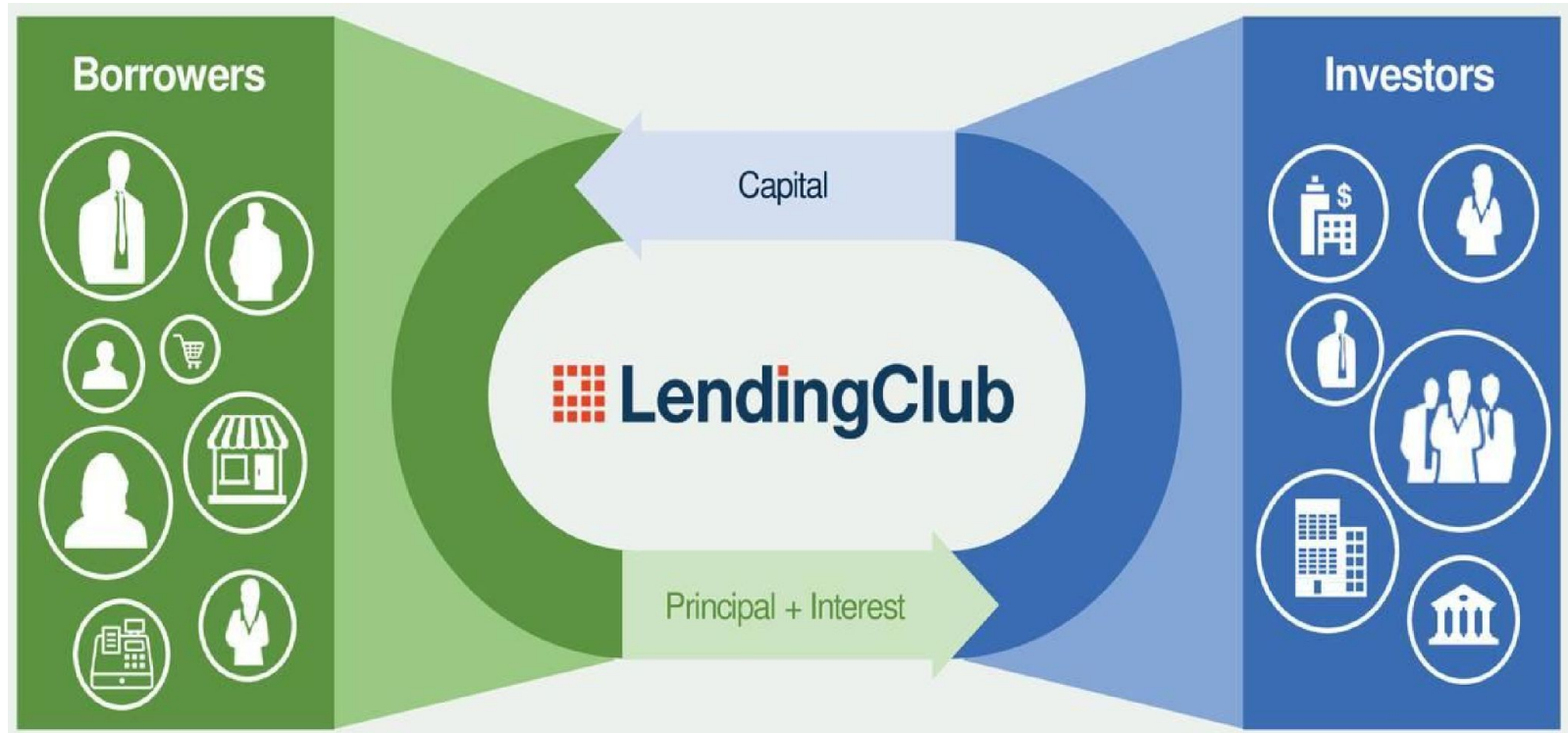




# Lending Club

## Loan Approval Optimization

# Lending Club Business Operation



# Problem

- **Charged Off Loans**
- **Default Loans**
- **Delayed Payments**



# Solution

- **Identify “Bad Loans”**



# Stakeholders



**Investors**



# Data Information

- Source: Lending Club Loan Data 2007-11 - dataset by jaypeedevlin
- Data pertaining to **2007-2011**
- **Loan Status** and relevant financial information
- Number of entries: **42,538**
- Number of features: **115**

# Features

- Loan Amount
- Term
- Installment
- Grade
- Employment Length
- Home Ownership
- Annual Income
- Verification Status
- Loan Status
- Purpose
- Loan Title
- Address State
- Debt/Income Ratio
- Delinquency
- Earliest Credit Line
- Inquiry
- # of Open Credit Lines
- Public Records
- Total Credit Revolving Balance
- ...

# Feature Engineering

## 1. Data Cleaning

- Removed empty columns/null values

## 2. Feature Selection

- Removed redundancy/data leakage
- Retained relevant/useful features

## 3. Conversion to Numerical Dtype

- Numerical: **revol\_util**
- Ordinal: **grade/emp\_length**

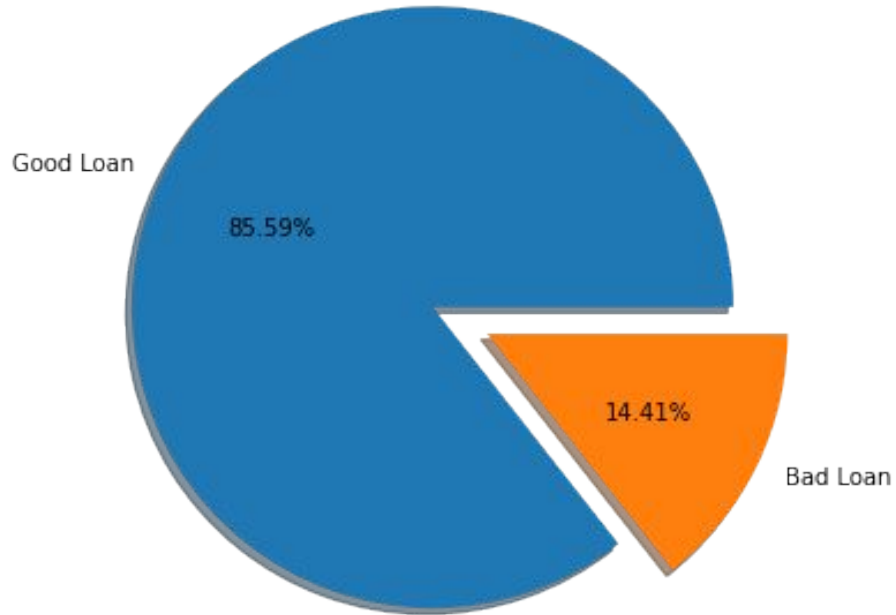
## 4. New Feature Added

- **fico\_range\_avg**

## 5. Target Feature Engineering

- Excluded loans in-progress
- New classification: **loan\_type**
- Binarization

Percentage of Each Loan Type



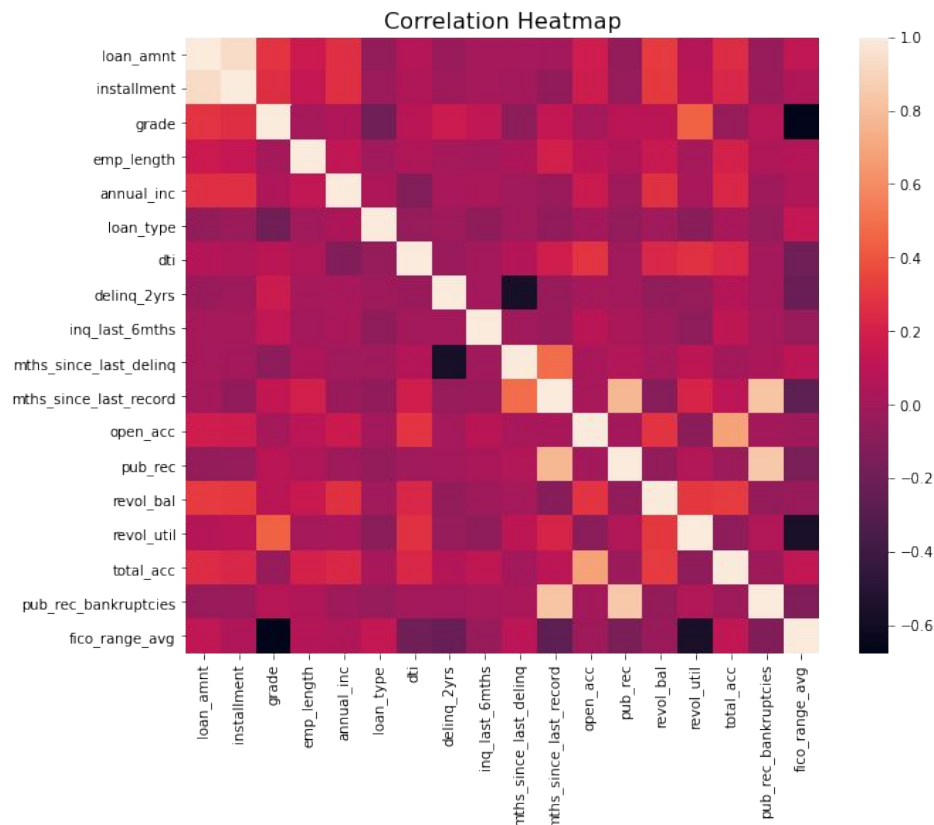
## Good Loan

- Fully Paid

## Bad Loan

- Charged Off
- Does Not Meet the Credit Policy



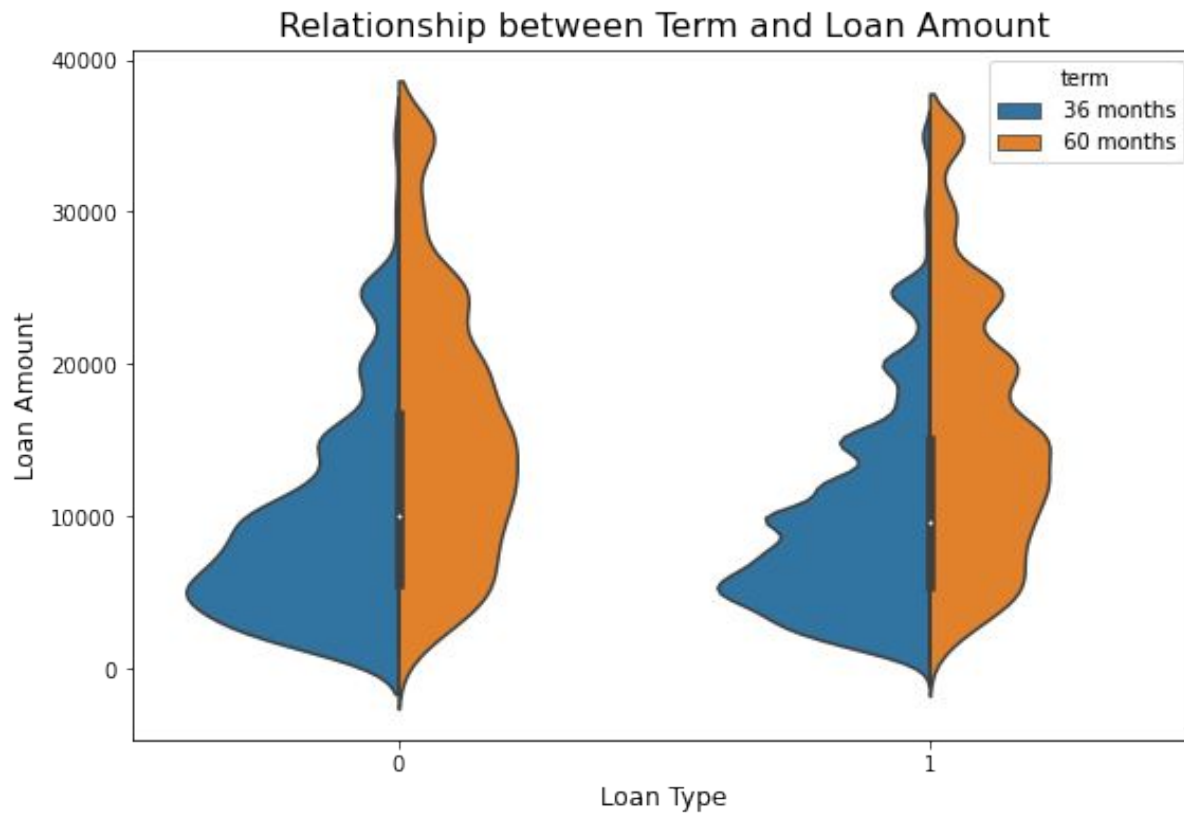


- Notable **loan\_type** correlations

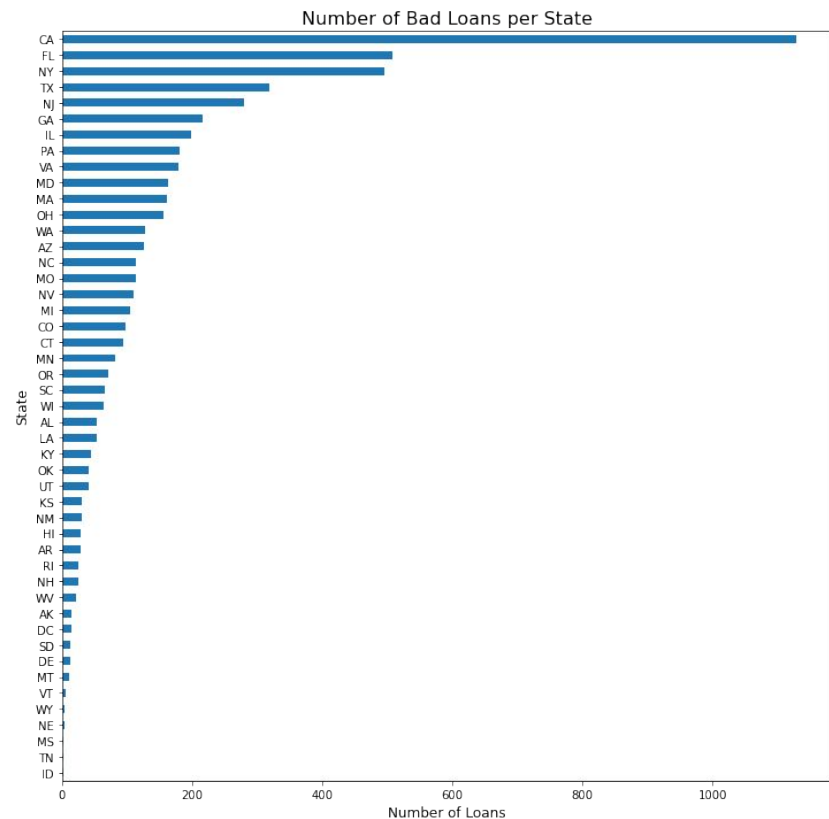
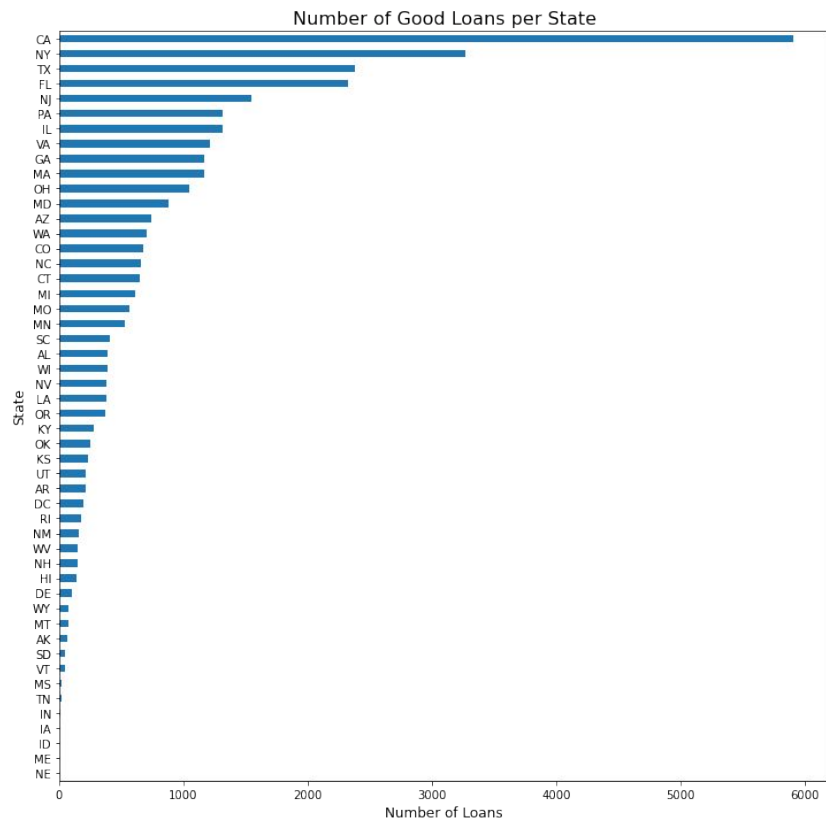
- **grade**
- **fico\_range\_avg**
- **revol\_util**
- **Inq\_last\_6mths**

- Intrinsically linked features

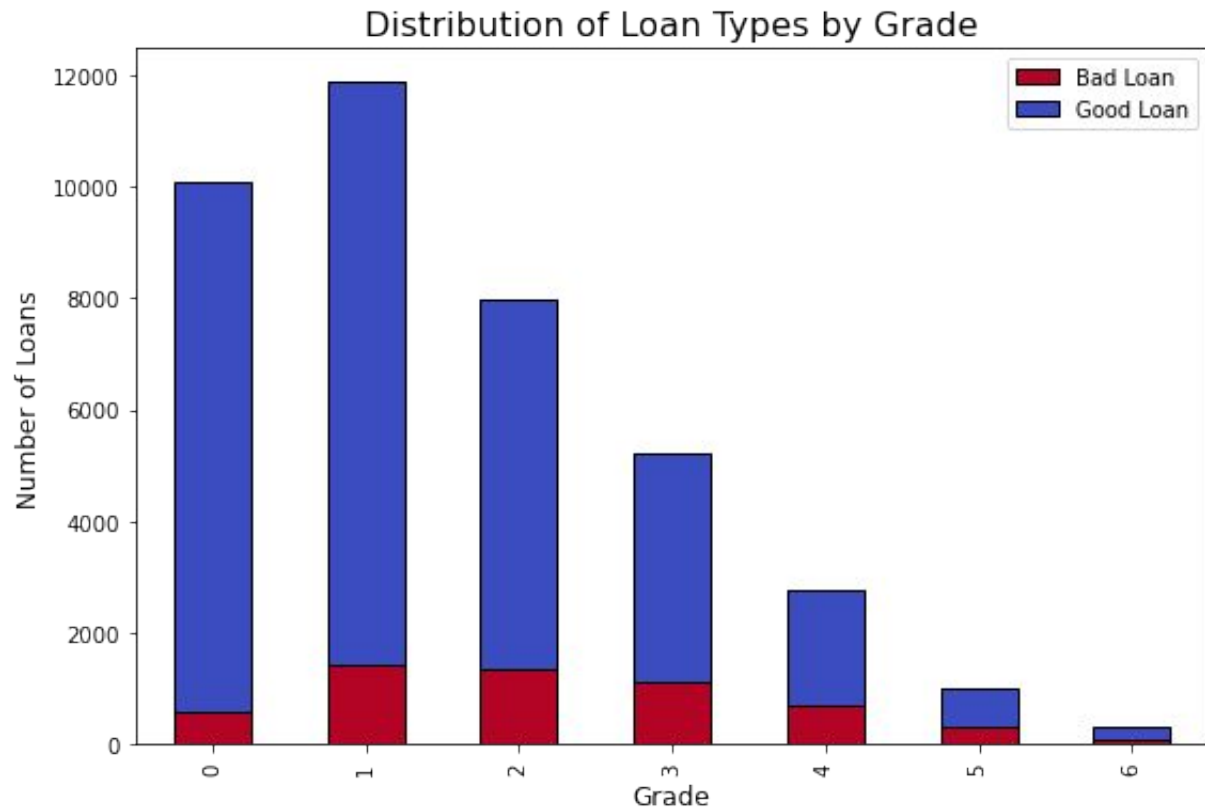
- **Installment/loan\_amnt**
- **pub\_rec/pub\_rec\_bankruptcies**
- **pub\_rec/mths\_since\_last\_record**
- **open\_acc/total\_acc**
- ...



Loan amount tends to increase with longer term regardless of the loan type



States with highest # of good loans  $\approx$  States with highest # of bad loans



Grade Description: 0 as worst, 6 as best

# Modeling Overview

- Supervised Learning
- Binary Classification
  - 1: Good Loan
  - 0: Bad Loan
- Highly Imbalanced Data
- Machine Learning Tools: Scikit-Learn, Imbalanced-Learn, XGBoost

# Modeling Procedure

## I. Data Preprocessing

1. One-Hot Encoding
2. Datetime Objects to Ordinal Numeric
3. Training Test Split (80%: 20%)
4. Feature Standardization
5. Minority Class Oversampling

## II. Randomized Search

- Stratified K-Fold with  $cv = 5$
- $n\_iter = 30$
- $scoring = "roc\_auc"$

## III. Training with Tuned Parameters

## IV. Performance Evaluation

- Evaluation Metric: F1, ROC AUC

# Techniques/Algorithms Used

## Resampling Technique

- SMOTE from Imbalanced-Learn

## Weighting Technique

- `class_weight = "balanced"` in Scikit-Learn classifiers

## Hyperparameter Tuning Technique

- Randomized Search from Scikit-Learn

## Classification Algorithms

- 1) Logistic Regression
- 2) Random Forest
- 3) Support Vector Machine
- 4) XGBoost

# Model Comparison

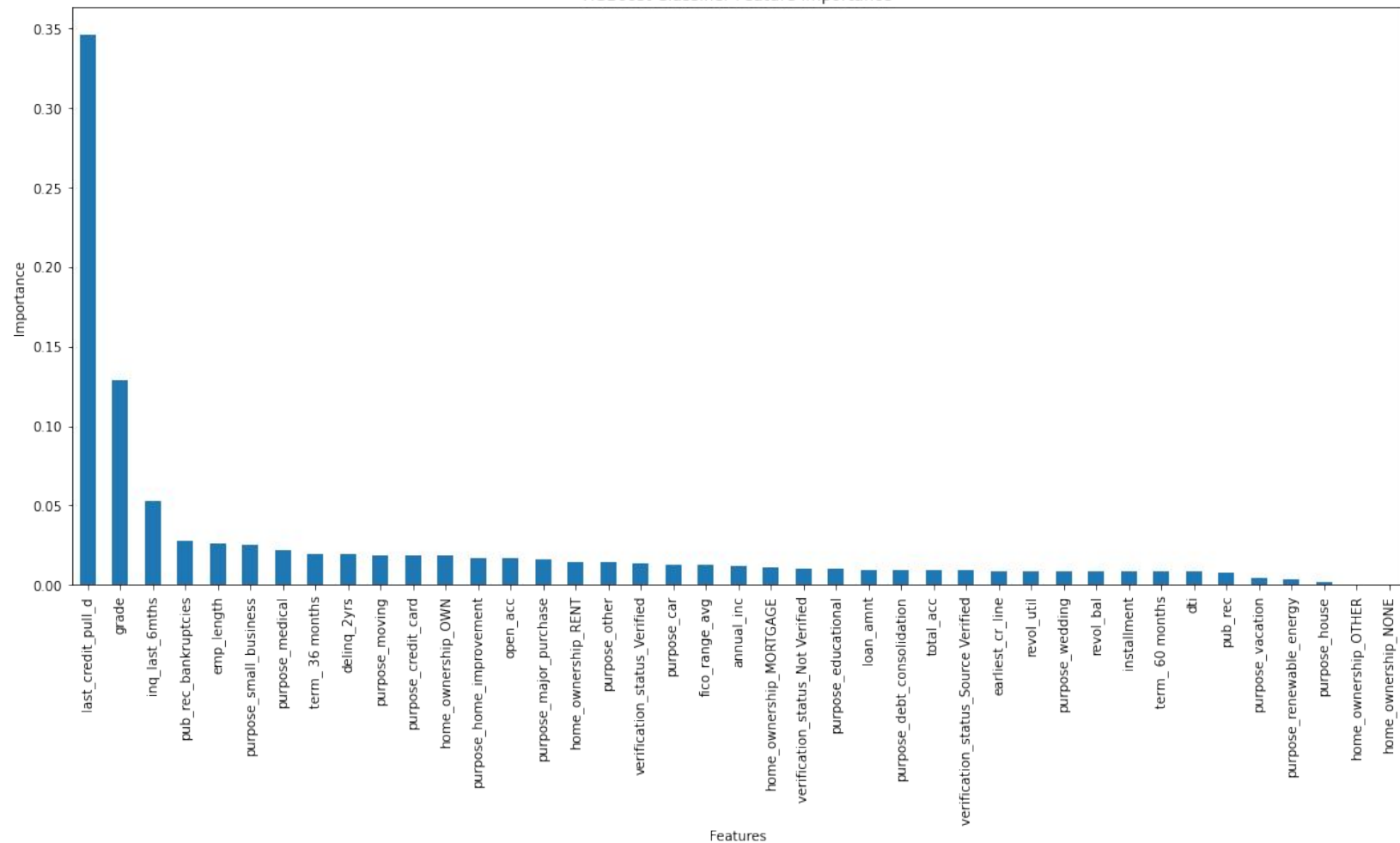
Model	ROC AUC	Minority F1	Majority F1
Logistic Regression	0.78	0.47	0.80
Random Forest	0.74	0.53	0.91
Support Vector Machine	0.67	0.41	0.87
XGBoost	0.80	0.56	0.89

Best: XGBoost

Worst: Logistic Regression



# XGBoost Classifier Feature Importance



# Assumptions/Limitations

- Primary Assumption
  - Validity of the data
- Limitations
  - Outdated (year 2007-2011)
  - Large volume of missing features
  - Limited entries (only 39239 entries)

# Conclusion

- Only 24 out of 115 features used
- Best model: XGBoost
- Primary features of importance: **last\_credit\_pull\_d, grade, inq\_last\_6mths**
- Prospective improvement
  - Up-to-date dataset
  - Hyperparameter tuning with Grid Search
  - Alternative classifier algorithms: Neural Network, Deep Learning, etc.