

Loan Approval Prediction – Checkpoint 1

Group: Predictive Minds

Team Members:

- **Suraj Ravindra Rao (ASU ID: 1233990435)**
- **Om Gaurav Patel (ASU ID: 1233379110)**
- **Darren Alastair Ferrao (ASU ID: 1237189124)**
- **Mohammad Asgar Khan Dalwai (ASU ID: 1237671684)**

Background & Problem

Loan approval prediction is essential for financial institutions.

Problem: Binary classification → Approved (0) vs Default (1).

Importance:

Risk Reduction

Reduces loan defaults and financial risk.

Process Automation

Automates decisions for faster processing.

Fair & Scalable

Ensures fairness and scalability.

Dataset Overview

Dataset Evolution

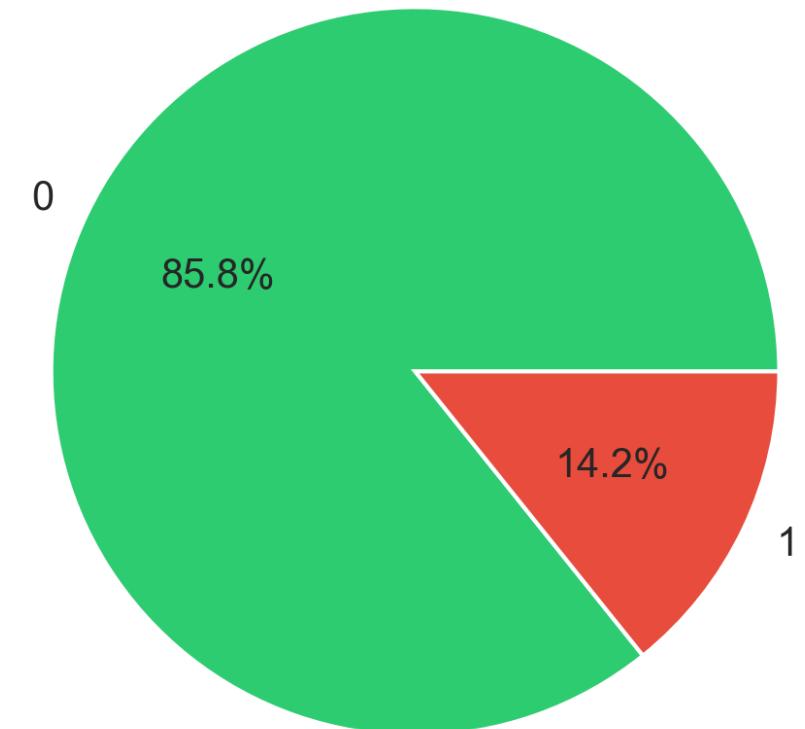
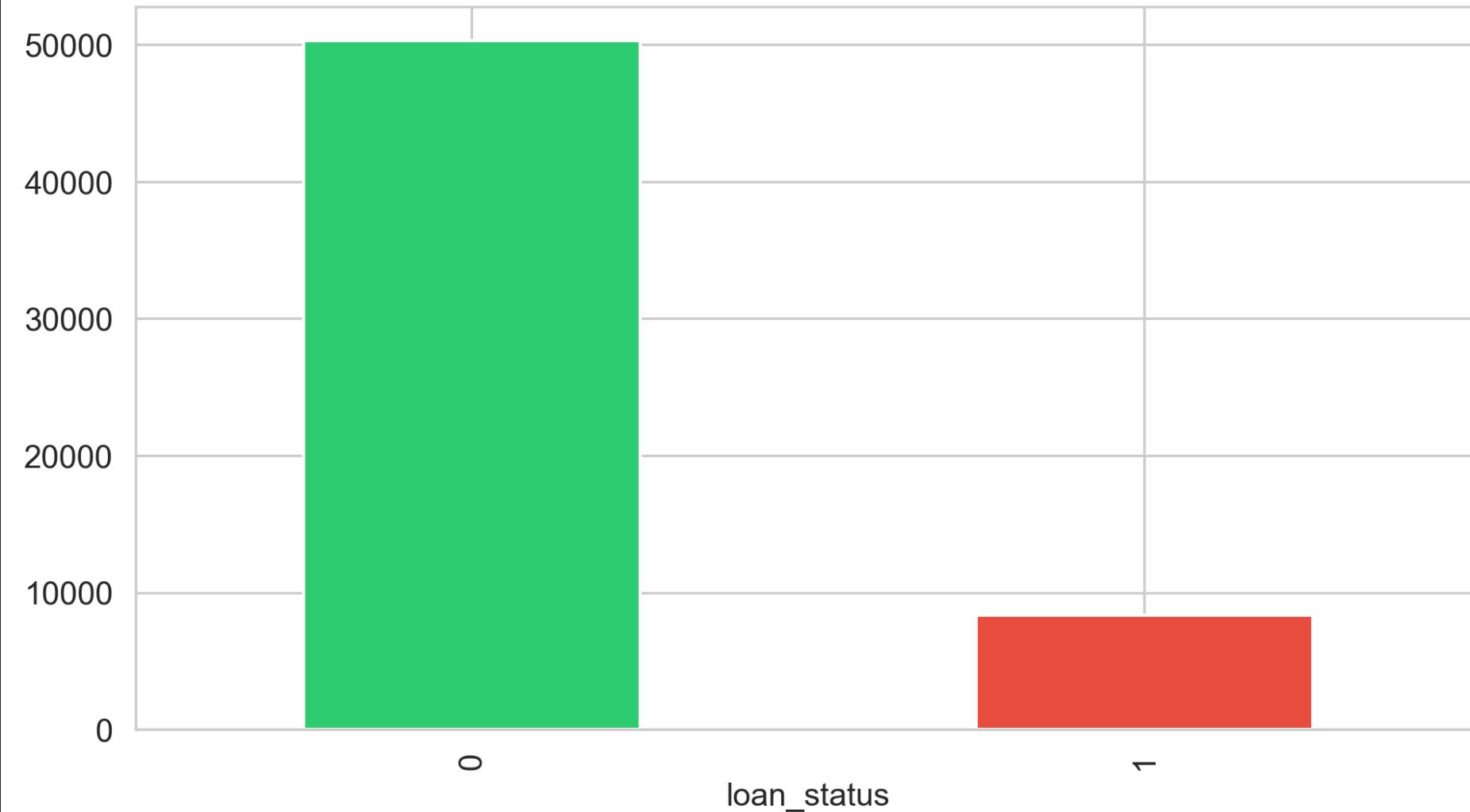
Pitch dataset: 4,269 samples (initial plan).

Updated dataset (Checkpoint 1): Kaggle Playground dataset, 58,645 samples × 13 features.

Features: demographic, employment, loan details, credit history.

Train/Test split: 80/20.

Loan Status Distribution



Numerical Feature Exploration

Age

Concentrated in 20s–30s.

Income

Skewed with large outliers.

Loan Amount

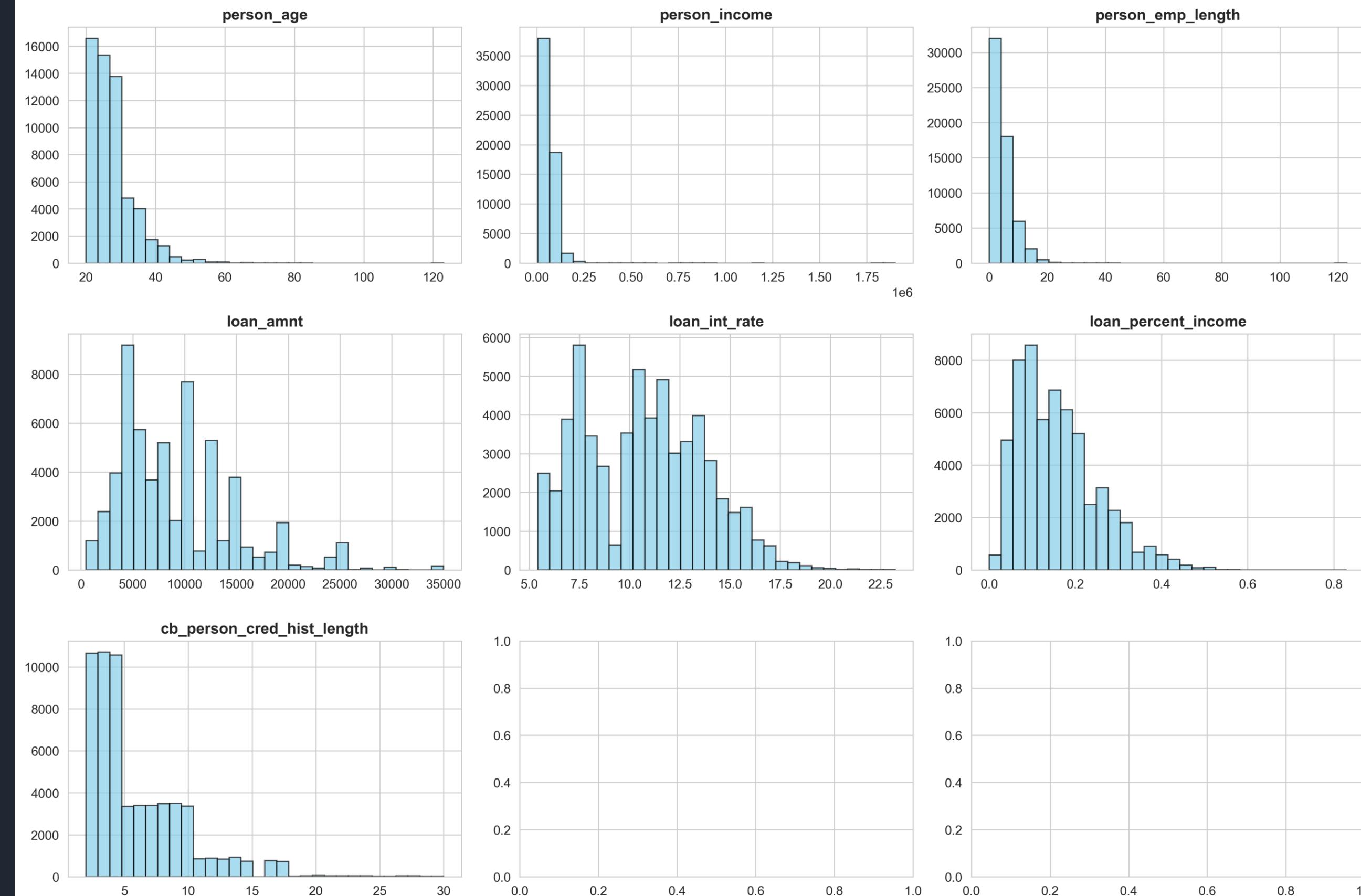
Most between 5k–15k.

Interest Rate

7–15% range dominates.

Credit History

Majority < 10 years.



Categorical Feature Exploration



Home Ownership

Rent and Mortgage are common.



Loan Intent

Education and Medical dominate.



Loan Grade

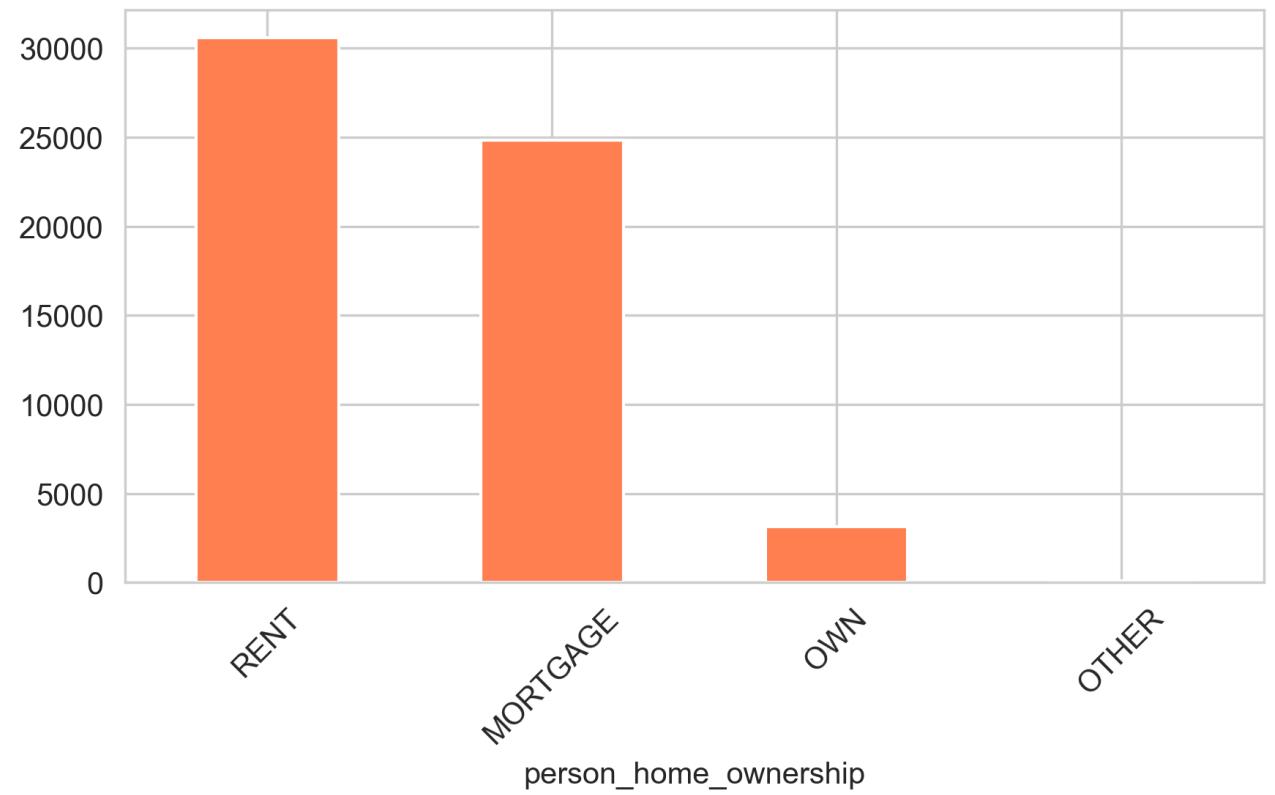
Lower grades (D–G) have higher default risk.



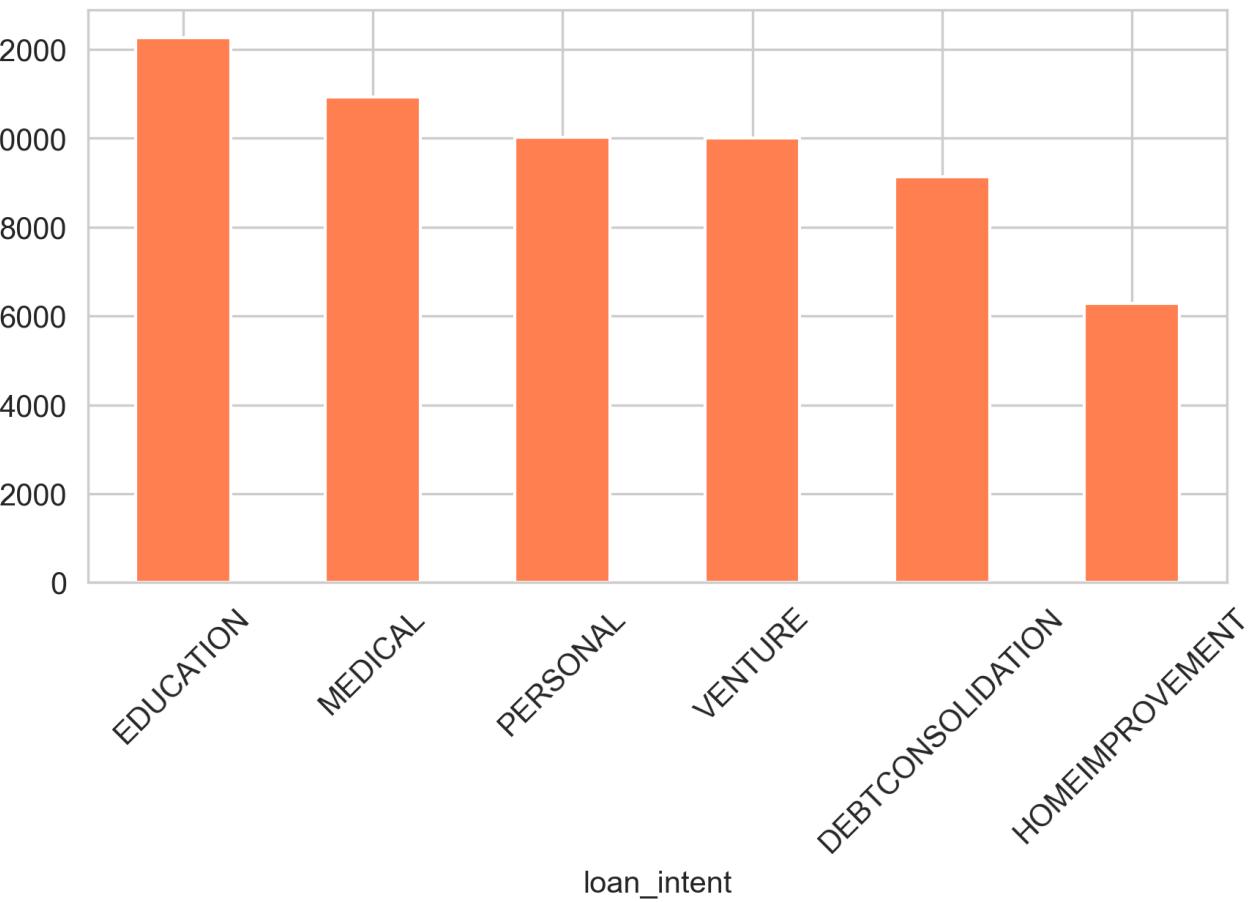
Default History

Strong predictor of loan status.

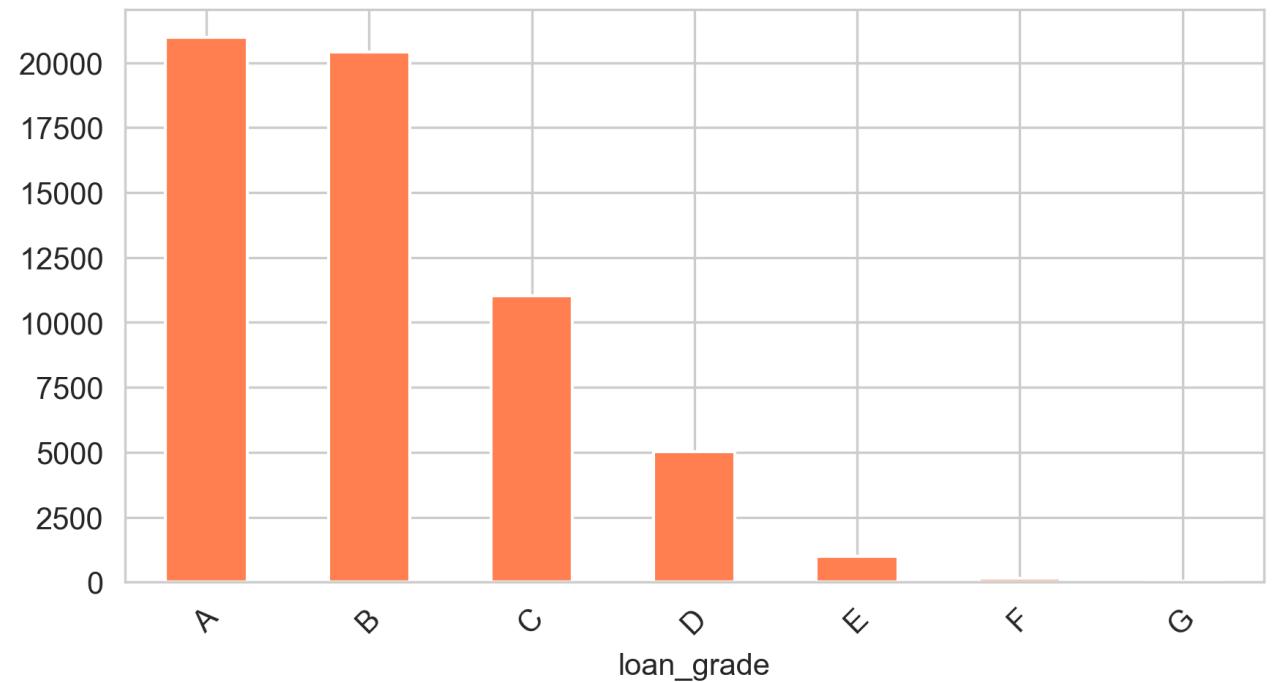
person_home_ownership



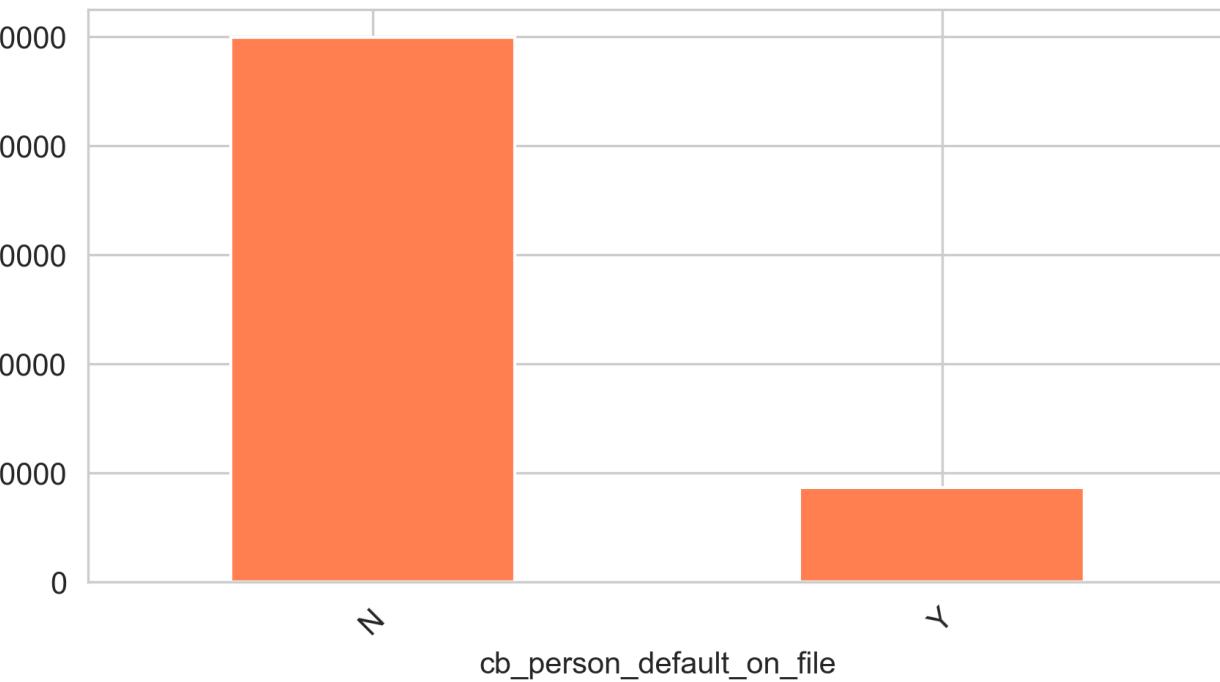
loan_intent



loan_grade



cb_person_default_on_file



Correlations & Relationships

1

Strong Correlation

Loan Percent Income ↔ Loan Amount.

2

Age Relationship

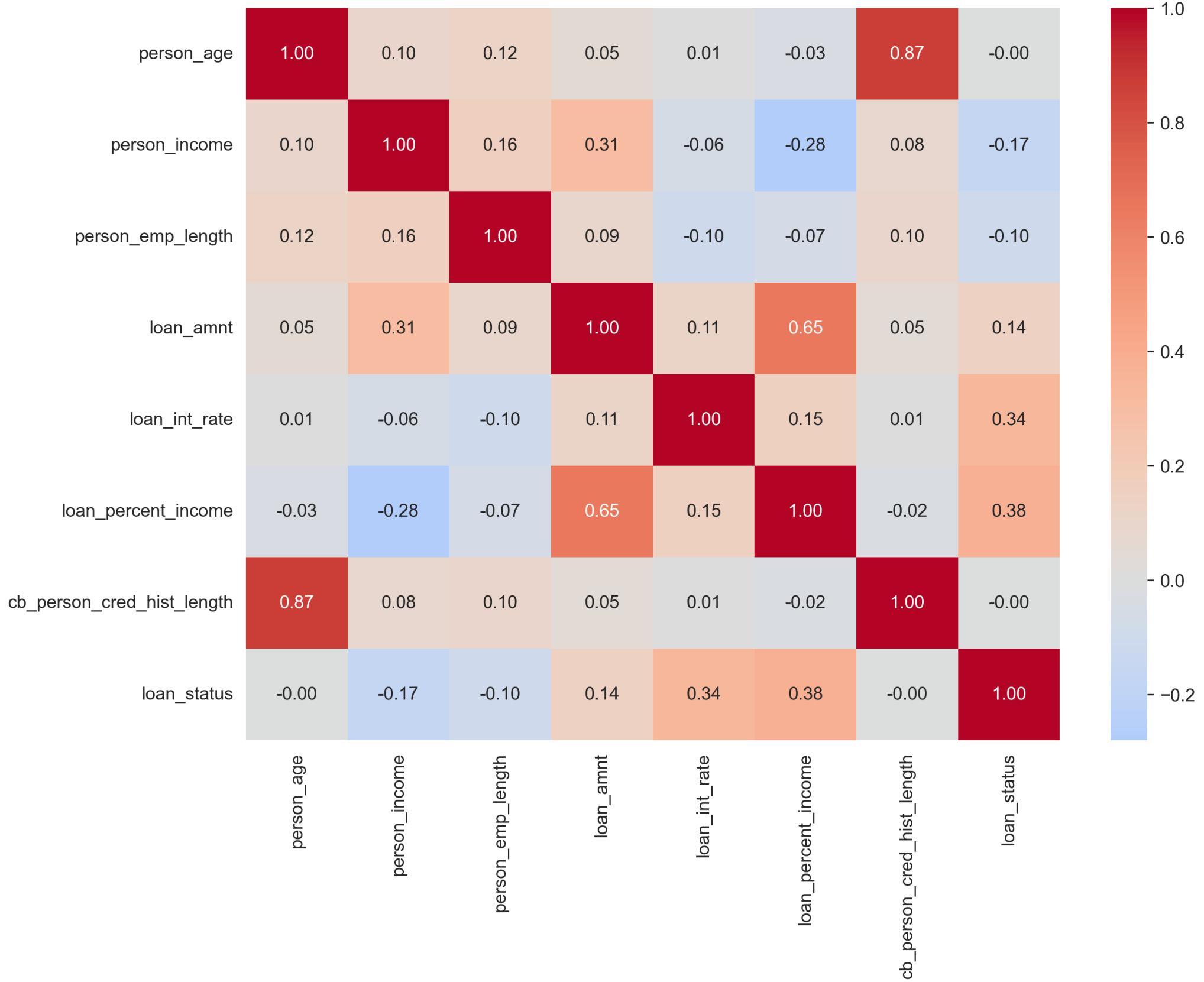
Age ↔ Credit history length also related.

3

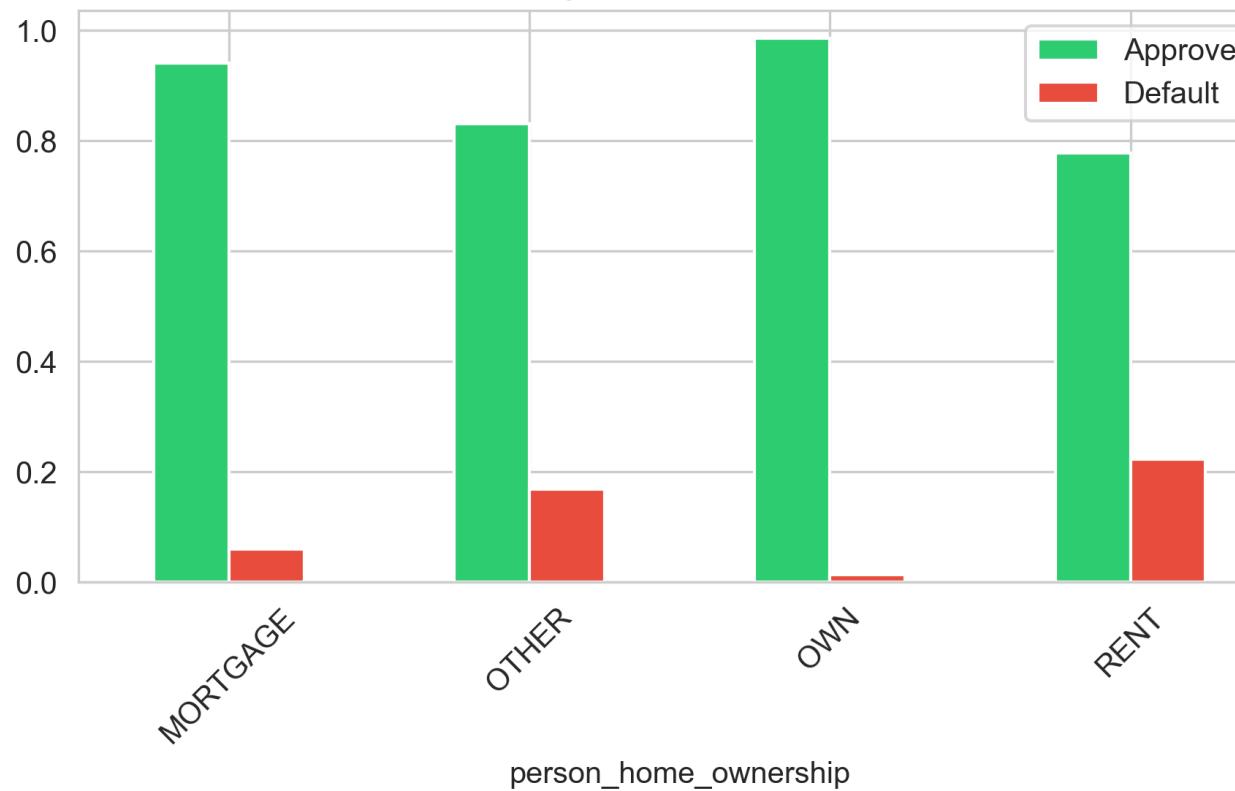
Status Influence

Loan Status influenced by loan grade & percent income.

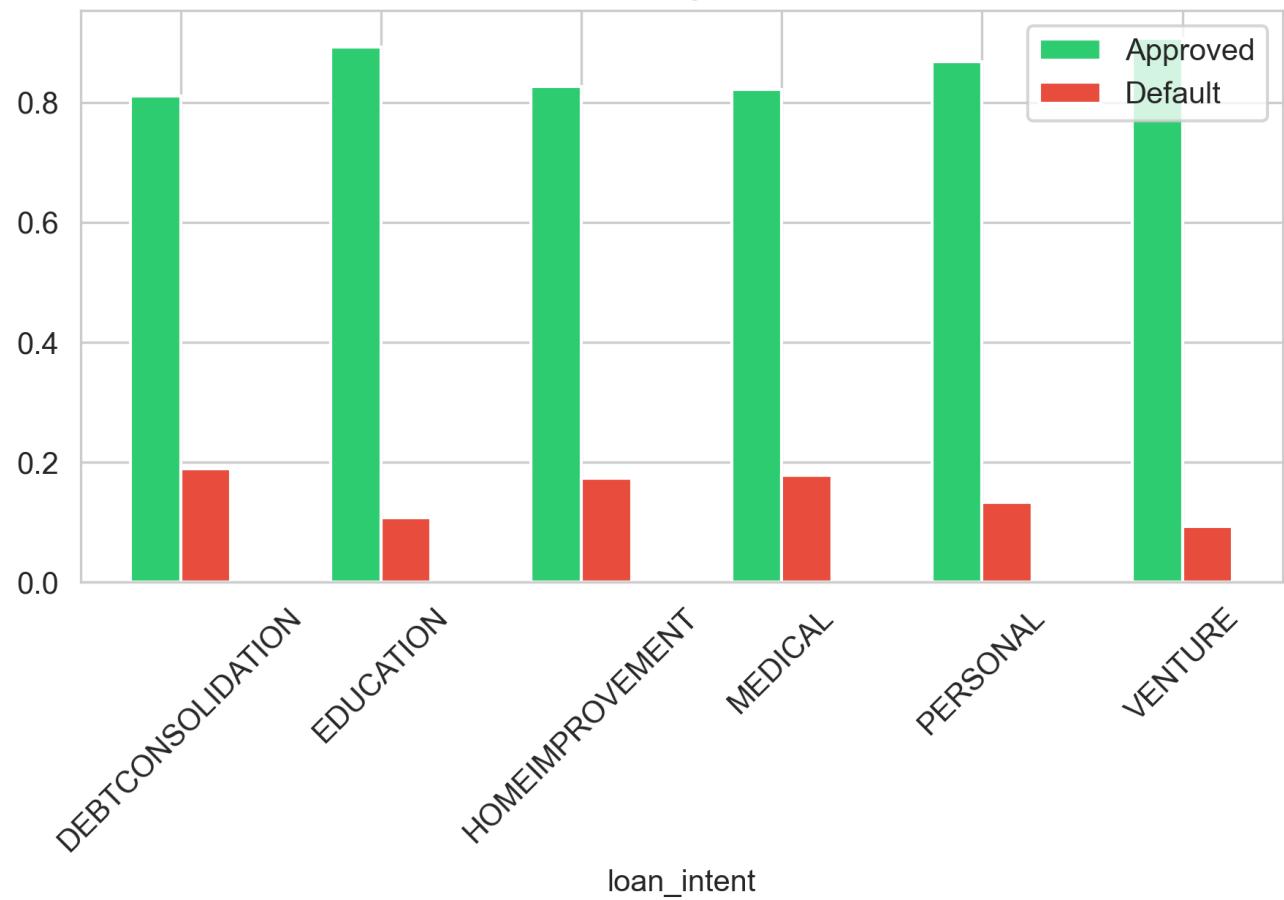
Correlation Matrix



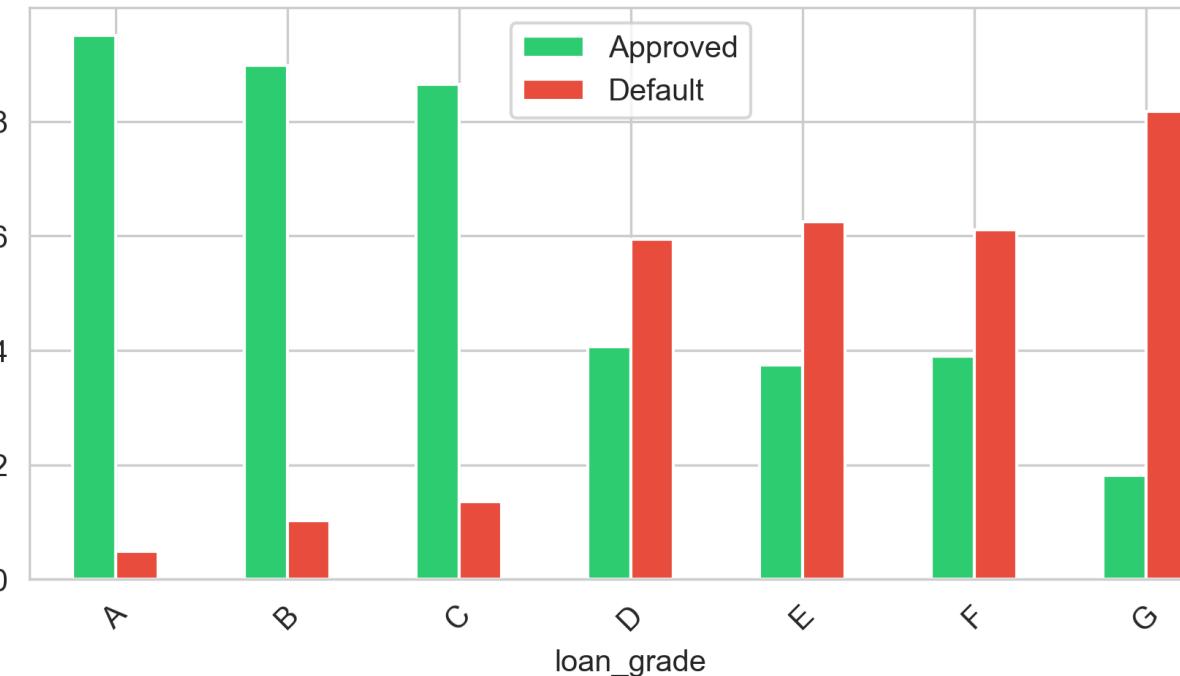
Loan Status by person_home_ownership



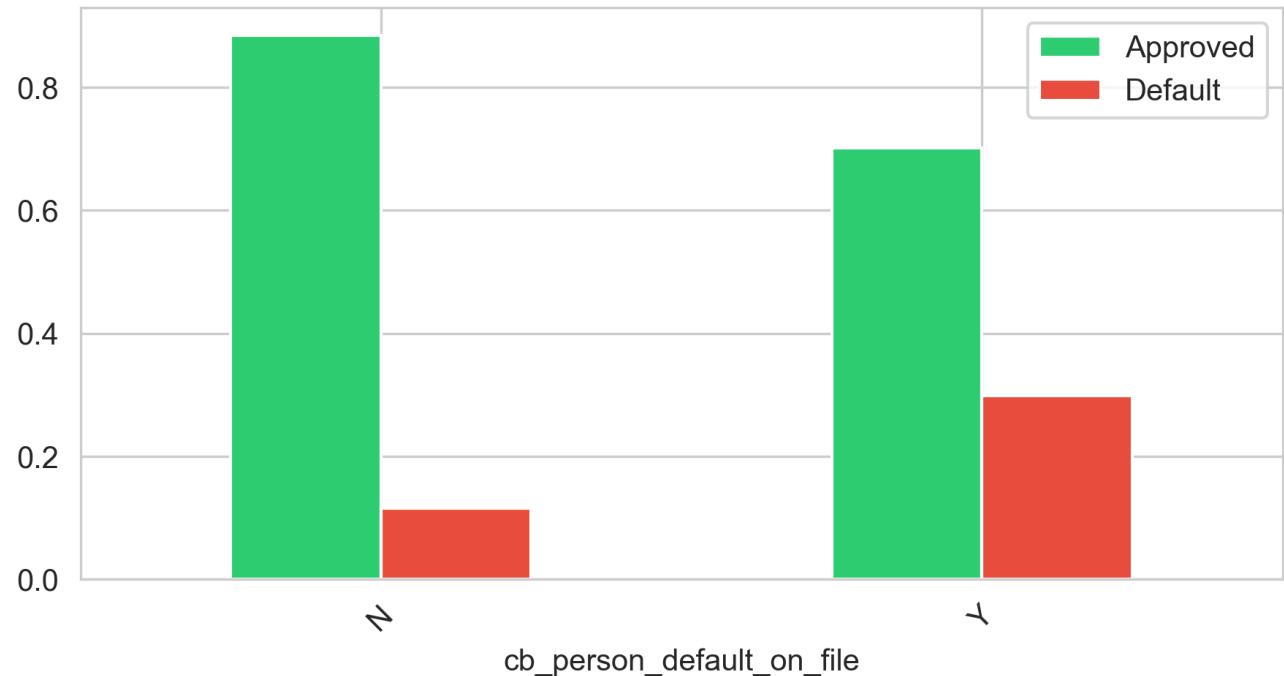
Loan Status by loan_intent



Loan Status by loan_grade



Loan Status by cb_person_default_on_file



Data Preprocessing

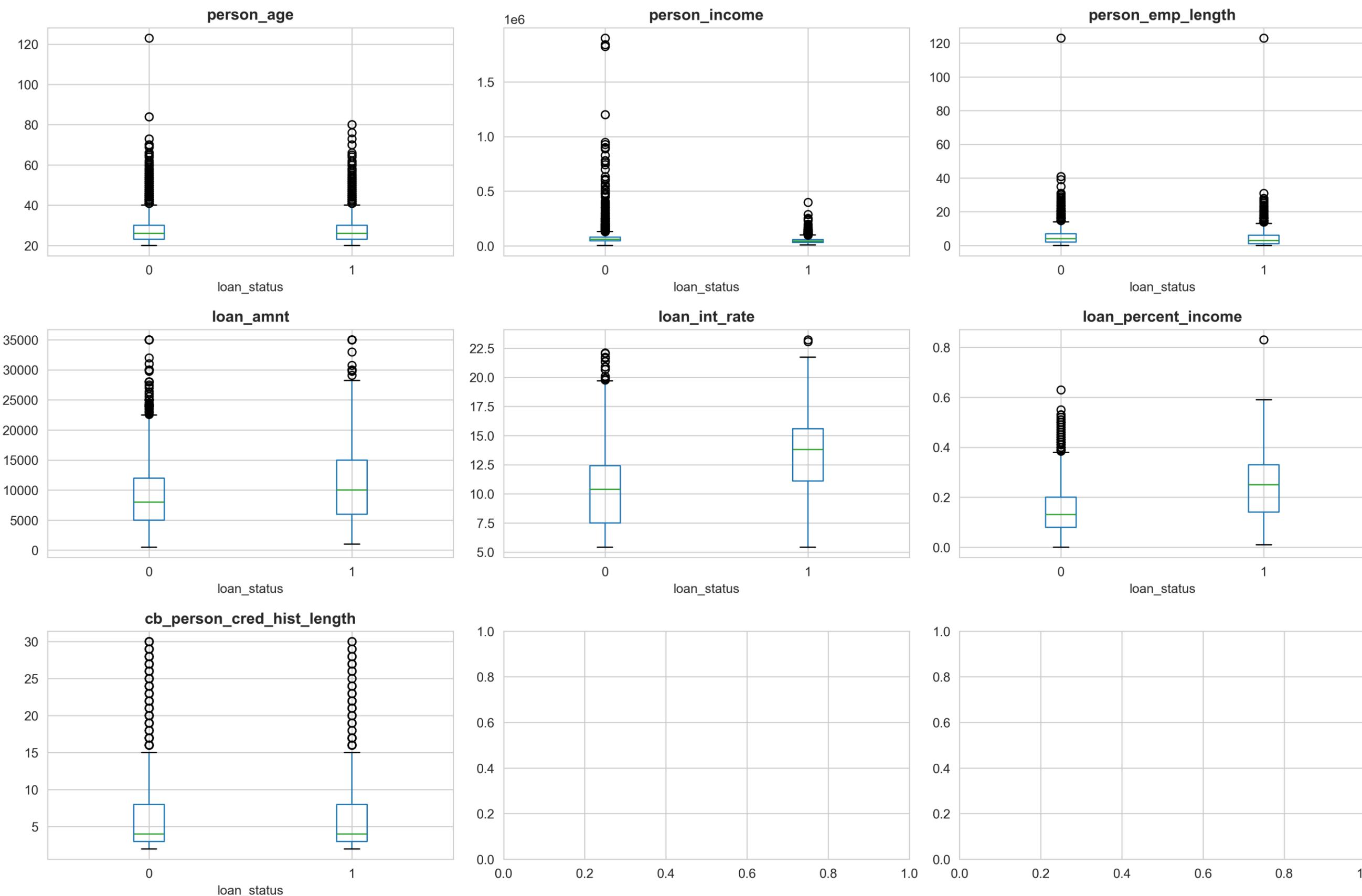
Data Cleaning & Transformation

- Dropped IDs and irrelevant fields.
- Label encoded categorical features.
- Standardized numerical features.

Feature Engineering:

Income-to-loan ratio

Age-to-employment ratio



Model Training & Evaluation

Models: Logistic Regression, Decision Tree, Random Forest.

Metrics: Accuracy, Precision, Recall, F1, ROC-AUC.

Results:

Logistic Regression

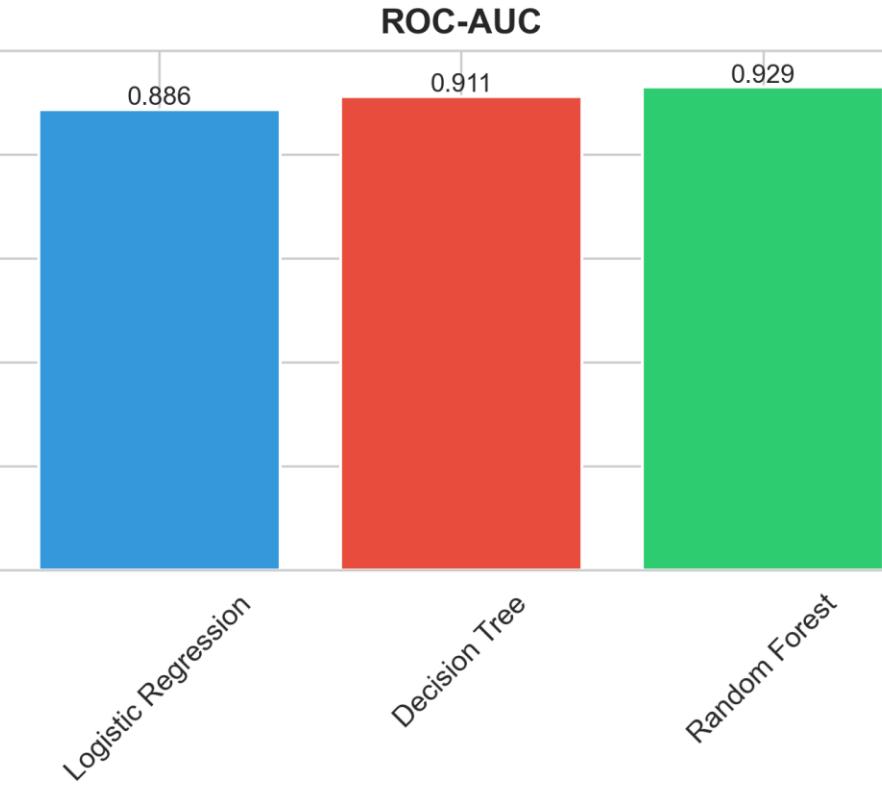
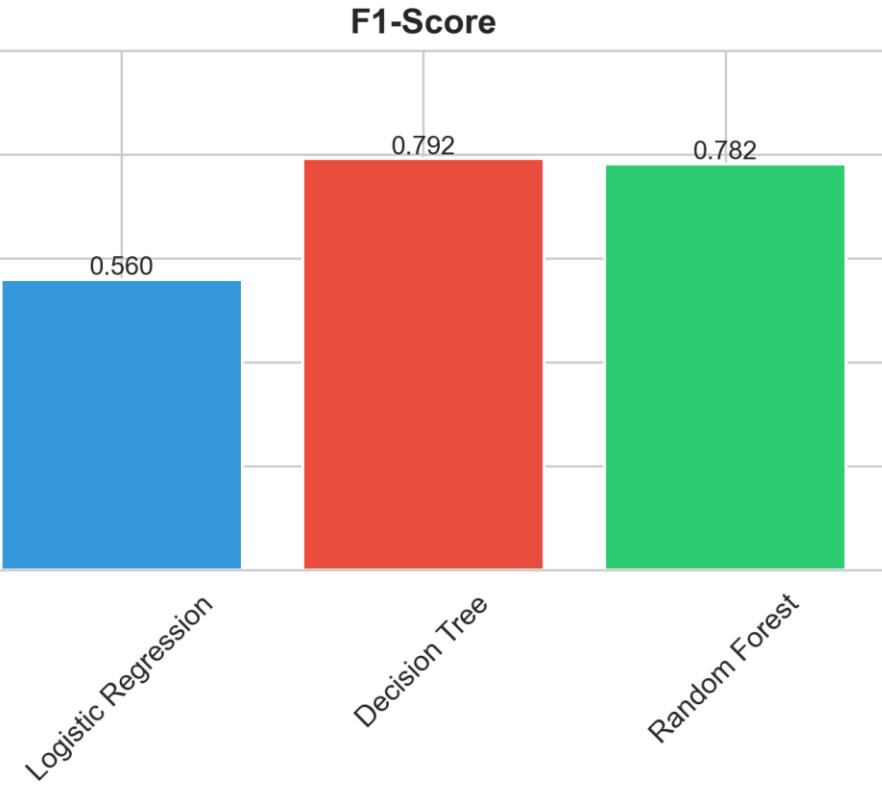
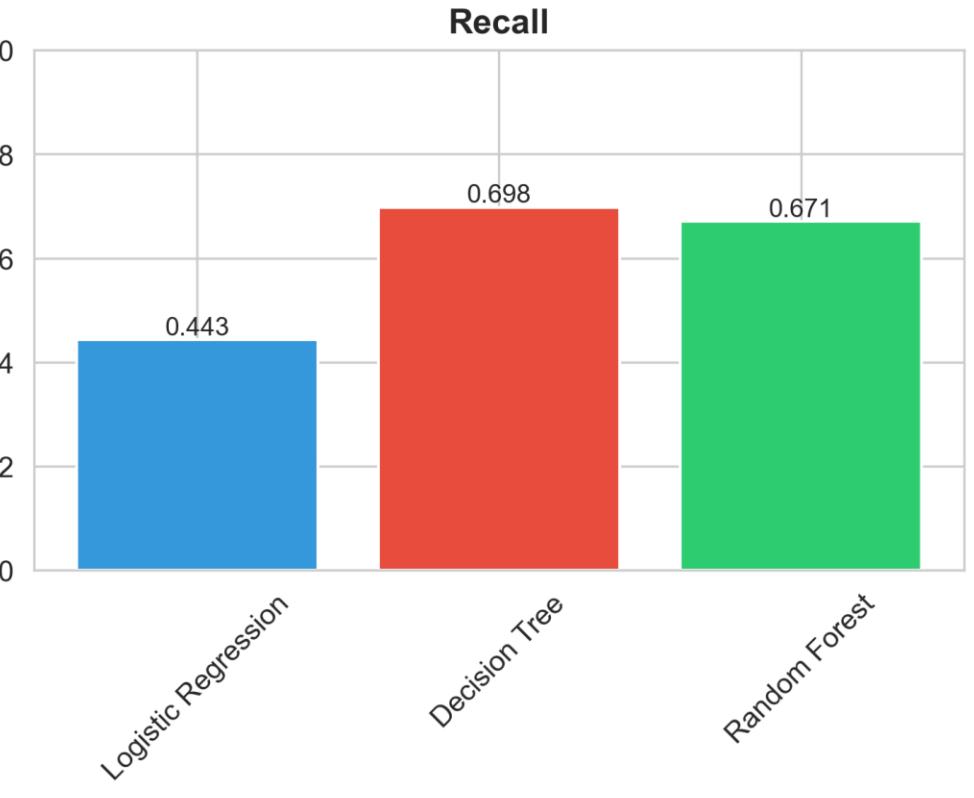
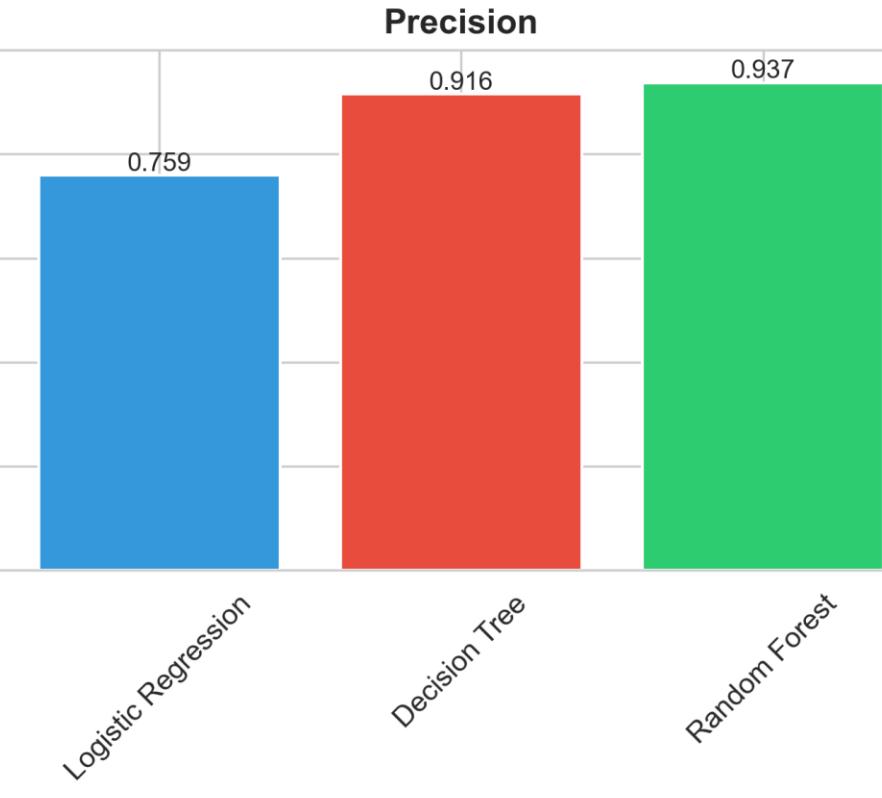
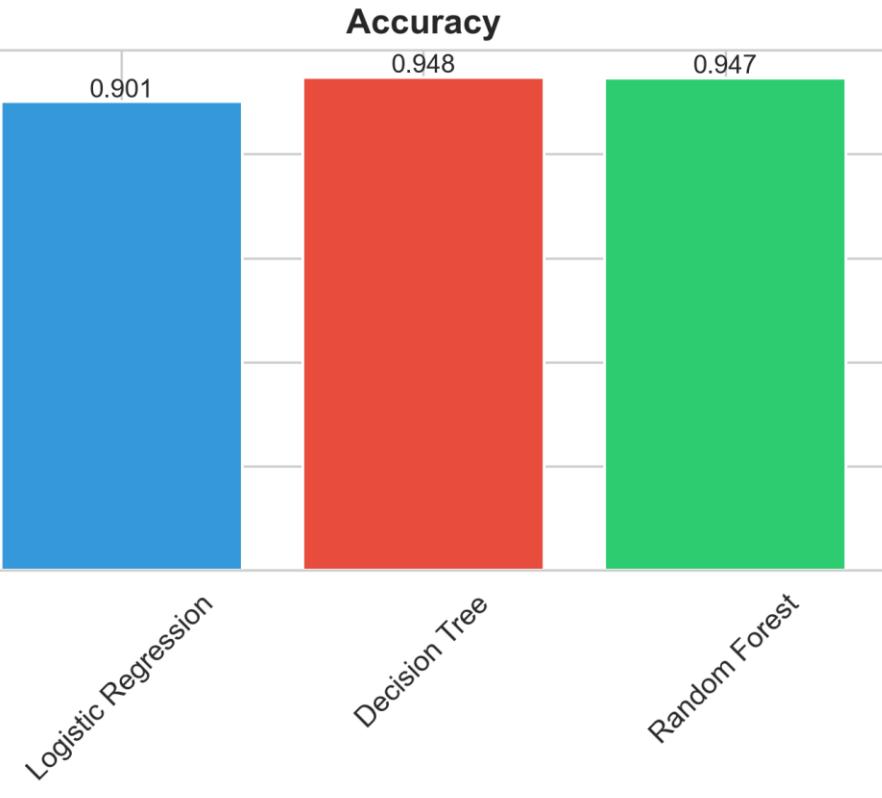
Accuracy 0.90, Recall 0.44.

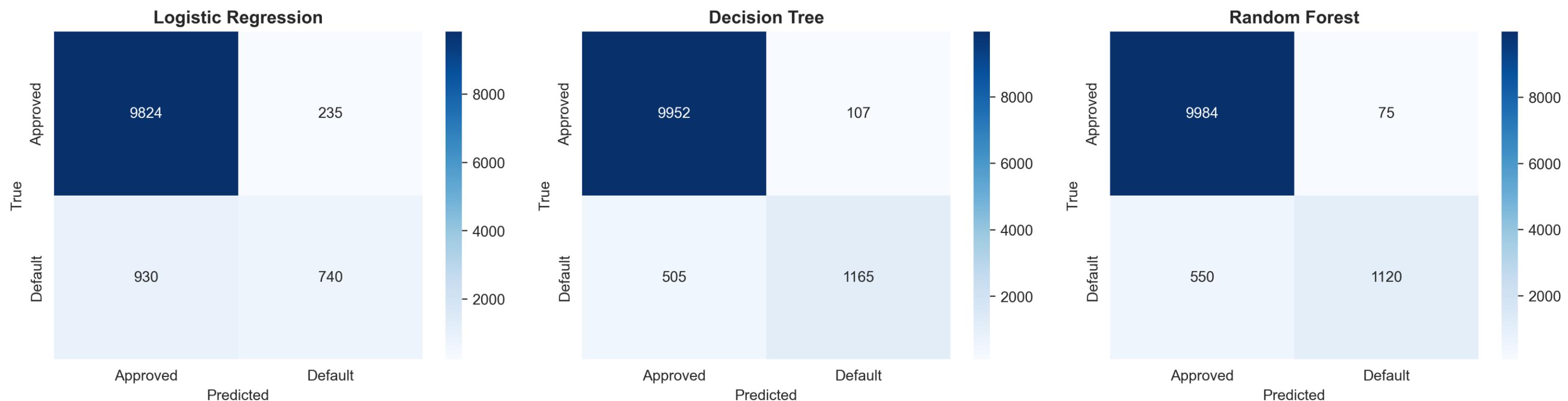
Decision Tree

Accuracy 0.95, Recall 0.70.

Random Forest

Accuracy 0.95, ROC-AUC 0.93.





Feature Importance

Top predictors:

01

Loan Percent Income

02

Loan Grade

03

Income-to-loan ratio

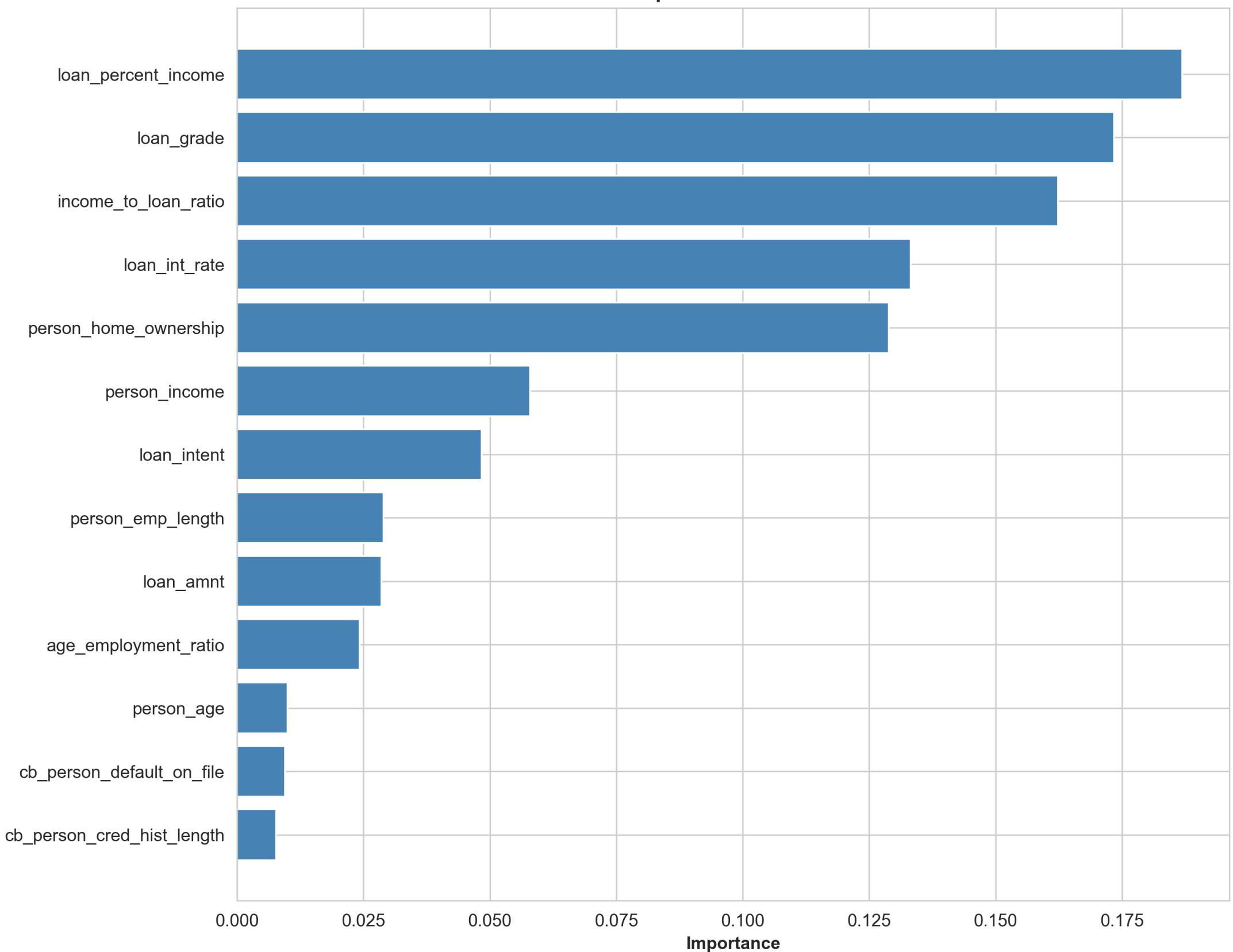
04

Interest rate

05

Home ownership

Feature Importance - Random Forest



Insights & Next Steps

Insights:

- Tree-based models perform best.
- Class imbalance hurts recall for defaults.
- Key predictors: income ratios, loan grade.

Next Steps:

- Apply SMOTE/undersampling for balance.
- Hyperparameter tuning.
- Try XGBoost/LightGBM.
- Add interpretability (SHAP, LIME).

Thank You

Video Link: <https://drive.google.com/file/d/1DM88fetLoI78E1IDGRsYyeiD0g5knKzl/view?usp=sharing>