



# LEAD SCORING CASE STUDY

LE THANH DAT





# DATA PREPROCESSING, FEATURE SELECTION, MODEL OPTIMIZATION

# NECESSARY LIBRARIES

- pandas, numpy
- seaborn, matplotlib
- sklearn.metrics
- optbinning
- xgboost

```
import pandas as pd
from optbinning import BinningProcess
```

```
from sklearn.linear_model import LogisticRegression
from sklearn.model_selection import train_test_split, RandomizedSearchCV
from sklearn.metrics import accuracy_score, classification_report
```

```
from xgboost import XGBClassifier
from sklearn.model_selection import RandomizedSearchCV
from sklearn.metrics import accuracy_score, classification_report, roc_curve, roc_auc_score
import matplotlib.pyplot as plt
```

# DATA EXPLORATION

- Check duplicated/missing values
- Handle outliers

## Numeric & categorical variables

```
def choose_numeric_categorical(data):  
    numeric_vars = data.select_dtypes(include=['number']).columns.tolist()  
    categorical_vars = data.select_dtypes(include=['object', 'category']).columns.tolist()  
    return numeric_vars, categorical_vars  
numeric_vars, categorical_vars = choose_numeric_categorical(data)
```

## Check missing values

```
missing_values_percentage = round(100 * (data.isna().sum() / len(data)), 2)  
sorted_missing_values_percentage = missing_values_percentage.sort_values(ascending=False)  
sorted_missing_values_percentage
```

## Data exploration

```
import matplotlib.pyplot as plt  
import seaborn as sns  
import numpy as np  
  
def plot_numeric_variable_analysis(numeric_columns, dataframe):  
    for numeric_column in numeric_columns:  
        plt.figure(figsize=(12, 6))  
        print(numeric_column)  
  
        # Calculate quartiles  
        quartiles = np.percentile(dataframe[numeric_column].dropna(), [25, 50, 75, 90, 99])  
        print(f'25th percentile: {quartiles[0]}')  
        print(f'50th percentile: {quartiles[1]}')  
        print(f'75th percentile: {quartiles[2]}')  
        print(f'90th percentile: {quartiles[3]}')  
        print(f'99th percentile: {quartiles[4]}')  
  
        # Plot Violin plot  
        plt.subplot(1, 2, 1)  
        sns.violinplot(y=dataframe[numeric_column])  
        plt.title(f'Violin Plot of {numeric_column}')  
        plt.ylabel(numeric_column)  
  
        # Plot CDF (Cumulative Distribution Function)  
        plt.subplot(1, 2, 2)  
        sorted_values = np.sort(dataframe[numeric_column].dropna())  
        cdf_values = np.arange(1, len(sorted_values) + 1) / len(sorted_values)  
        plt.plot(sorted_values, cdf_values, marker='.', linestyle='none')  
        plt.title(f'CDF of {numeric_column}')  
        plt.xlabel(numeric_column)  
        plt.ylabel('CDF')  
  
        plt.tight_layout()  
        plt.show()
```

```
plot_numeric_variable_analysis(numeric_vars, data)
```

# DATA EXLORATION

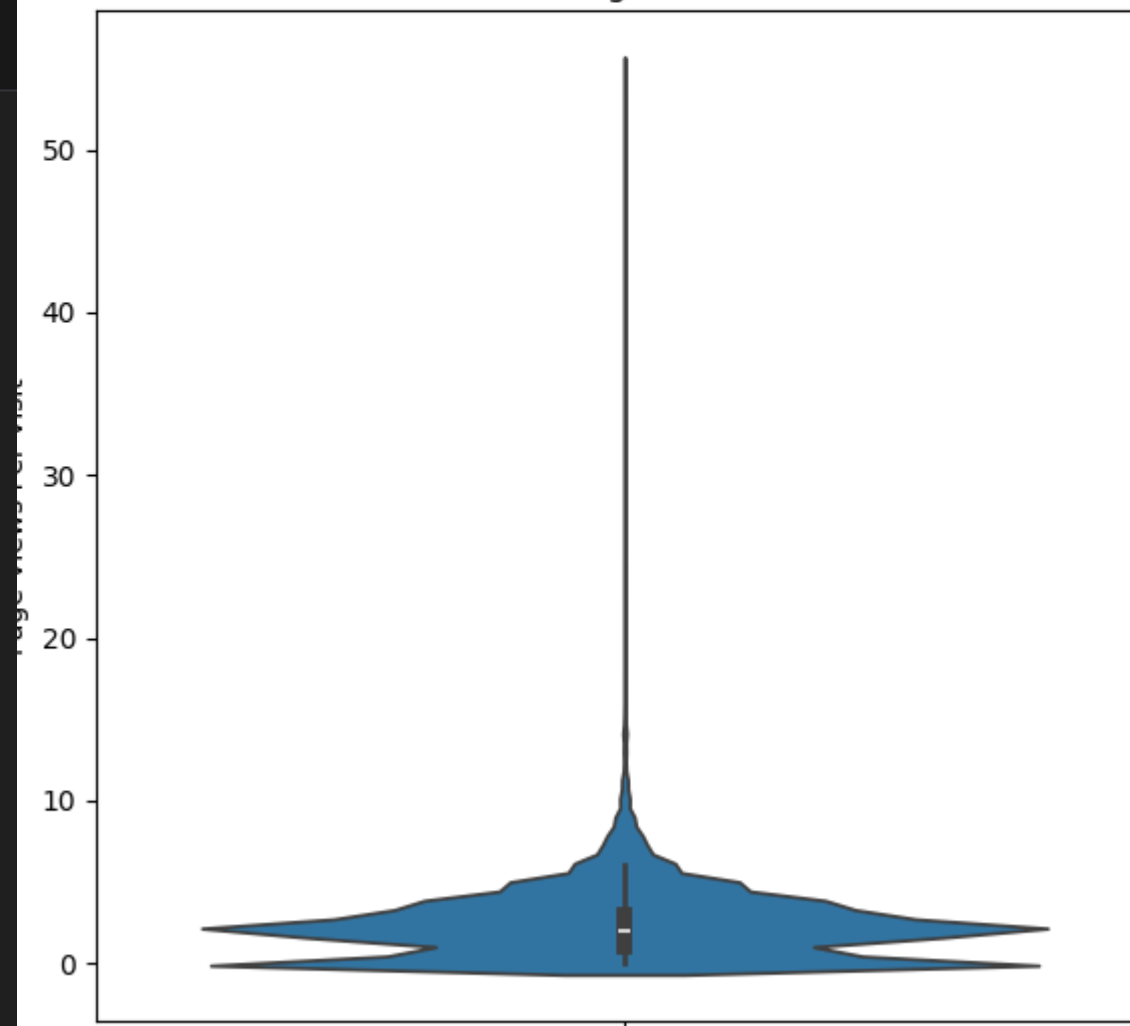
The images show that there are numerous NaN values across the columns. Notably, the Page Views Per Visit column contains some outliers.

```
Page Views Per Visit
25th percentile: 1.0
50th percentile: 2.0
75th percentile: 3.0
90th percentile: 5.0
99th percentile: 9.0
```

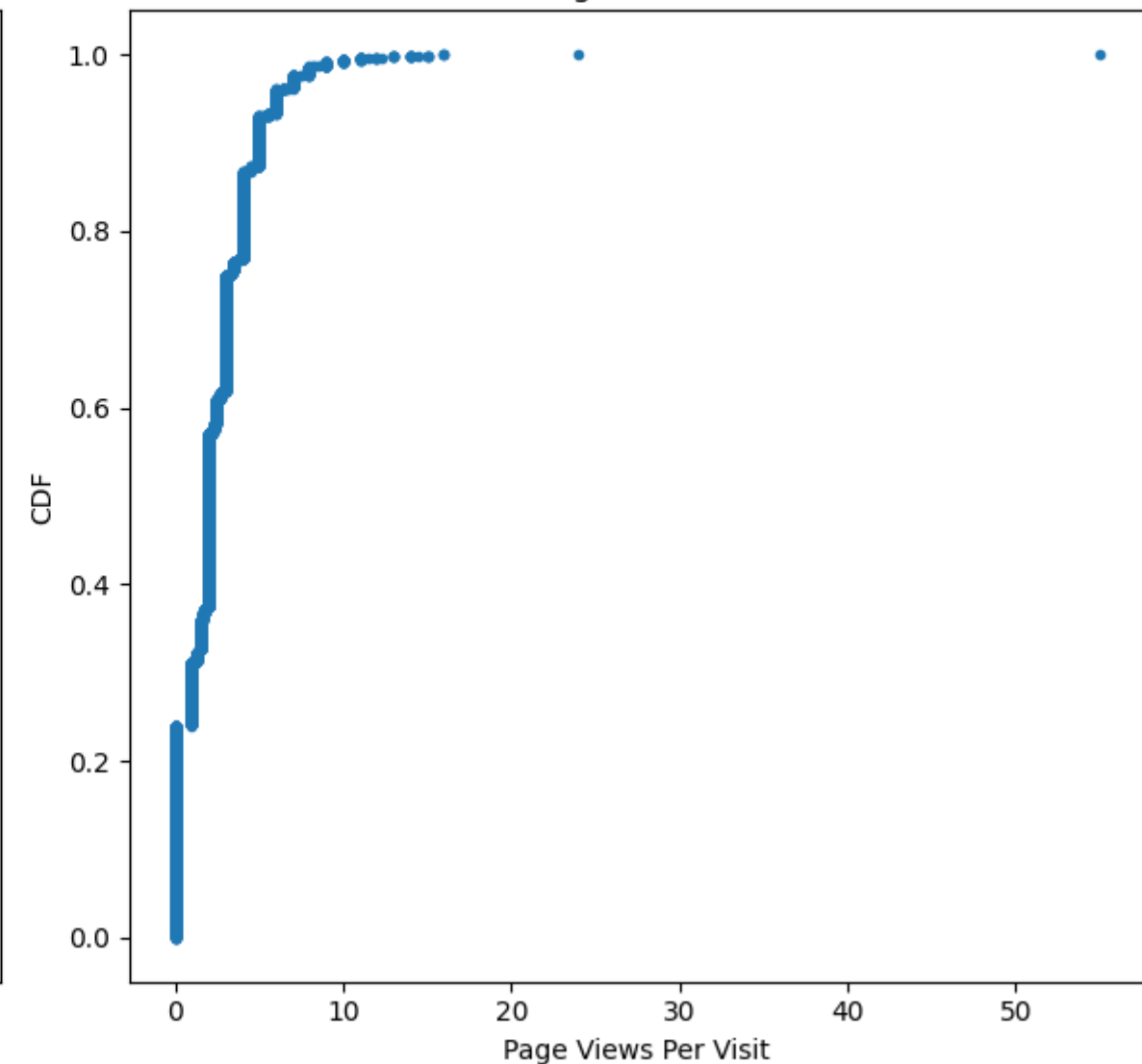
```
missing_values_percentage = round(100 * (data.isna().sum() / len(data)), 2)
sorted_missing_values_percentage = missing_values_percentage.sort_values(ascending=False)
sorted_missing_values_percentage
```

|   |       |
|---|-------|
| Lead Quality                                  | 51.59 |
| Asymmetrique Activity Index                   | 45.65 |
| Asymmetrique Profile Score                    | 45.65 |
| Asymmetrique Profile Index                    | 45.65 |
| Asymmetrique Activity Score                   | 45.65 |
| Tags  | 36.29 |
| Lead Profile                                  | 29.32 |
| What matters most to you in choosing a course | 29.32 |
| What is your current occupation               | 29.11 |
| Country                                       | 26.63 |
| How did you hear about X Education            | 23.89 |
| Specialization                                | 15.56 |
| City  | 15.37 |
| Page Views Per Visit                          | 1.48  |
| TotalVisits                                   | 1.48  |
| Last Activity                                 | 1.11  |
| Lead Source                                   | 0.39  |
| I agree to pay the amount through cheque      | 0.00  |
| A free copy of Mastering The Interview        | 0.00  |
| Get updates on DM Content                     | 0.00  |
| Update me on Supply Chain Content             | 0.00  |
| Lead Origin                                   | 0.00  |
| X Education Forums                            | 0.00  |
| Receive More Updates About Our Courses        | 0.00  |
| Through Recommendations                       | 0.00  |
| ...   |       |
| Converted                                     | 0.00  |
| Do Not Call                                   | 0.00  |
| Do Not Email                                  | 0.00  |
| Last Notable Activity                         | 0.00  |

Violin Plot of Page Views Per Visit



CDF of Page Views Per Visit



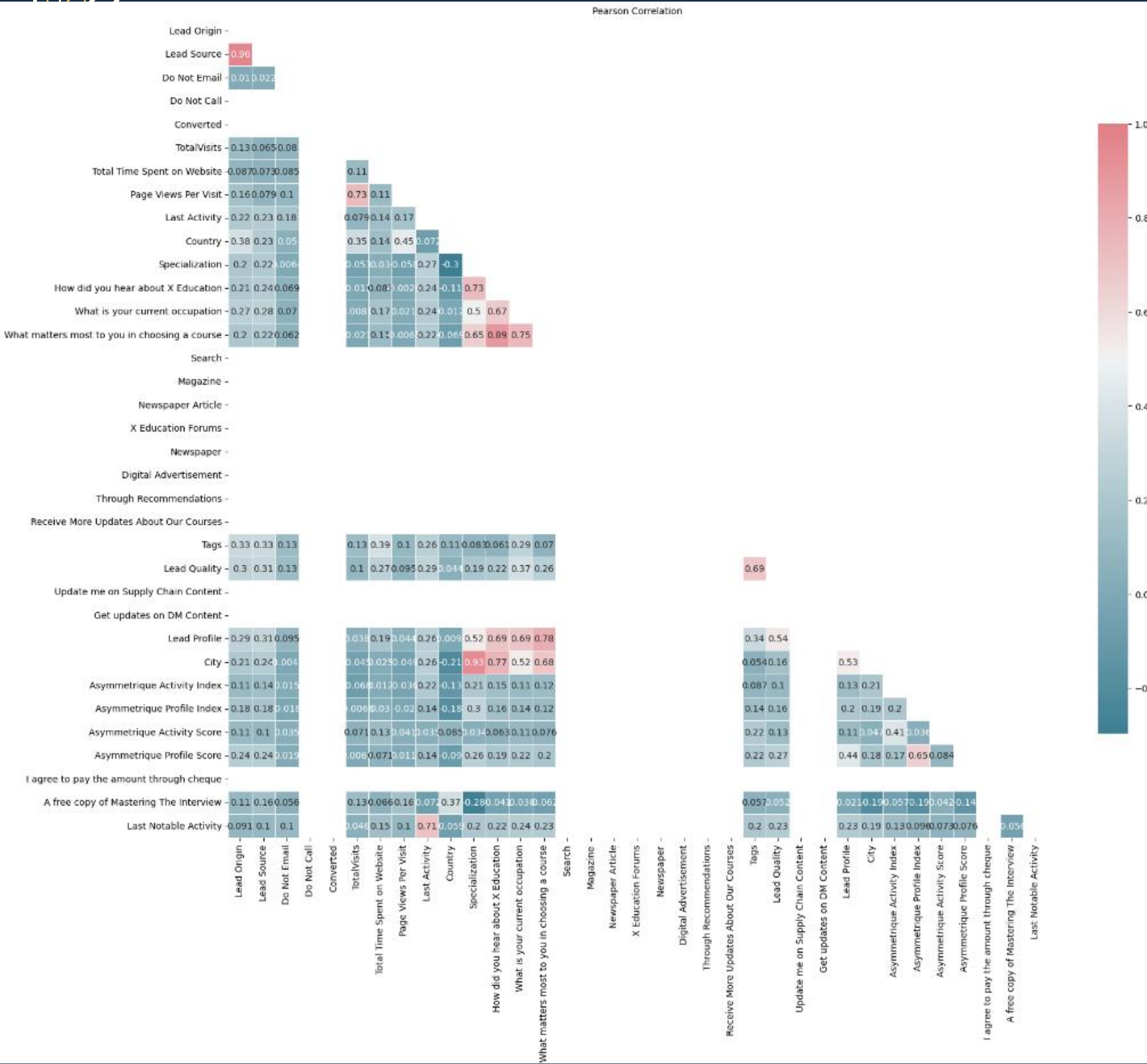


- Pearson correlation to check multi-collinearity: from the Lead data, export to iv\_table.csv and binning\_table.csv to identify which variables have the highest IV and should be included

1



# FEATURE SELECTION USING CORRELATION AND IV



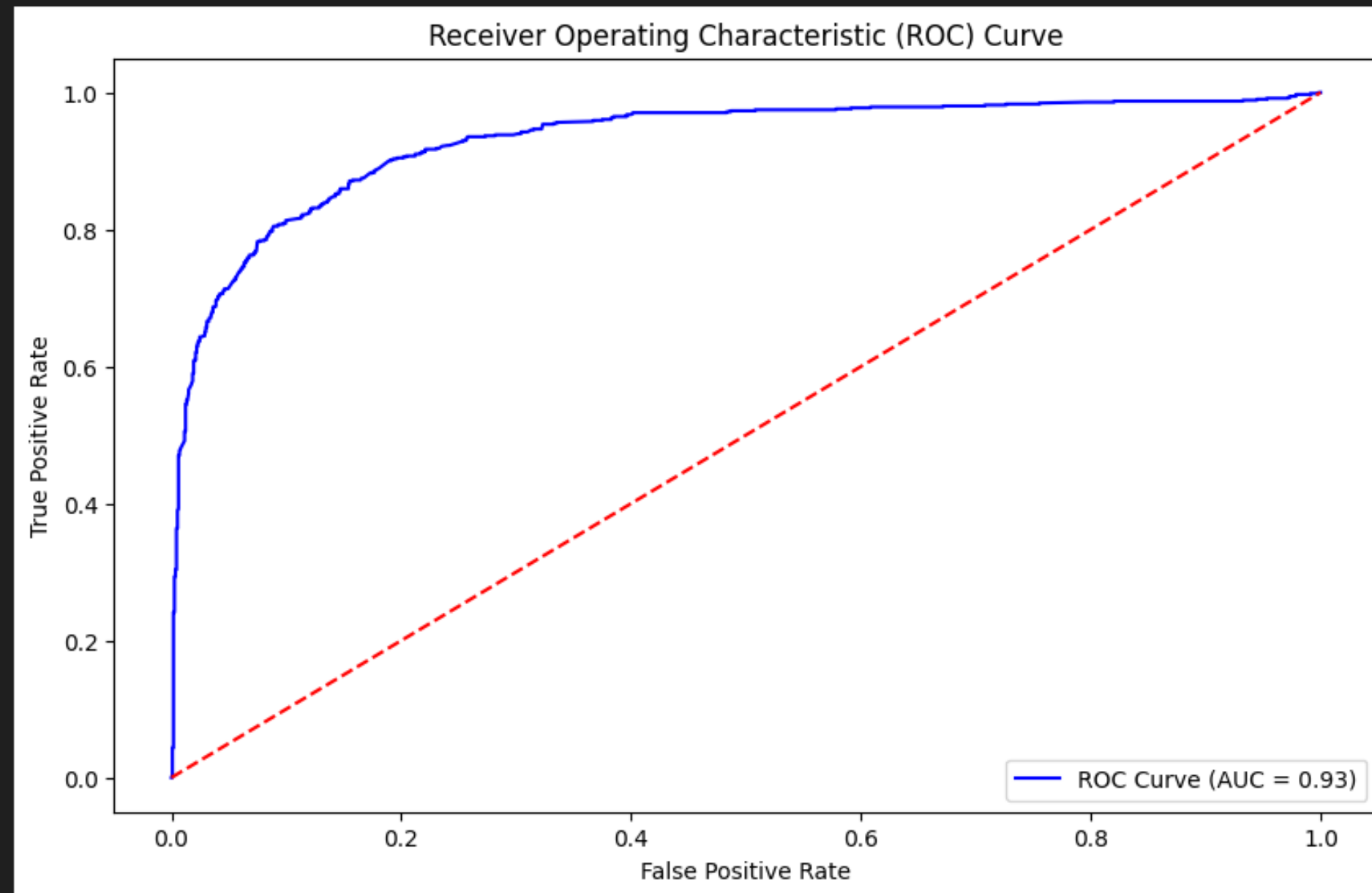
Using heatmap visualization and remove highly correlated with low IV feature values:

Certain conditions must be applied to filter for useful variables: correlation should be less than 0.6, and if it exceeds 0.6, the variable with the lower IV should be removed. Additionally, only variables with an IV greater than 0.6 should be selected.

# MODEL PERFORMANCE SUMMARY

- Optimized Logistic Regression
  - Accuracy: 0.86851
  - AUC Score: 0.93027

Accuracy: 0.8685064935064936  
AUC Score: 0.9302710575249249

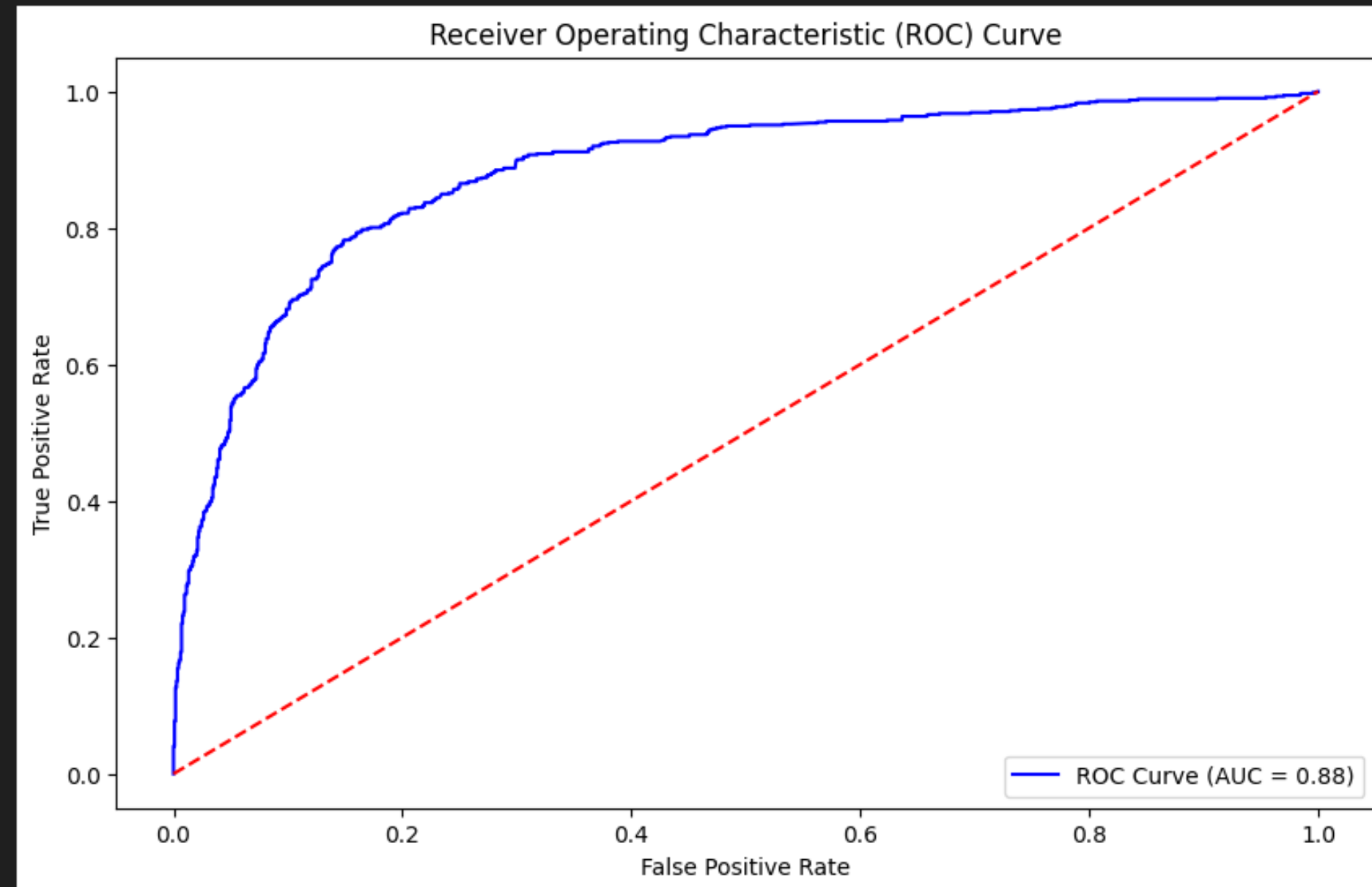




# MODEL PERFORMANCE SUMMARY

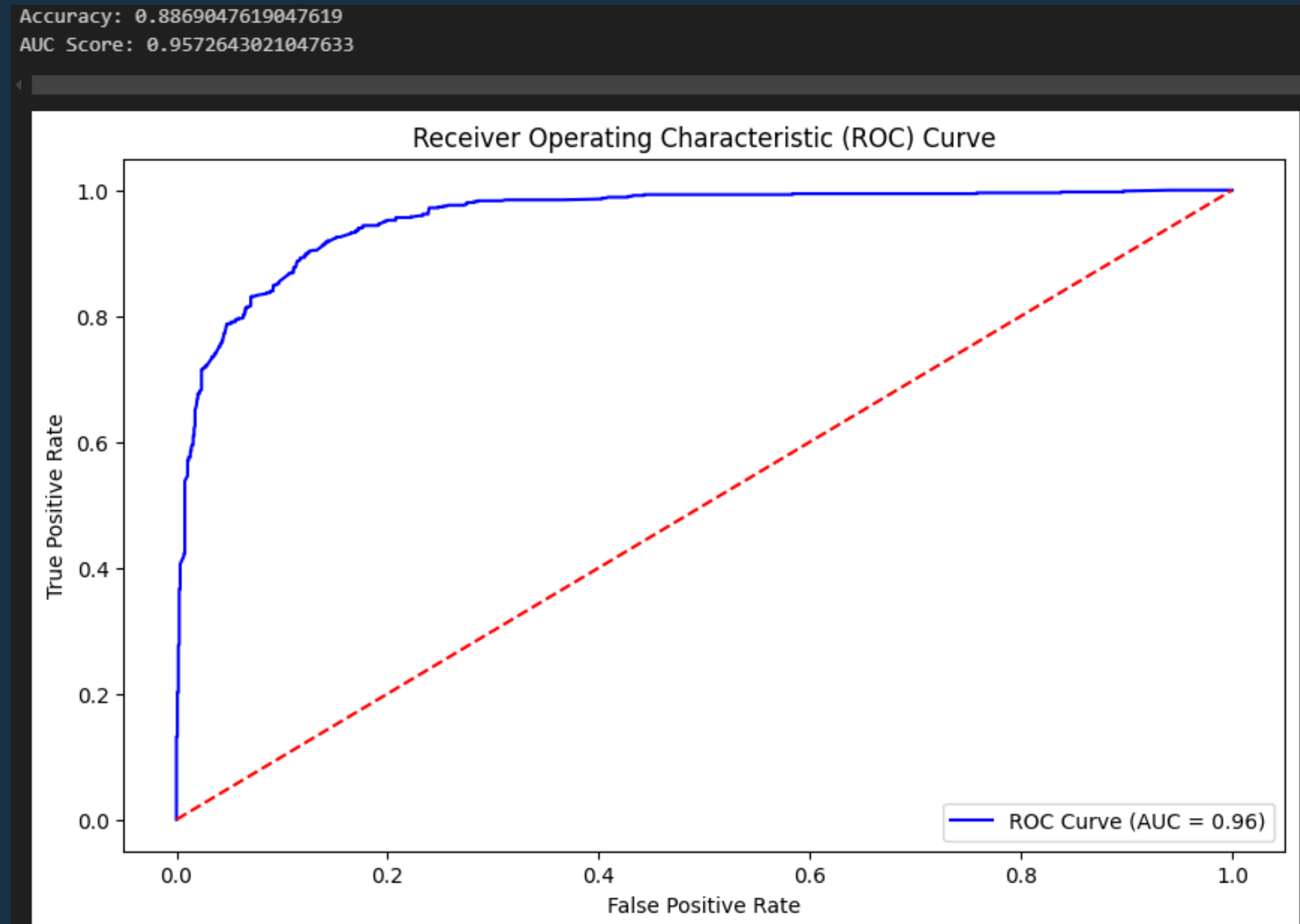
- Baseline Logistic Regression
  - Accuracy: 0.81818
  - AUC Score: 0.88224

Accuracy: 0.81818181818182  
AUC Score: 0.8822468695600569



# MODEL PERFORMANCE SUMMARY

- Optimized XGBoost
  - Accuracy: 0.88690
  - AUC Score: 0.95726



# XGBOOST THRESHOLD OPTIMIZATION

Update threshold for accuracy optimization

- Highest optimal threshold: 0.649999
- Highest accuracy: 0.888528

```
▶ # Define the range of thresholds to evaluate
threshold_values = np.arange(0.4, 0.9, 0.05)

# Get the predicted probabilities from the best model
y_test_pred_probabilities = best_xgb_model.predict_proba(x_test_data)[: , 1]

# Initialize variables to store the best threshold and the corresponding accuracy
optimal_threshold = 0.0
highest_accuracy = 0.0

# Iterate through each threshold value
for threshold in threshold_values:
    # Apply the threshold to get binary predictions
    y_threshold_predictions = (y_test_pred_probabilities >= threshold).astype(int)

    # Calculate accuracy for the current threshold
    current_accuracy = accuracy_score(y_test_labels, y_threshold_predictions)

    # Update the best threshold if the current accuracy is higher
    if current_accuracy > highest_accuracy:
        highest_accuracy = current_accuracy
        optimal_threshold = threshold

# Output the best threshold and corresponding accuracy
print(f"Optimal Threshold: {optimal_threshold}")
print(f"Highest Accuracy: {highest_accuracy}")
```

[27]

```
... Optimal Threshold: 0.6499999999999999
Highest Accuracy: 0.8885281385281385
```