# Loan Dataset Case-study Report

- **Rahul B**
- **Alhad Parashtekar**

# Aim

- **The aim of the study is to analyze the risks associated with the Bank's decision for loan approval when it receives a loan application using borrower information. The analysis will help us understand the underlying pattern connecting 'loan_status' column in the database and rest of the columns.**

- **The aim of exploratory analysis is to figure out columns which are most influential or highly correlated with the 'loan_status' and gain additional insights.**

*Each slide of the presentation would contain a figure showing impact of each variable present in the dataset on the loan status.*

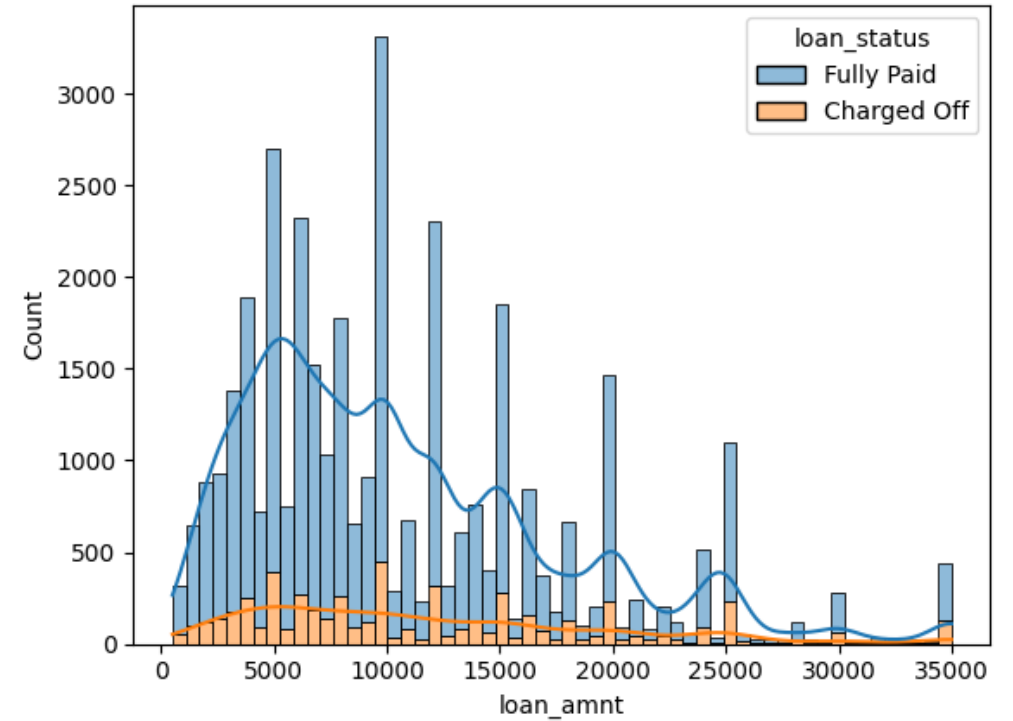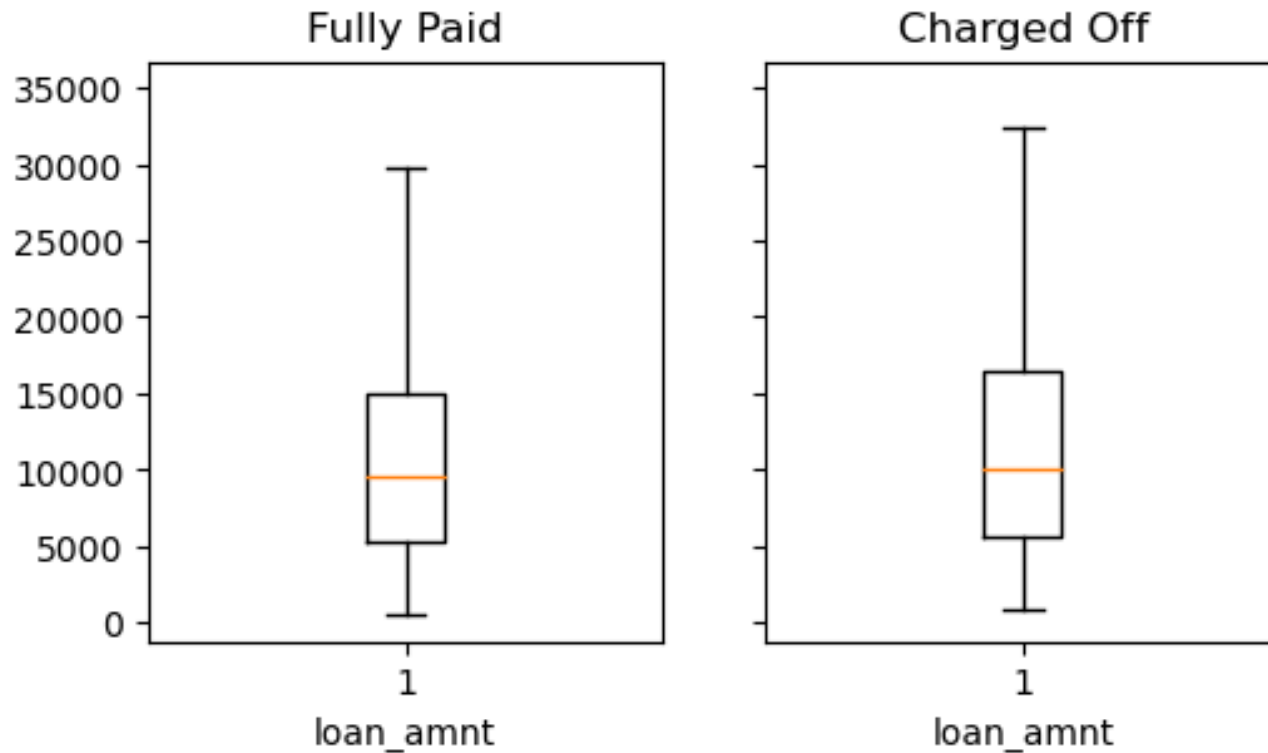# *Columns from the database considered for exploratory analysis*

- *The database contains 111 columns, out of which 54 columns contains only NaN values, and 2 columns contain only '0' as values. These 56 columns are not considered for the analysis (EDA).*

- *The 'member_id' and 'id' are not necessary to understand the correlation, hence not considered in EDA.*

- Rest of the columns are selected for analysis:
  *'loan_amnt', 'funded_amnt', 'int_rate', 'installment', 'grade', 'sub_grade', 'addr_state', 'all_util', 'annual_inc', 'chargeoff_within_12_mths', 'delinq_2yrs', 'funded_amnt_inv', 'loan_status', 'mort_acc', 'num_actv_bc_tl', 'num_actv_rev_tl', 'num_il_tl', 'num_sats', 'pymnt_plan', 'revol_bal', 'revol_util', 'tax_liens', 'tot_coll_amt', 'zip_code', 'term', 'purpose'*
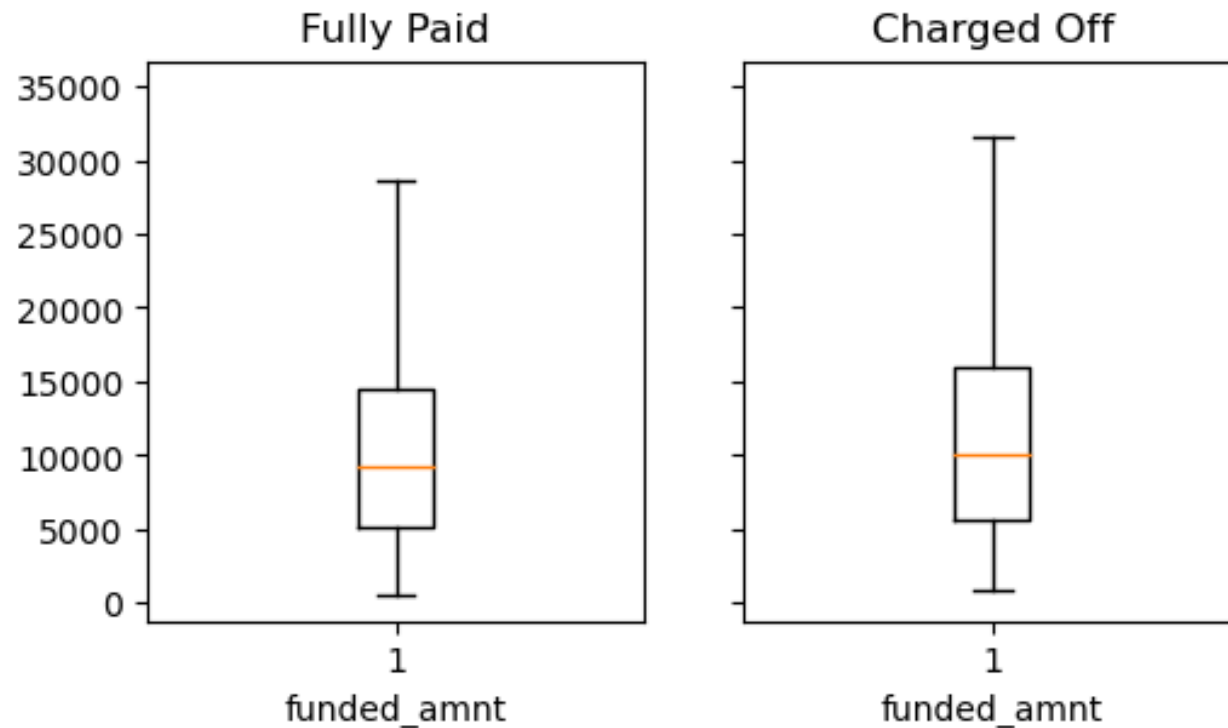
# Univariate Analysis

# Impact of *Loan Amount*

Though the median and other quartiles are higher in charged-off loans, no significant difference in distribution can be observed
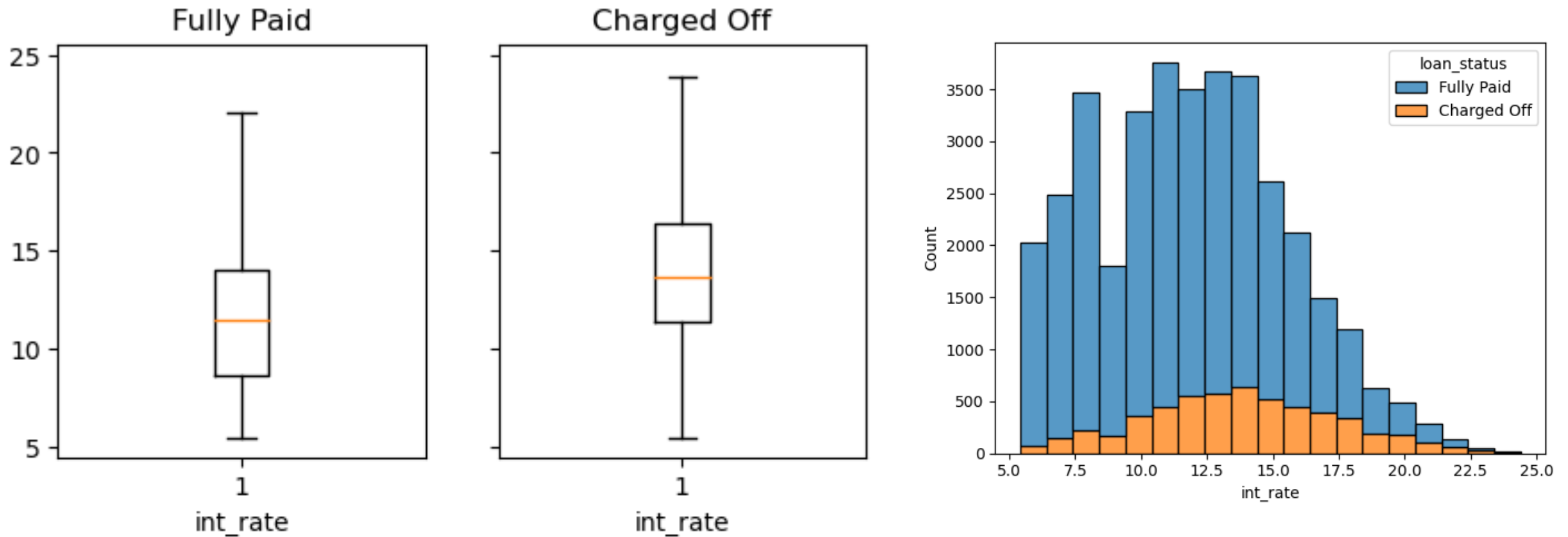
# Impact of *Funded Amount*

**Though the spread of 'funded_amnt' is slightly larger for 'charged_off', no significant difference in distribution can be observed.**
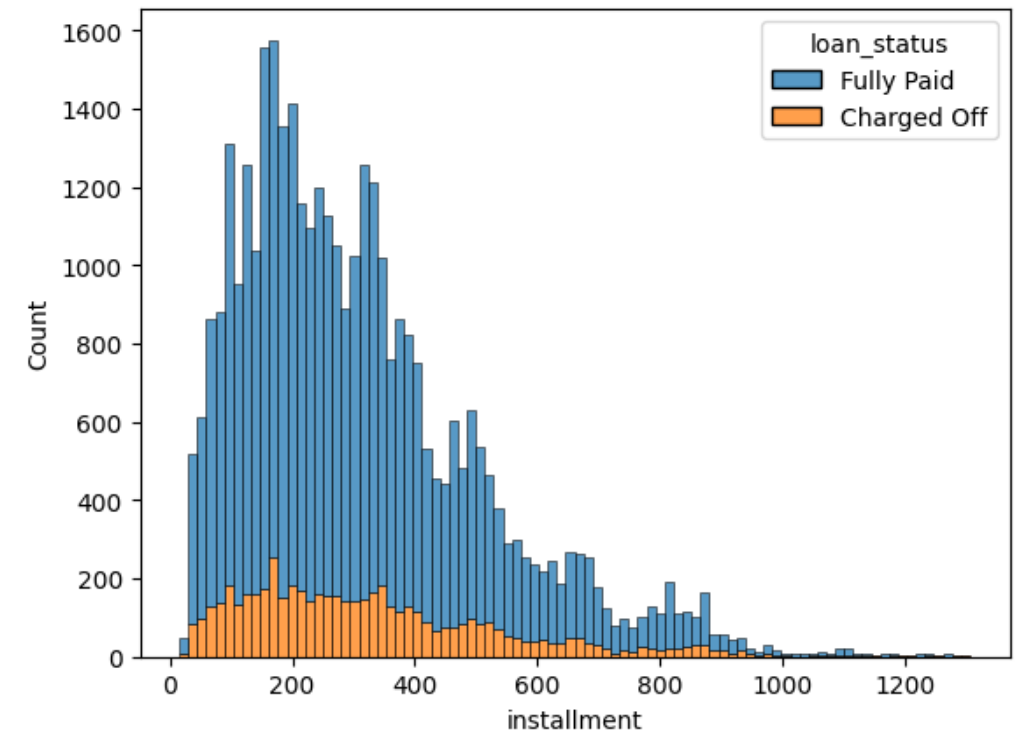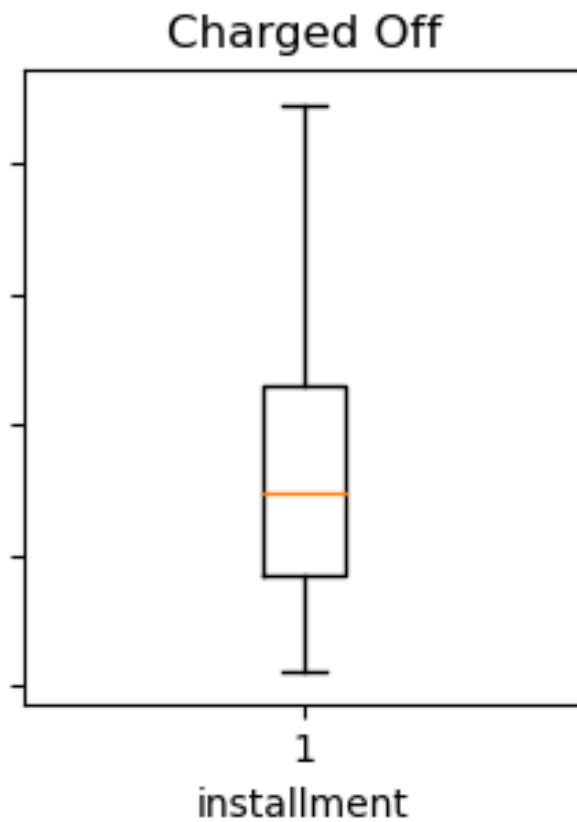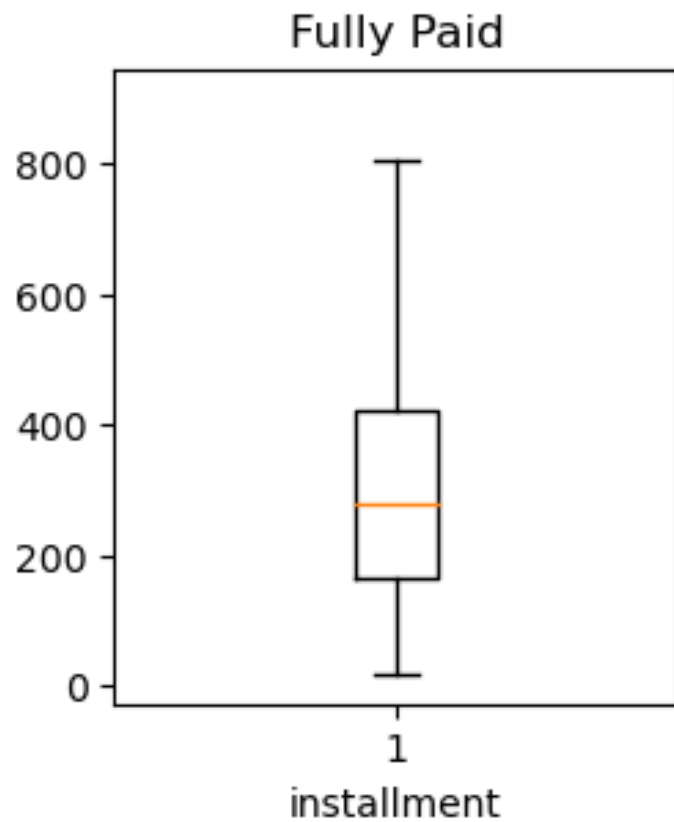
# Impact of *Interest Rate*

**The median and other quartiles (25 and 75) are considerably higher in case of 'charged_off'**

**loans than 'fully_paid' loans.** **This variable seems to influence the 'loan_status'.**
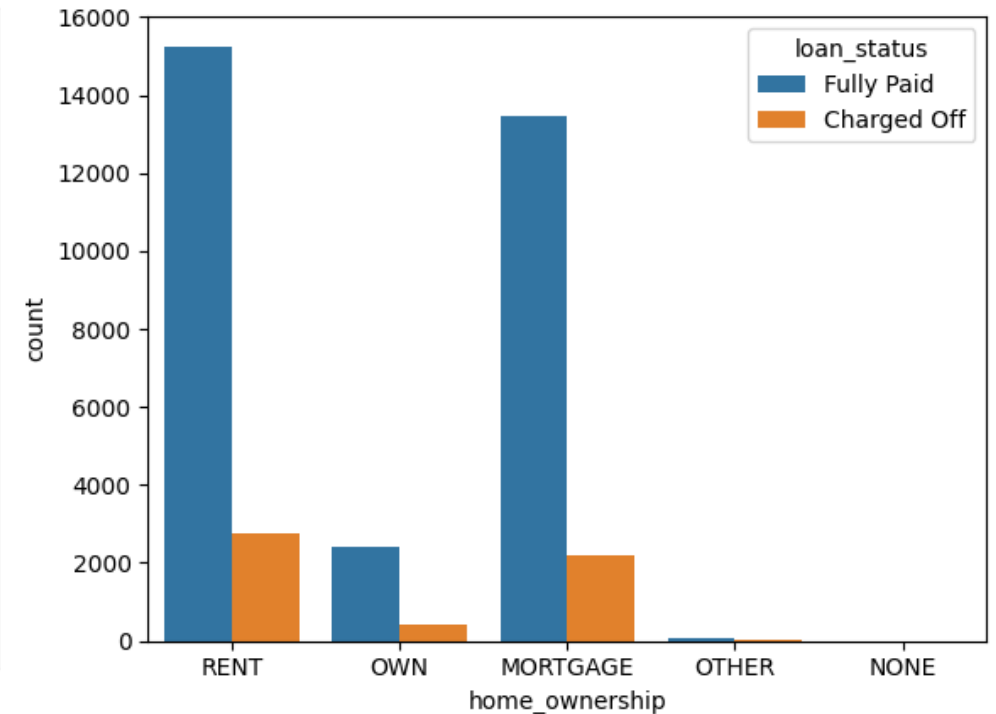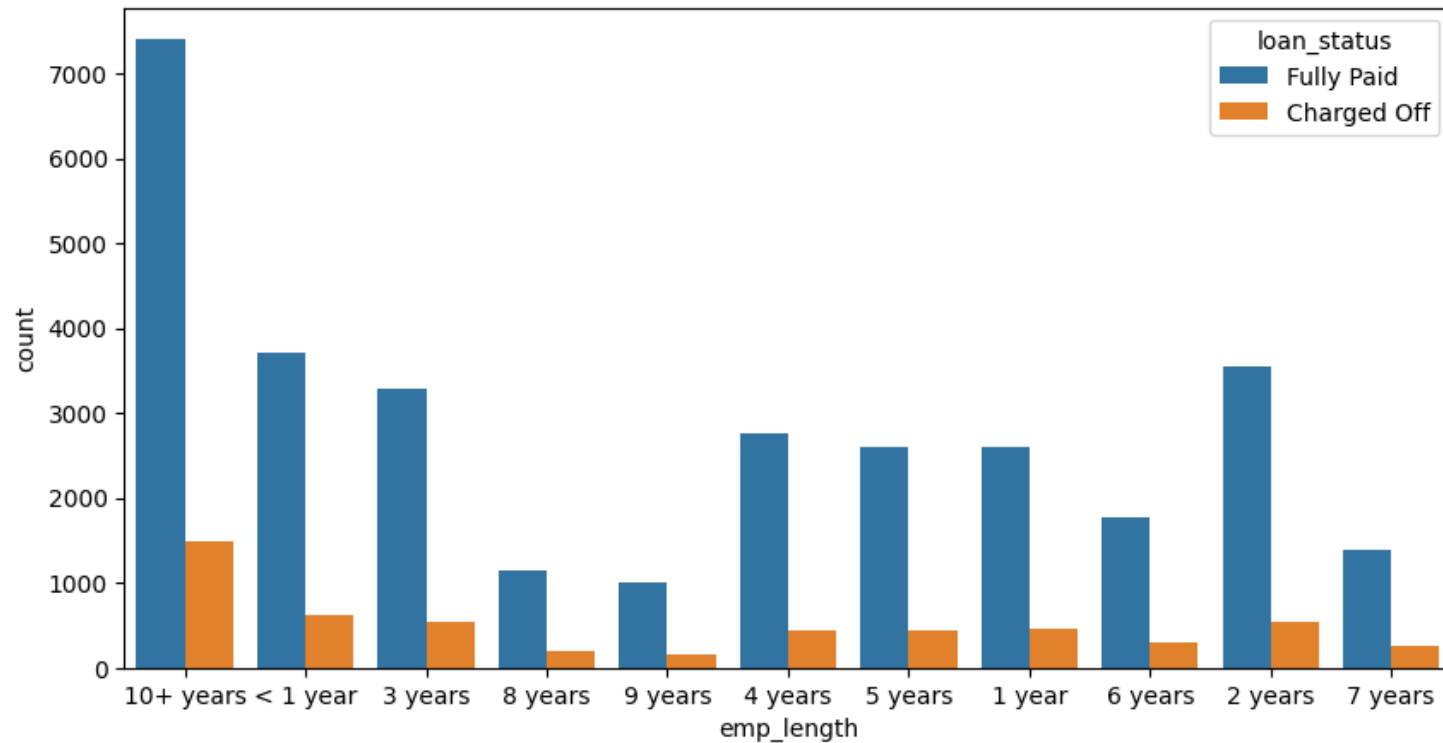
# Impact of *Interest Rate*

**No significant difference can be observed here in distribution for 'loan_amnt'.**
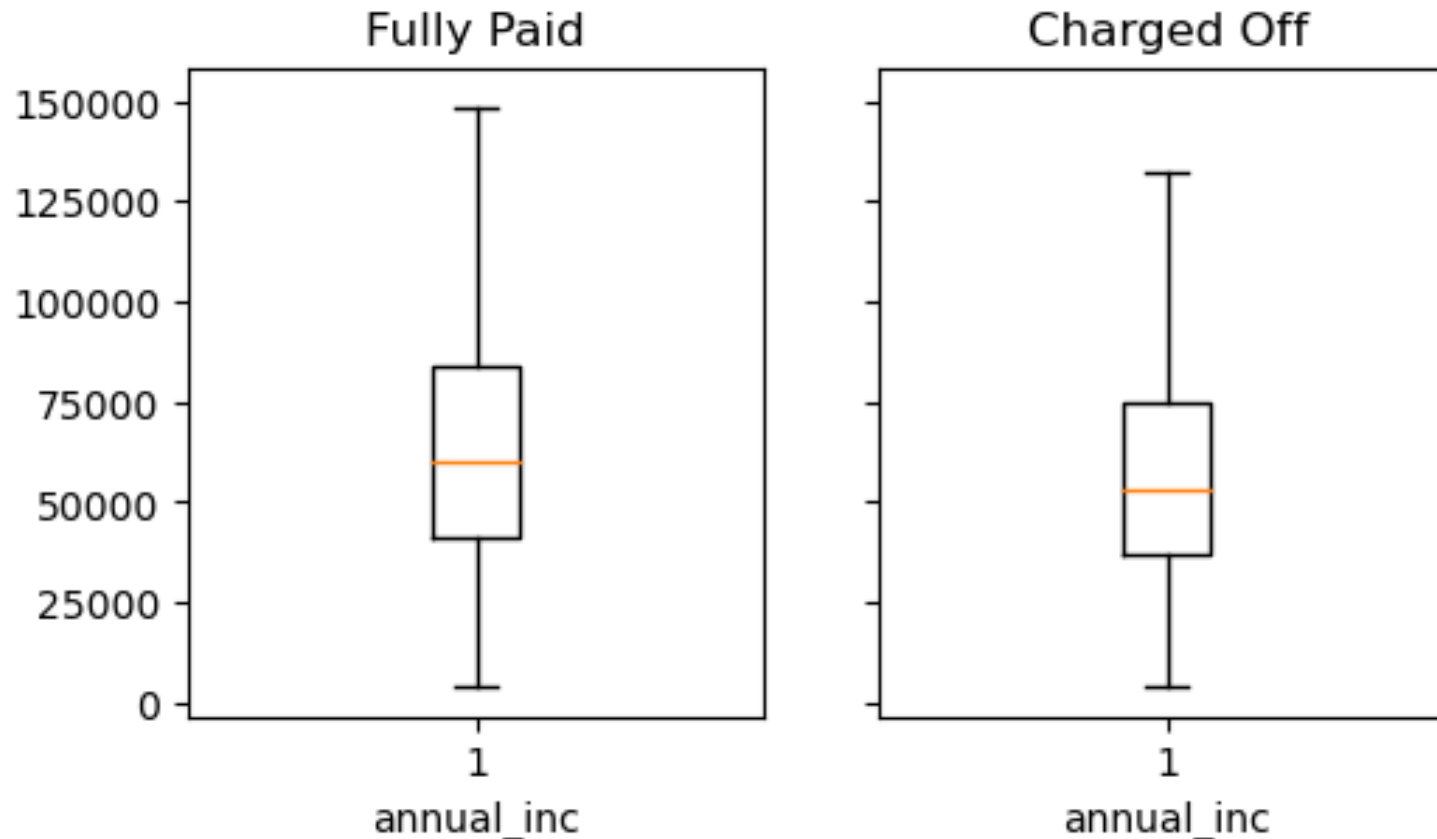
# *Impact of **Employment Length***

**No significant difference can be observed here in distribution for 'emp_length' and**
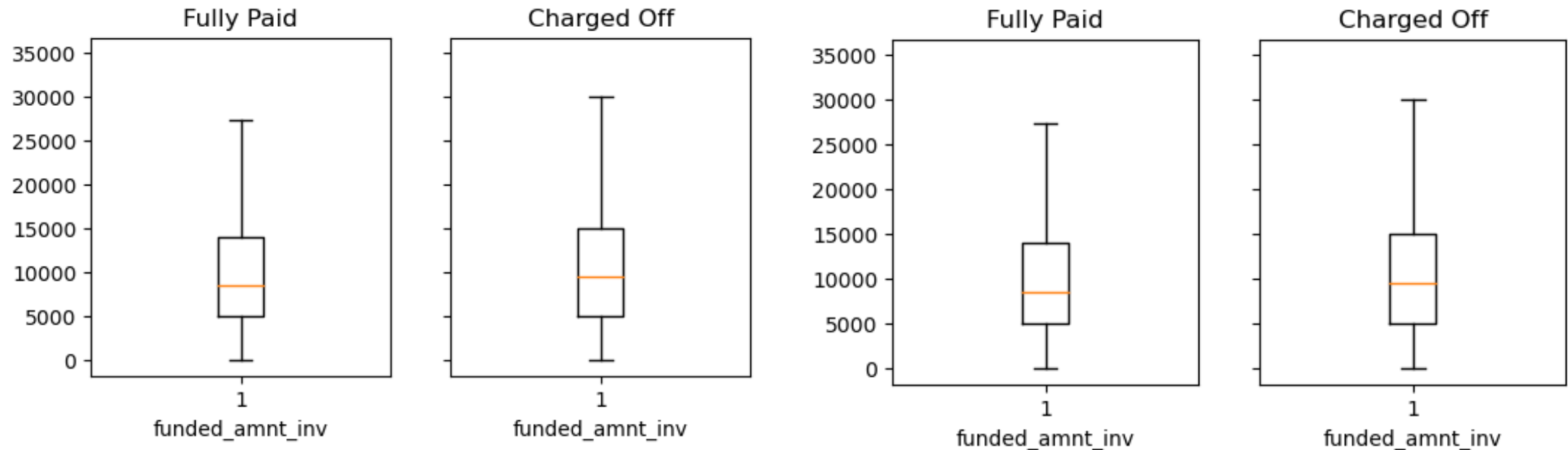
**'home_ownership'.**

# *Impact of **Annual Income***

**The fully paid loans have higher median and other quartiles for 'annual_income'. The variable**

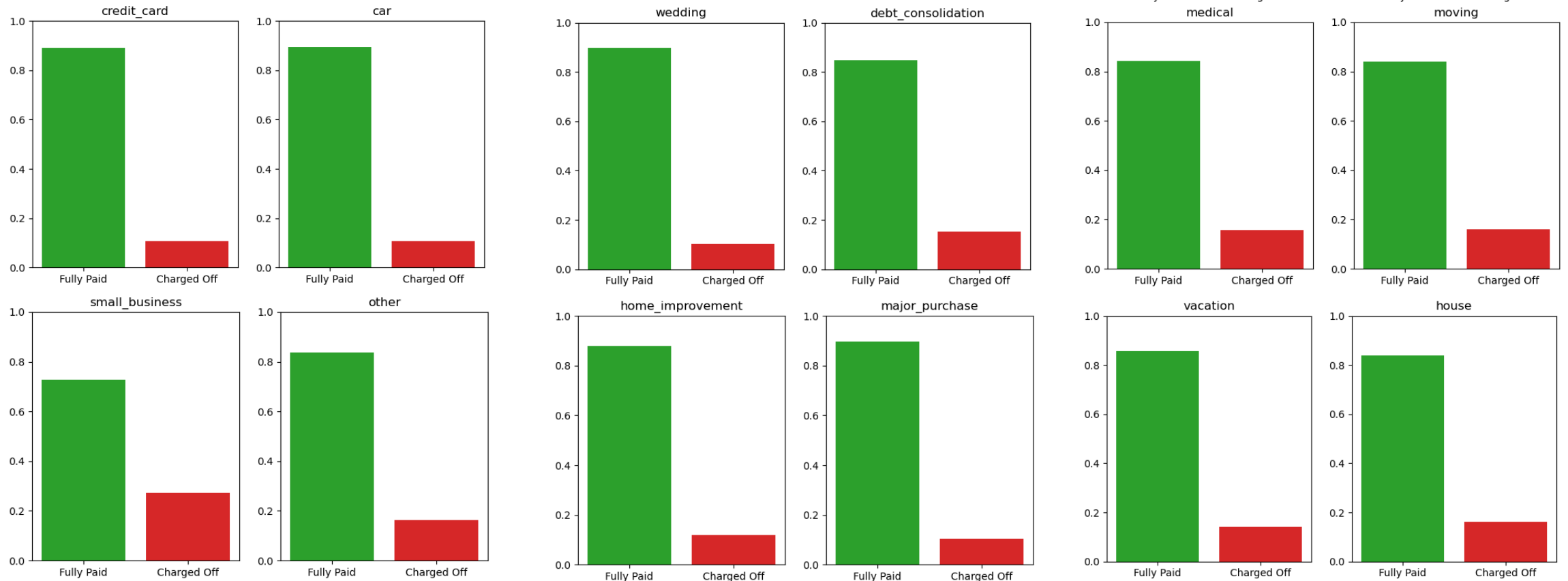**seems to have some impact on the loan status.**

# *Impact of Interest Rate*

**No significant difference can be observed here in distribution for 'funded_amnt_inv' and 'revol_bal'.**
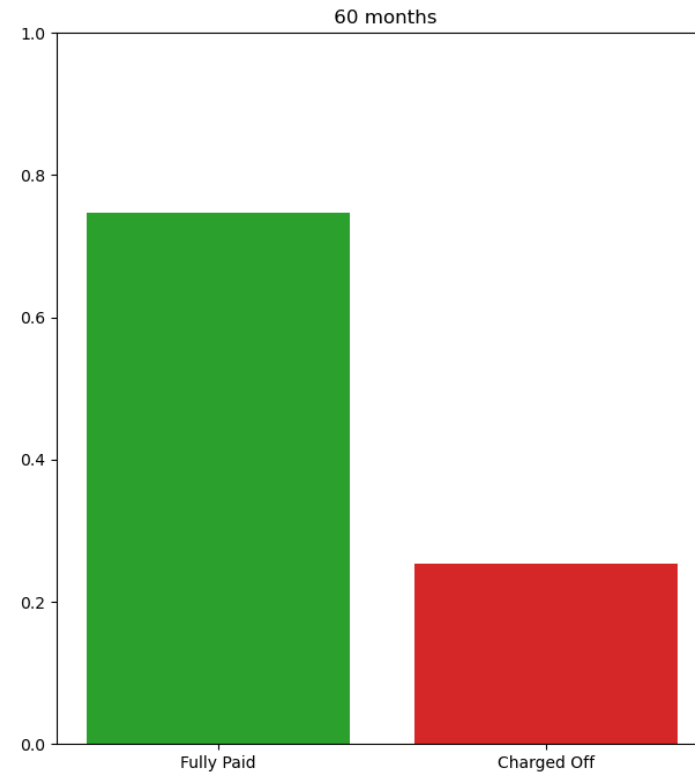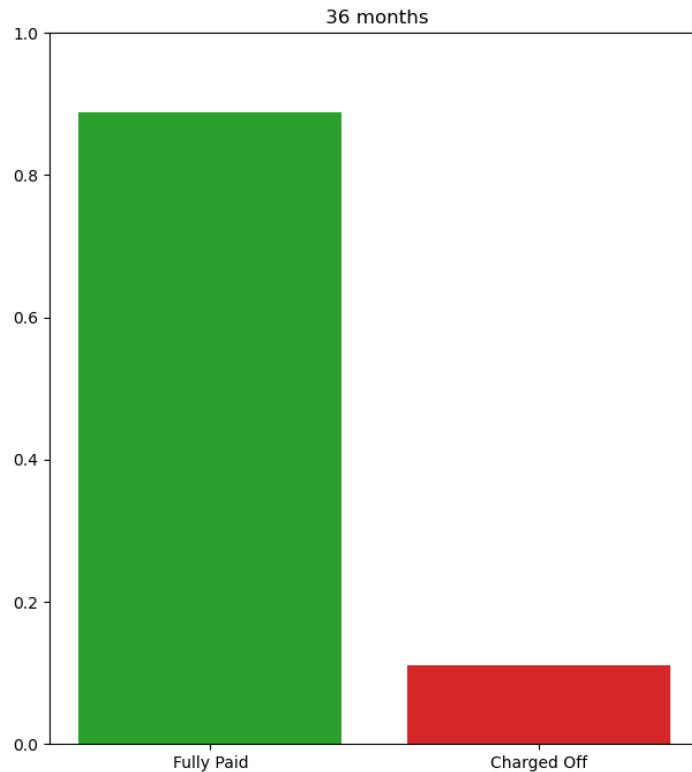
# Impact of *Purpose*

*The purpose of the loan seems to impact the 'loan_status'*

*variables. Some purpose values have higher percentage of*
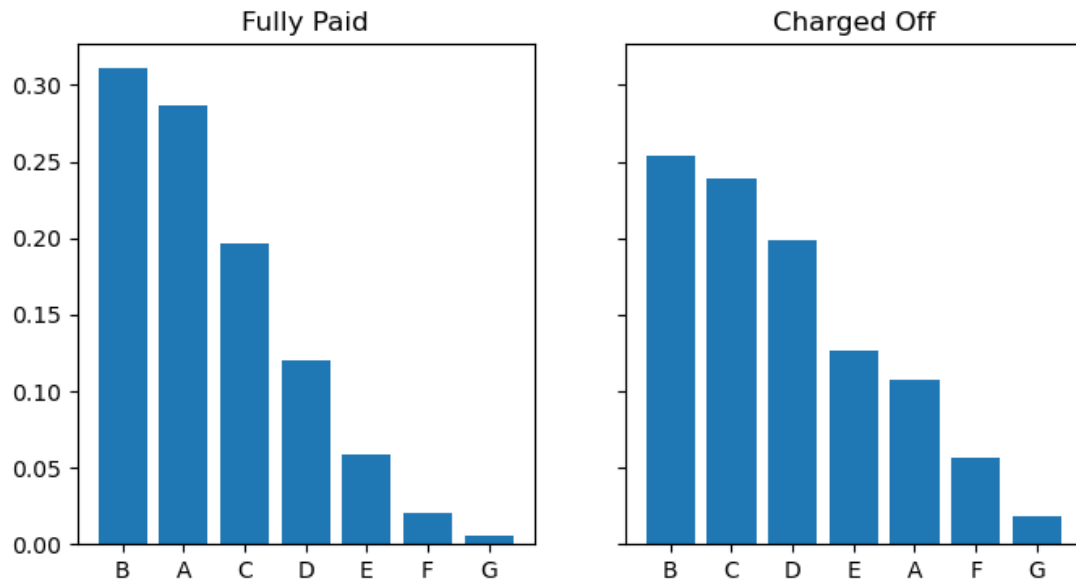
*the charged-off loans compared to others.*

# *Impact of **Term***

**Peercentage of 'charged_off' loans is almost double in case of higher terms, i.e., '60 months',**

**when compared to shorted months '36 months'. This variable seems to have larger impact on**
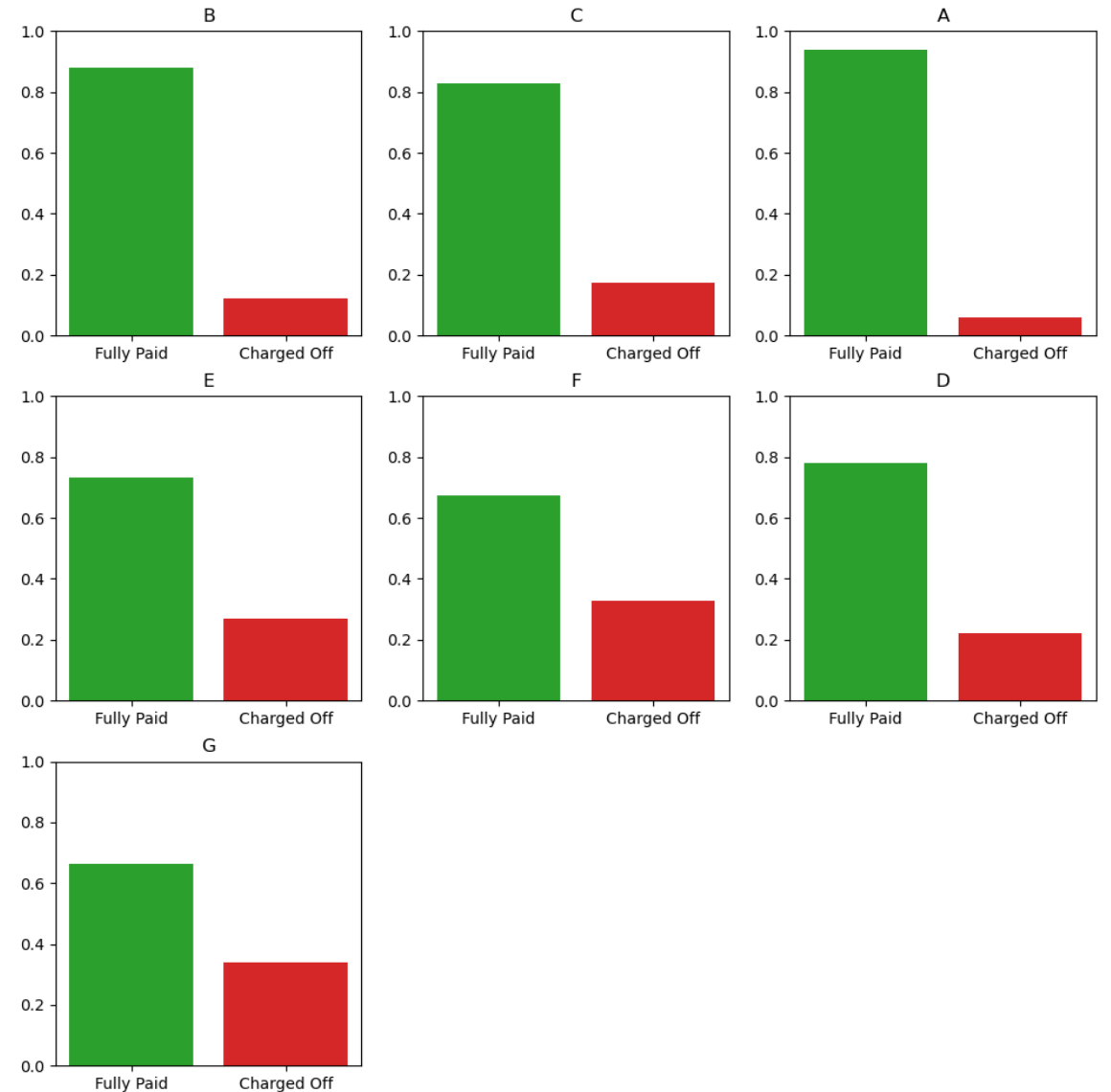
**the loan status.**

# Impact of *Grade*

As the grade quality is declining (A to G), we see higher percentage of the 'charged_off' loans. Hence, 'grade' is an influential variable.
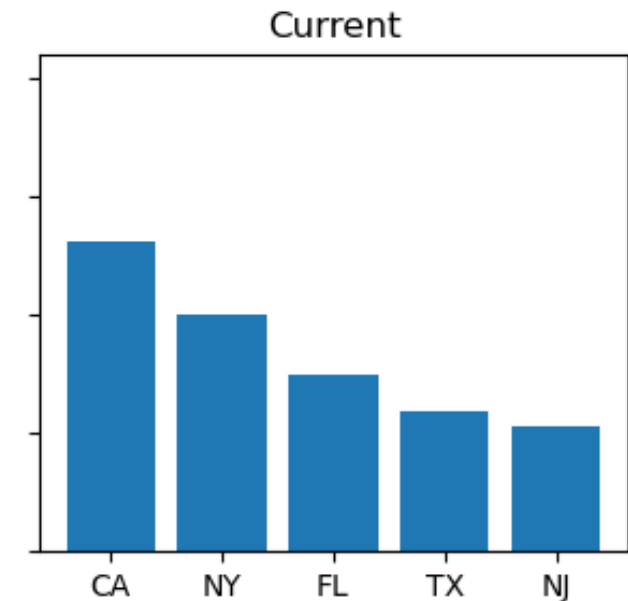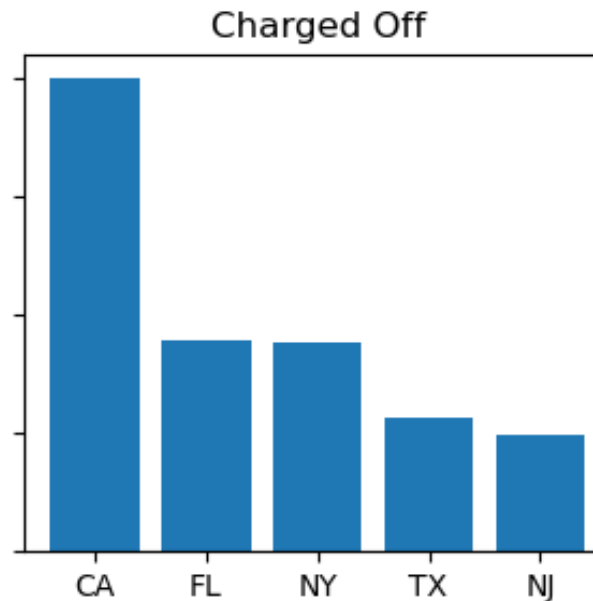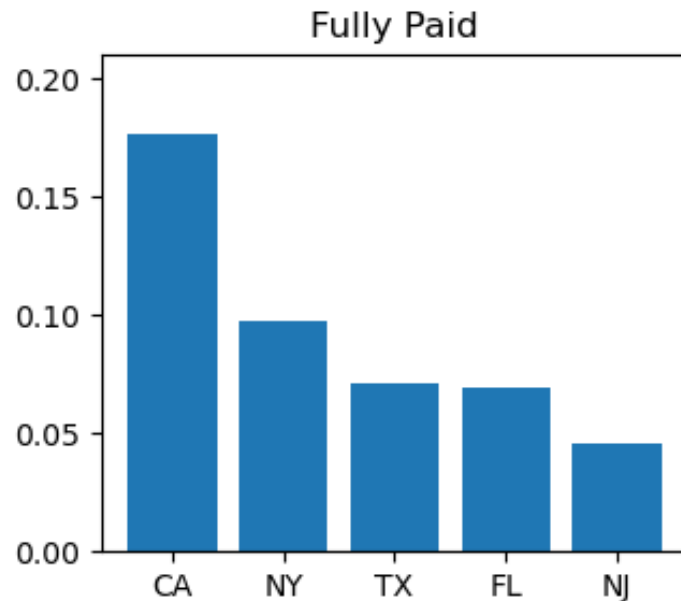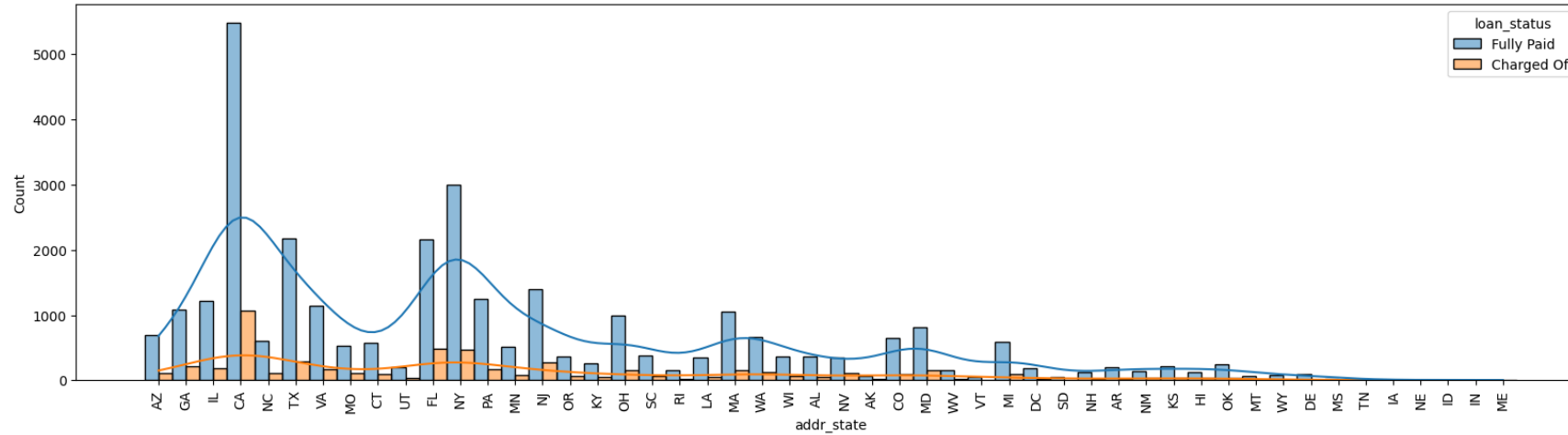


Percentage of each grade in the 'fully_paid' and 'charged_off' loans



Percentage of 'fully_paid' and 'charged_off' loans in each grade.

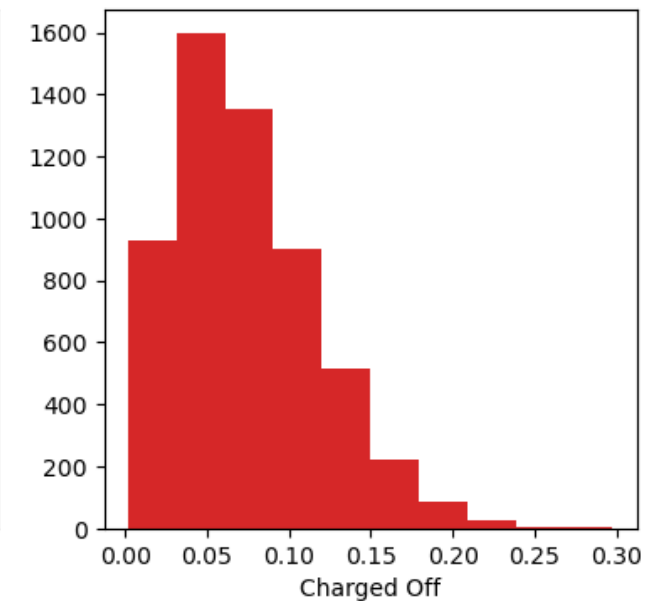**No particular impact of the city, i.e., 'addr_state' can be observed.**

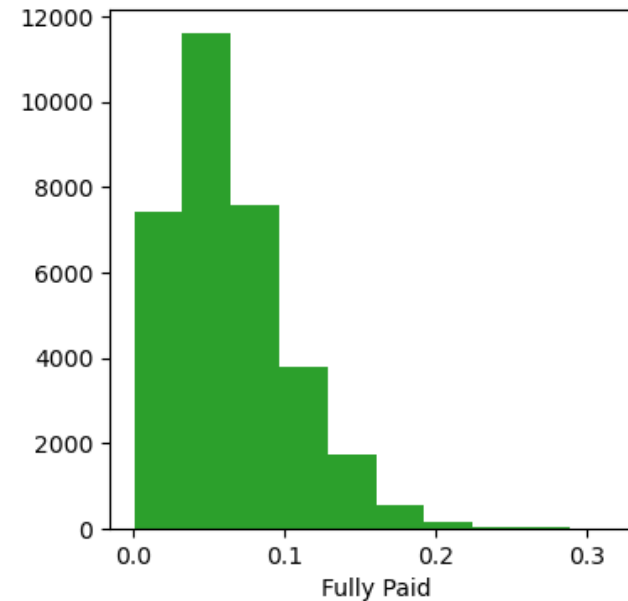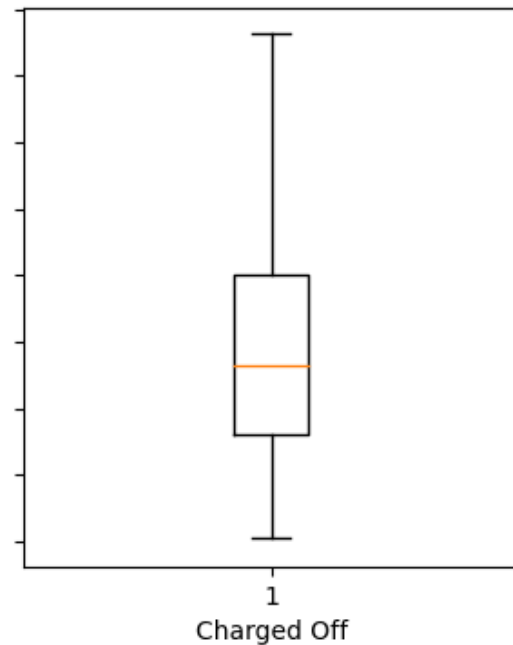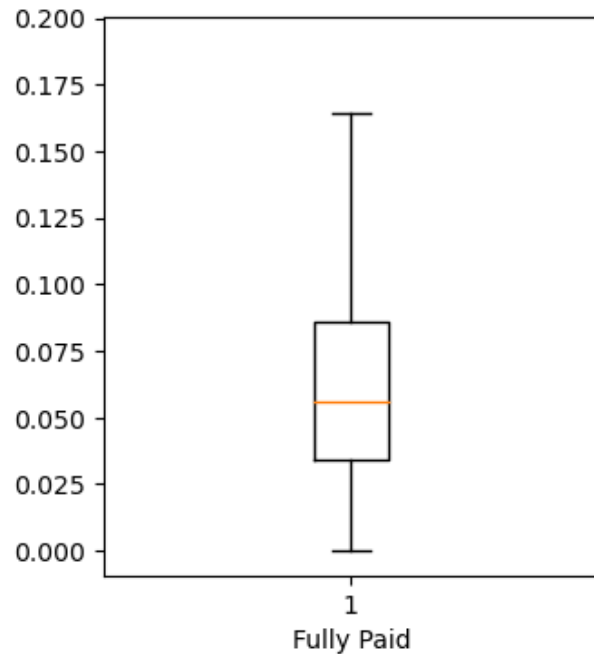# Impact of *Installment to Income ratio*

We created a new column with following formula

**Ratio = installment*12/annual_inc**

**The charged-off loans have slightly higher quartiles for the 'ratio'. However, no significant difference can be observed here in distribution for 'ratio'.**
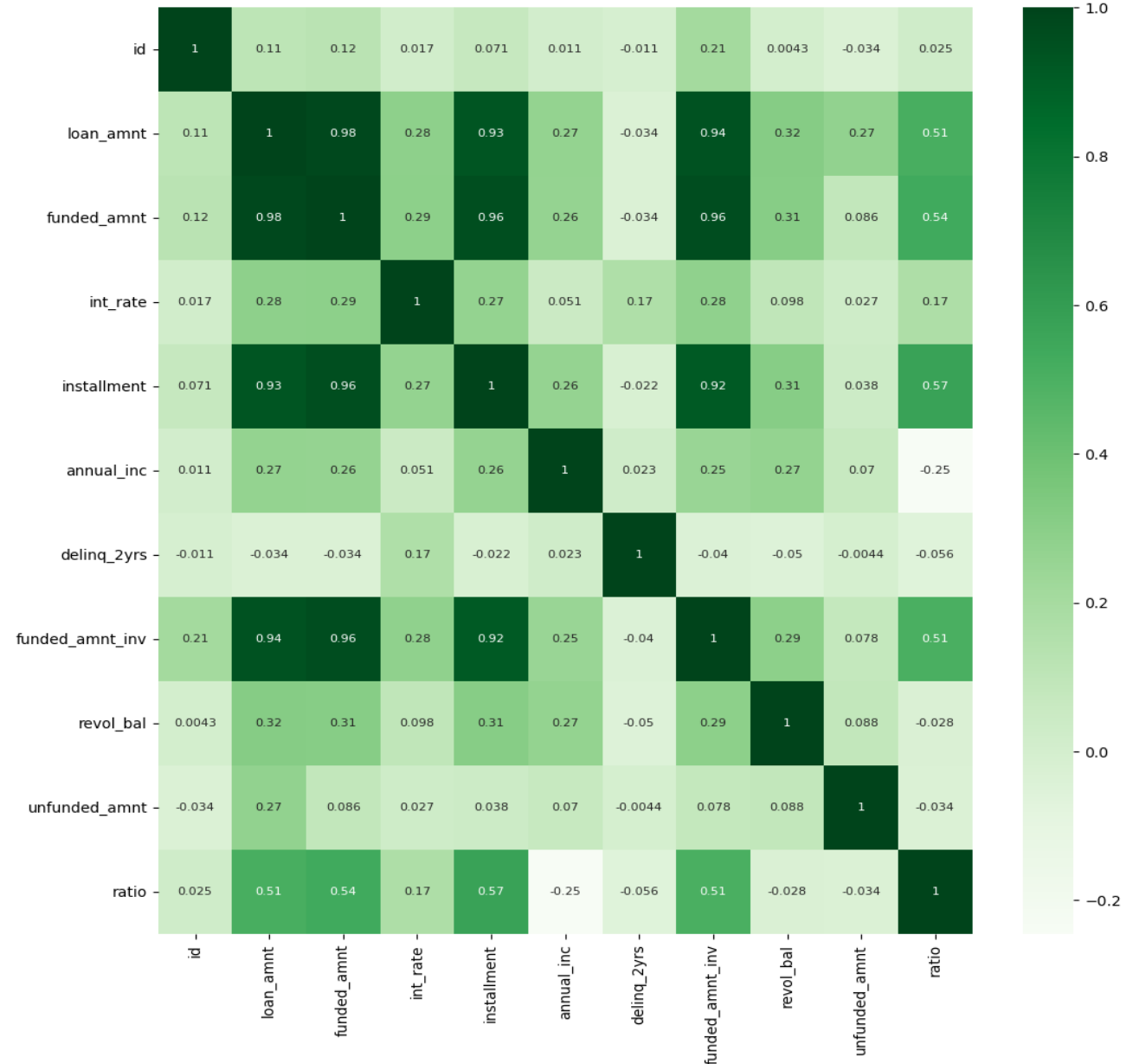
# Bi-variate Analysis

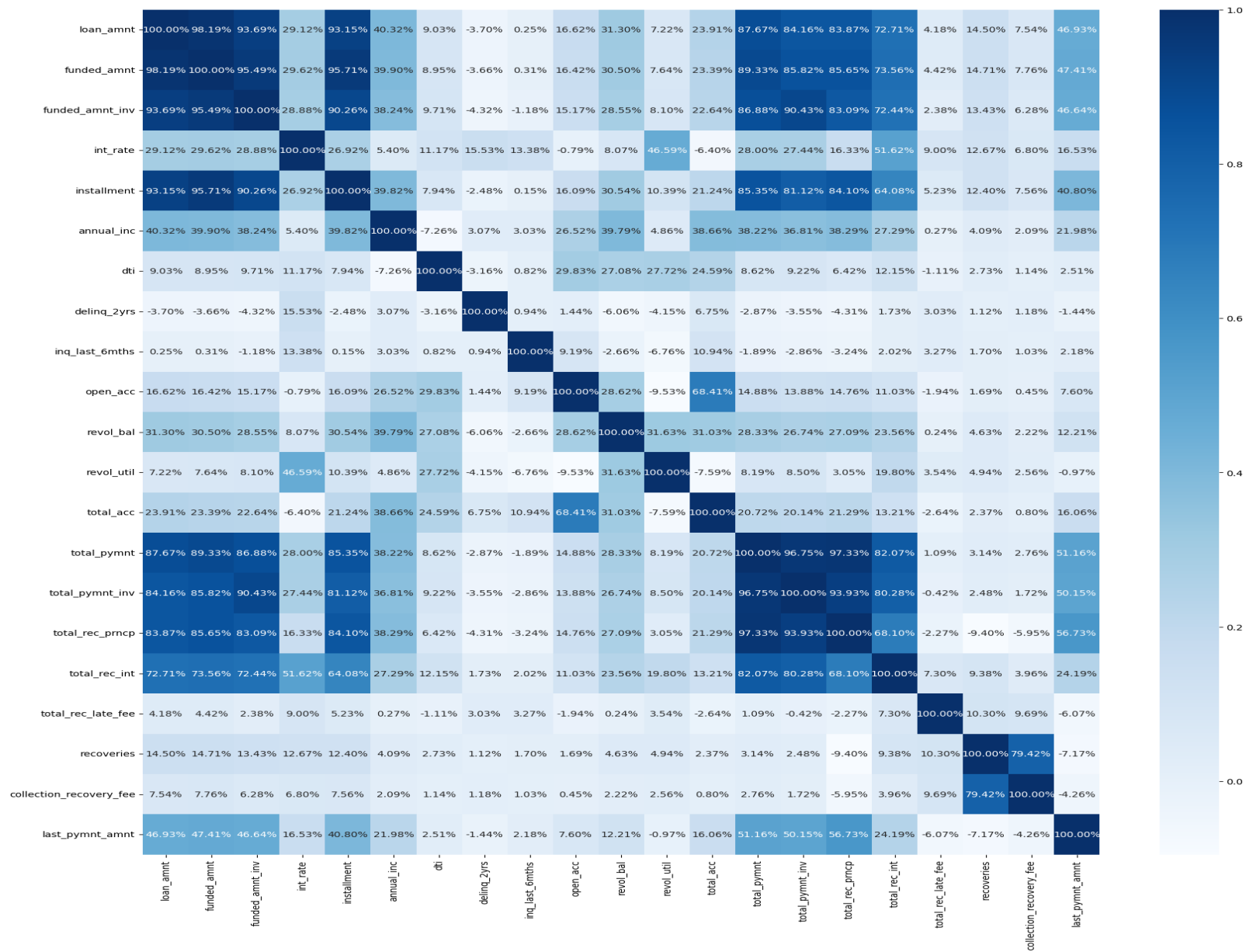# Correlation Analysis

*Some variables show very high co-relation.*

*However, the correlation is very-much expected.*

*For example, it is obvious that the installment will be higher in the fraction of the loan_amnt.*

*From the highly co-related variables, i.e., 'loan_amnt', 'funded_amnt', 'installment', 'funded_amnt_inv', only one variable can be considered for model as the univariate analysis already showed that these variables do not have any significant impact.*

# Correlation Analysis

# *Pair Plots*

**Pair plots provide same information as that we find in the univariate analysis.**

# Summary and Observations

- *The variables 'loan_amnt', 'funded_amnt', 'installment', and 'funded_amnt_inv', are highly correlated. Only one variable can be considered for model as the univariate analysis already showed that these variables do not have any significant impact.*
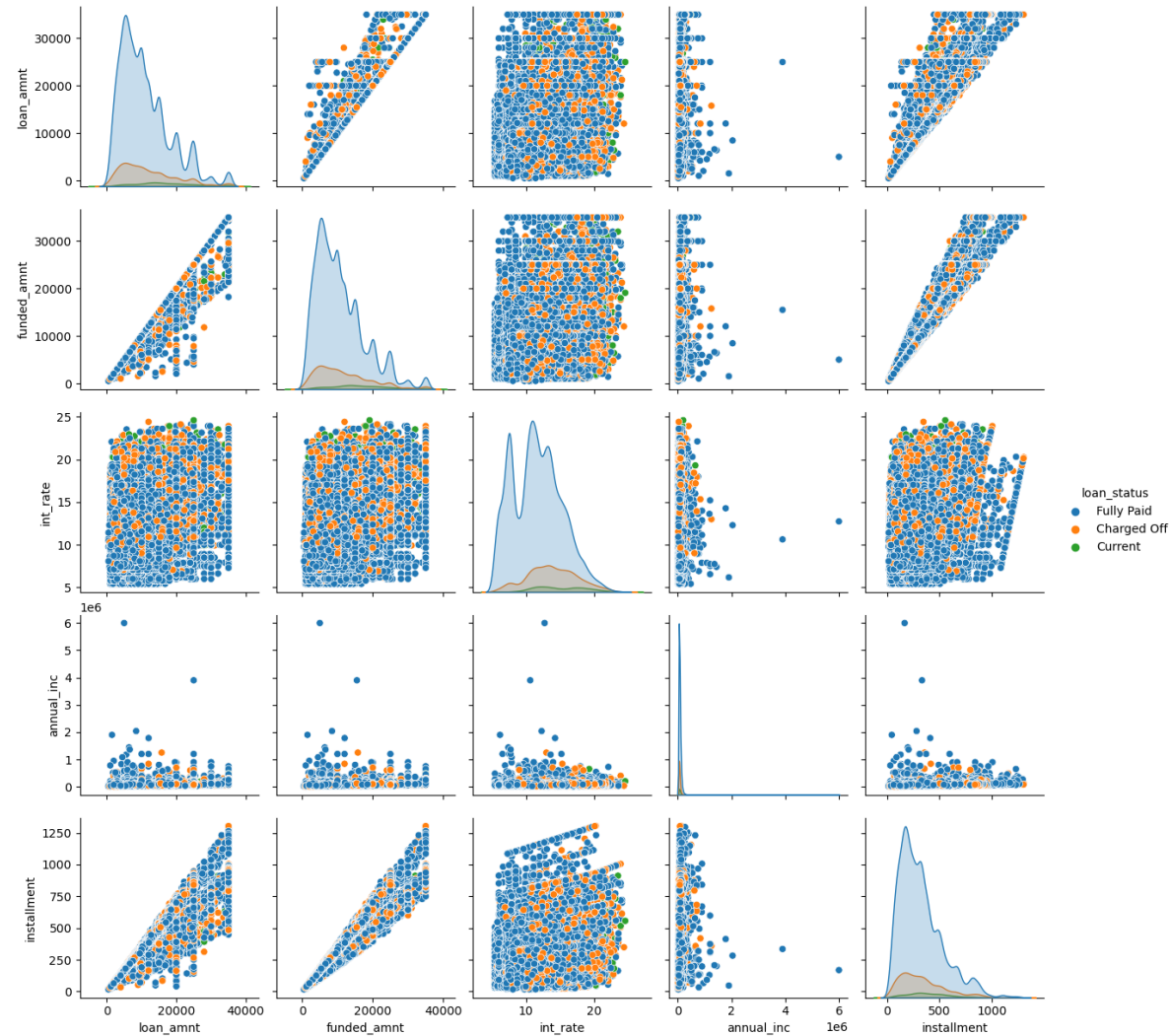
- *The variables 'grade', 'term', 'purpose', and 'int rate' seems to have larger influence on the 'loan_status'.*

- *Rest of the variables have middling to low impact on the 'loan_status'.*