc_449: 0.9507 - loss: 0.1863
Epoch 124/250
45/45 0s 382us/step -
auc_449: 0.9509 - loss: 0.1859
Epoch 125/250
45/45 0s 388us/step -
auc_449: 0.9512 - loss: 0.1854
Epoch 126/250
45/45 0s 393us/step -
auc_449: 0.9515 - loss: 0.1850
Epoch 127/250
45/45 0s 378us/step -
auc_449: 0.9518 - loss: 0.1845
Epoch 128/250
45/45 0s 378us/step -
auc_449: 0.9520 - loss: 0.1841
Epoch 129/250
45/45 0s 385us/step -
auc_449: 0.9522 - loss: 0.1836
Epoch 130/250
45/45 0s 418us/step -
auc_449: 0.9524 - loss: 0.1832
Epoch 131/250
45/45 0s 396us/step -
auc_449: 0.9528 - loss: 0.1828
Epoch 132/250
45/45 0s 372us/step -
auc_449: 0.9530 - loss: 0.1823
Epoch 133/250
45/45 0s 372us/step -
auc_449: 0.9532 - loss: 0.1819
Epoch 134/250
45/45 0s 374us/step -
auc_449: 0.9534 - loss: 0.1815
Epoch 135/250
45/45 0s 436us/step -
auc_449: 0.9535 - loss: 0.1810

Epoch 136/250
45/45 0s 373us/step -
auc_449: 0.9537 - loss: 0.1806
Epoch 137/250
45/45 0s 379us/step -
auc_449: 0.9539 - loss: 0.1802
Epoch 138/250
45/45 0s 374us/step -
auc_449: 0.9541 - loss: 0.1798
Epoch 139/250
45/45 0s 373us/step -
auc_449: 0.9543 - loss: 0.1794
Epoch 140/250
45/45 0s 377us/step -
auc_449: 0.9546 - loss: 0.1790
Epoch 141/250
45/45 0s 380us/step -
auc_449: 0.9548 - loss: 0.1786
Epoch 142/250
45/45 0s 376us/step -
auc_449: 0.9552 - loss: 0.1782
Epoch 143/250
45/45 0s 400us/step -
auc_449: 0.9553 - loss: 0.1778
Epoch 144/250
45/45 0s 383us/step -
auc_449: 0.9555 - loss: 0.1774
Epoch 145/250
45/45 0s 377us/step -
auc_449: 0.9557 - loss: 0.1770
Epoch 146/250
45/45 0s 402us/step -
auc_449: 0.9559 - loss: 0.1766
Epoch 147/250
45/45 0s 385us/step -
auc_449: 0.9560 - loss: 0.1762
Epoch 148/250
45/45 0s 369us/step -
auc_449: 0.9562 - loss: 0.1759
Epoch 149/250
45/45 0s 375us/step -
auc_449: 0.9564 - loss: 0.1755
Epoch 150/250
45/45 0s 382us/step -
auc_449: 0.9567 - loss: 0.1751
Epoch 151/250
45/45 0s 380us/step -
auc_449: 0.9569 - loss: 0.1747

Epoch 152/250
45/45 0s 383us/step -
auc_449: 0.9571 - loss: 0.1743
Epoch 153/250
45/45 0s 391us/step -
auc_449: 0.9574 - loss: 0.1740
Epoch 154/250
45/45 0s 374us/step -
auc_449: 0.9576 - loss: 0.1736
Epoch 155/250
45/45 0s 371us/step -
auc_449: 0.9579 - loss: 0.1732
Epoch 156/250
45/45 0s 379us/step -
auc_449: 0.9580 - loss: 0.1728
Epoch 157/250
45/45 0s 373us/step -
auc_449: 0.9582 - loss: 0.1724
Epoch 158/250
45/45 0s 373us/step -
auc_449: 0.9584 - loss: 0.1721
Epoch 159/250
45/45 0s 379us/step -
auc_449: 0.9586 - loss: 0.1717
Epoch 160/250
45/45 0s 389us/step -
auc_449: 0.9587 - loss: 0.1714
Epoch 161/250
45/45 0s 405us/step -
auc_449: 0.9588 - loss: 0.1710
Epoch 162/250
45/45 0s 376us/step -
auc_449: 0.9588 - loss: 0.1706
Epoch 163/250
45/45 0s 378us/step -
auc_449: 0.9589 - loss: 0.1703
Epoch 164/250
45/45 0s 381us/step -
auc_449: 0.9591 - loss: 0.1699
Epoch 165/250
45/45 0s 379us/step -
auc_449: 0.9592 - loss: 0.1695
Epoch 166/250
45/45 0s 380us/step -
auc_449: 0.9594 - loss: 0.1692
Epoch 167/250
45/45 0s 378us/step -
auc_449: 0.9596 - loss: 0.1688

Epoch 168/250
45/45 0s 372us/step -
auc_449: 0.9597 - loss: 0.1685
Epoch 169/250
45/45 0s 372us/step -
auc_449: 0.9599 - loss: 0.1681
Epoch 170/250
45/45 0s 396us/step -
auc_449: 0.9601 - loss: 0.1678
Epoch 171/250
45/45 0s 375us/step -
auc_449: 0.9603 - loss: 0.1674
Epoch 172/250
45/45 0s 375us/step -
auc_449: 0.9605 - loss: 0.1671
Epoch 173/250
45/45 0s 374us/step -
auc_449: 0.9608 - loss: 0.1668
Epoch 174/250
45/45 0s 373us/step -
auc_449: 0.9609 - loss: 0.1664
Epoch 175/250
45/45 0s 369us/step -
auc_449: 0.9611 - loss: 0.1661
Epoch 176/250
45/45 0s 372us/step -
auc_449: 0.9613 - loss: 0.1658
Epoch 177/250
45/45 0s 371us/step -
auc_449: 0.9614 - loss: 0.1654
Epoch 178/250
45/45 0s 373us/step -
auc_449: 0.9615 - loss: 0.1651
Epoch 179/250
45/45 0s 369us/step -
auc_449: 0.9617 - loss: 0.1648
Epoch 180/250
45/45 0s 371us/step -
auc_449: 0.9619 - loss: 0.1644
Epoch 181/250
45/45 0s 420us/step -
auc_449: 0.9621 - loss: 0.1641
Epoch 182/250
45/45 0s 379us/step -
auc_449: 0.9623 - loss: 0.1638
Epoch 183/250
45/45 0s 375us/step -
auc_449: 0.9624 - loss: 0.1634

Epoch 184/250
45/45 0s 368us/step -
auc_449: 0.9625 - loss: 0.1631
Epoch 185/250
45/45 0s 382us/step -
auc_449: 0.9627 - loss: 0.1628
Epoch 186/250
45/45 0s 371us/step -
auc_449: 0.9628 - loss: 0.1625
Epoch 187/250
45/45 0s 379us/step -
auc_449: 0.9630 - loss: 0.1622
Epoch 188/250
45/45 0s 375us/step -
auc_449: 0.9631 - loss: 0.1618
Epoch 189/250
45/45 0s 374us/step -
auc_449: 0.9633 - loss: 0.1615
Epoch 190/250
45/45 0s 375us/step -
auc_449: 0.9634 - loss: 0.1612
Epoch 191/250
45/45 0s 371us/step -
auc_449: 0.9634 - loss: 0.1609
Epoch 192/250
45/45 0s 374us/step -
auc_449: 0.9636 - loss: 0.1606
Epoch 193/250
45/45 0s 370us/step -
auc_449: 0.9637 - loss: 0.1603
Epoch 194/250
45/45 0s 370us/step -
auc_449: 0.9639 - loss: 0.1600
Epoch 195/250
45/45 0s 406us/step -
auc_449: 0.9640 - loss: 0.1598
Epoch 196/250
45/45 0s 377us/step -
auc_449: 0.9641 - loss: 0.1595
Epoch 197/250
45/45 0s 373us/step -
auc_449: 0.9642 - loss: 0.1592
Epoch 198/250
45/45 0s 372us/step -
auc_449: 0.9643 - loss: 0.1589
Epoch 199/250
45/45 0s 373us/step -
auc_449: 0.9644 - loss: 0.1586

Epoch 200/250
45/45 0s 367us/step -
auc_449: 0.9646 - loss: 0.1583
Epoch 201/250
45/45 0s 369us/step -
auc_449: 0.9647 - loss: 0.1581
Epoch 202/250
45/45 0s 371us/step -
auc_449: 0.9649 - loss: 0.1578
Epoch 203/250
45/45 0s 370us/step -
auc_449: 0.9650 - loss: 0.1575
Epoch 204/250
45/45 0s 370us/step -
auc_449: 0.9651 - loss: 0.1572
Epoch 205/250
45/45 0s 368us/step -
auc_449: 0.9653 - loss: 0.1570
Epoch 206/250
45/45 0s 371us/step -
auc_449: 0.9654 - loss: 0.1567
Epoch 207/250
45/45 0s 371us/step -
auc_449: 0.9656 - loss: 0.1564
Epoch 208/250
45/45 0s 369us/step -
auc_449: 0.9656 - loss: 0.1561
Epoch 209/250
45/45 0s 370us/step -
auc_449: 0.9658 - loss: 0.1558
Epoch 210/250
45/45 0s 372us/step -
auc_449: 0.9659 - loss: 0.1556
Epoch 211/250
45/45 0s 372us/step -
auc_449: 0.9661 - loss: 0.1553
Epoch 212/250
45/45 0s 371us/step -
auc_449: 0.9662 - loss: 0.1550
Epoch 213/250
45/45 0s 370us/step -
auc_449: 0.9663 - loss: 0.1547
Epoch 214/250
45/45 0s 370us/step -
auc_449: 0.9664 - loss: 0.1545
Epoch 215/250
45/45 0s 404us/step -
auc_449: 0.9666 - loss: 0.1542

Epoch 216/250
45/45 0s 378us/step -
auc_449: 0.9667 - loss: 0.1539
Epoch 217/250
45/45 0s 374us/step -
auc_449: 0.9668 - loss: 0.1536
Epoch 218/250
45/45 0s 370us/step -
auc_449: 0.9669 - loss: 0.1534
Epoch 219/250
45/45 0s 373us/step -
auc_449: 0.9671 - loss: 0.1531
Epoch 220/250
45/45 0s 375us/step -
auc_449: 0.9672 - loss: 0.1528
Epoch 221/250
45/45 0s 376us/step -
auc_449: 0.9673 - loss: 0.1526
Epoch 222/250
45/45 0s 381us/step -
auc_449: 0.9674 - loss: 0.1523
Epoch 223/250
45/45 0s 376us/step -
auc_449: 0.9675 - loss: 0.1520
Epoch 224/250
45/45 0s 376us/step -
auc_449: 0.9676 - loss: 0.1517
Epoch 225/250
45/45 0s 375us/step -
auc_449: 0.9678 - loss: 0.1515
Epoch 226/250
45/45 0s 371us/step -
auc_449: 0.9679 - loss: 0.1512
Epoch 227/250
45/45 0s 373us/step -
auc_449: 0.9679 - loss: 0.1509
Epoch 228/250
45/45 0s 373us/step -
auc_449: 0.9681 - loss: 0.1506
Epoch 229/250
45/45 0s 376us/step -
auc_449: 0.9682 - loss: 0.1503
Epoch 230/250
45/45 0s 372us/step -
auc_449: 0.9683 - loss: 0.1501
Epoch 231/250
45/45 0s 374us/step -
auc_449: 0.9683 - loss: 0.1498

Epoch 232/250
45/45 0s 363us/step -
auc_449: 0.9685 - loss: 0.1495
Epoch 233/250
45/45 0s 364us/step -
auc_449: 0.9687 - loss: 0.1492
Epoch 234/250
45/45 0s 370us/step -
auc_449: 0.9688 - loss: 0.1490
Epoch 235/250
45/45 0s 377us/step -
auc_449: 0.9689 - loss: 0.1487
Epoch 236/250
45/45 0s 372us/step -
auc_449: 0.9690 - loss: 0.1484
Epoch 237/250
45/45 0s 374us/step -
auc_449: 0.9692 - loss: 0.1481
Epoch 238/250
45/45 0s 370us/step -
auc_449: 0.9692 - loss: 0.1479
Epoch 239/250
45/45 0s 369us/step -
auc_449: 0.9694 - loss: 0.1476
Epoch 240/250
45/45 0s 369us/step -
auc_449: 0.9696 - loss: 0.1473
Epoch 241/250
45/45 0s 370us/step -
auc_449: 0.9698 - loss: 0.1471
Epoch 242/250
45/45 0s 369us/step -
auc_449: 0.9699 - loss: 0.1468
Epoch 243/250
45/45 0s 413us/step -
auc_449: 0.9700 - loss: 0.1465
Epoch 244/250
45/45 0s 421us/step -
auc_449: 0.9701 - loss: 0.1463
Epoch 245/250
45/45 0s 396us/step -
auc_449: 0.9702 - loss: 0.1460
Epoch 246/250
45/45 0s 385us/step -
auc_449: 0.9704 - loss: 0.1457
Epoch 247/250
45/45 0s 377us/step -
auc_449: 0.9705 - loss: 0.1455

```
Epoch 248/250
45/45          0s 374us/step -
auc_449: 0.9706 - loss: 0.1452
Epoch 249/250
45/45          0s 376us/step -
auc_449: 0.9708 - loss: 0.1449
Epoch 250/250
45/45          0s 375us/step -
auc_449: 0.9709 - loss: 0.1447
```

[299]: <keras.src.callbacks.history.History at 0x34bf25e90>

```
[300]: # Generate the predicted probability of y
y_score_mlp = mlp.predict(X_test_preproc)

# Convert probability to y value
y_pred_mlp = (y_score_mlp > 0.5).astype(int)

# Generate the test loss and AUC
test_loss, auc_mlp=mlp.evaluate(X_test_preproc, y_test)
```

```
45/45          0s 549us/step
45/45          0s 341us/step -
auc_449: 0.9546 - loss: 0.1823
```

```
[301]: #Generate performance metrics
accuracy_mlp = accuracy_score(y_test, y_pred_mlp)
conf_matrix_mlp = confusion_matrix(y_test, y_pred_mlp)
classification_rep_mlp = classification_report(y_test, y_pred_mlp)
# auc_mlp=roc_auc_score(y_test, y_score_mlp)

print(f"Test accuracy: \n {accuracy_mlp:.2%}")
print ("Confusion Matrix : \n", conf_matrix_mlp)
print("\nClassification Report: : \n",classification_rep_mlp)
print(f"AUC: \n {auc_mlp:.2%}")
```

Test accuracy:
92.03%

Confusion Matrix :
[[1161 44]
[69 143]]

Classification Report: :

	precision	recall	f1-score	support
0	0.94	0.96	0.95	1205
1	0.76	0.67	0.72	212
accuracy			0.92	1417

macro avg	0.85	0.82	0.84	1417
weighted avg	0.92	0.92	0.92	1417

AUC:
94.38%

```
[302]: data = {'Models': ['Random Forest', 'XGBoost', 'MLP'],
              'F1 score on test set': [0.76, 0.80, 0.72]}
data= pd.DataFrame(data = data)
data.head()
```

```
[302]:
```

	Models	F1 score on test set
0	Random Forest	0.76
1	XGBoost	0.80
2	MLP	0.72

6 Conclusion

I'll select XGBoost as predictive model for this business problem due to its high f1 score(0.8) and AUC (97.12%).

This model can effectively identify 82% of customers who will close their credit cards. However, within the customers identified as attrited, 21% actually do not close their cards. A trade-off relationship is observed between the precision and recall of this model. Therefore, further evaluation is needed to compare the cost of customer churn with the reduction in profit from offering promotions to customers who will not churn. If the cost of customer churn is greater, the model's scale_pos_weight should be increased to improve its recall. Conversely, if the reduction in profit from offering promotions to customers who will not churn is greater, the model's scale_pos_weight should be adjusted downwards to improve precision.