

Lending Club Case Study

Overview

1.Data Understanding and Quality check

2.Data Cleaning

3.Data Analysis

- Uni Variate analysis results
- Bi Variate analysis results

4.Recommendations from analysis

5.Summary:

The Lending Club case study aims to analyze loan default patterns, identify factors influencing default, and identify customer and loan attributes to reduce defaults in future loan grants.

Overview

Loan Dataset has below customers:

1. Full Paid -> Completed Loan
2. Current -> Loan In Progress
3. Charged off -> Defaulted loan

In the above list we will only consider “Charged Off” customers for this analysis as the interest is mainly on defaulted loan attributes and customer attributes

Expected Analysis Outcomes

Factors influencing loan default means as part of this analysis identify variables which are strong indicators of default loan

Data Understanding and Quality check

- Checking the loan dataset shape shows 39717 rows and 111 columns
<class 'pandas.core.frame.DataFrame'> RangeIndex: 39717 entries, 0 to 39716 Columns: 111 entries, id to total_il_high_credit_limit dtypes: float64(74), int64(13), object(24) memory usage: 33.6+ MB
- Checking for null values shows there are many columns with more null/na/NaN values
id 0 member_id 0 loan_amnt 0 funded_amnt 0 funded_amnt_inv 0 ... tax_liens 39 tot_hi_cred_lim 39717 total_bal_ex_mort 39717 total_bc_limit 39717 total_il_high_credit_limit 39717 Length: 111, dtype: int64
- Checking for these loan attributes in given Data Dictionary shows these are not very much important attributes needed for our analysis and so cleaning these will be better

DATA CLEANING

Step1:

Dropping all rows with missing values – Couldn't find any rows with complete missing values and so nothing to be dropped

Step2:

Dropping columns with at least 1 or max null values -> after this step the loan dataset shape becomes 39717 rows and 43 columns

Step3:

Making sure no null values present after these cleaning

Step4:

Listing the columns remaining in the loan dataset after cleaning

```
Index(['id', 'member_id', 'loan_amnt', 'funded_amnt', 'funded_amnt_inv', 'term', 'int_rate', 'installment', 'grade', 'sub_grade',  
'home_ownership', 'annual_inc', 'verification_status', 'issue_d', 'loan_status', 'pymnt_plan', 'url', 'purpose', 'zip_code',  
'addr_state', 'dti', 'delinq_2yrs', 'earliest_cr_line', 'inq_last_6mths', 'open_acc', 'pub_rec', 'revol_bal', 'total_acc',  
'initial_list_status', 'out_prncp', 'out_prncp_inv', 'total_pymnt', 'total_pymnt_inv', 'total_rec_prncp', 'total_rec_int',  
'total_rec_late_fee', 'recoveries', 'collection_recovery_fee', 'last_pymnt_amnt', 'policy_code', 'application_type',  
'acc_now_delinq', 'delinq_amnt'], dtype='object')
```

Step5:

Filtering for only defaulted loan customers (with loan_status as "Charged Off") resulted in dataset shape of

5627 rows and 43 columns

DATA ANALYSIS – UNI VARIATE ANALYSIS

Considering below 8 variables for univariate analysis from this defaulted loan dataset

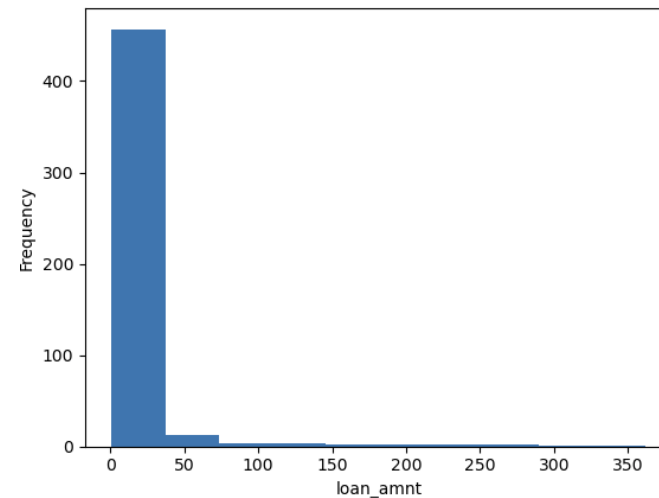
1. loan_amnt
2. term
3. int_rate
4. grade
5. home_ownership
6. annual_inc
7. verification_status
8. purpose

Variable 1: Loan_amnt

Analysis Outcome: Defaulted

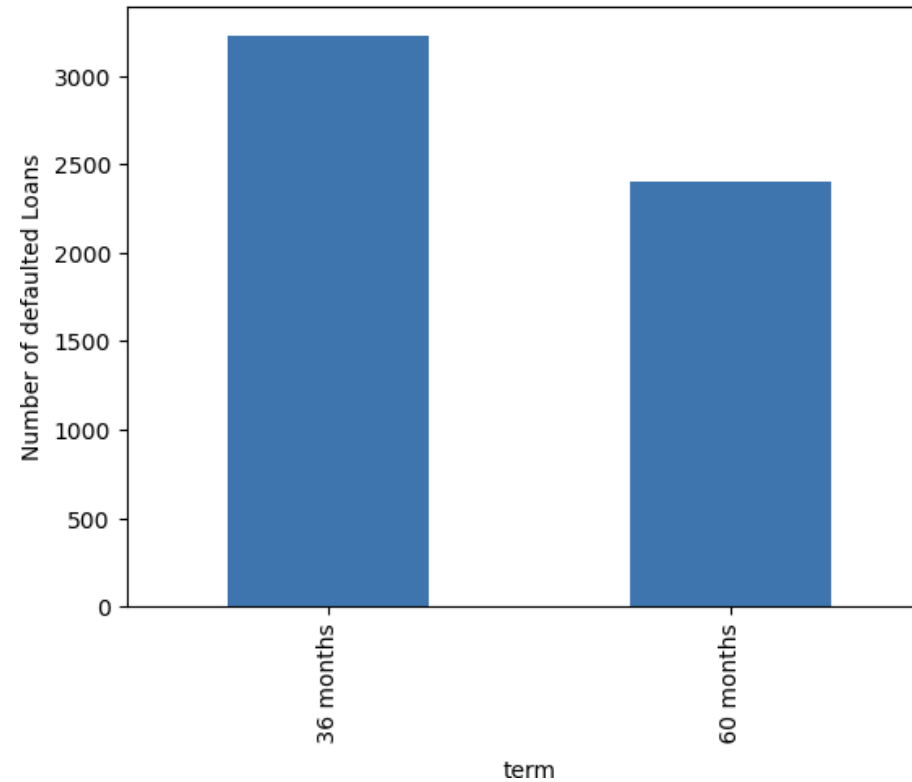
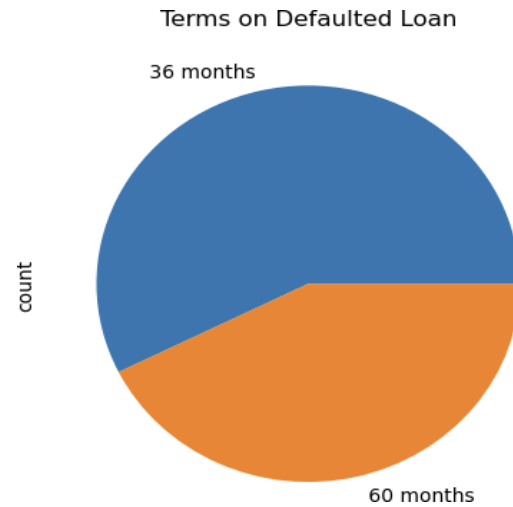
loan customers took Loan

amount < 50,000rs



DATA ANALYSIS – UNI VARIATE ANALYSIS

Variable 2: Term

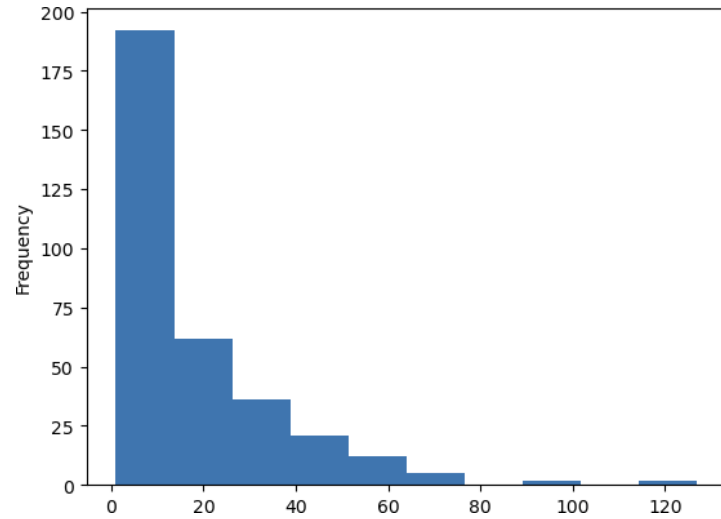


Analysis Outcome:

Loans with tenure of 36 months defaulted more than the loans with tenure of 60 months might be due to HIGH EMI per month due to less tenure

DATA ANALYSIS – UNI VARIATE ANALYSIS

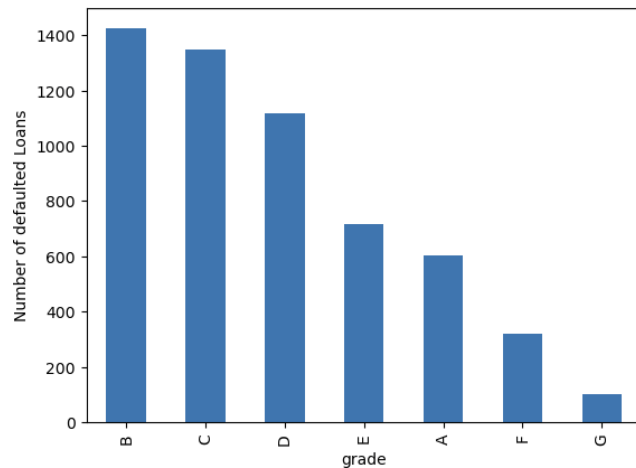
Variable 3: int_rate



Analysis Outcome:

Interest rates for most of the defaulted loan customers are less than 20%

Variable 4: grade

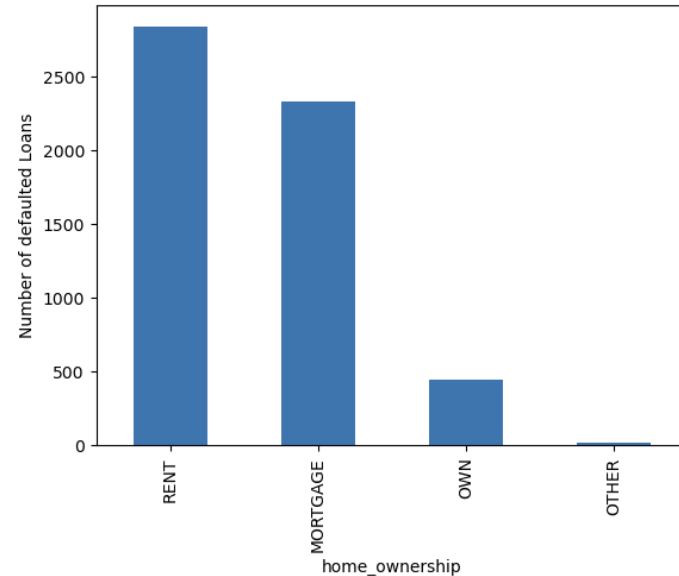


Analysis Outcome:

Grades - B, C, D are the top 3 loan grades which got defaulted

DATA ANALYSIS – UNI VARIATE ANALYSIS

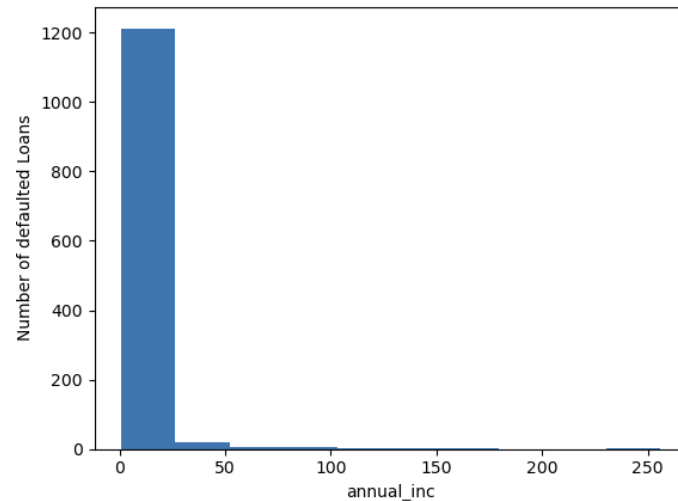
Variable 5: home_ownership



Analysis Outcome:

Customers who are either in RENTED or MORTGAGED Homes defaulted loans more times

Variable 6: annual_inc



Analysis Outcome:

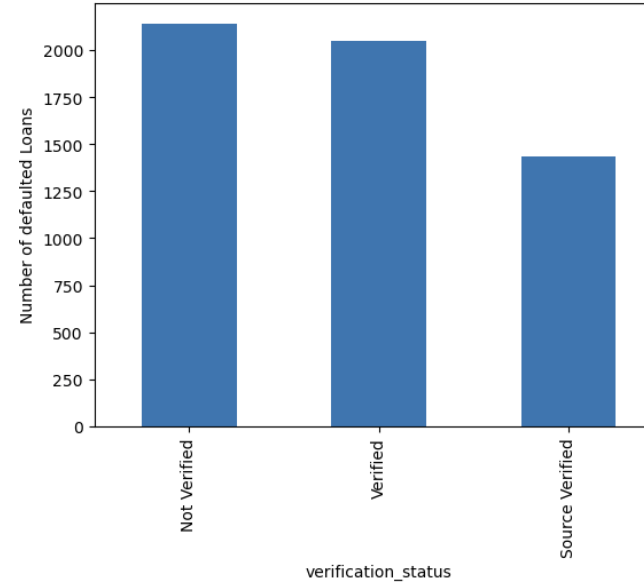
Customers who defaulted loans have annual income < 50k

DATA ANALYSIS – UNI VARIATE ANALYSIS

Variable 7: Verification Status

Analysis Outcome:

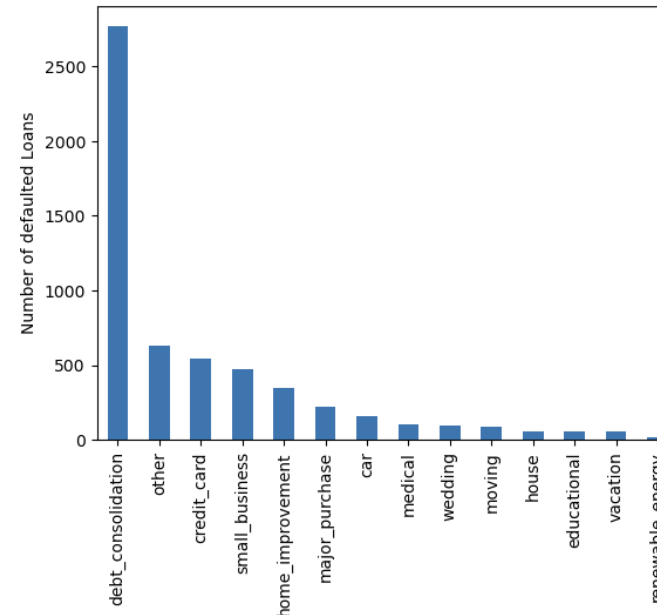
Sum of Verified and Source Verified > Not Verified which means more verified loans defaulted than not verified. This looks weird and need to find the root cause for the same, it could be lapse in the process of verification



Variable 8: Purpose

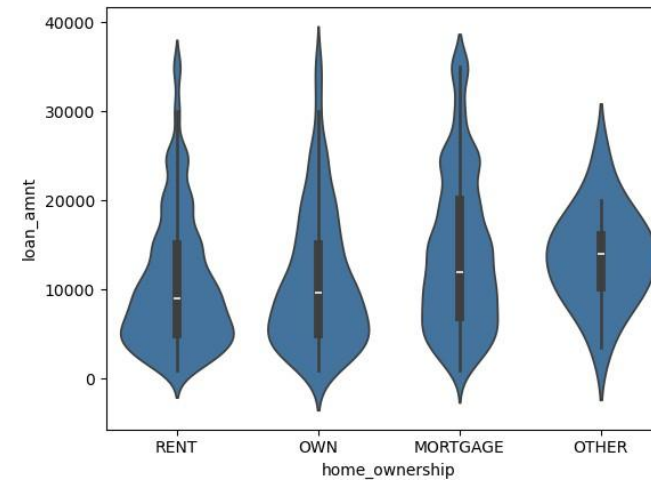
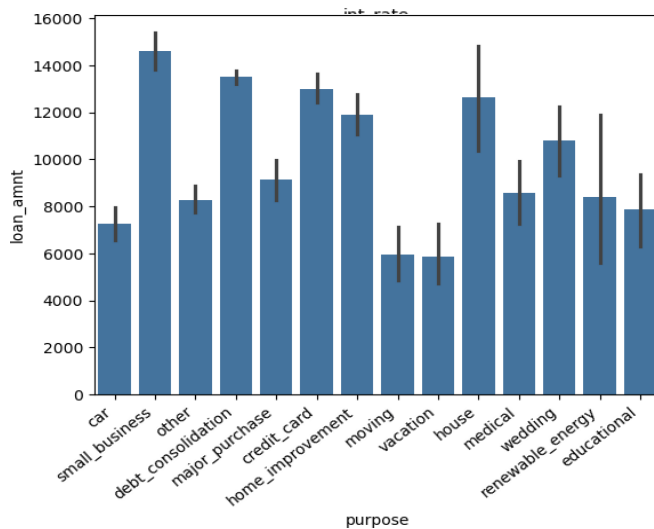
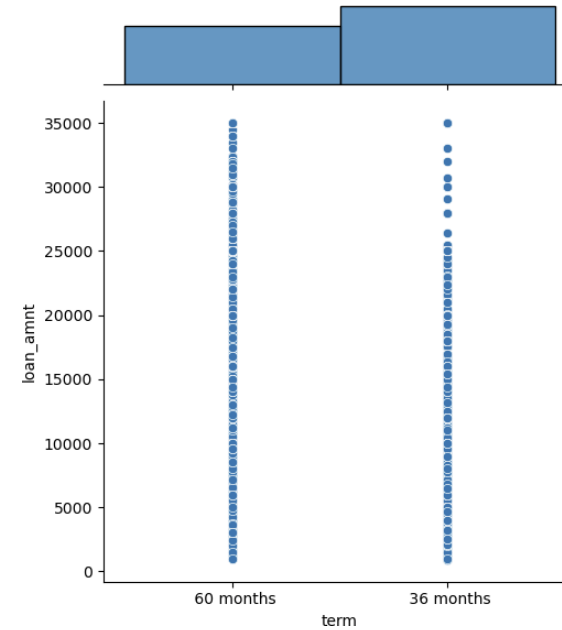
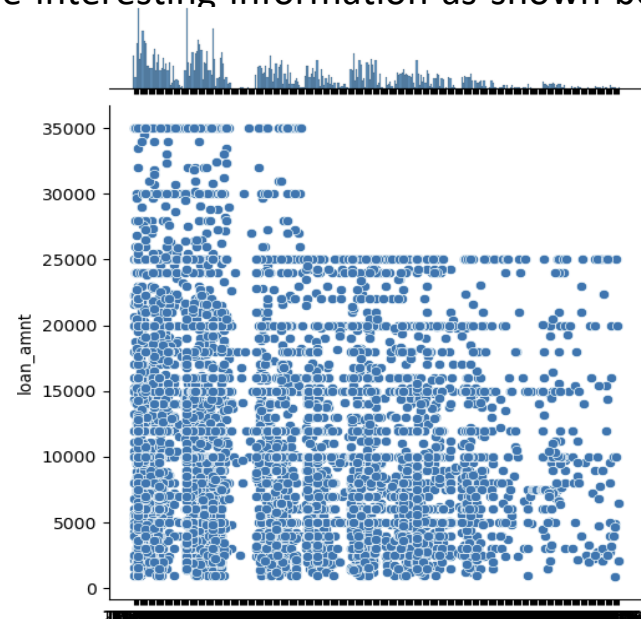
Analysis Outcome:

Customer who took loan for the purpose of debt consolidation defaulted more than any other purpose



DATA ANALYSIS – BI VARIATE ANALYSIS

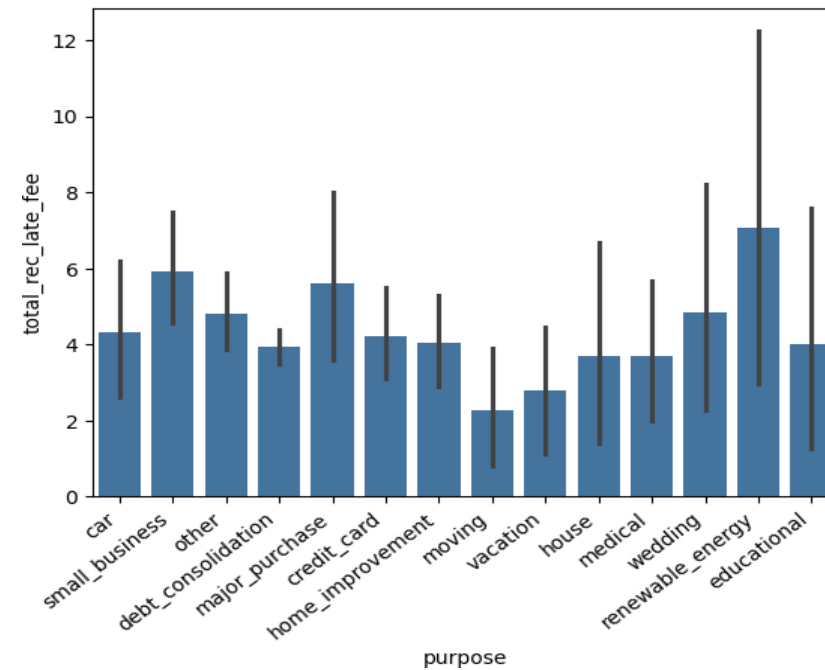
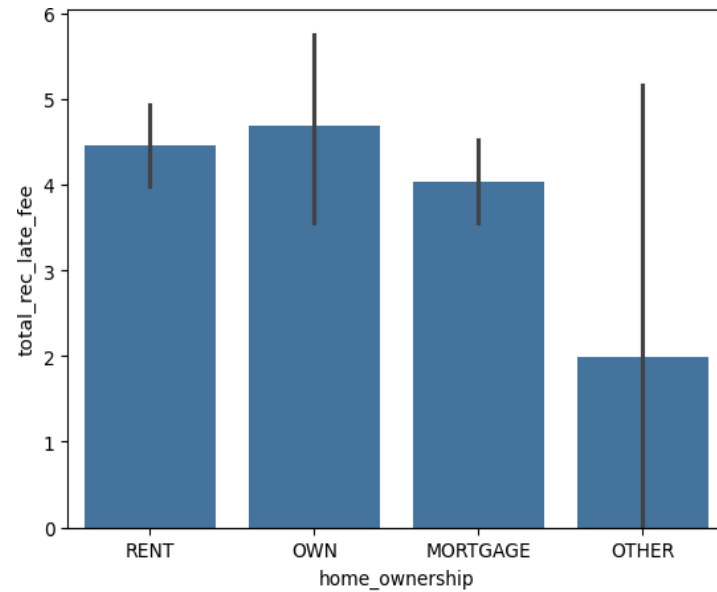
Comparison of Loan amnt with different other variables like term, home_ownership, int_rate, purpose shows some interesting information as shown below



DATA ANALYSIS – BI VARIATE ANALYSIS

Other few Bi Variate analysis results shown below

1. Homeowner ship vs Total recovery late fee
2. Purpose vs Total recovery late fee



Recommendations based on Analysis

Some Interesting outcomes of this analysis are

Default Risk Factors in Loans

- Loans with less than 36 months tenure are more likely to default due to high monthly EMI.
- Customers who have rented or mortgaged their homes also tend to default.
- Annual income less than 50k increases the risk of default.
- Verified customers default more often than non-verified ones, suggesting potential lapses or corruption in verification processes.
- Customers consolidating multiple loans may default heavily, indicating the need for cautious lending practices.

SUMMARY

Future Loan Grant Recommendations

- Consideration of recommendations for loan verification process.
- Further investigation needed to reduce loan defaults.

Team Members:

1. PALANI ELLAPPAN
2. PABAN DAS

Please refer below GitHub Link:

<https://github.com/palaniellappank/LendingClubCaseStudy>