# Midterm project presentation

**Loan Default Dataset Analysis**

GitHub repository:
https://github.com/Bail111/Loan
_Default-Analysis.git

REPORTER：Qinjunjie Pu

October 25th

01
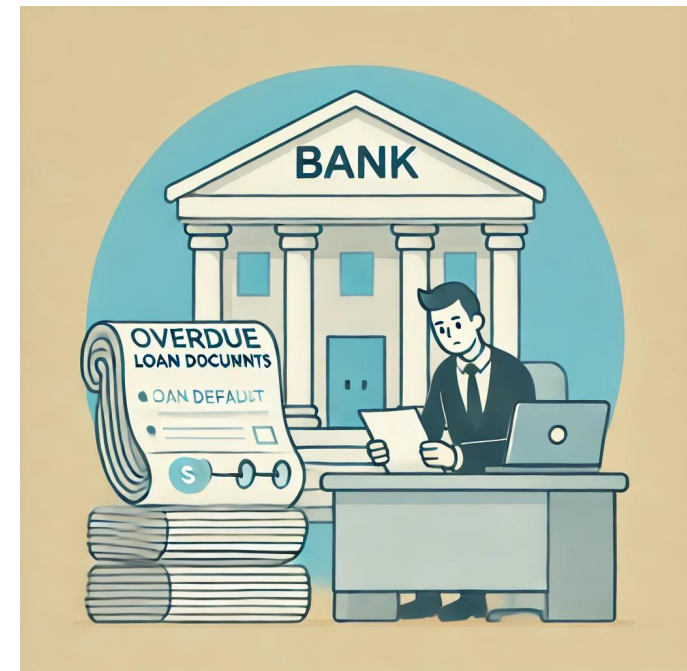
troduction

# Introduction

This project aims to predict whether a loan will default.

Minimizing financial losses and optimizing lending decisions.

## Classification problem

Data Source: kaggle

Data Collection: past data on the loan borrowers

# Exploratory Data Analysis

# EDA

Brief description of dataset

```
row: 148670
col: 34

12 continuous features
21 categorical features
```
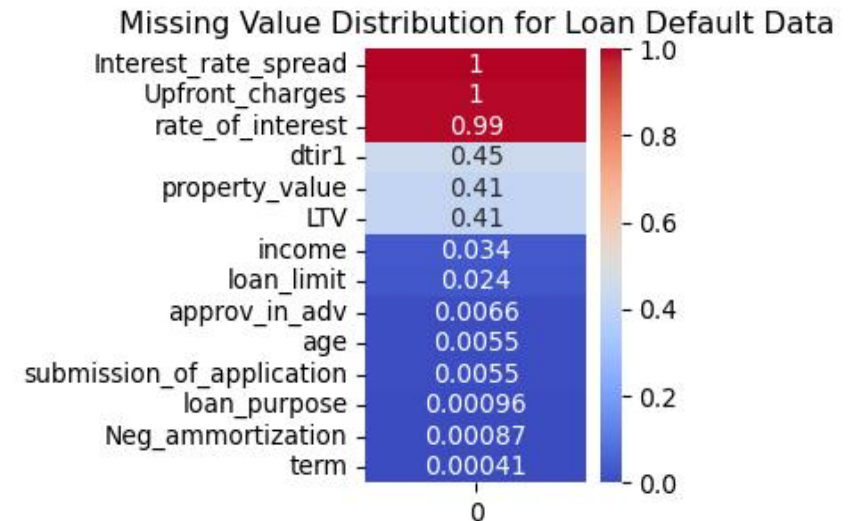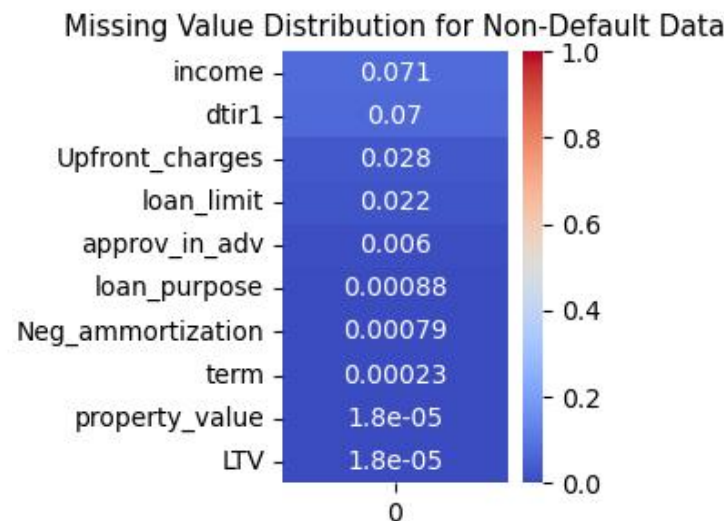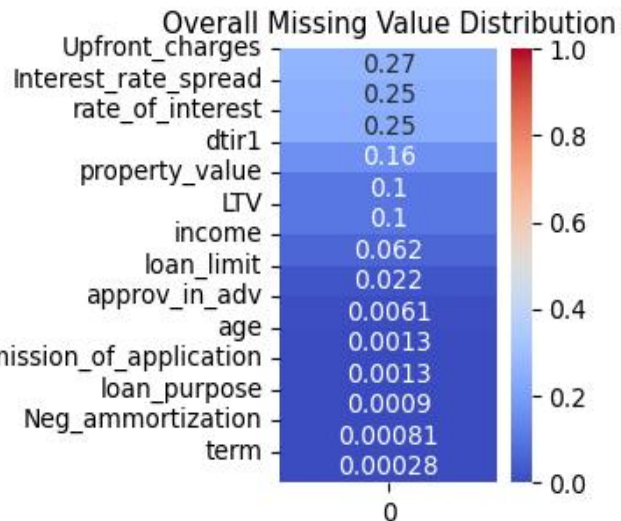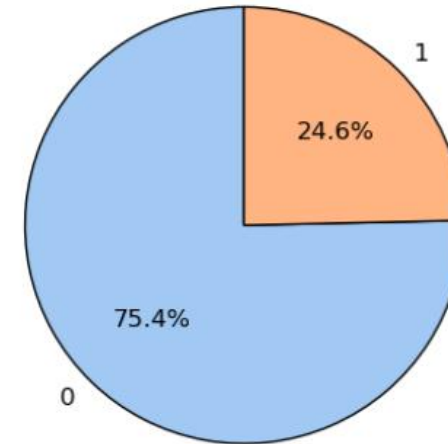
**Loan Default Status Distribution**



**Overall Missing Value Distribution**

| feature | value |
| --- | --- |
| Upfront_charges | 0.27 |
| Interest_rate_spread | 0.25 |
| rate_of_interest | 0.25 |
| dtir1 | 0.16 |
| property_value | 0.1 |
| LTV | 0.1 |
| income | 0.062 |
| loan_limit | 0.022 |
| approv_in_adv | 0.0061 |
| age | 0.0013 |
| submission_of_application | 0.0013 |
| loan_purpose | 0.0009 |
| Neg_ammortization | 0.00081 |
| term | 0.00028 |

**Missing Value Distribution for Non-Default Data**

| feature | value |
| --- | --- |
| income | 0.071 |
| dtir1 | 0.07 |
| Upfront_charges | 0.028 |
| loan_limit | 0.022 |
| approv_in_adv | 0.006 |
| loan_purpose | 0.00088 |
| Neg_ammortization | 0.00079 |
| term | 0.00023 |
| property_value | 1.8e-05 |
| LTV | 1.8e-05 |

**Missing Value Distribution for Loan Default Data**

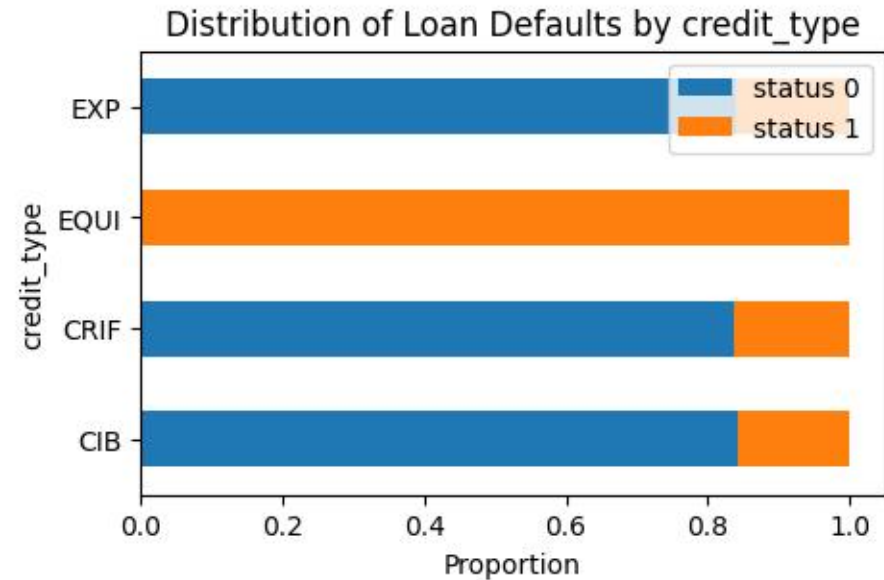| feature | value |
| --- | --- |
| Interest_rate_spread | 1 |
| Upfront_charges | 1 |
| rate_of_interest | 0.99 |
| dtir1 | 0.45 |
| property_value | 0.41 |
| LTV | 0.41 |
| income | 0.034 |
| loan_limit | 0.024 |
| approv_in_adv | 0.0066 |
| age | 0.0055 |
| submission_of_application | 0.0055 |
| loan_purpose | 0.00096 |
| Neg_ammortization | 0.00087 |
| term | 0.00041 |

# EDA

Analysis for categorical features

**credit_type :** applicant's type of credit


Distribution of credit_type
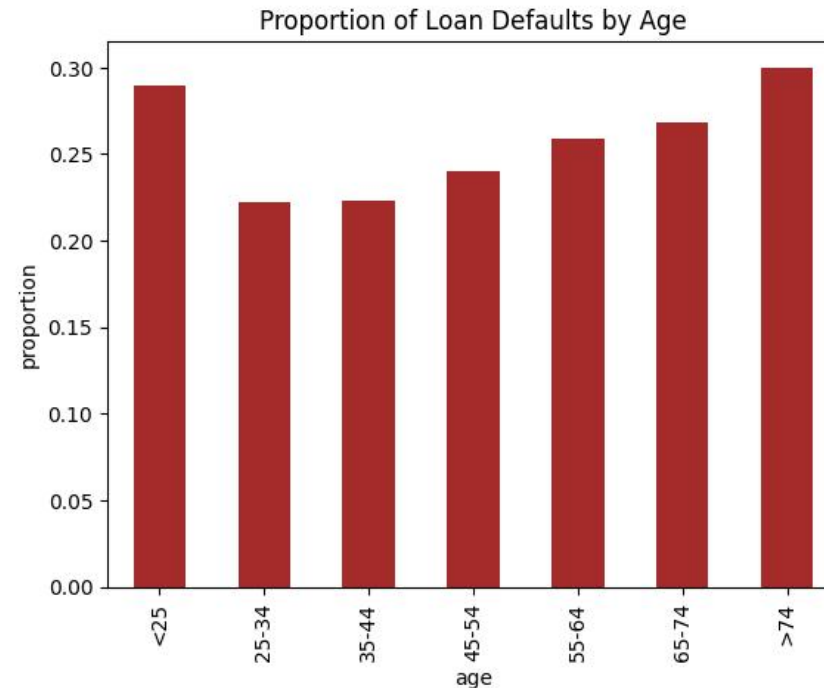

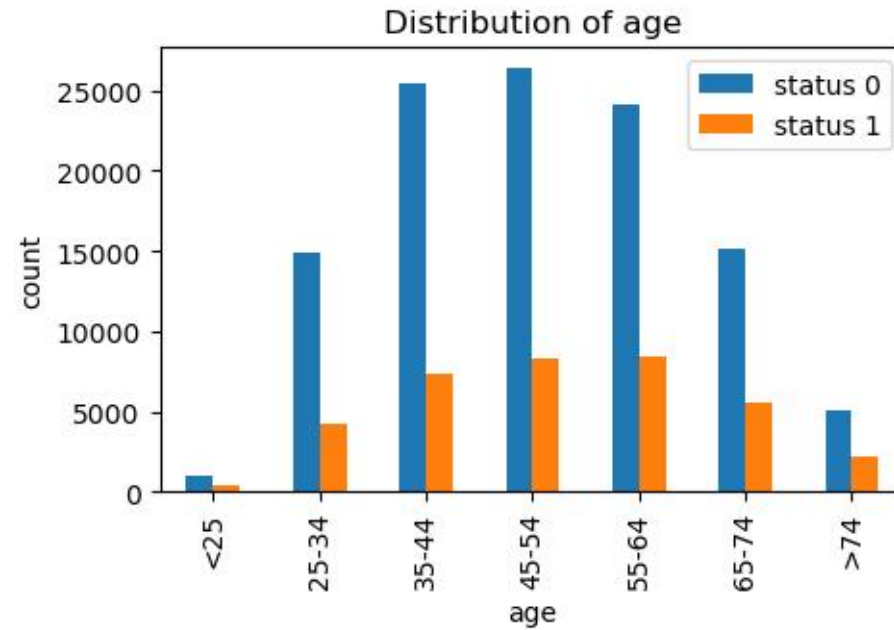Distribution of Loan Defaults by credit_type

# EDA



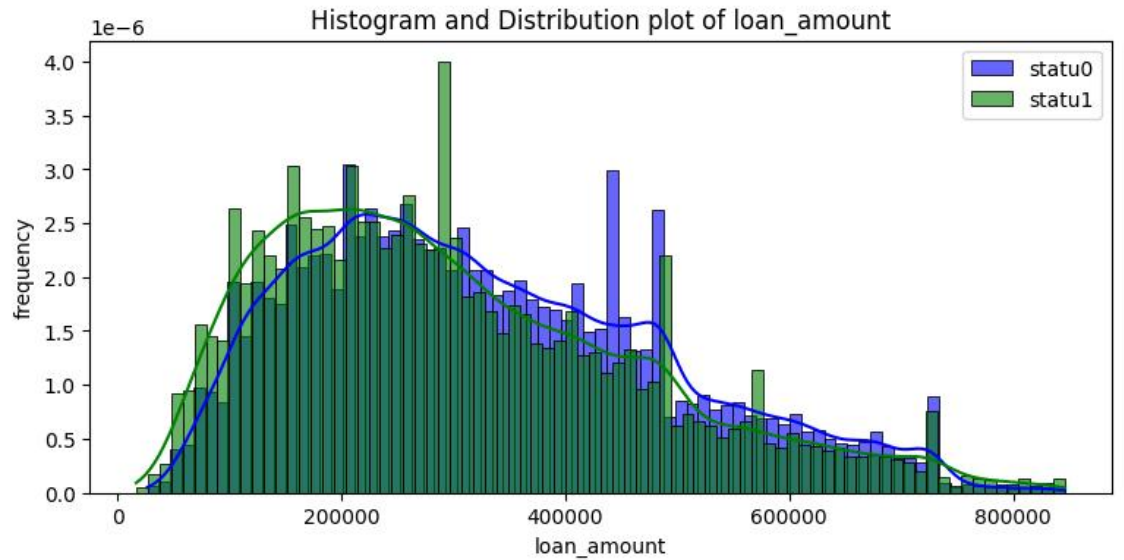🌐 Analysis for categorical features
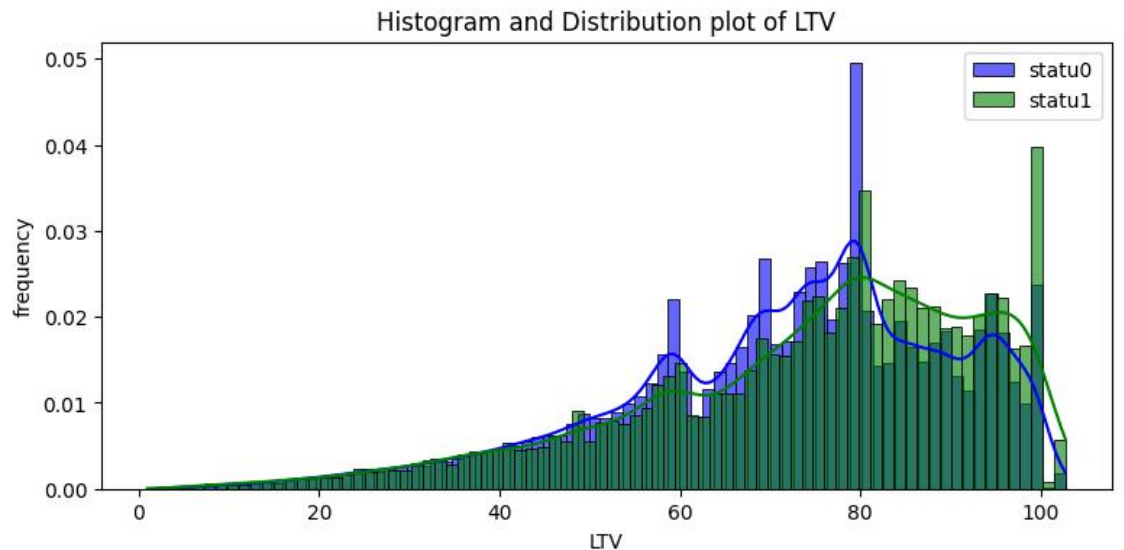
**age:**   the age of the applicant

# EDA

Analysis for continuous features
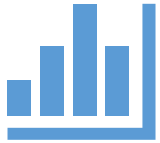
**loan_amount:** amount of money being borrowed

**LTV:** loan-to-value ratio, calculated as the loan amount divided by the property value

# EDA

Analysis for continuous features



Scatter Plot of Loan Amount vs Income by Loan Default Status (Log-Scaled)
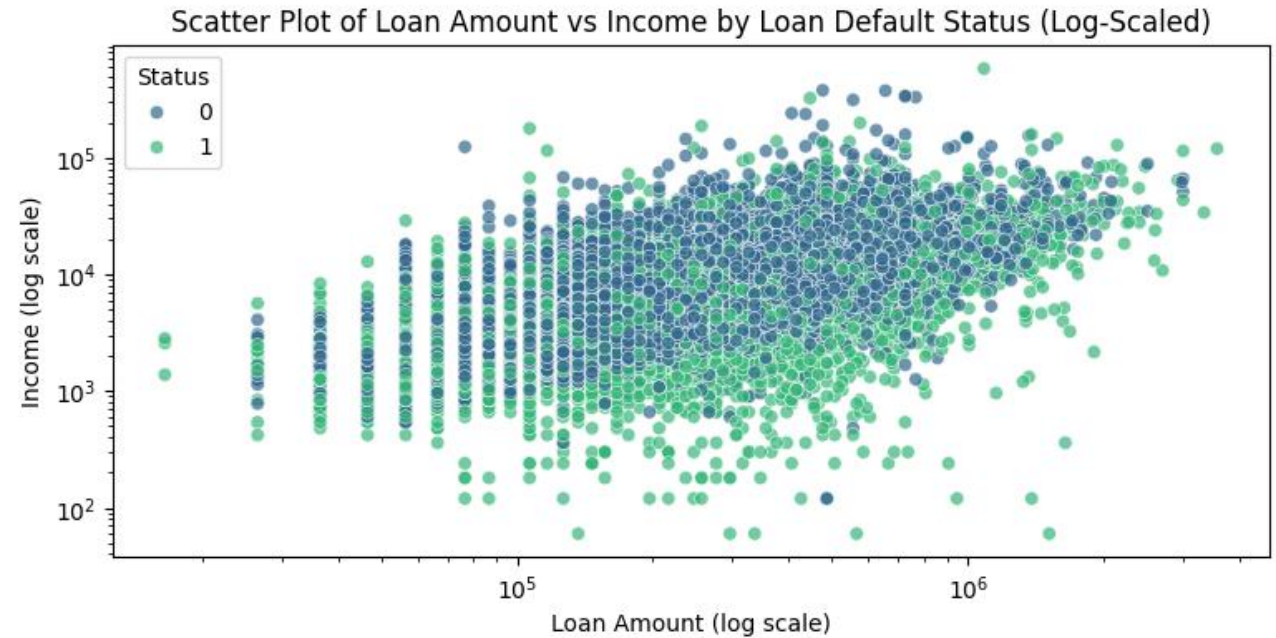
**Income:** applicant's annual income

**loan_amount:** amount of money being borrowed

# Splitting and preprocessing

# Splitting and preprocessing

Splitting method
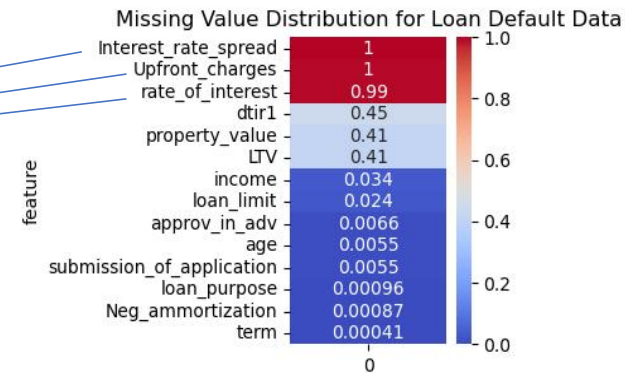
20% test data

Loan default
Dataset

80% other data

Using StratifiedKFold to 10-fold training and validation data

# Splitting and preprocessing

## Preprocessing

ID, year (2019 for all data)

1. Delete features

Missing Value Distribution for Loan Default Data

| feature | value |
|---|---|
| Interest_rate_spread | 1 |
| Upfront_charges | 1 |
| rate_of_interest | 0.99 |
| dtir1 | 0.45 |
| property_value | 0.41 |
| LTV | 0.41 |
| income | 0.034 |
| loan_limit | 0.024 |
| approv_in_adv | 0.0066 |
| age | 0.0055 |
| submission_of_application | 0.0055 |
| loan_purpose | 0.00096 |
| Neg_ammortization | 0.00087 |
| term | 0.00041 |

2. For categorical and ordinal features:
            Impute missing value with 'missing'

Binary encoding        e.g. business_or_commercial  ('nob/c' to 0,  'b/c' to 1)
Onehot encoding        e.g. region
Ordinal encoding       e.g. age  ('<25', '25-34', '35-44', etc.)

3. For continuous features:
Minmax scalar        e.g. credit_score
Standard scalr        e.g. income

Before preprocessing:  33 features
After preprocessing: 53 features

**Thanks for listening**

end