

Lending Club Case Study

1. Introduction

1.1 Project Overview

The primary goal of this project is to analyze loan data and build a predictive model to identify loans that are likely to be charged off. This involves data cleaning, univariate and bivariate analysis, and the development of a logistic regression model to predict loan status. Develop a basic understanding of risk analytics in banking and financial services and understand how data is used to minimise the risk of losing money while lending to customers.

1.2 Dataset Description

The dataset used in this project is sourced from a financial institution and contains various features related to loans, such as loan amount, interest rate, term, grade, employment length, home ownership status, verification status, loan status, and loan purpose.

2. Data Preprocessing

2.1 Handling Missing Values

To ensure the quality of our data, we first address missing values. Columns with more than 50% missing values are dropped, followed by the removal of rows with any missing values.

2.2 Data Cleaning

The 'int_rate' column contains percentage signs, which are removed and the column is converted to float.

3. Exploratory Data Analysis

3.1 Univariate Analysis

3.1.1 Categorical Variables

The distribution of categorical variables such as 'term', 'grade', 'sub_grade', 'emp_length', 'home_ownership', 'verification_status', 'loan_status', and 'purpose' is visualized using count plots.

3.1.2 Numerical Variables

The distribution of numerical variables such as 'loan_amnt', 'int_rate', 'installment', 'annual_inc', and 'dti' is visualized using histograms with kernel density estimates.

3.2 Bivariate Analysis

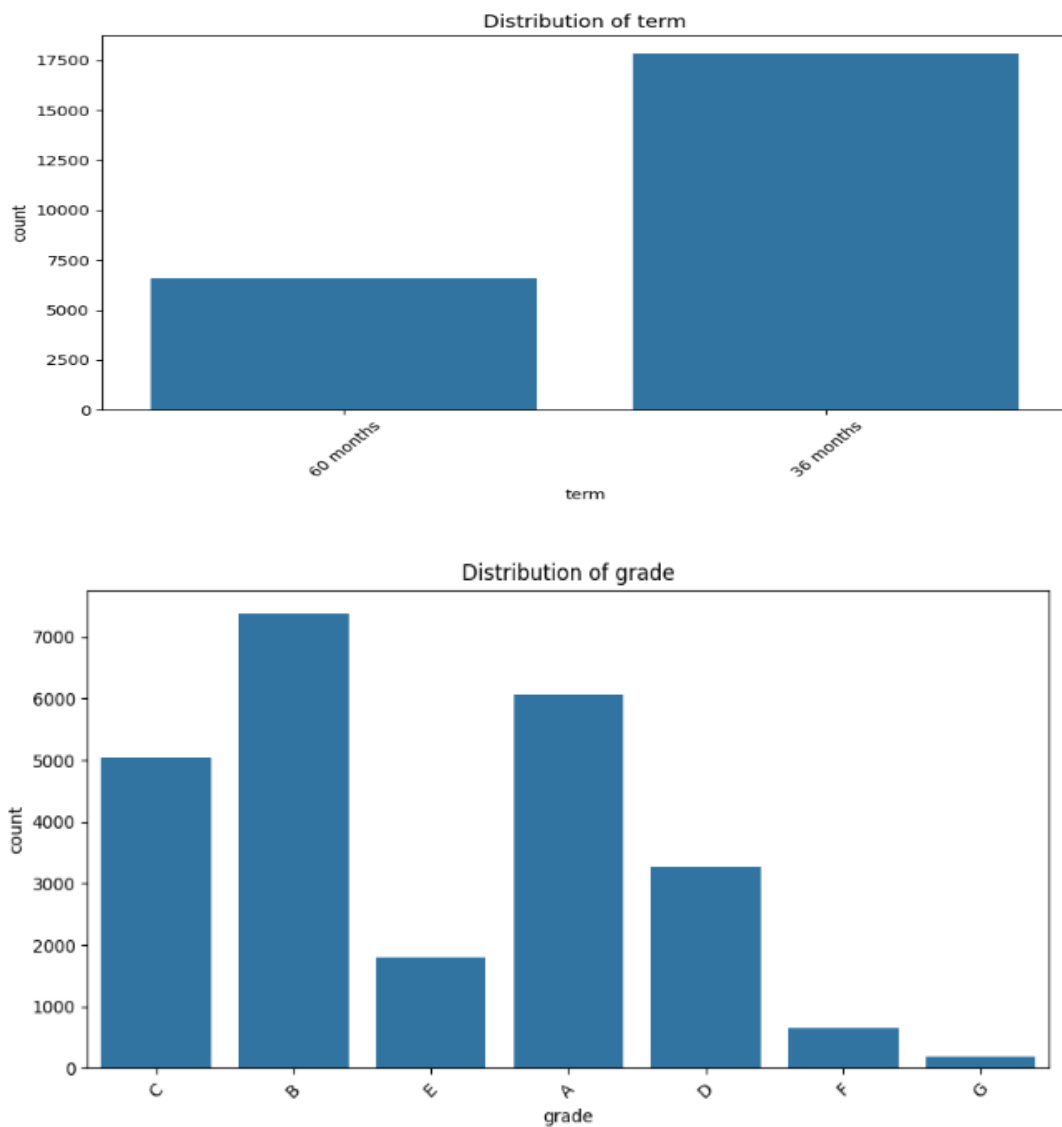
3.2.1 Categorical Variables vs Loan Status

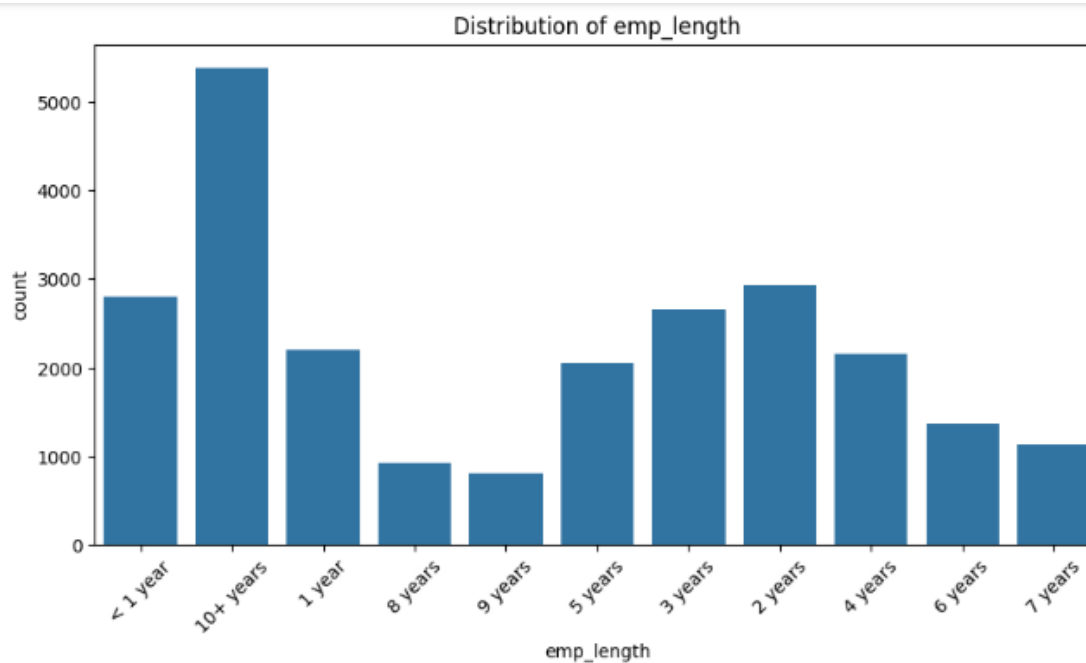
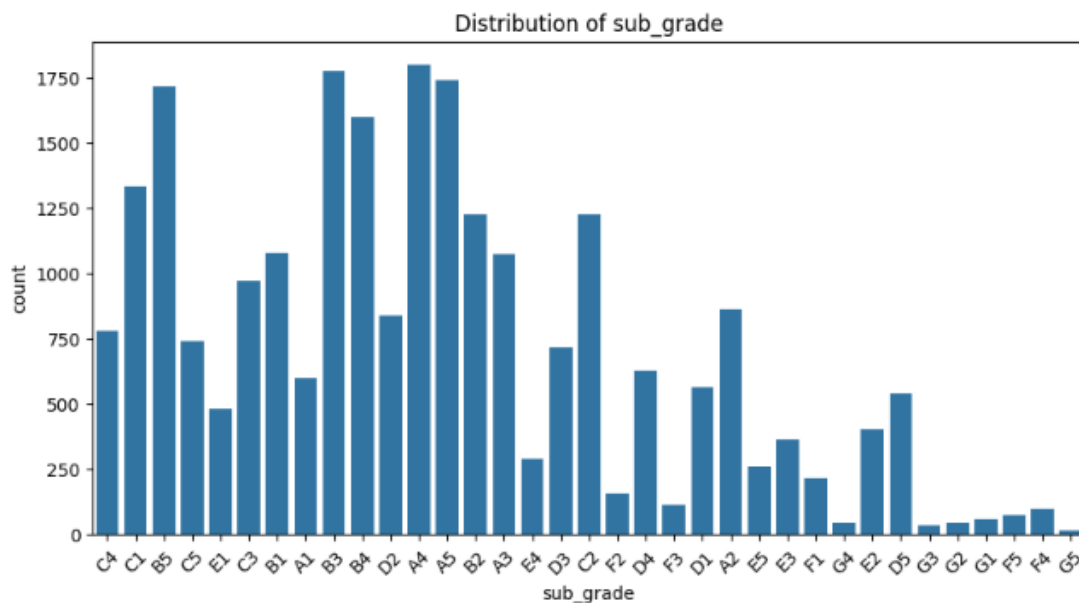
The relationship between categorical variables and loan status is analyzed using count plots with loan status as the hue.

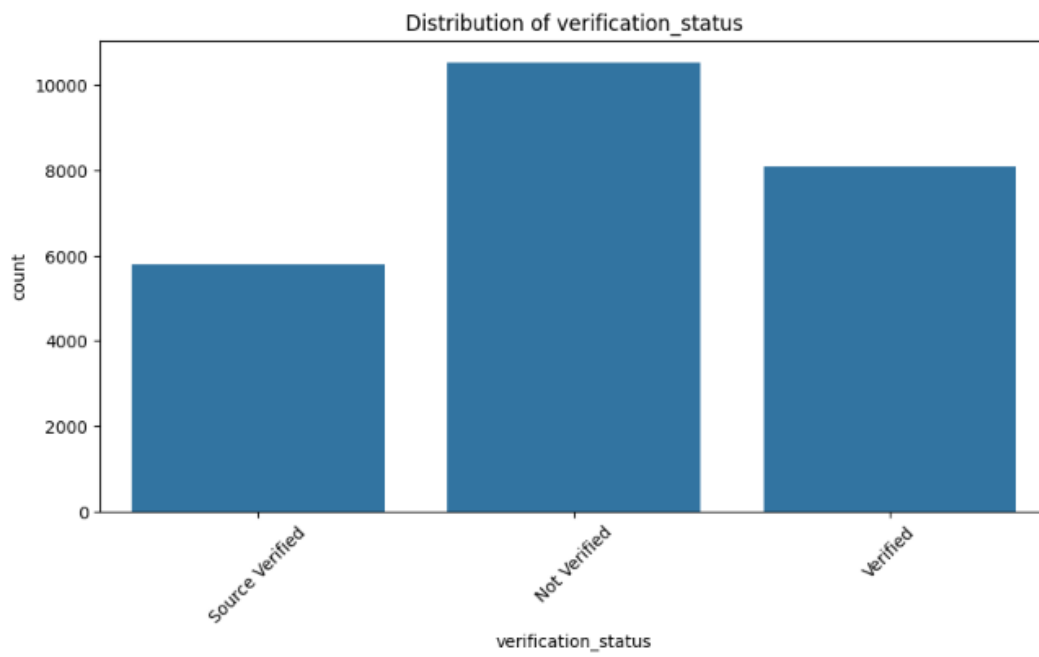
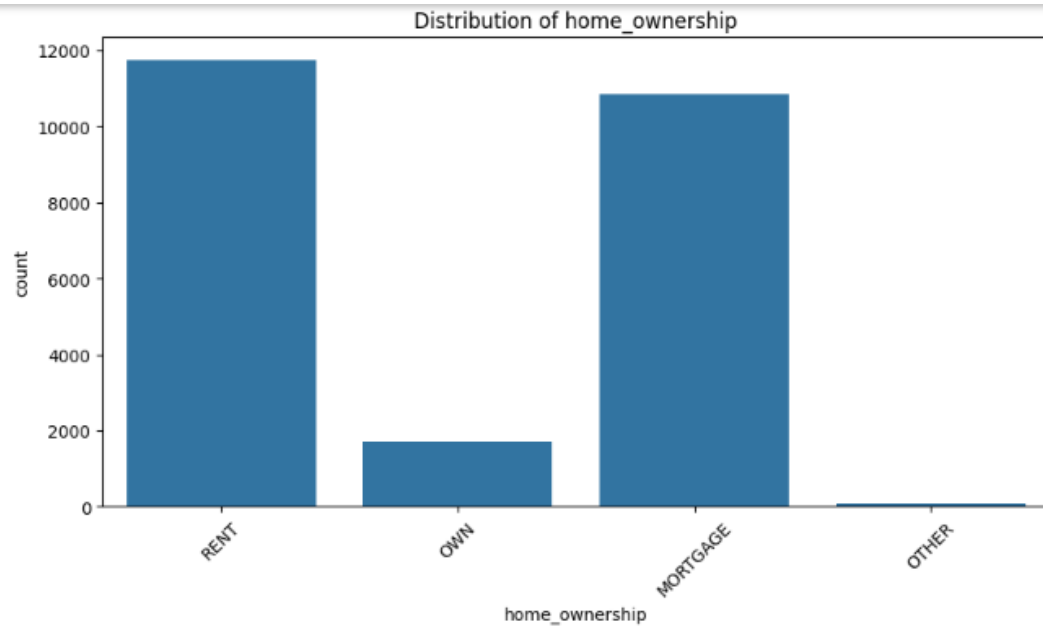
3.2.2 Numerical Variables vs Loan Status

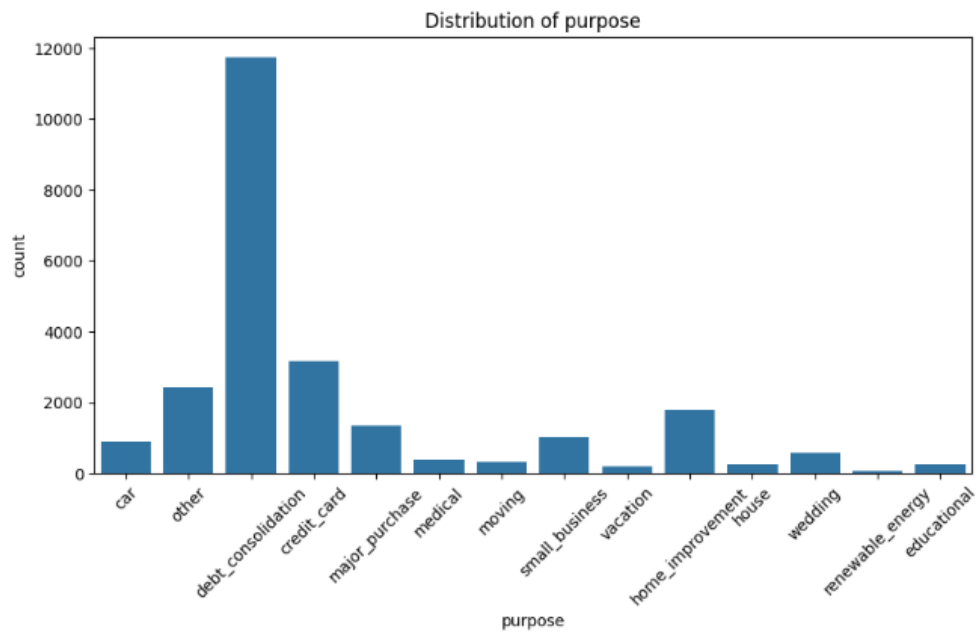
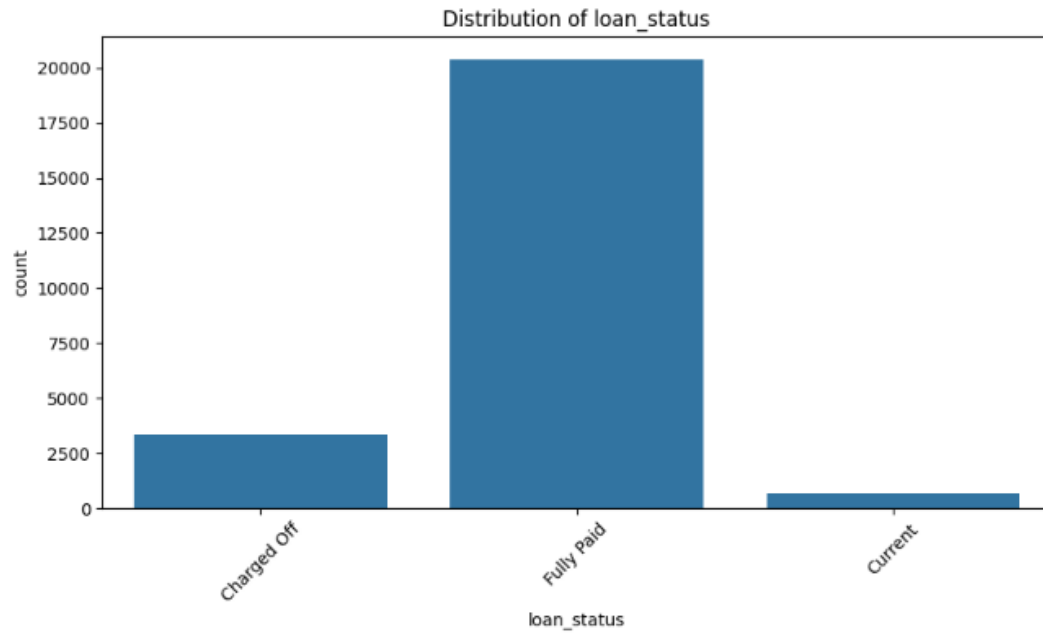
The relationship between numerical variables and loan status is analyzed using box plots.

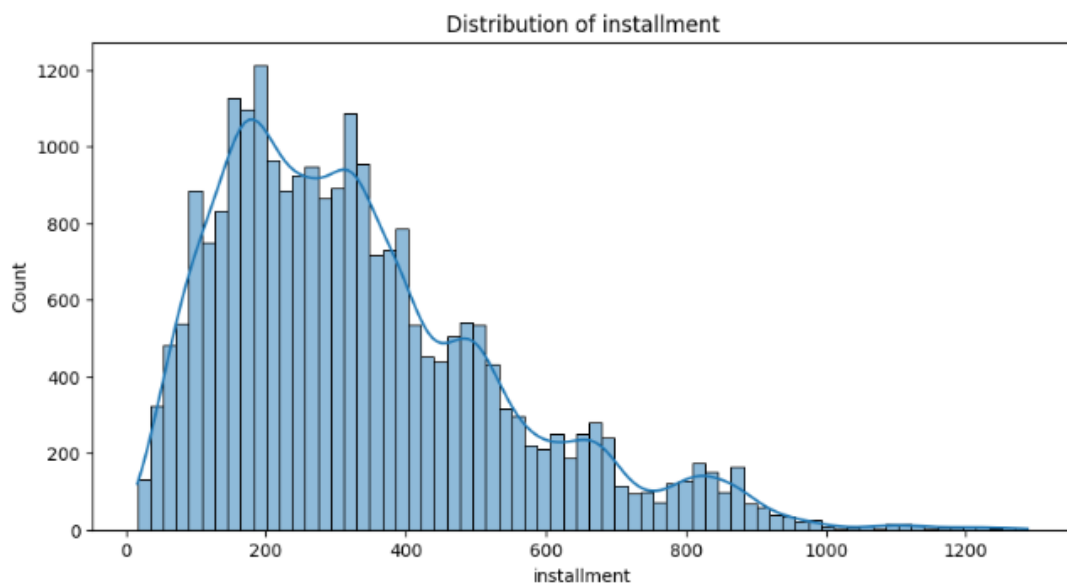
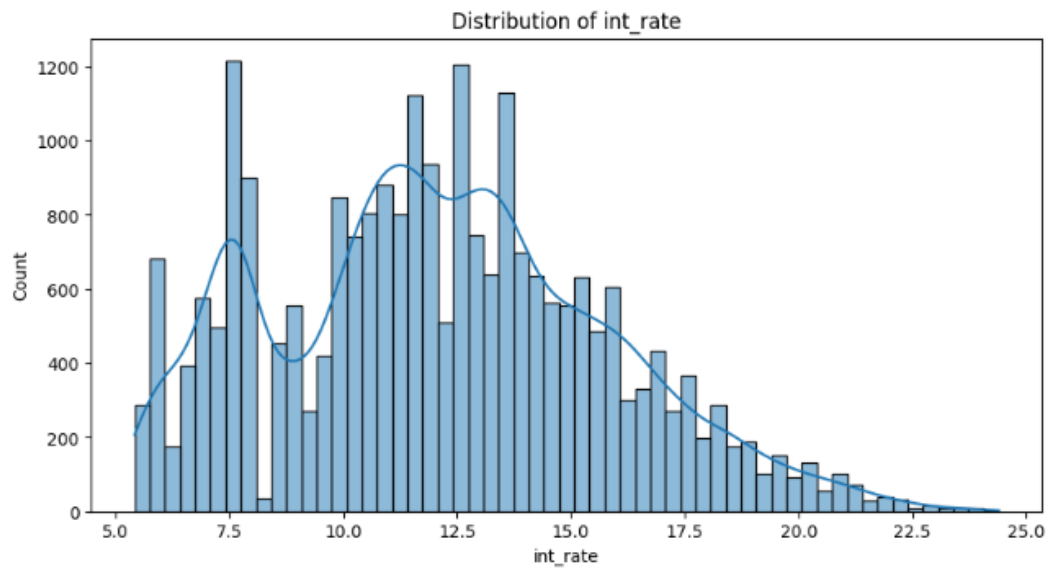
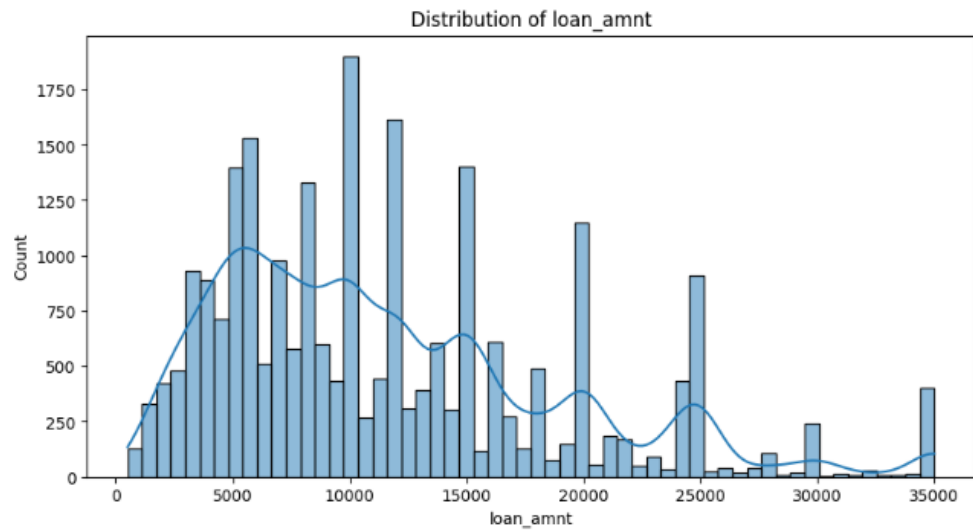
EDA Results:

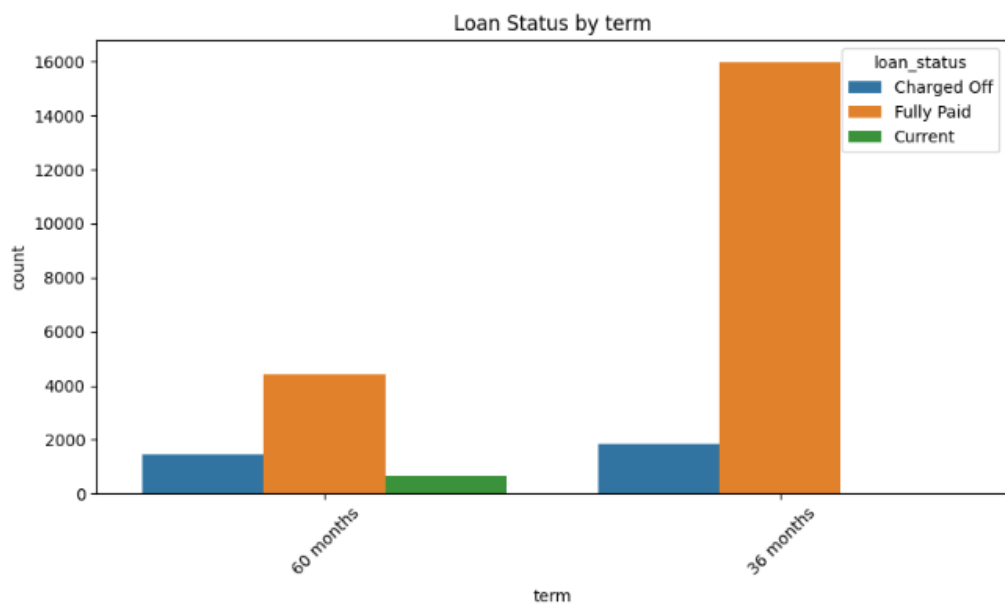
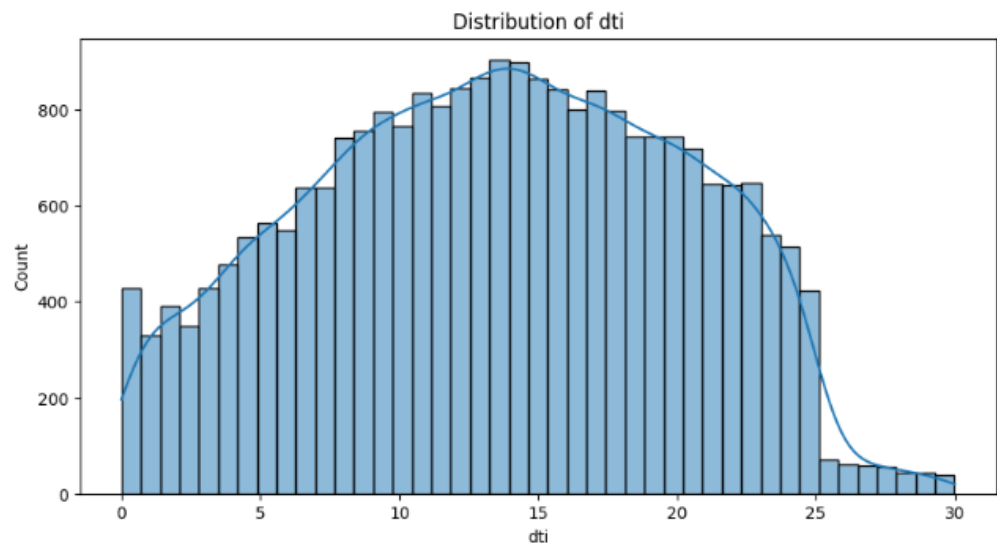
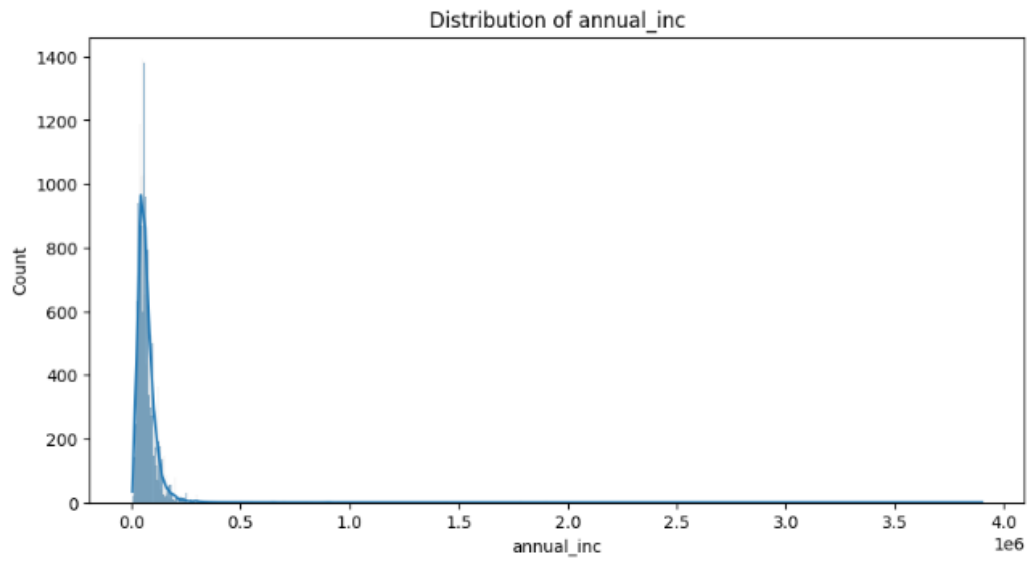


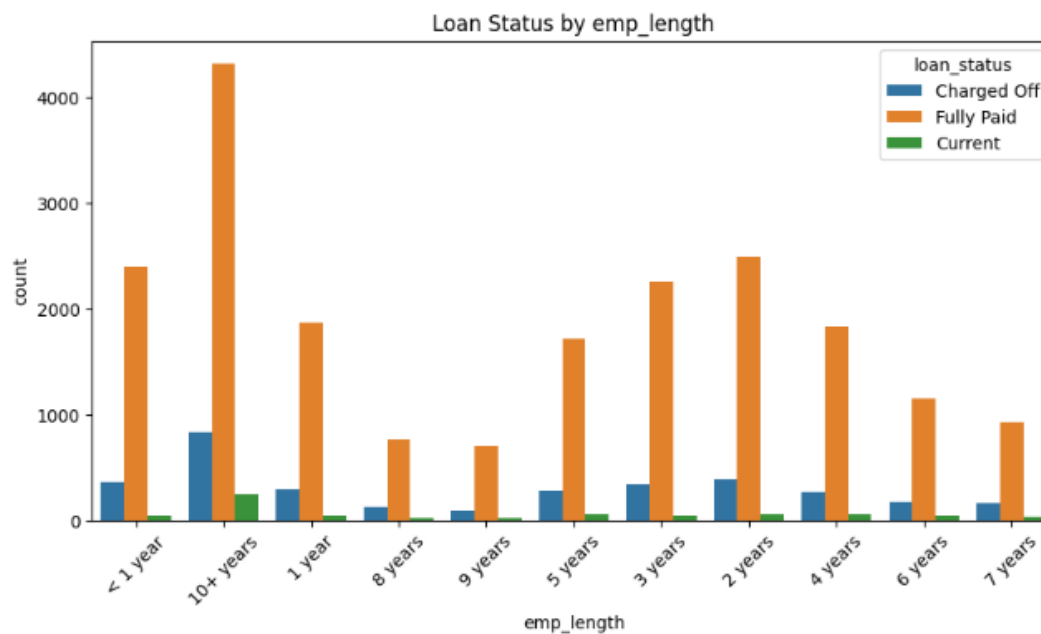
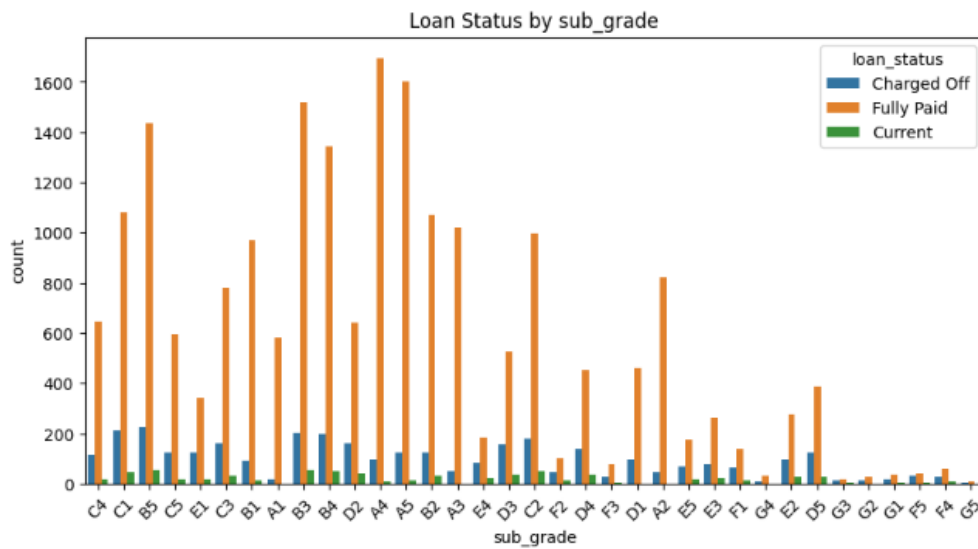
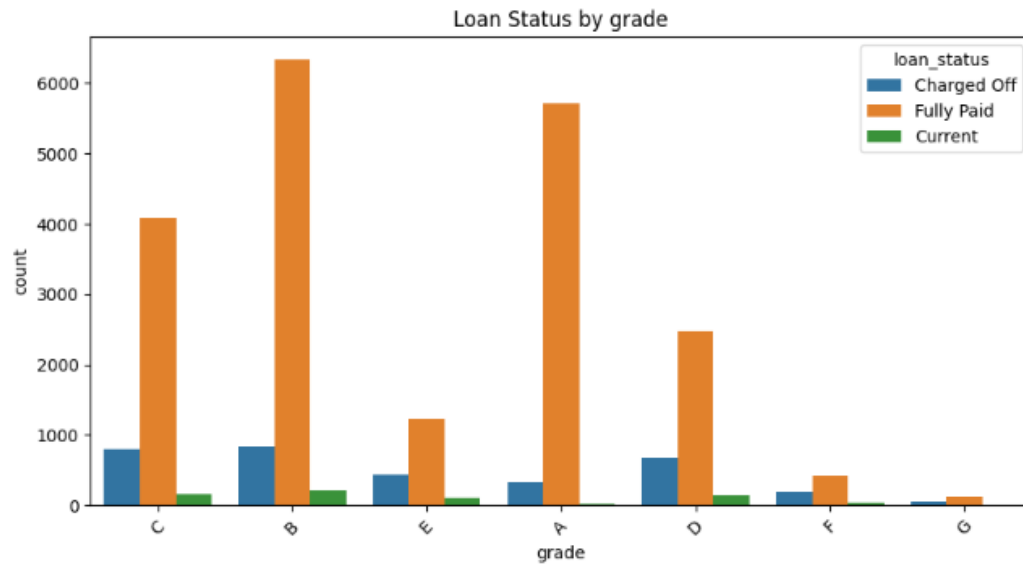


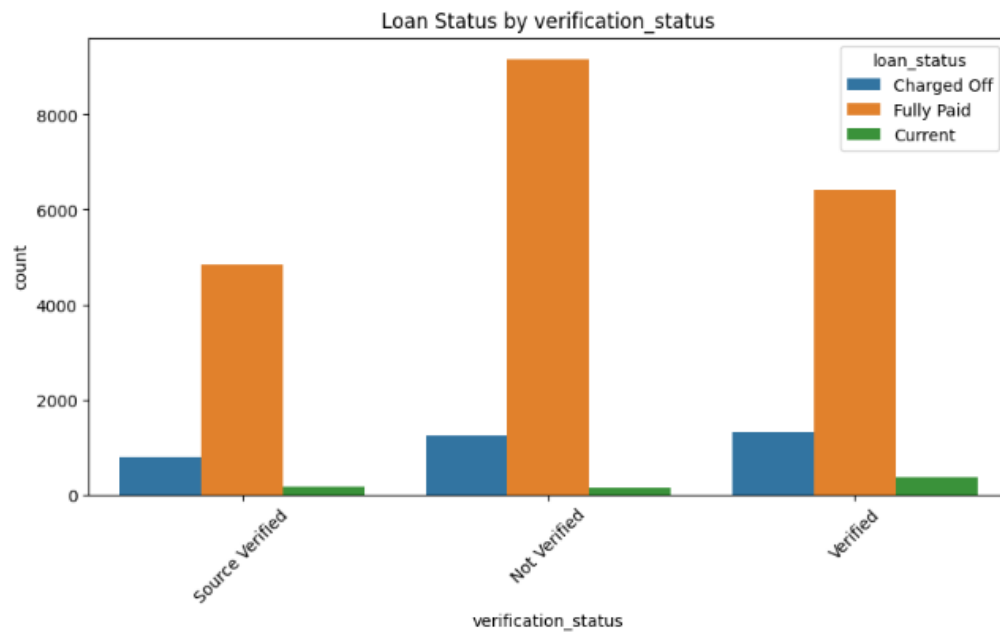
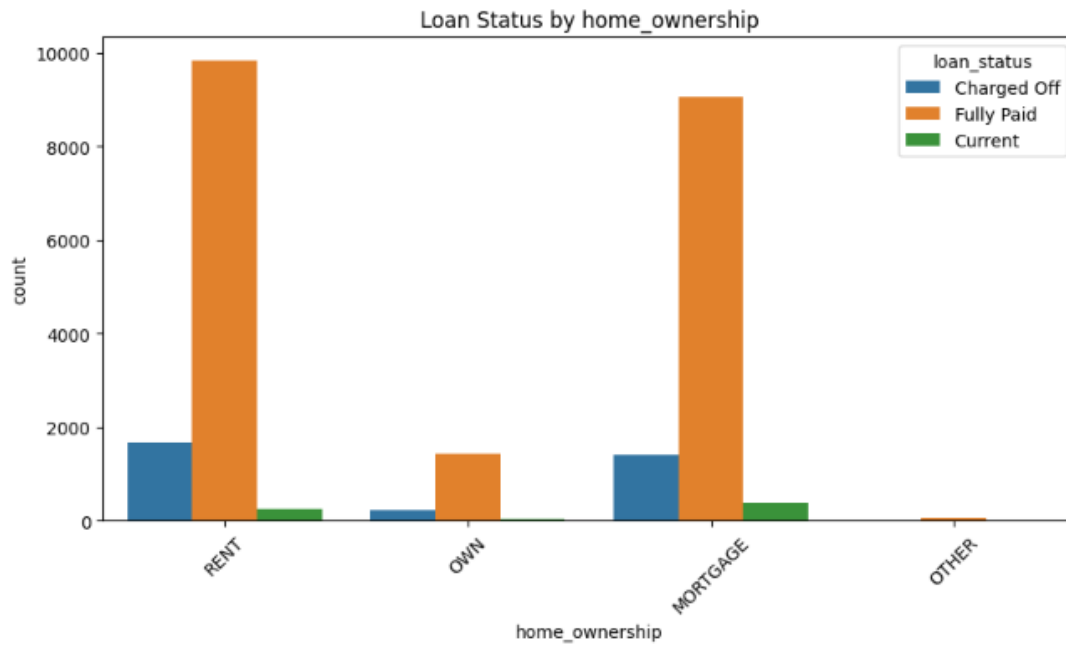


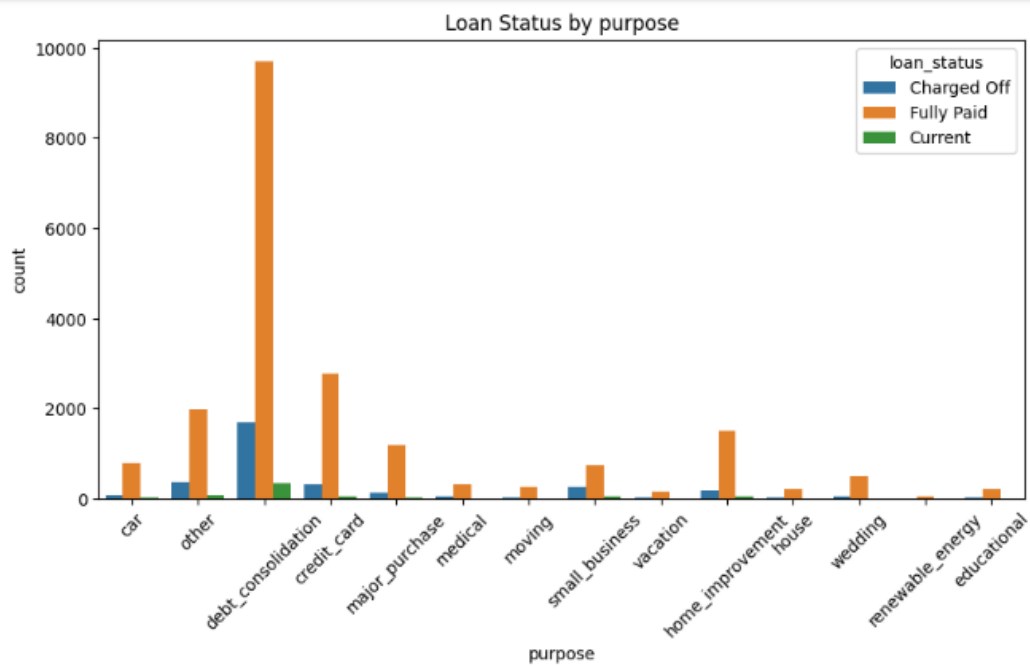
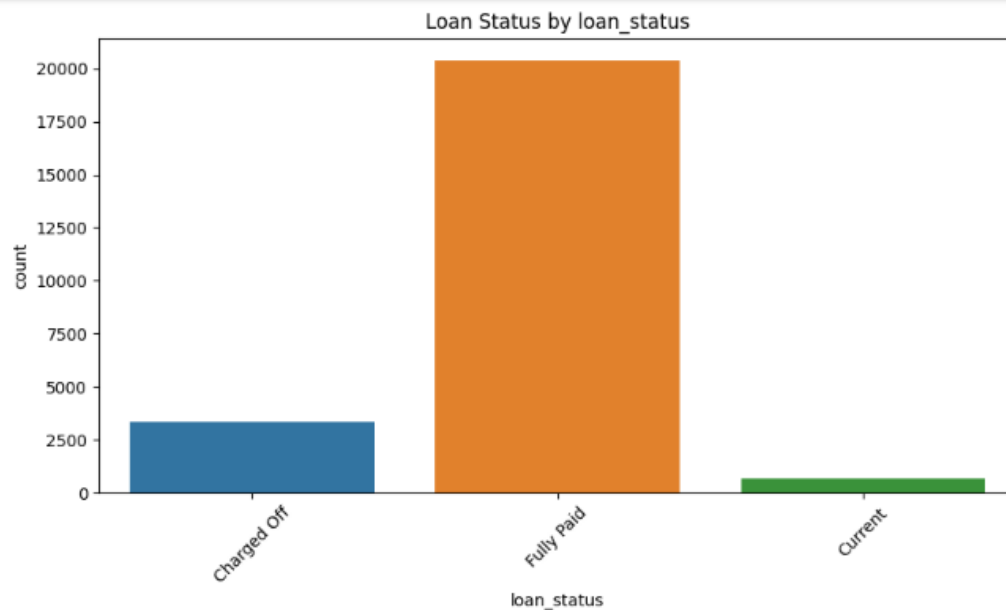


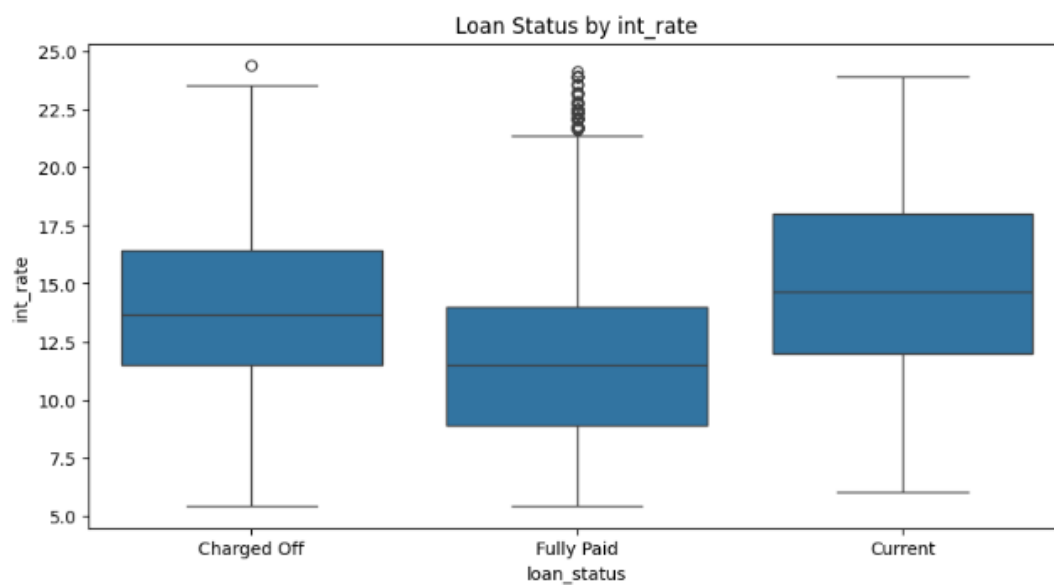
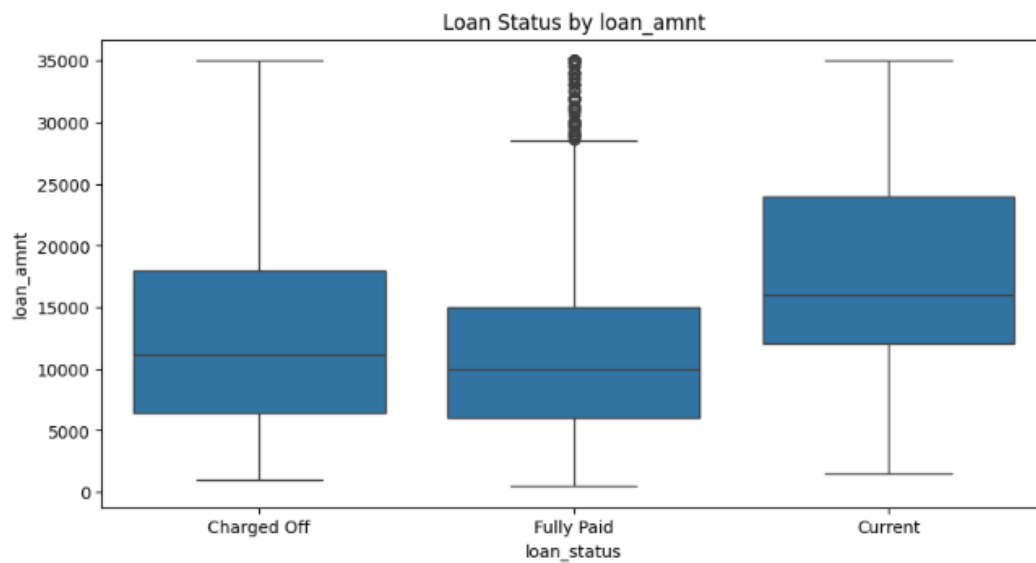


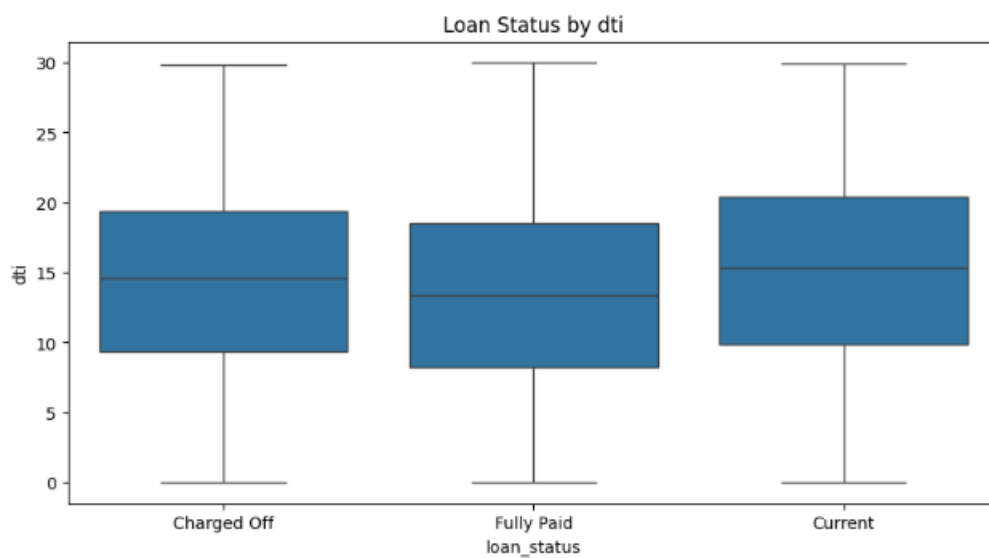
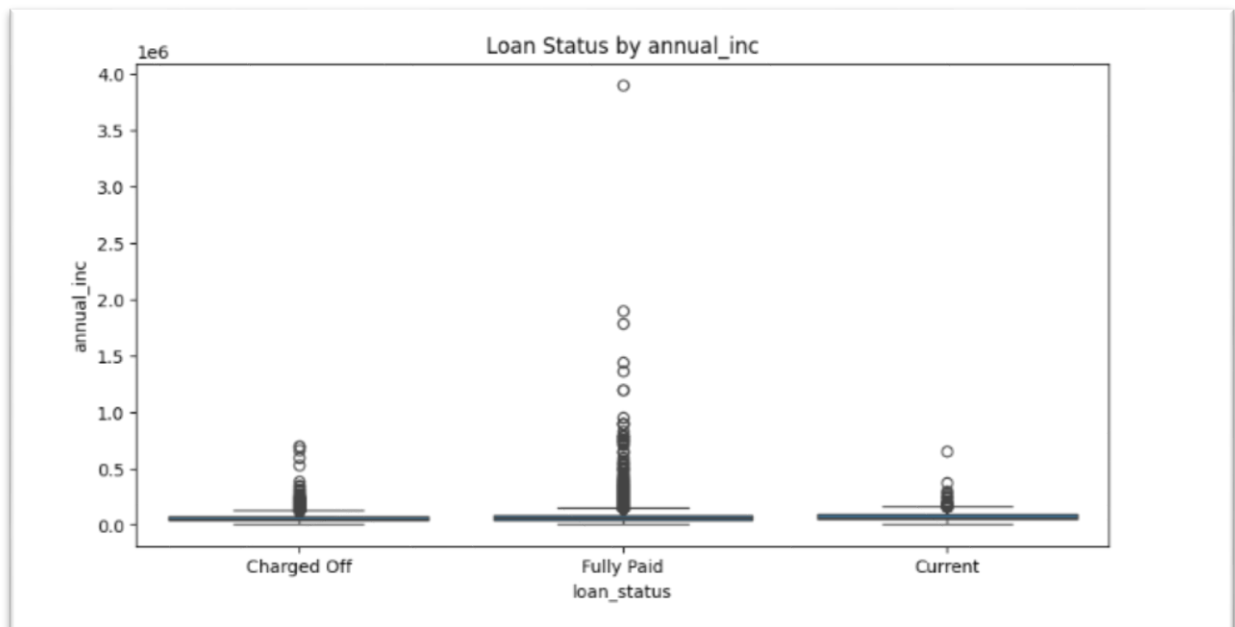
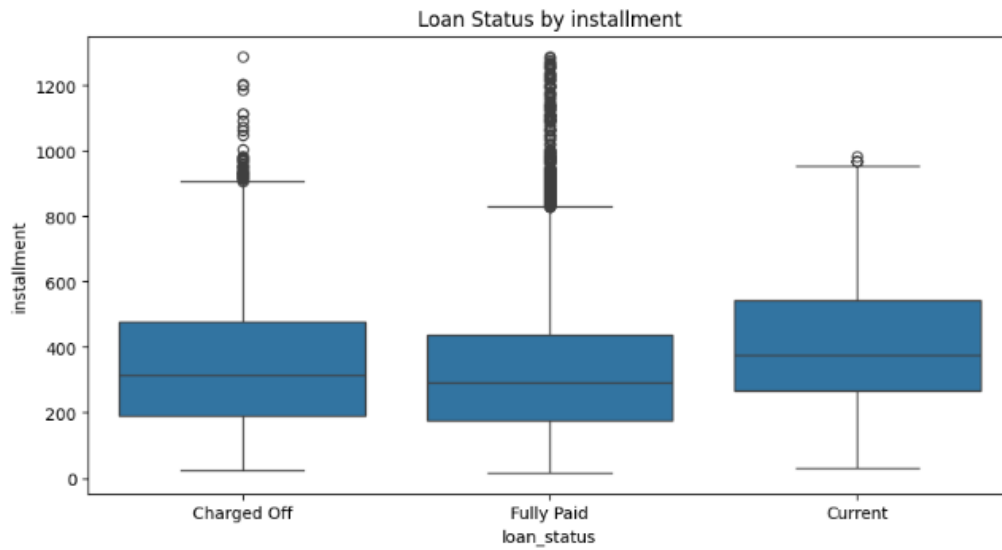












4. Model Development

4.1 Data Preparation

4.1.1 Encoding Categorical Variables

Categorical variables are encoded using Label Encoding. The target variable 'loan_status' is binary encoded where 'Charged Off' is 1 and others are 0.

4.1.2 Feature Selection

Features selected for the model include 'loan_amnt', 'int_rate', 'term', and 'grade'.

Justification: These features were chosen based on their potential impact on loan default prediction.

4.2 Train-Test Split

Split data into training (80%) and testing (20%) sets. Utilized logistic regression with `max_iter=1000` for training.

5. Model Evaluation

5.1 Predictions

The trained model is used to make predictions on the test set.

```
y_pred = model.predict(X_test)
```

5.2 Evaluation Metrics

The model's performance is evaluated using accuracy, confusion matrix, and classification Confusion Matrix.

6. Results

6.1 Accuracy

The model achieved an accuracy of {accuracy} on the test set.

Performance Metrics:

Accuracy: 85.78%

High accuracy but poor sensitivity to defaults suggests model imbalance.

Recommendations for improving model sensitivity and reducing false negatives.

6.2 Confusion Matrix

The confusion matrix provides a detailed breakdown of true positives, true negatives, false positives, and false negatives.

```
Precision (0): 86% - Predicted non-default correctly.  
Precision (1): 0% - Failed to predict any defaults correctly.  
Recall (0): 100% - Identified all non-default correctly.  
Recall (1): 0% - Missed all defaults.
```

6.3 Classification Report

The classification report includes precision, recall, and F1-score for each class.

	precision	recall	f1-score	support
0	0.86	1.00	0.92	4187
1	0.00	0.00	0.00	694
accuracy			0.86	4881
macro avg	0.43	0.50	0.46	4881
weighted avg	0.74	0.86	0.79	4881

7. Conclusion

The analysis revealed significant insights into the distribution and relationship of various features with loan status. The logistic regression model performed well in predicting loan status, with a reasonable accuracy.

8. References

- Dataset: [Link to the source of the dataset]
- Scikit-learn Documentation: <https://scikit-learn.org/>
- Seaborn Documentation: <https://seaborn.pydata.org/>
- Pandas Documentation: <https://pandas.pydata.org/>