# Predicting Loan Defaults: A Data-Driven Approach to Credit Risk Analysis

**BEE2041 - Data Science in Economics**

Student Number - 720017170

## Table of contents

## 1. Introduction

Access to credit is a important driver of economic growth, allowing households or buisnesses to invest, expand and smooth consumption. However, credit risk remains a fundemental challenge for financial institutions, as loan defaulting can lead to substantial financial losses for both the company and for stakeholders. The ability to predict these defaults is vital for lending institutions to mitigate their risk and make more informed lending predictions. Recent advancements in machine learning (ML) have aided in the development of robust predictive models that outperform traditional credit-scoring methods (Yang, 2024)

Ensemble methods such as Random Forest (RF), XGBoost, and Light Gradient Boosting Machines (LGBM), have shown significant promise in improving classification accuracy over traditional statistical methods (Yadav, 2025). These models offer enhanced predictive capacity due to their ability to capture non-linear relationships in borrower data, providing financial institutions with more reliable risk assessment (Roy, 2025)

This study aims to explore a data-driven approach to credit risk analysis by using ML methods to predict loan defaulting. Logistic regression (LR), RF, XGBoost and LGBM have all been implemented and compared using standard performance metrics such as accuracy, precision, recall, F1-score and area under the curve (AUC). Moreover, exploratory data analysis will be conducted to examine the distribution of important financial variables, identify correlations and allow for optimised feature selection to improve model performance.

Due to the increasing reliance on alternative data sources and advanced computational methods in the financial sector, the results of this study may have significant practical implications. Improved credit risk analysis can help lenders reduce default rates, minimise losses and promote more inclusive access to credit (Ellsworth, 2025). By leveraging the latest ML methods, this project aims to contribute to the growing body of research on predictive analytics in finance and support more robust lending practices (Khoshkhoy Nilash & Esmaeilpour, 2025).

## 2. Data

Prior to conducting the analysis of credit risk, we need to understand and organise the data. For this analysis we will be using a loan defaulting dataset from Kaggle (reference), consisting of 12 variables/columns and 28,501 observations.

```
PersonAge            0
PersonIncome         0
PersonHomeOwnership  0
PersonEmpLength      0
LoanIntent           0
LoanGrade            0
LoanAmnt             0
```

```
LoanIntRate            0
LoanStatus             0
LoanPercentIncome      0
PreviousDefault        0
CredHistory            0
dtype: int64
```

Table 1: Variable Information

| Variable | Data Type | Definition |
| --- | --- | --- |
| PersonAge | int64 | Age of the borrower |
| PersonIncome | int64 | Income of the borrower |
| PersonHomeOwnership | object | Home ownership of the borrower |
| PersonEmpLength | float64 | Employment length of the borrower |
| LoanIntent | object | Intention of the loan |
| LoanGrade | int64 | Loan grade |
| LoanAmnt | int64 | Amount of the loan (USD) |
| LoanIntRate | float64 | Loan interest rate |
| LoanStatus | int64 | Loan status (0 - not defaulted, 1 - defaulted) |
| LoanPercentIncome | float64 | Loan percentage of income |
| PreviousDefault | object | If the borrower has defaulted before |
| CredHistory | int64 | Credit history length |

## 2.1 Preparing the Data

Table 2: Missing Values in Each Column

| Variable | Missing Values |
| --- | --- |
| PersonAge | 0 |
| PersonIncome | 0 |
| PersonHomeOwnership | 0 |
| PersonEmpLength | 887 |
| LoanIntent | 0 |
| LoanGrade | 0 |
| LoanAmnt | 0 |
| LoanIntRate | 3095 |
| LoanStatus | 0 |
| LoanPercentIncome | 0 |
| PreviousDefault | 0 |
| CredHistory | 0 |

Talk about how missing values were handled.

## 2.2 Descriptive Statistics

Table 3: Summary Statistics of Numeric Variables

| Variable | N | Mean | Median | SD | Min | Max |
|---|---|---|---|---|---|---|
| PersonAge | 32415.0 | 27.7 | 26.0 | 6.3 | 20.0 | 144.0 |
| PersonIncome | 32415.0 | 65908.6 | 55000.0 | 52533.0 | 4000.0 | 2039784.0 |
| PersonEmpLength | 32415.0 | 4.8 | 4.0 | 4.1 | 0.0 | 123.0 |
| LoanGrade | 32415.0 | 1.2 | 1.0 | 1.2 | 0.0 | 6.0 |
| LoanAmnt | 32415.0 | 9594.0 | 8000.0 | 6322.8 | 500.0 | 35000.0 |
| LoanIntRate | 32415.0 | 11.0 | 11.0 | 3.2 | 5.4 | 23.4 |
| LoanStatus | 32415.0 | 0.2 | 0.0 | 0.4 | 0.0 | 1.0 |
| LoanPercentIncome | 32415.0 | 0.2 | 0.2 | 0.1 | 0.0 | 0.8 |
| CredHistory | 32415.0 | 5.8 | 4.0 | 4.1 | 2.0 | 30.0 |

## 2.3 Distribution Analysis



Figure ?: Histograms of all Numeric Variables

4

Figure ?: Box Plots of All Variables Before Normalisation



Figure ?: Box Plots of All Variables After Normalisation

Figure ?: Distribution of Default Before and After Downsampling

Downsampled the dataset to ensure that the models didn't get affected by the magnitude of the majority class = can affect performance metrics. Allows a higher recall score
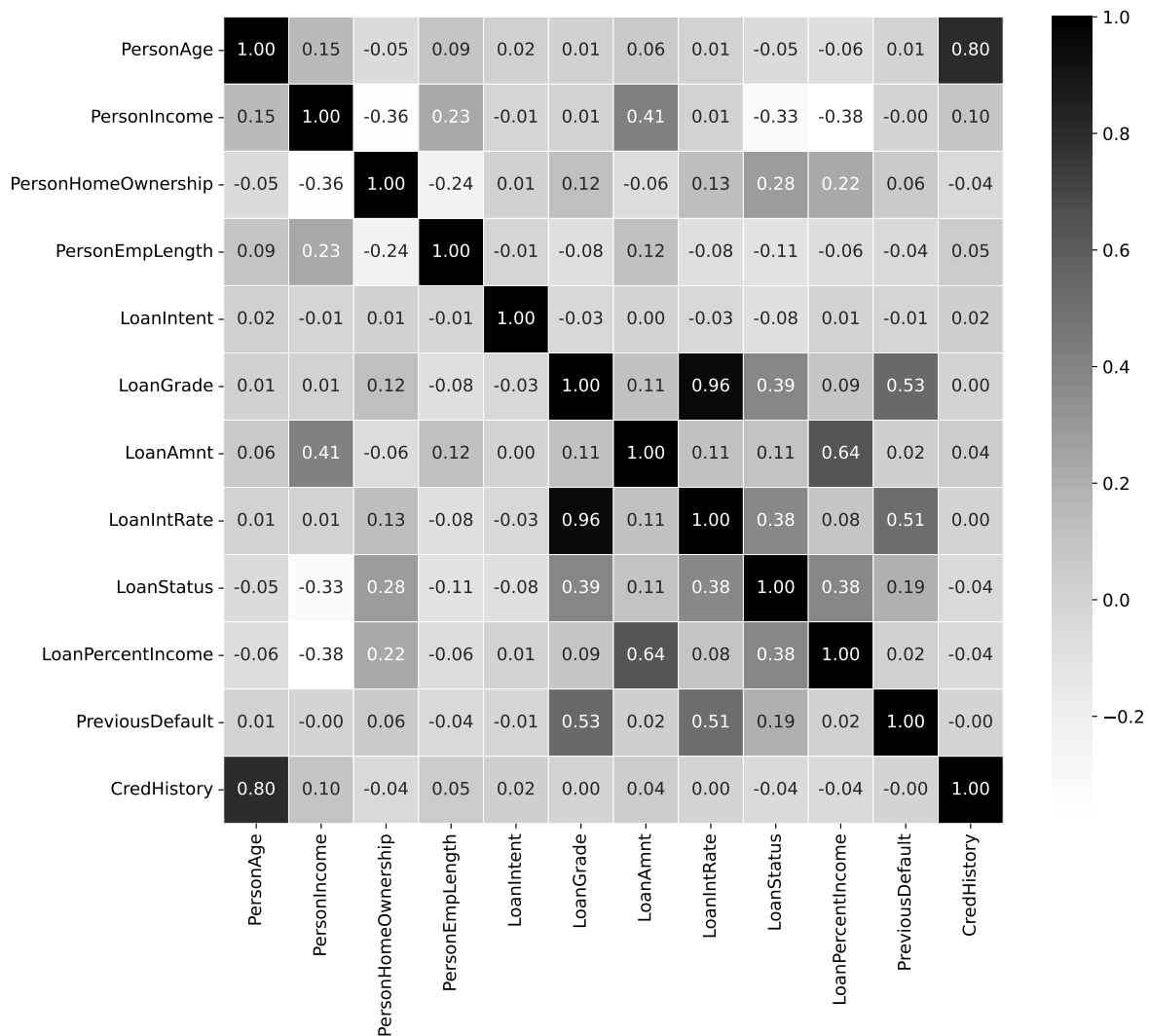
## 2.4 Correlation Analysis



Figure ?: Correlation Plot of All Variables

Table 4: Variance Inflation Factor (VIF) Values

| Feature | VIF |
|---|---|
| PersonAge | 1.493000 |
| PersonIncome | 9.698000 |
| PersonHomeOwnership | 1.200000 |
| PersonEmpLength | 1.059000 |
| LoanIntent | 1.001000 |
| LoanGrade | 2.990000 |
| LoanAmnt | 12.973000 |
| LoanIntRate | 3.117000 |
| LoanPercentIncome | 12.429000 |
| PreviousDefault | 1.254000 |
| CredHistory | 1.461000 |

Due to high multicollinearity between some variables, when using logistic regresion ridge and lasso regression are implemented to reduce effects of multicollinearity. Other models handle multicollinearity

# 3. Results and Discussion

## 3.1 Logistic Regression



Figure ?: ROC Curve for Logistic Regression Model

Figure ?: Confusion Matrix for Logistic Regression Model



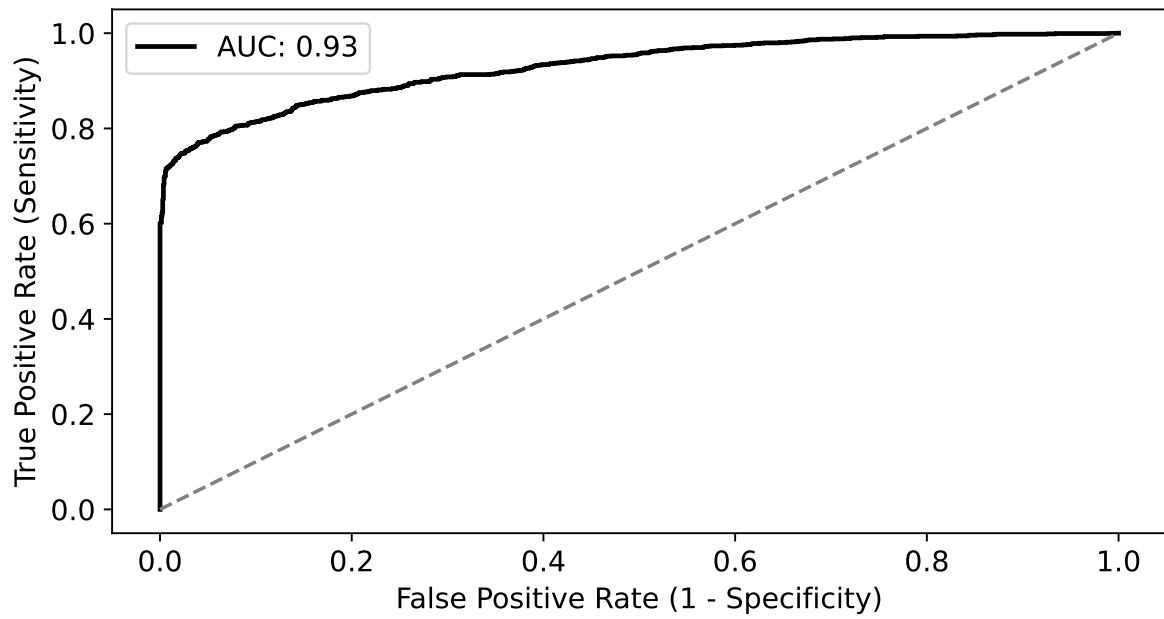Figure ?: Odds Ratios for Logistic Regression Model

## 3.2 Random Forest
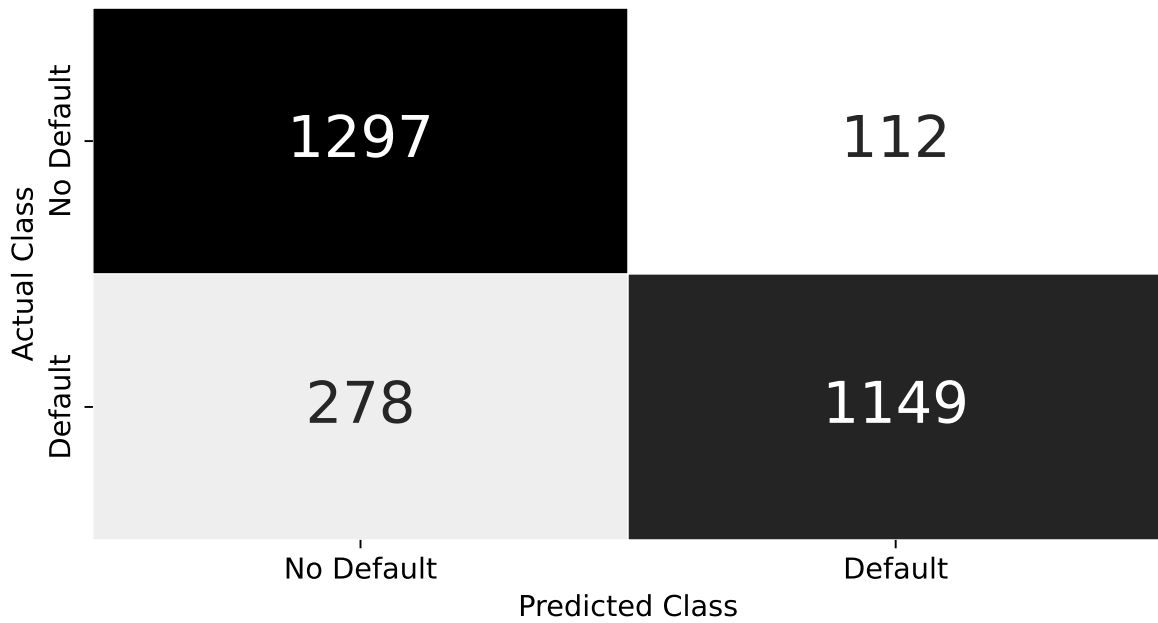


Figure ?: ROC Curve for Random Forest Model


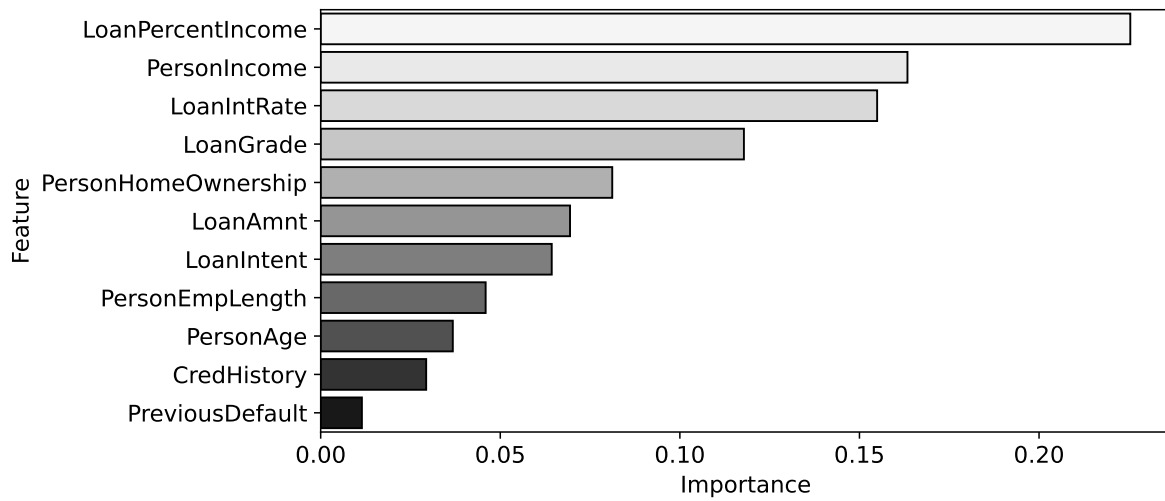
Figure ?: Confusion Matrix for Random Forest Model

Figure ?: Feature Importances from Random Forest Model
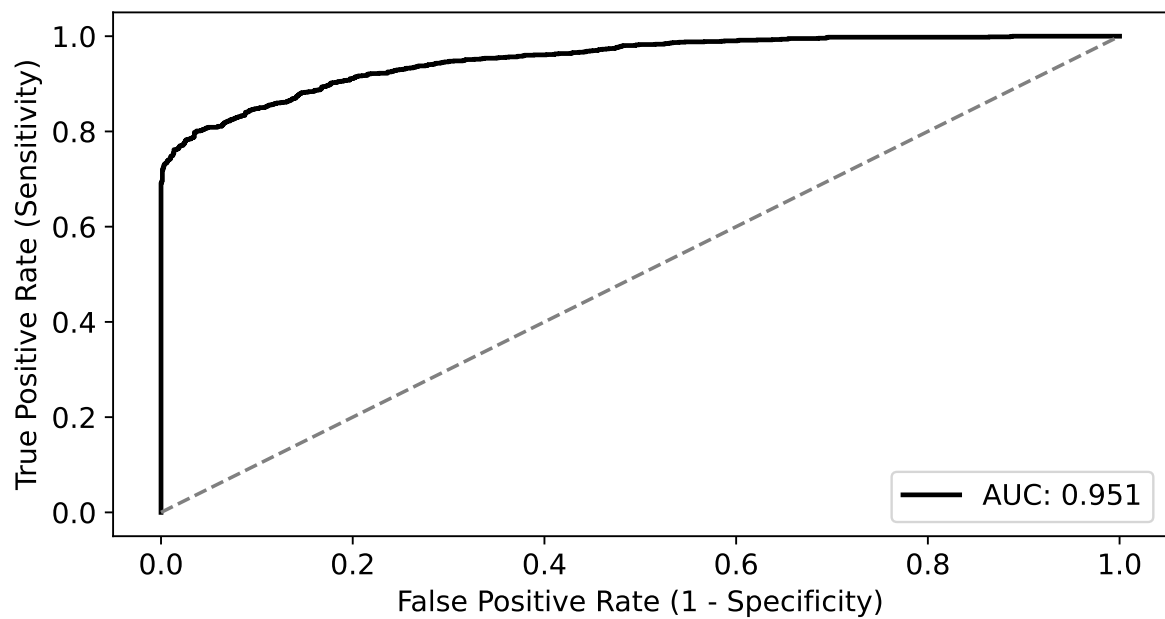
## 3.3 XGBoost
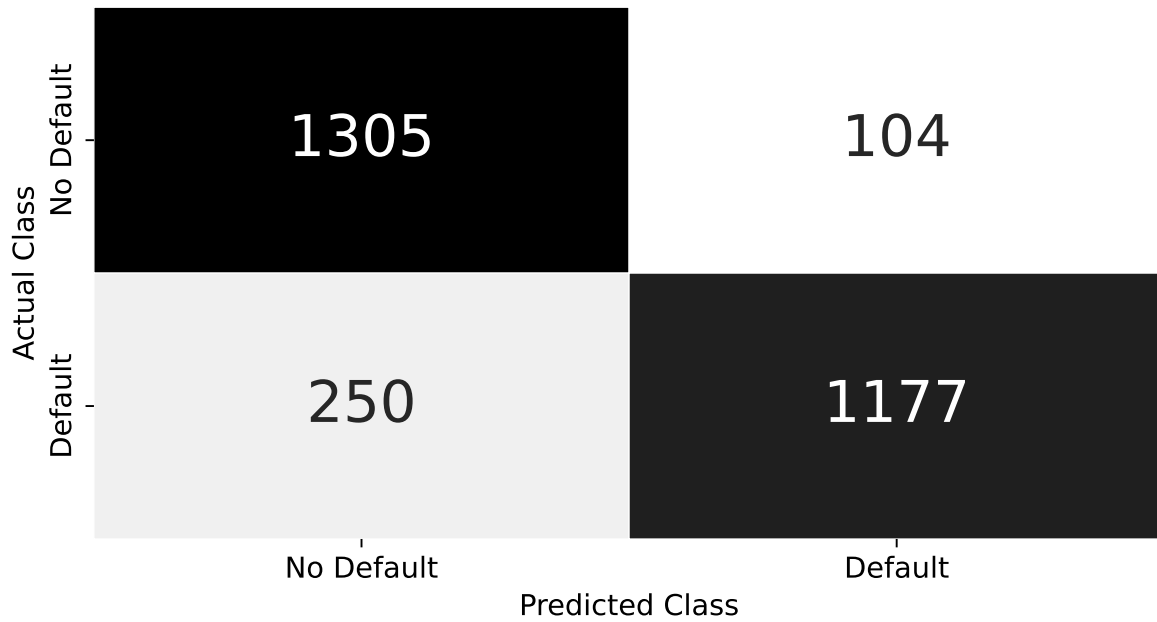


Figure ?: ROC Curve for XGBoost Model

Figure ?: Confusion Matrix for XGBoost Model
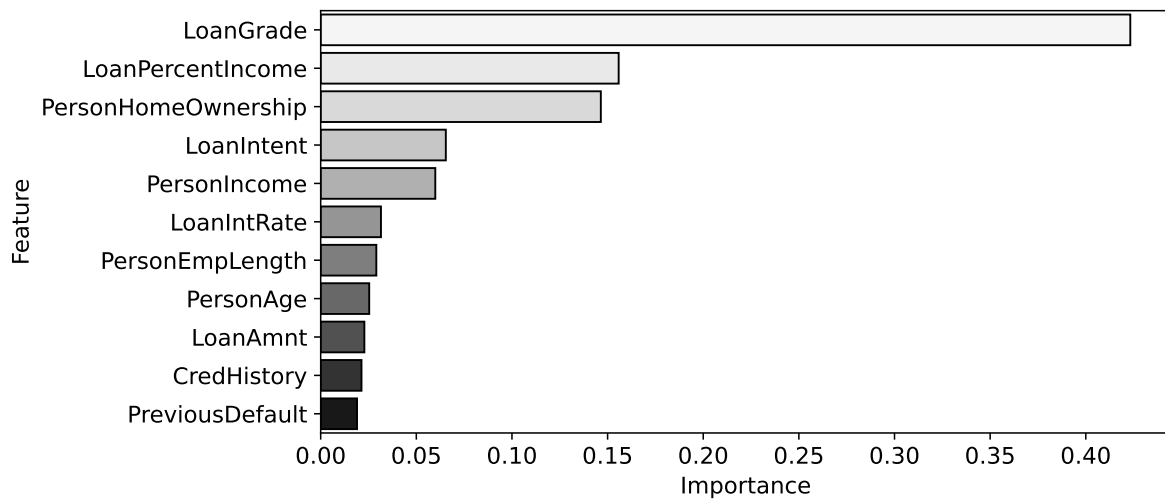


Figure ?: Feature Importances from XGBoost Model
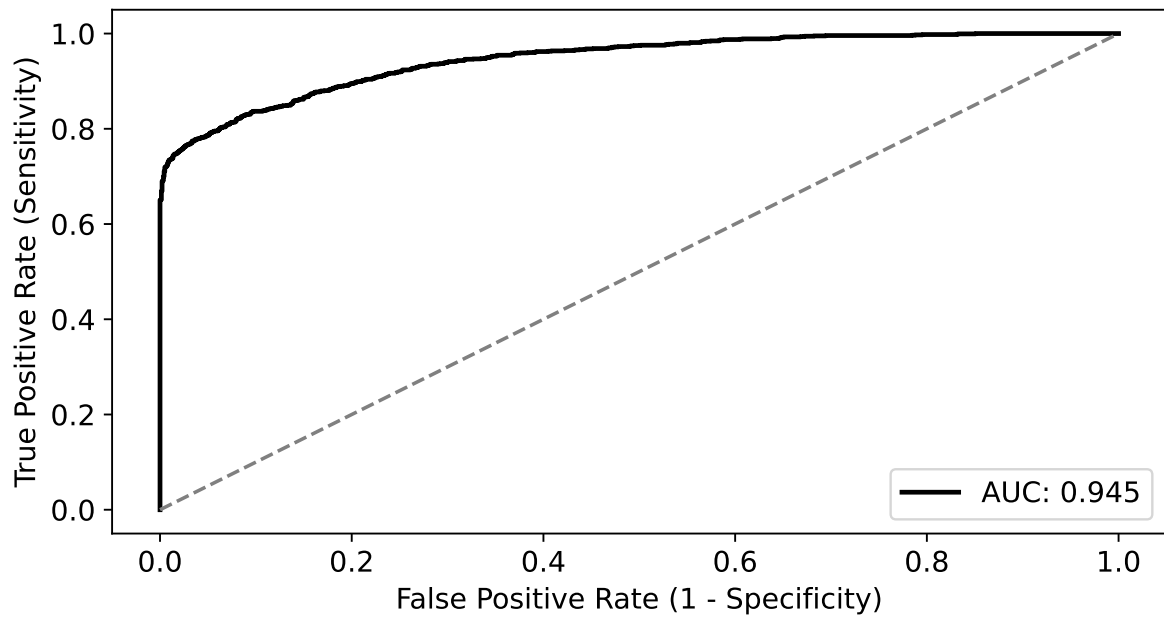
## 3.4 Light Gradient Boosted Machine
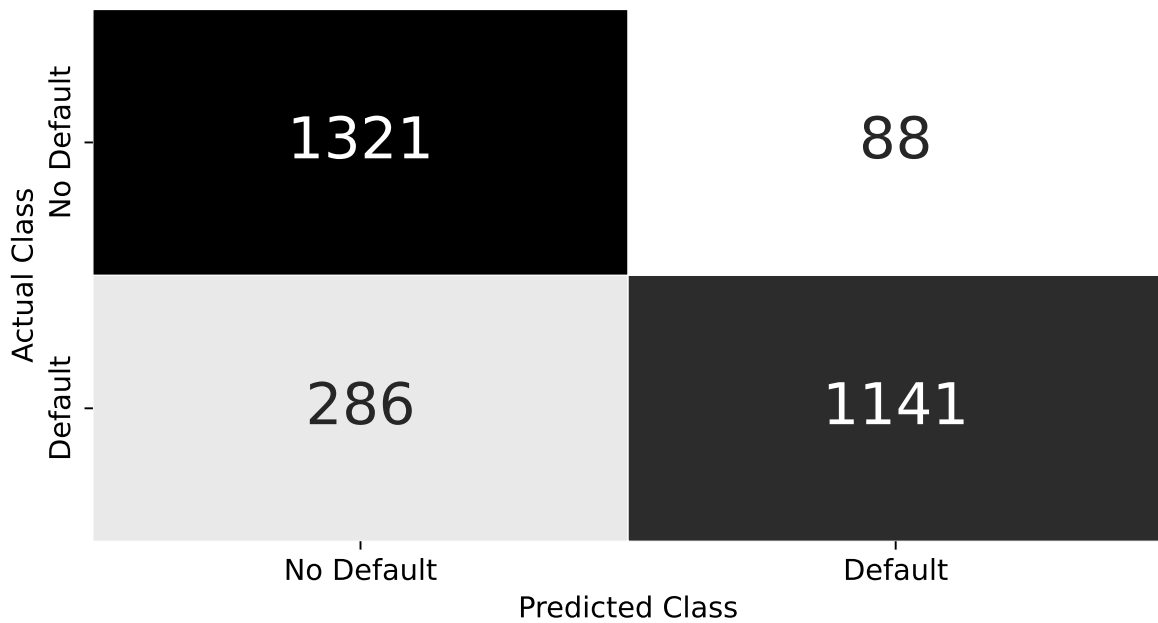


Figure ?: ROC Curve for LightGBM Model



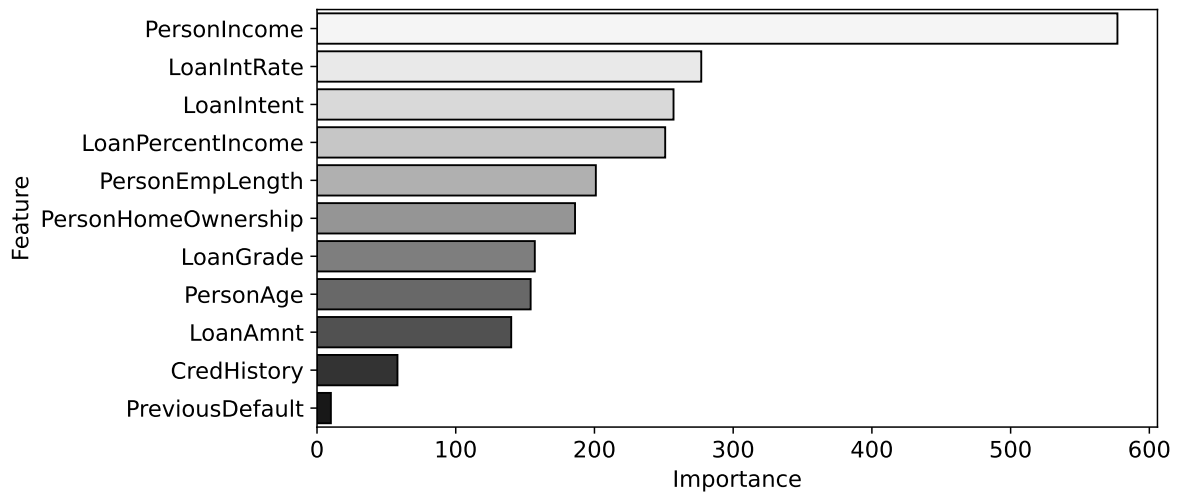Figure ?: Confusion Matrix for LightGBM Model

Figure ?: Feature Importances from LightGBM Model
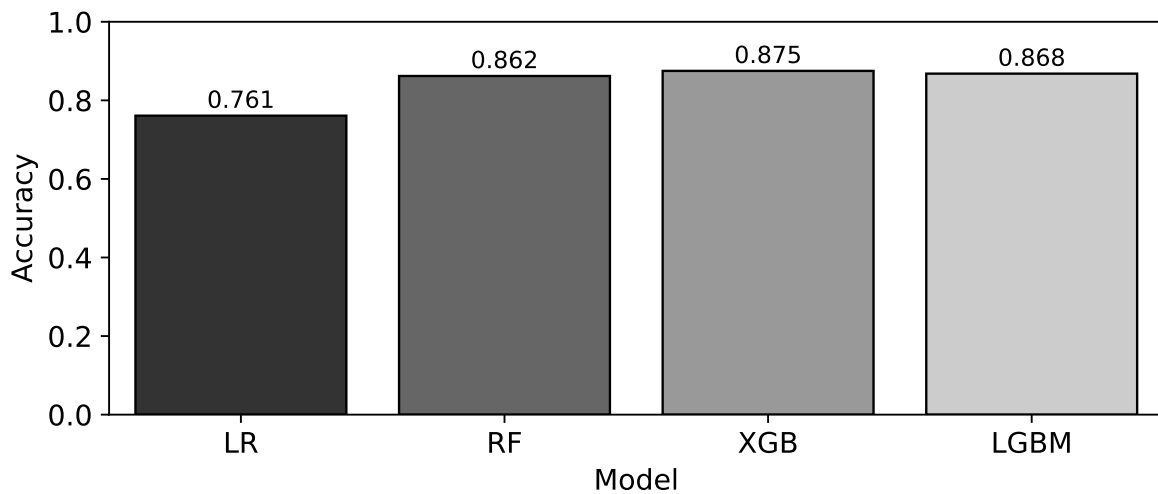
## 3.5 Model Evaluation and Comparisons



Figure ?: Accuracy for Each Model

Table 5: Performance Metrics for Each Model

| Model | Accuracy | Precision | Recall | F1 Score | AUC | Log Loss | Brier Score |
|-------|----------|-----------|--------|----------|------|----------|-------------|
| LR | 0.761 | 0.756 | 0.776 | 0.766 | 0.839 | 8.604 | 0.239 |
| RF | 0.862 | 0.911 | 0.805 | 0.855 | 0.93 | 4.957 | 0.138 |
| XGB | 0.875 | 0.919 | 0.825 | 0.869 | 0.951 | 4.499 | 0.125 |
| LGBM | 0.868 | 0.928 | 0.8 | 0.859 | 0.945 | 4.753 | 0.132 |

## 4. Conclusion

Link to Github Repository = https://github.com/JoshLG18/DSE-EMP-Project