❖ **Presentation on: Lending Club Data (Loan Predictions)**

**An Overview**

**Prepared by – Shubham Chaudhari**

**Amol Shinde**

**Imarticus Learning Institute, Pune**

# ❖ Table of Content

✓ **Lending Club is a US peer-to-peer lending company, headquartered in San Francisco, California. It was the first peer-to-peer lender to register its offerings as securities with the Securities and Exchange Commission (SEC), and to offer loan trading on a secondary market. Lending Club is the world's largest peer-to-peer lending platform.**

✓ **Understanding the Dataset**

These files contain complete loan data for all loans issued through the June 2007- Dec 2015,

including the current loan status and latest payment information. The file containing loan data

through the "present" contains complete loan data for all loans issued through the previous

completed calendar quarter. Additional features include credit scores, number of finance

inquiries, address including zip codes, and state, and collections among others. The file is a

matrix of about 855969 observations and 73 variables

✓ **Checking number of Columns and Rows**

```
print(loan_df)
(855969, 73)
```

✓ **Describe the Dataset**

| | Count | Percent |
|---|---|---|
| emp_title | 51462 | 5.799326 |
| emp_length | 44825 | 5.051393 |
| annual_inc | 4 | 0.000451 |
| desc | 761351 | 85.797726 |
| title | 152 | 0.017129 |
| delinq_2yrs | 29 | 0.003268 |
| mths_since_last_record | 750326 | 84.555303 |
| open_acc | 29 | 0.003268 |

✓ **Concluding X-Variables and Y-Variables**

✓ **Checking NA's, NULL's, ?, BLANKS**
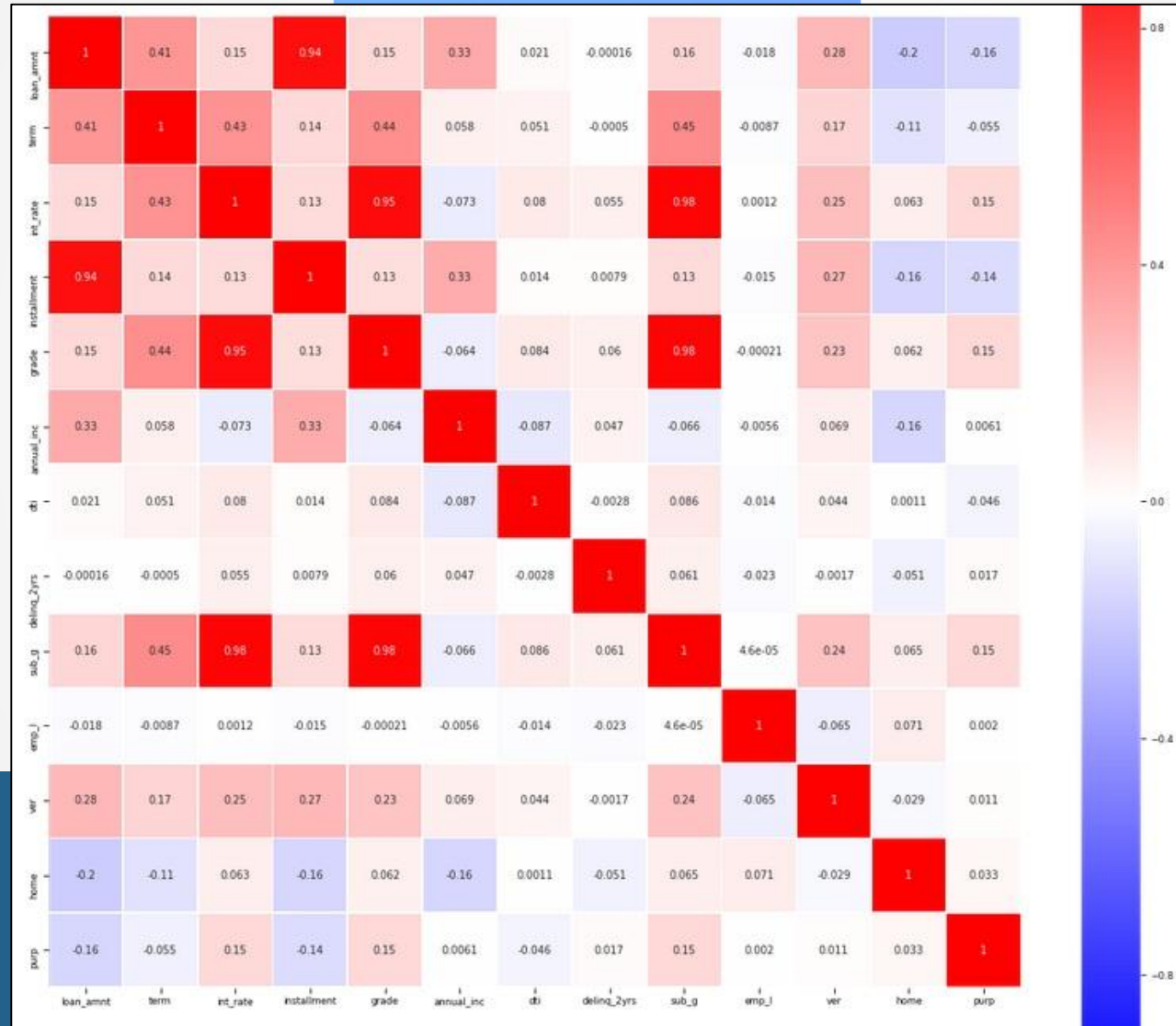
Checked for Nulls and Blanks in the dataset

✓ **Dealing with missing values**

Columns having missing values greater than 80% have been dropped.

✓ **Dropping Unnecessary Columns**

## 1. Correlation between the features

✓ **Finding the correlation between the variables**

It can be seen from the plot above that loan amount and installment have a very high correlation amongst each other (0.94). This is intuitive since a person who takes a large sum of loan would require extra time to repay it back. Also, interest rate, sub grade and grade have a very high correlation between them. This is obvious since interest rate is decided by grades once the grades are decided, a subgrade is assigned to that loan (leading to high correlation).

✓ **Selecting important variables**

```
df_LC = df1.filter(['loan_amnt','term','int_rate','installment','grade','sub_grade','emp_length','home_
ownership',
                    'annual_inc','verification_status','purpose','dti','delinq_2yrs','loan_status'])
df_LC.dtypes
```

## ✓ Transformation

> **Before training the data, we would first transform the data to account for any skewness in the variable distribution.** *(Float)*

### Splitting Dataset into Training and Testing

> **We split the dataset based on 'issue_d'**
> **Train set - (June 2007 - May 2015)**
> **Test set – (June 2015 – Dec 2015)**

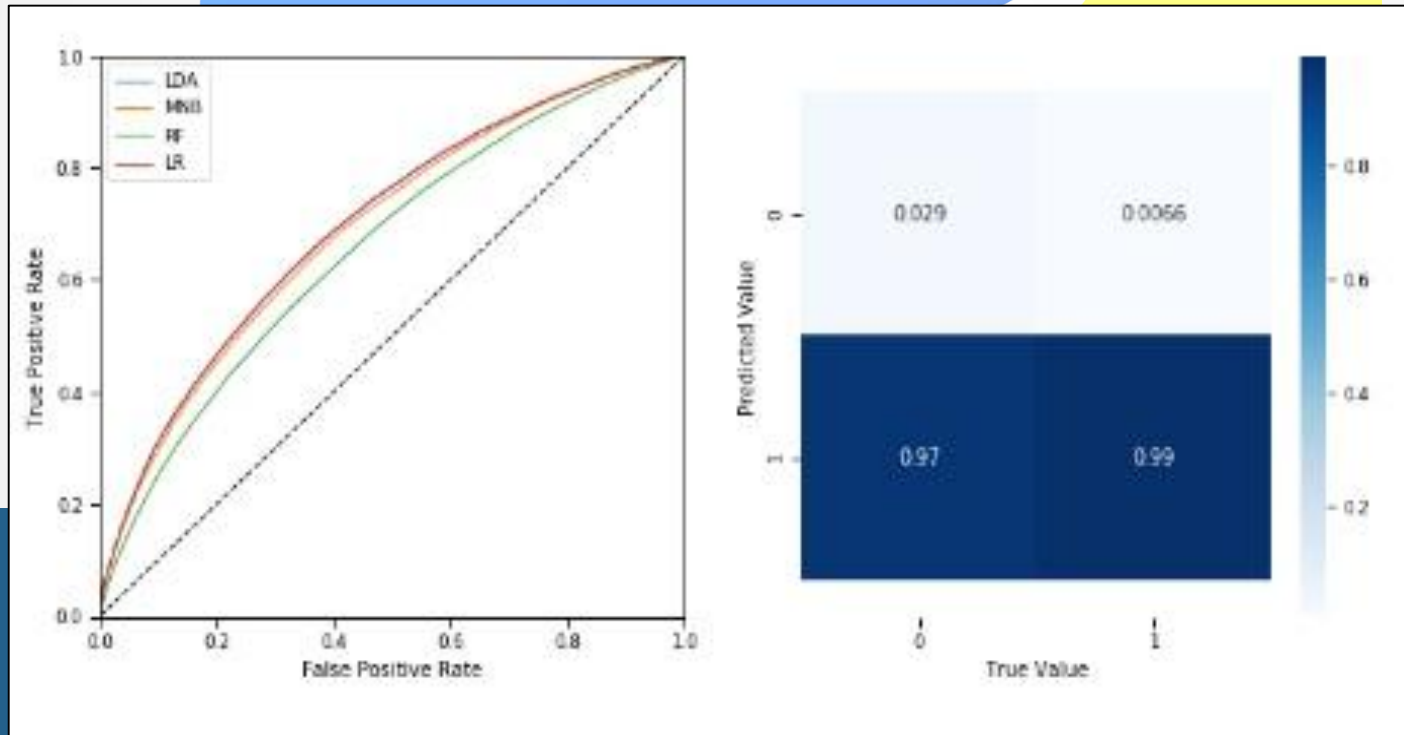## ✓ Model Selection

> **1. Logistic Regression**
> **2. Random Forest**

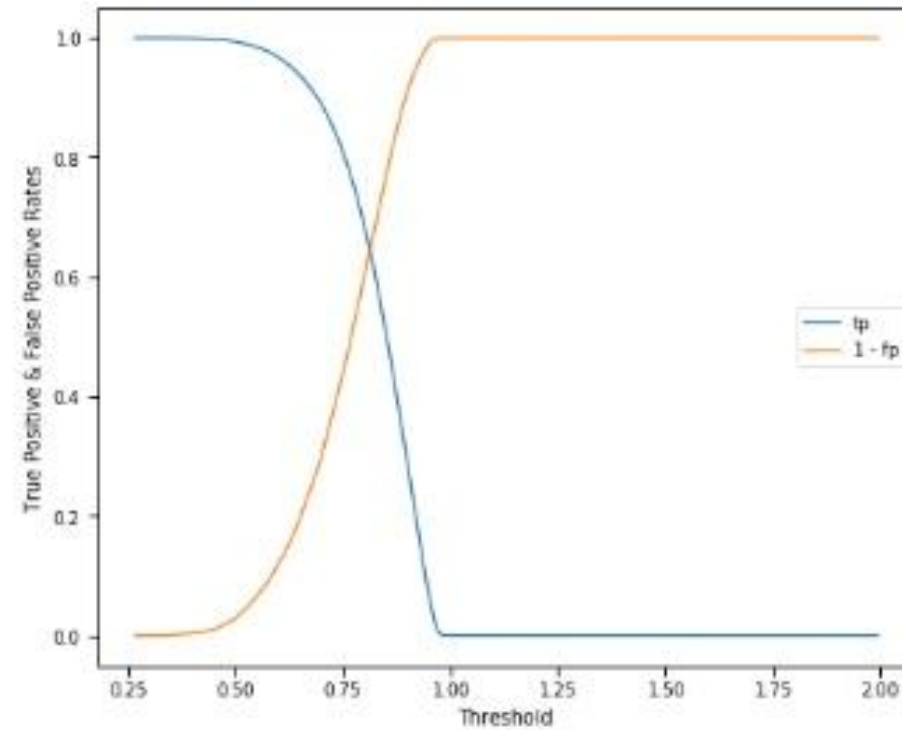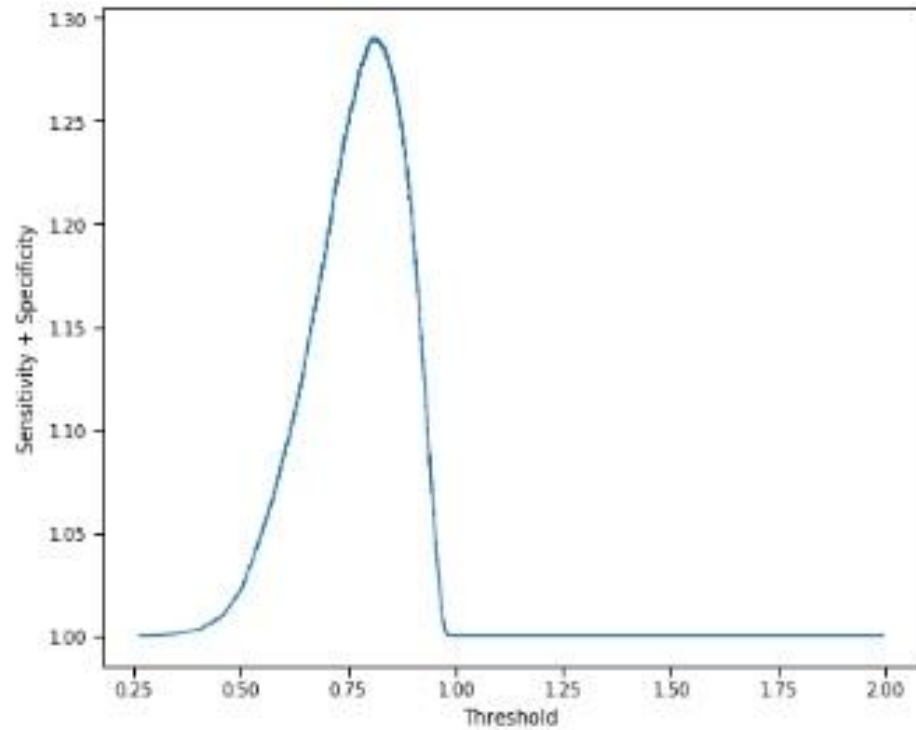## ✓ Cross- Validation (3-Folds)

# Results

- ✓ **Accuracy**
- ✓ **ROC Curve**
- ✓ **Confusion Matrix**

1. LR [0.92009416 0.92294362 0.92294362] 0.9219937927023743
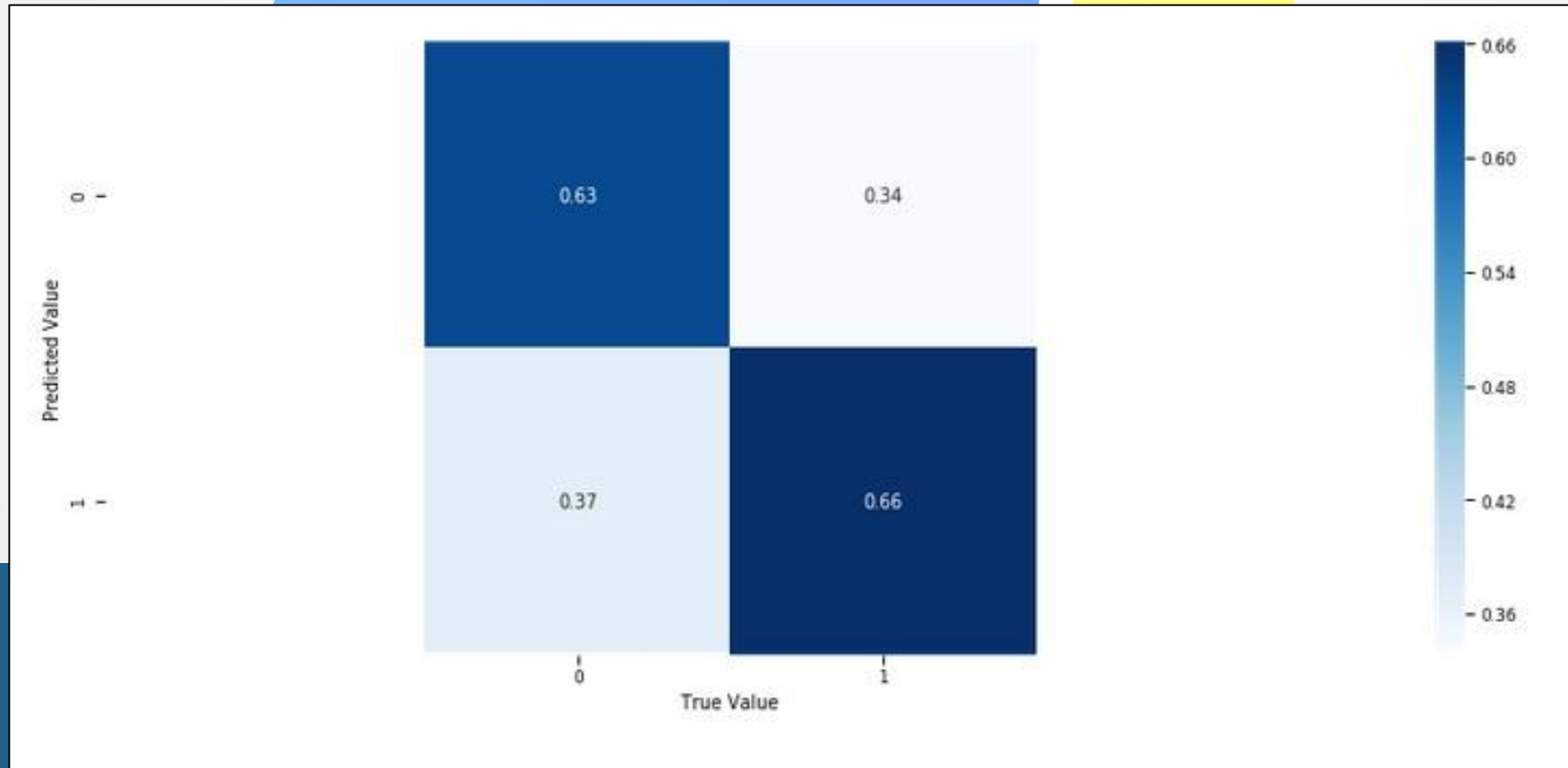2. RF [0.89259241 0.90865426 0.92266815] 0.9079716039306898
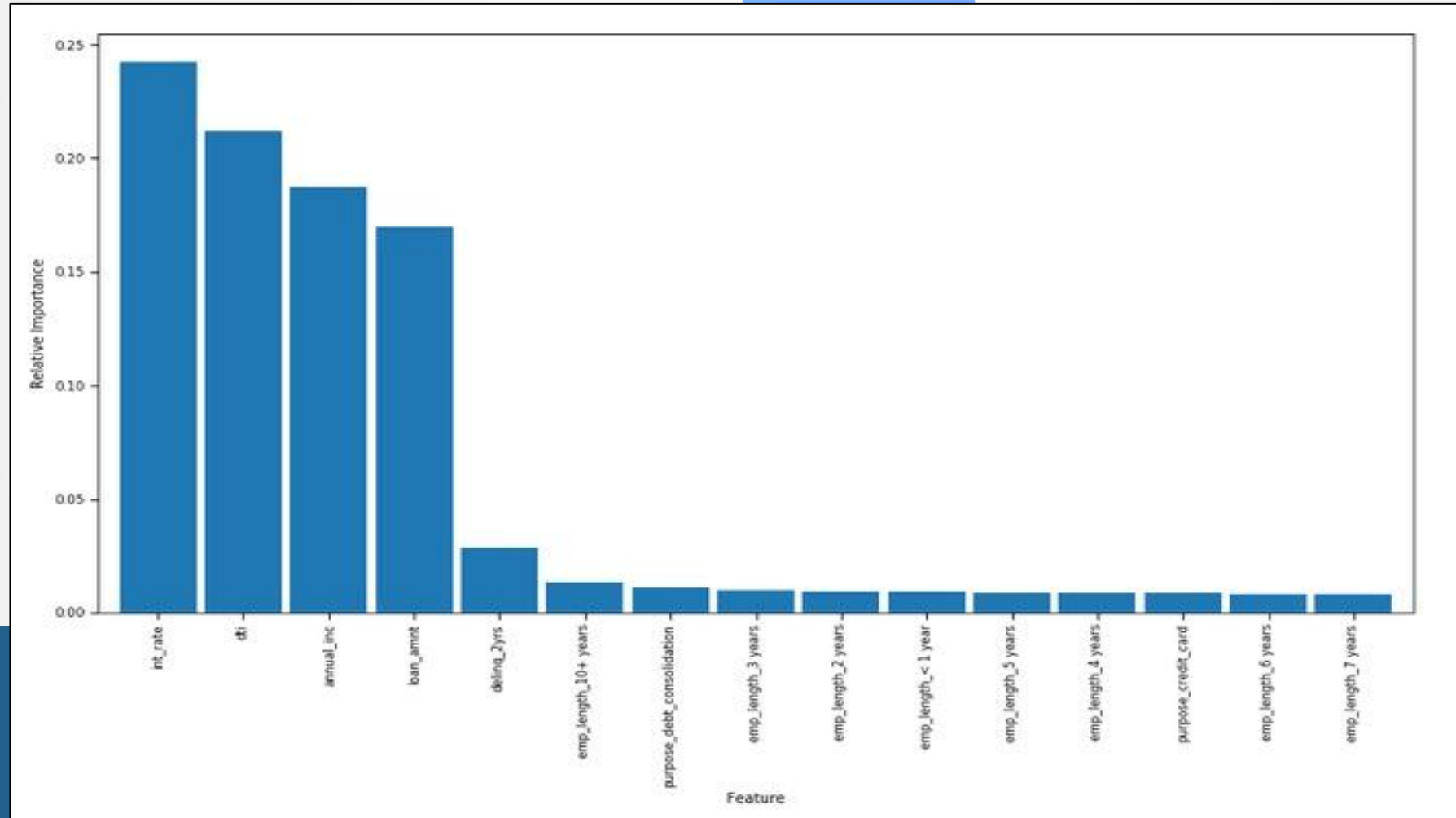
# Results

✓ **ROC Curve**

# Results

✓ **Optimal Threshold Values**

Optimal Threshold: 0.8079913653717011

✓ **Relative Important Variables**