

# Modeling of Prediction Loan Default with Machine Learning

SEAS 6401

Instructor: Dr. Benjamin Harvey

Group Member:

Rui Zhang

Ran Wei

# CONTENTS

1. Introduction
2. Feature Selection
3. Exploratory Data Analysis
4. Data Preprocessing
5. Model
6. Result
7. Conclusion

# INTRODUCTION

## ❑ Project Goal:

- ❑ Understand the applicant's profile
- ❑ Make a prediction

## ❑ Data Set:

- ❑ Kaggle – Lending Club 2007-2020(Q3)
- ❑ 2.9million loan records and 142 features

# Feature Selection

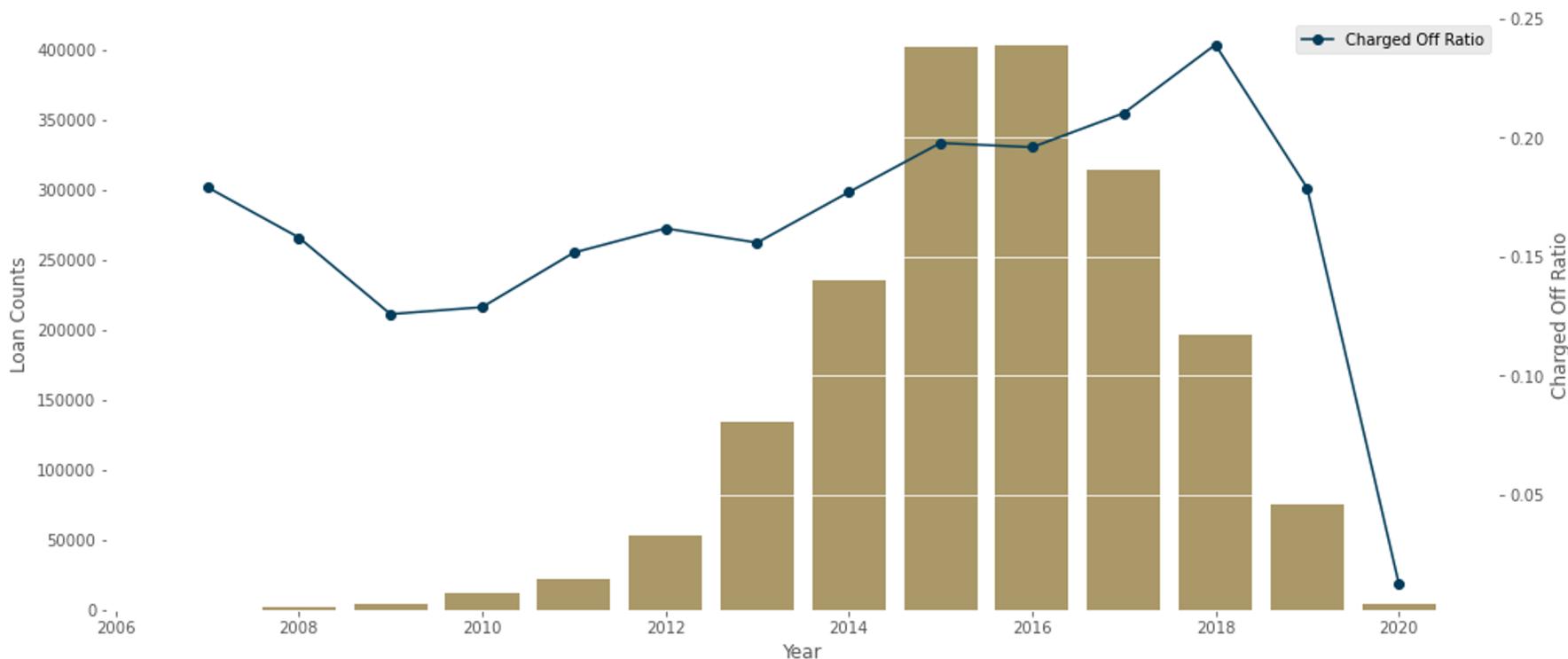
- Select 10000 rows
  - Set dependent variable
  - Filter loan status
  - Check correlation
  - Features would be available before funding the loan
- 

Dependent variable: 'loan\_status'  
Features:

- 1) 'loan\_amnt'
- 2) 'term'
- 3) 'int\_rate'
- 4) 'sub\_grade'
- 5) 'home\_ownership'
- 6) 'annual\_inc'
- 7) 'verification\_status'
- 8) 'purpose'
- 9) 'addr\_state'
- 10) 'dti' (highlighted)
- 11) 'open\_acc'
- 12) 'pub\_rec'
- 13) 'revol\_bal'
- 14) 'revol\_util' (highlighted)
- 15) 'initial\_list\_status'
- 16) 'application\_type'
- 17) 'mort\_acc' (highlighted)
- 18) 'pub\_rec\_bankruptcies'
- 19) 'loan\_status\_flag'
- 20) 'fico'
- 21) 'earliest\_cr\_line\_y'

# EXPLORATORY DATA ANALYSIS

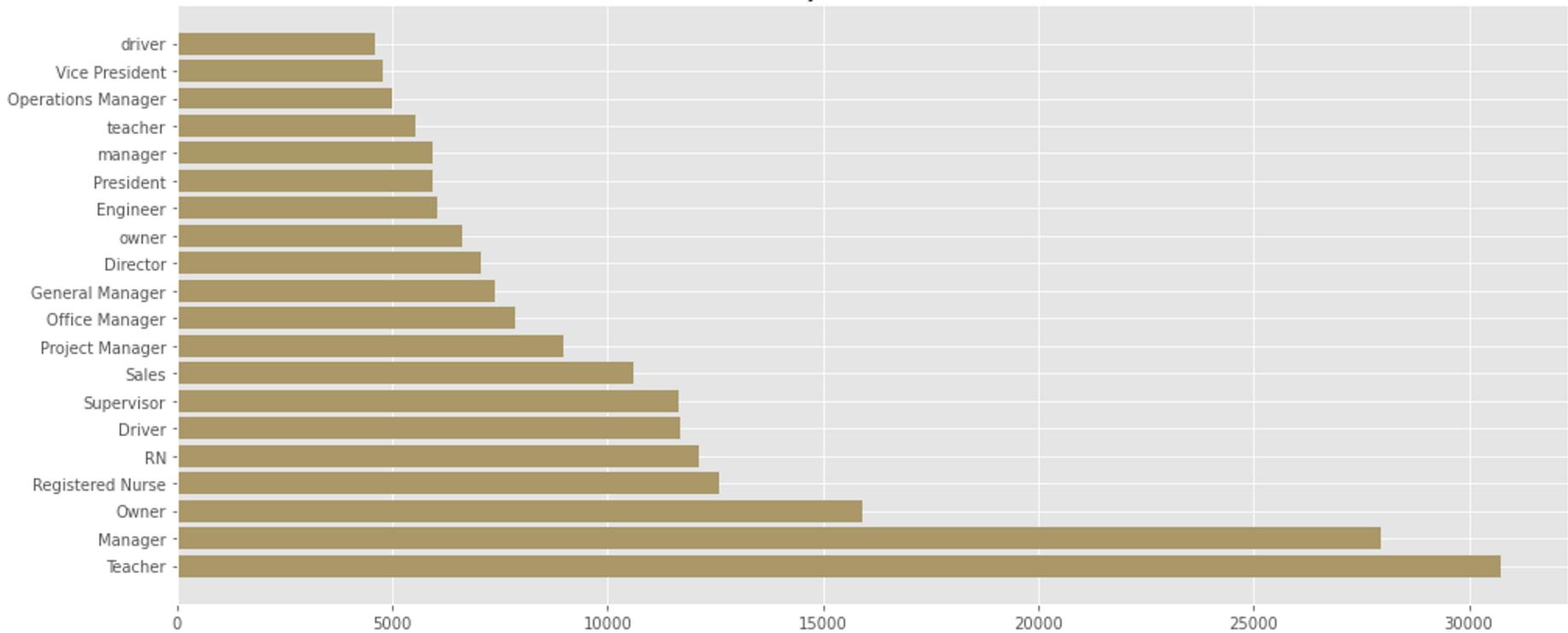
## □ Dependent Variable – Loan Status



# EXPLORATORY DATA ANALYSIS

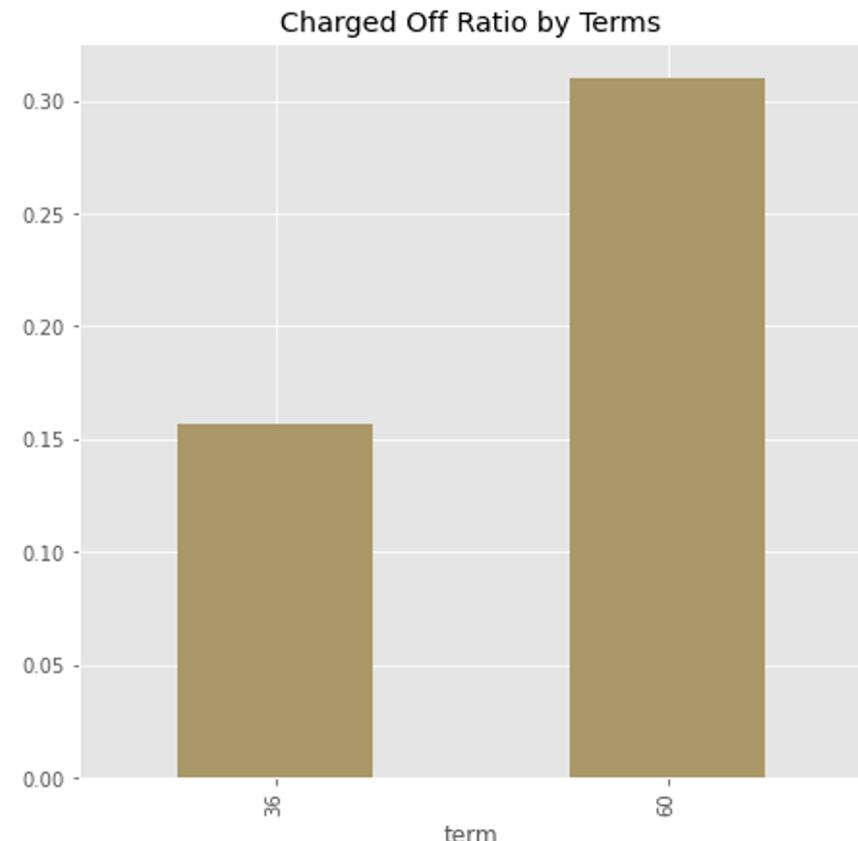
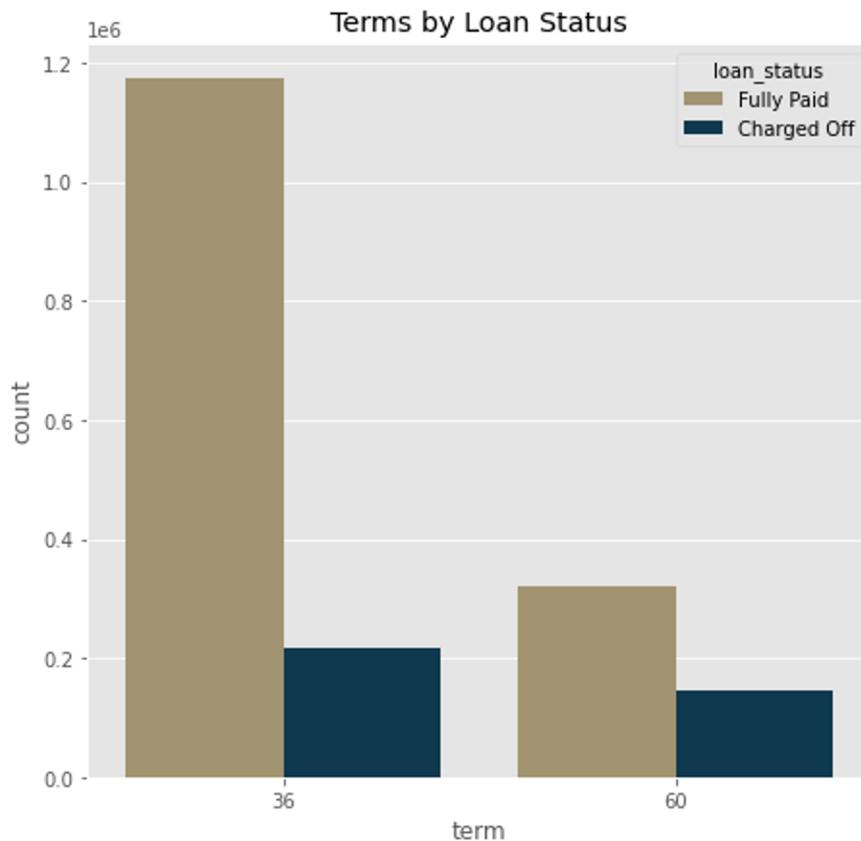
## □ Predictors – Job Title

The most 20 jobs title afforded a loan



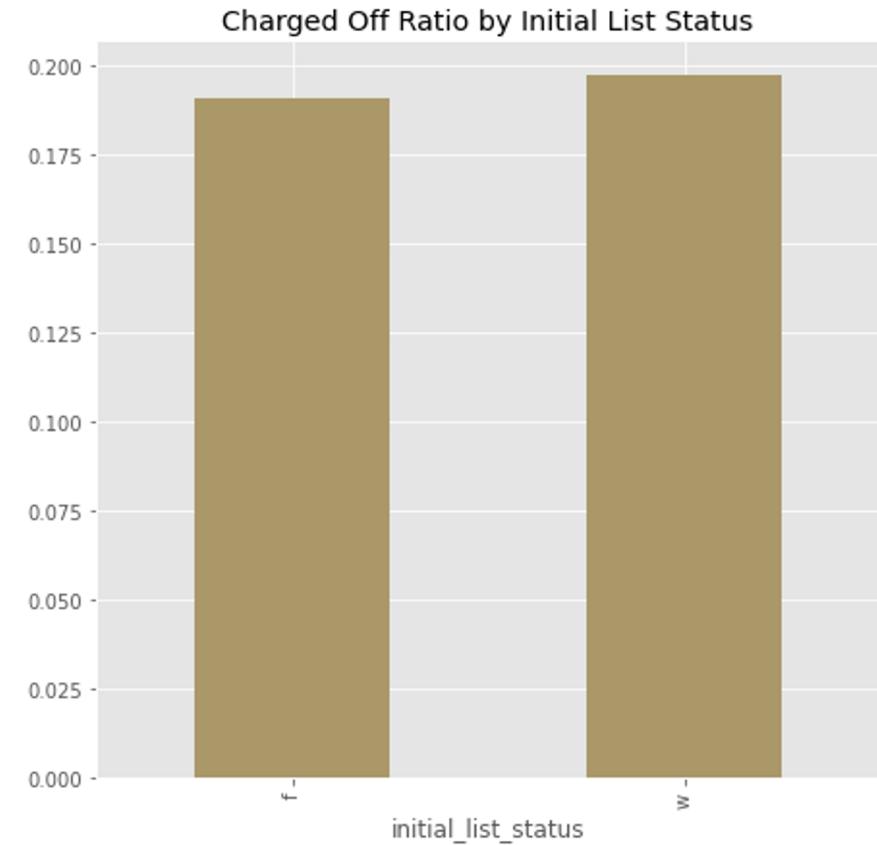
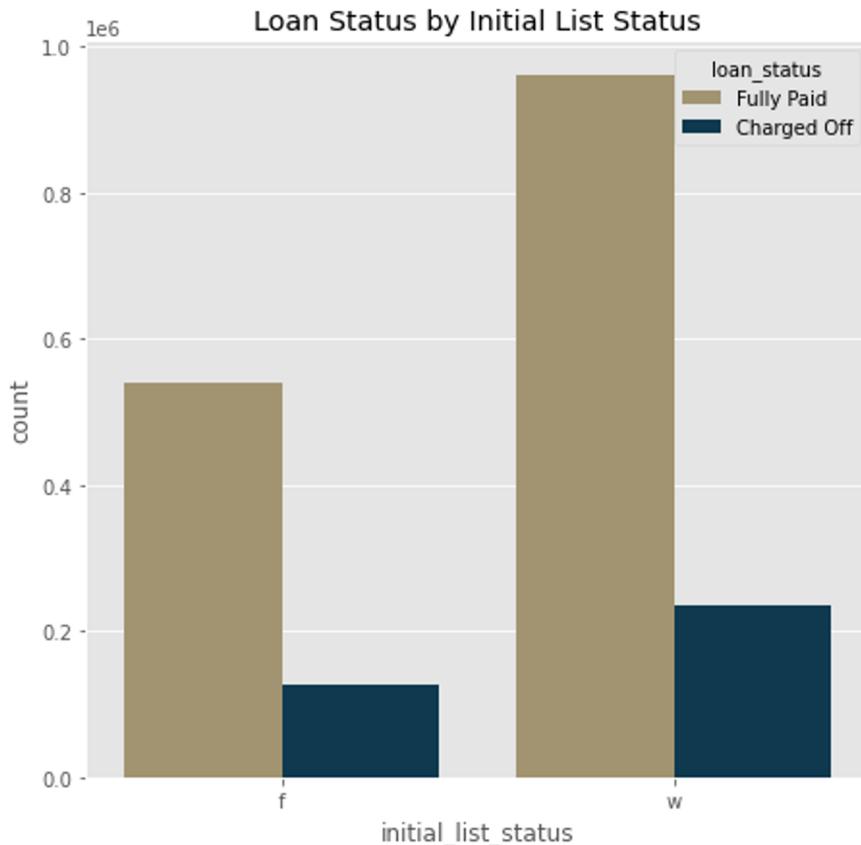
# EXPLORATORY DATA ANALYSIS

## □ Predictors – Loan Term



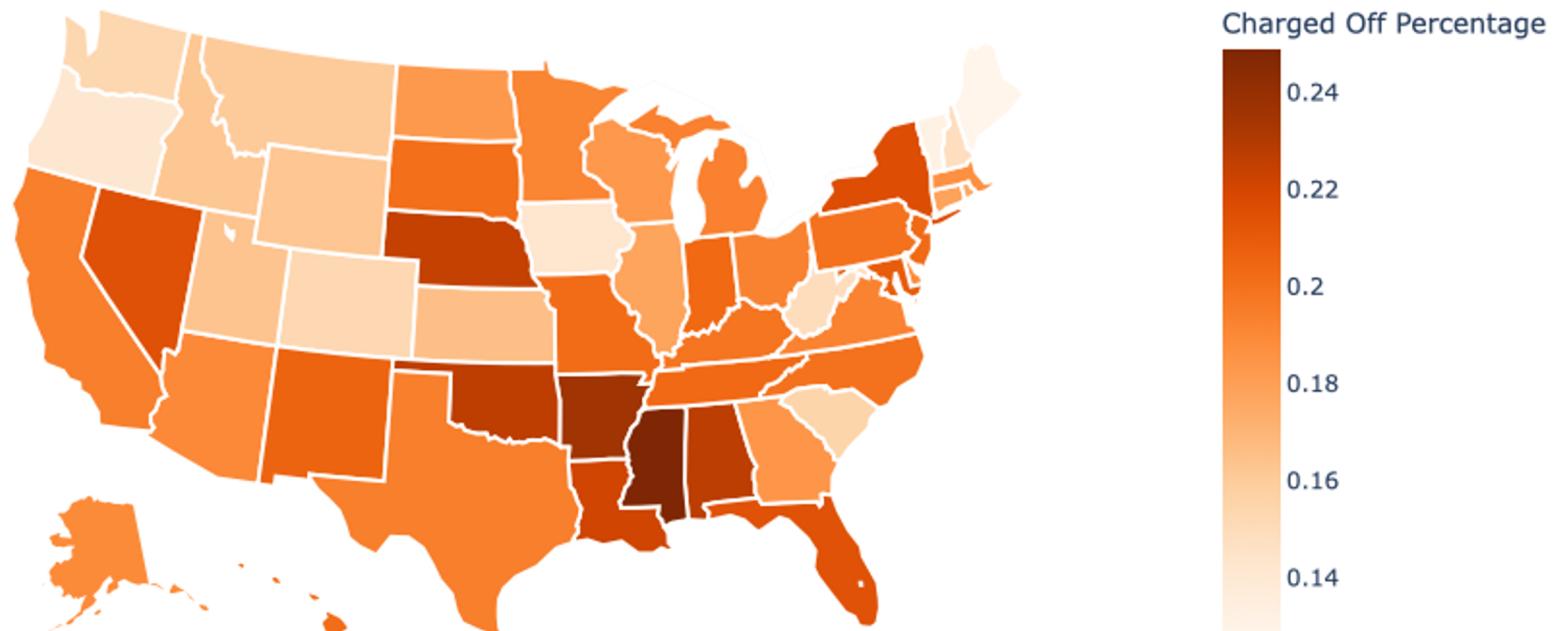
# EXPLORATORY DATA ANALYSIS

## ❑ Predictors – Application Type



# EXPLORATORY DATA ANALYSIS

## □ Predictors – State



# DATA PREPROCESSING

- ❑ Remove unnecessary features
- ❑ Missing Value

<code>mort_acc</code>	0.0254	→ Impute
<code>revol_util</code>	0.0008	
<code>dti</code>	0.0006	→ Drop
<code>pub_rec_bankruptcies</code>	0.0004	
.	.	.

- ❑ One-Hot Encoding

# MODEL

## ☐ Train Test Split

Train Data	80%
Test Data	20%

## ☐ Downsampling

Fully Paid	80%
Charged Off	20%



Fully Paid	50%
Charged Off	50%

# MODEL

## ❑ Model Selection

### ❑ Logistic Regression

### ❑ Random Forest

### ❑ XGBoost

## ❑ Hyperparameter Tuning

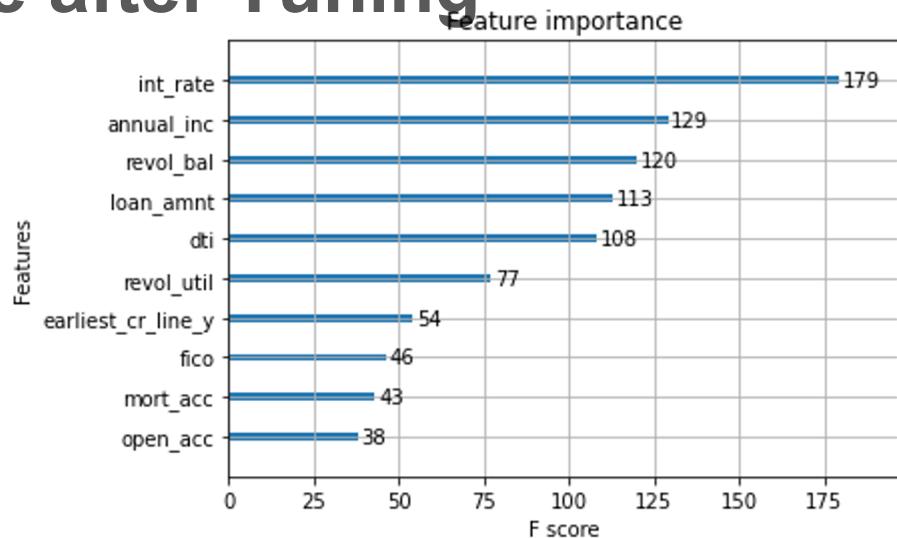
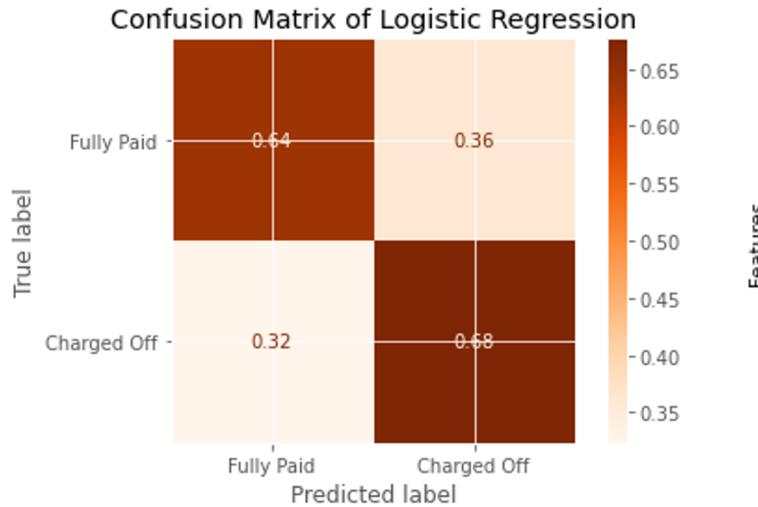
### ❑ Randomized Grid Search

*Summary of model performance on the test set*

	Logistic Regression	Random Forest	XGBoost
Precision	0.65	0.65	0.65
Recall	0.53	0.66	0.68
F1	0.58	0.65	0.66
AUC_ROC	0.68	0.71	0.72
Accuracy	0.62	0.65	0.66

# RESULT

## ❑ Model Performance after Tuning



## ❑ Web Application

## ❑ Streamlit Python Library

# CONCLUSION

- ❑ We built a predictive model to predict the loan status at the closing time.
- ❑ The performance of the model is not very good.
- ❑ We can choose a high threshold and use this model as a preventive tool for the investors to avoid risky loans.

Thank You  
&  
Questions?

---

THE GEORGE  
WASHINGTON  
UNIVERSITY

---

WASHINGTON, DC

---

THE GEORGE  
WASHINGTON  
UNIVERSITY

---

WASHINGTON, DC