



# Lending Club Case Study

Group Members:  
Mayank Bhardwaj  
Ritu Sharma





# Index

---

Objective

---

Data Understanding

---

Steps taken in python file for performing the EDA

---

Overall Analysis : All the customers

---

Bivariate Analysis b/w Fully paid and Charged off customer

---

Univariate Analysis for Charged off customer

---

Bivariate Analysis for Charged off customer

---

Multivariate Analysis for Charged off customer



## Objective

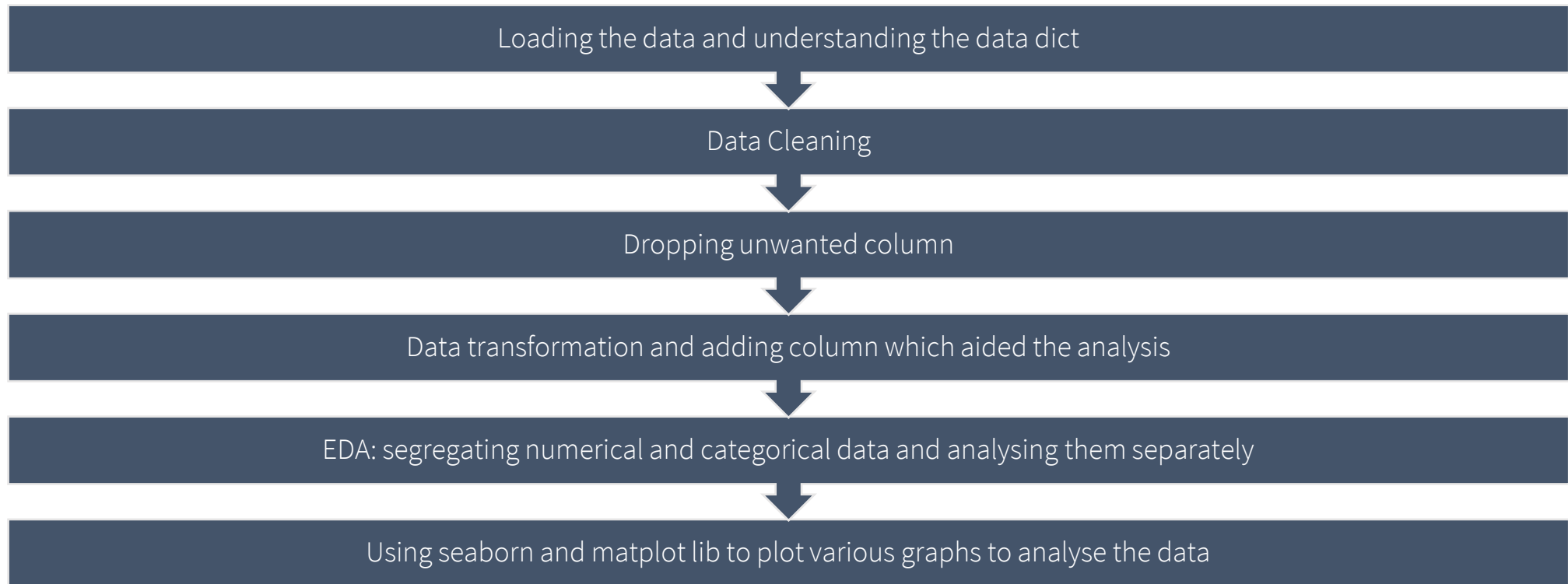
- Lending Club is America's largest online credit marketplace, and the first marketplace bank connecting borrowers and investors.
- When the company receives the loan application, it has to decide for loan approval on the basis of loan applicant
- The objective is to use EDA to understand how consumer attributes and loan attributes influence the tendency to default

## Data Understanding

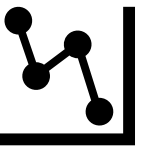
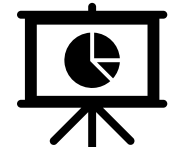
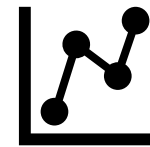
### Excluded Columns:

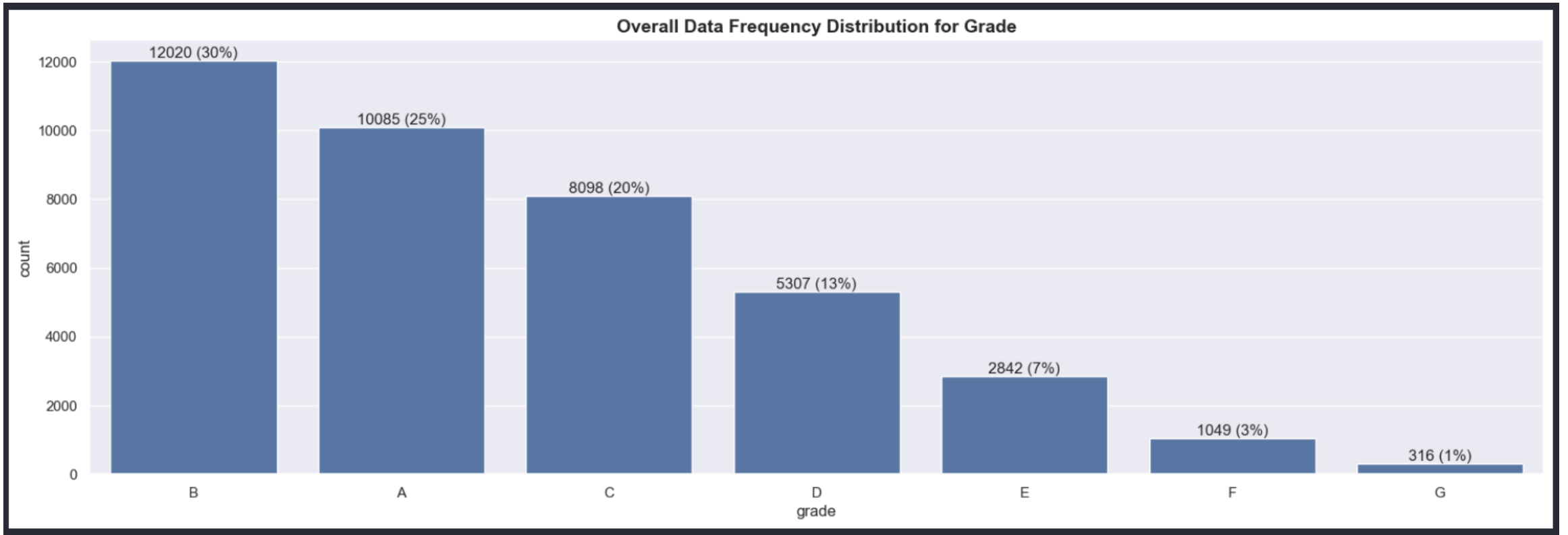
- In our analysis, we will not consider certain types of columns. It's important to note that this is a general categorization of the columns we will exclude from our approach, and it does not represent an exhaustive list.
- **Customer Behavior Columns**-Columns that describe customer behavior will not be factored into our analysis. The current analysis focuses on the loan application stage, while customer behavior variables pertain to post-approval actions. Consequently, these attributes will not influence the loan approval/rejection process.
- **Granular Data** -Columns providing highly detailed information that may not be necessary for our analysis will be omitted. For example, while the "grade" column may have relevance in creating business outcomes and visualizations, the "sub grade" column is excessively granular and will not be utilized in our analysis.
- 54 columns contain NA values only, and these columns will be removed namely acc\_open\_past\_24mths, all\_util, annual\_inc\_joint, avg\_cur\_bal, bc\_open\_to\_buy, bc\_util, dti\_joint, il\_util, inq\_fi, inq\_last\_12m, max\_bal\_bc, mo\_sin\_old\_il\_acct, mo\_sin\_old\_rev\_tl\_op, mo\_sin\_rcnt\_rev\_tl\_op, mo\_sin\_rcnt\_tl, mort\_acc, mths\_since\_last\_major\_derog, mths\_since\_rcnt\_il, mths\_since\_recent\_bc, mths\_since\_recent\_bc\_dlq, mths\_since\_recent\_inq, mths\_since\_recent\_revol\_delinq, num\_accts\_ever\_120\_pd, num\_actv\_bc\_tl, num\_actv\_rev\_tl, num\_bc\_sats, num\_bc\_tl, num\_il\_tl, num\_op\_rev\_tl, num\_rev\_accts, num\_rev\_tl\_bal\_gt\_0, num\_sats, num\_tl\_120dpd\_2m, num\_tl\_30dpd, num\_tl\_90g\_dpd\_24m, num\_tl\_op\_past\_12m, open\_acc\_6m, open\_il\_12m, open\_il\_24m, open\_il\_6m, open\_rv\_12m, open\_rv\_24m, pct\_tl\_nvr\_dlq, percent\_bc\_gt\_75, tot\_coll\_amt, tot\_cur\_bal, tot\_hi\_cred\_lim, total\_bal\_ex\_mort, total\_bal\_il, total\_bc\_limit, total\_cu\_tl, total\_il\_high\_credit\_limit, total\_rev\_hi\_lim, verification\_status\_joint
- Certain columns contain only 0 values, and these columns will also be dropped.
- 9 Columns with **single value** that do not contribute to the analysis will be removed.
- Columns with values that are single value but have other values as NA will be treated as constant and dropped.
- Columns with more than 65% of data being empty (mths\_since\_last\_delinq, mths\_since\_last\_record) will be dropped.
- Columns (id, member\_id) will be dropped as they are index variables with unique values and do not contribute to the analysis.
- Columns (emp\_title, desc, title) will be dropped as they contain descriptive text (nouns) and do not contribute to the analysis.
- The redundant column (url) will be dropped. Further analysis reveals that the URL is a static path with the loan ID appended as a query, making it redundant compared to the (id) column.
- These columns capture customer behavior recorded after loan approval and are not available at the time of loan approval. Thus, these variables will not be included in the analysis.
- Columns to be dropped: (delinq\_2yrs, earliest\_cr\_line, inq\_last\_6mths, open\_acc, pub\_rec, revol\_bal, revol\_util, total\_acc, out\_prncp, out\_prncp\_inv, total\_pymnt\_inv, total\_rec\_prncp, total\_rec\_int, total\_rec\_late\_fee, recoveries, collection\_recovery\_fee, last\_pymnt\_d, last\_pymnt\_amnt, last\_credit\_pull\_d, application\_type)

# Steps Taken to do analysis in python file:

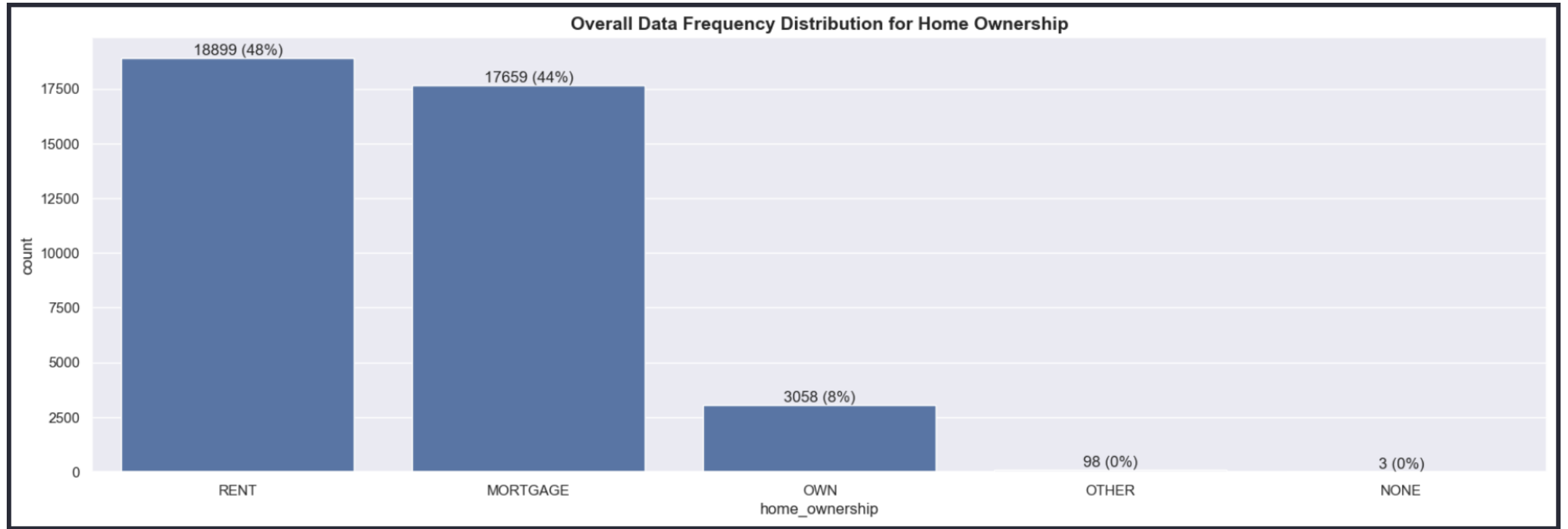


# Overall Analysis: Frequency Distribution of Grade, Home Ownership and Purpose



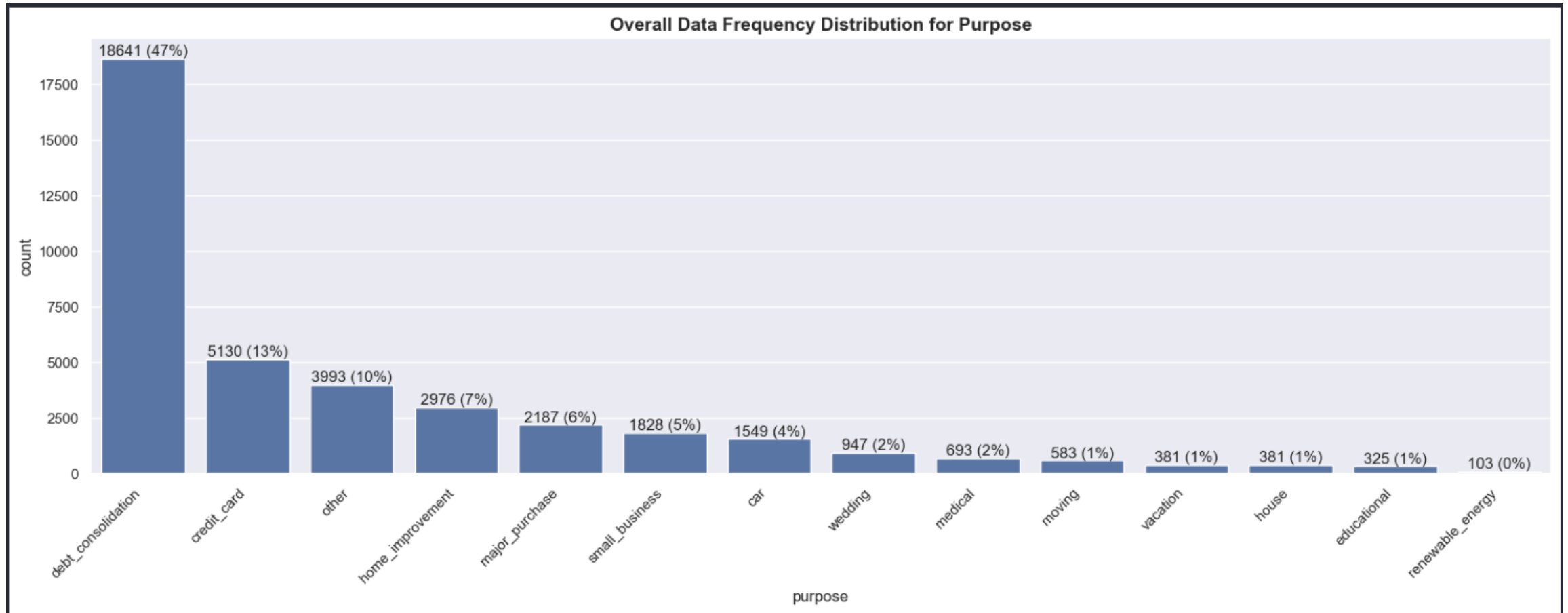


Overall Analysis based on grade show grade B customers are took highest Loan with a total count of 12020 which 30 % of overall customer count



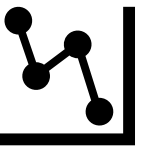
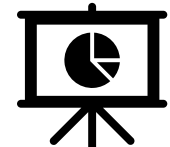
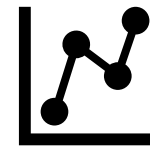
Overall Analysis based on Home Ownership show count of Rent and Mortgage are highest among all other ownership with a total of 18899 (48%) and 17659 (44%) respectively.

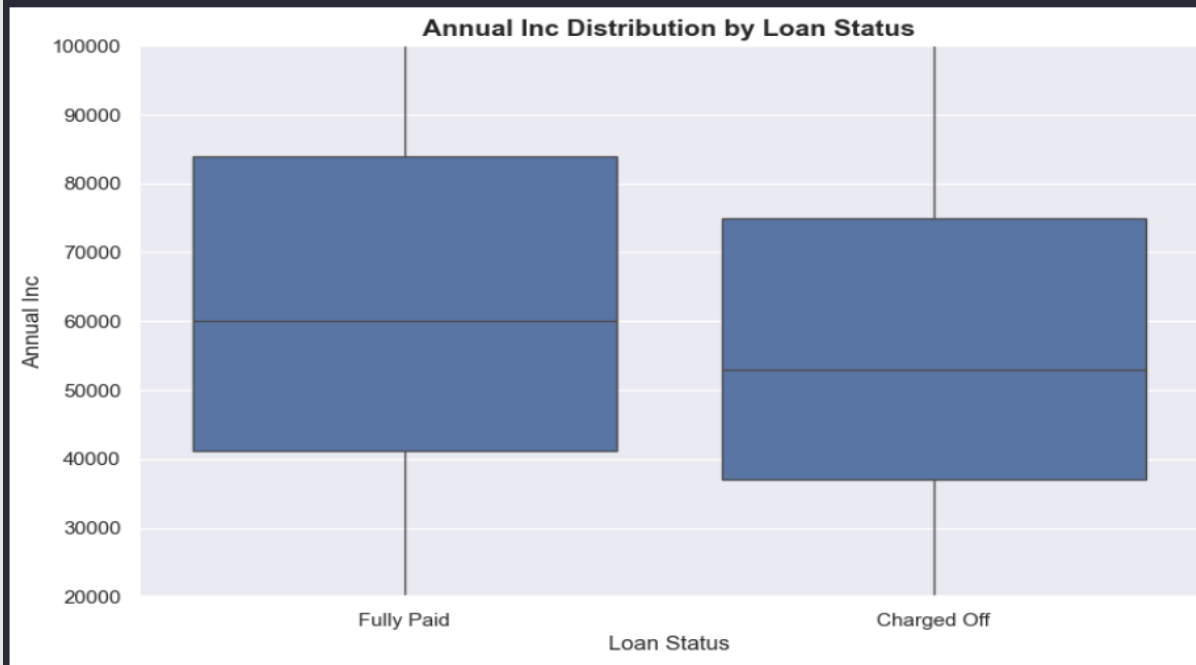
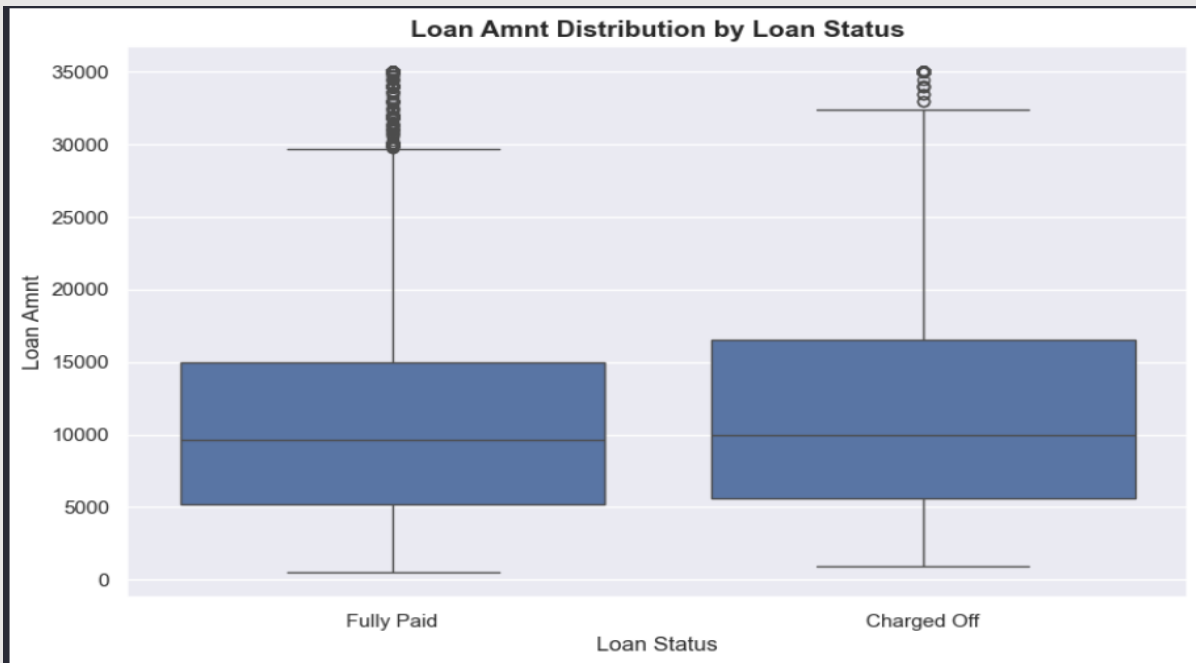




Overall Analysis based on Purpose show Debt Consolidation purpose grabs the highest count with a count of 18641 and with 47%

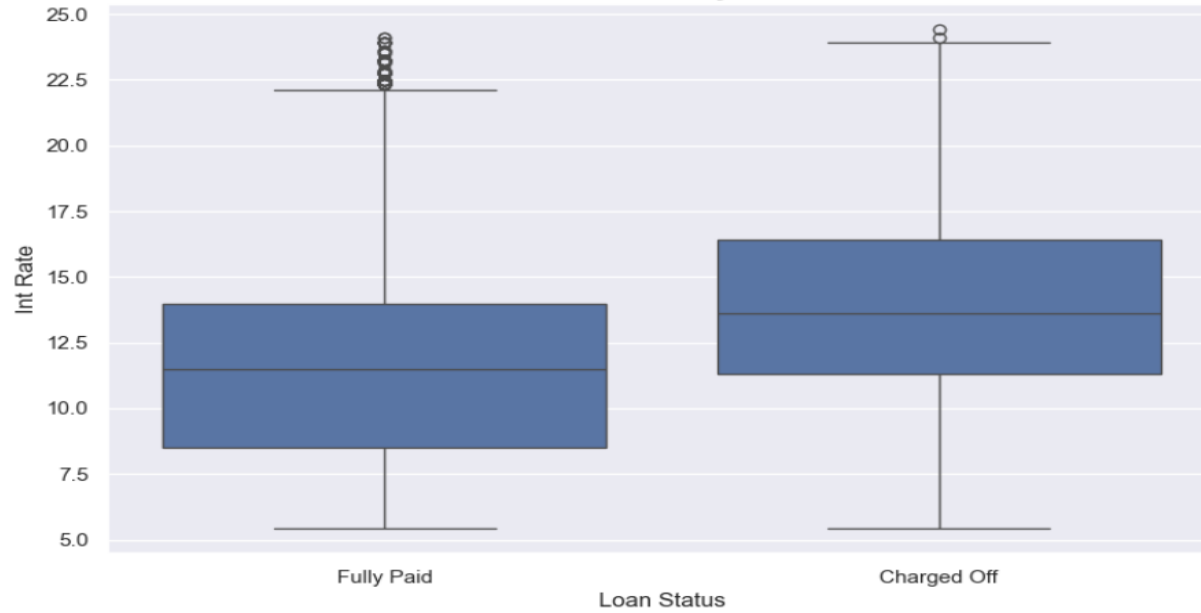
# Bivariate Analysis (between Fully Paid and Charged Off Customers) of Loan Amount, Interest Rate, Instalment with respect to Loan status



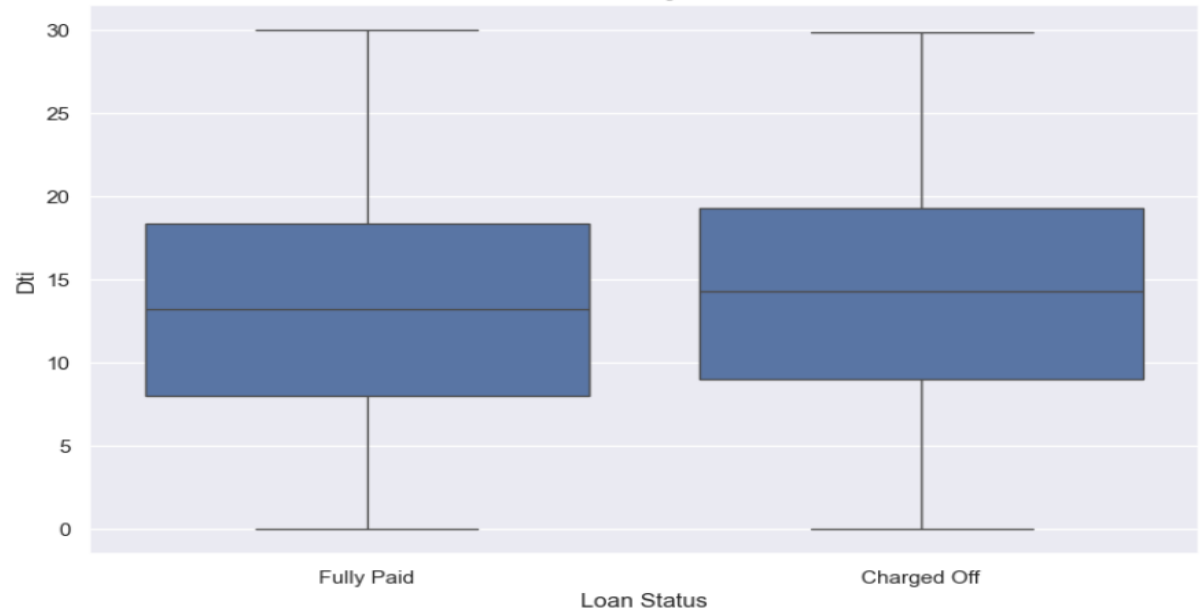


- There is a tendency to default when loan amount is higher as it is evident from the box plot.
- Lower income levels have higher tendency to default.

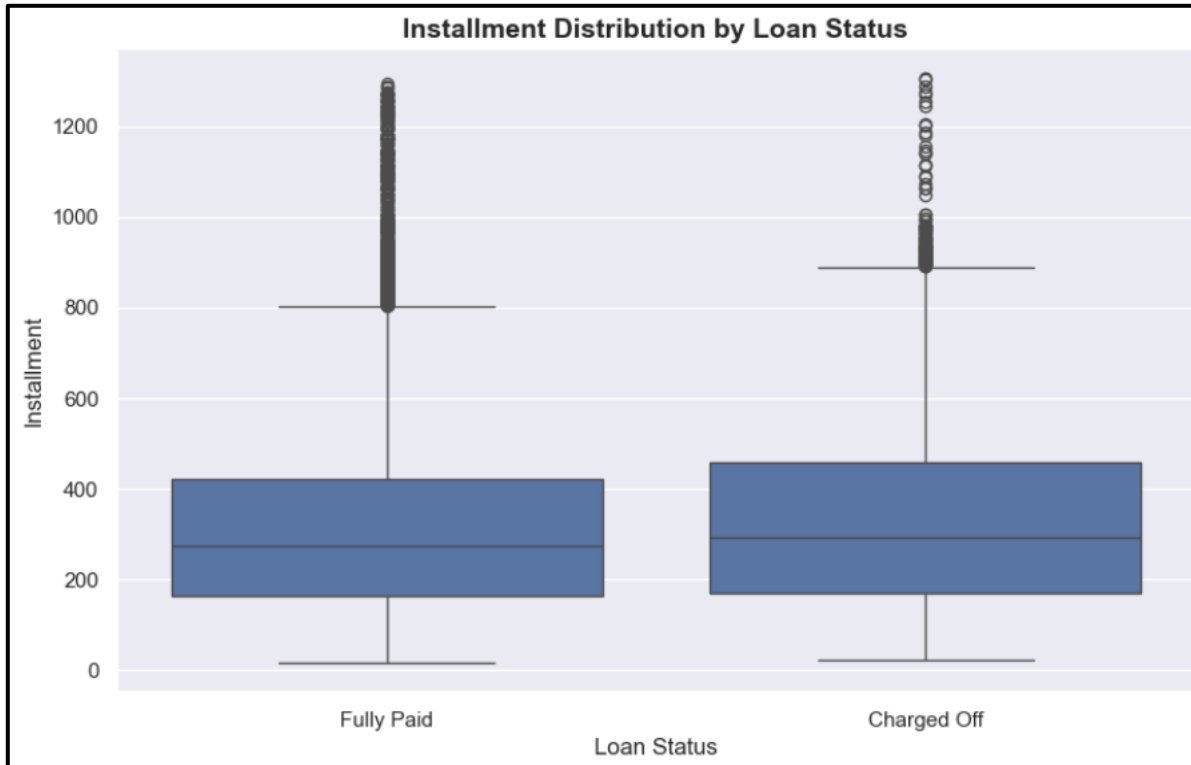
Int Rate Distribution by Loan Status



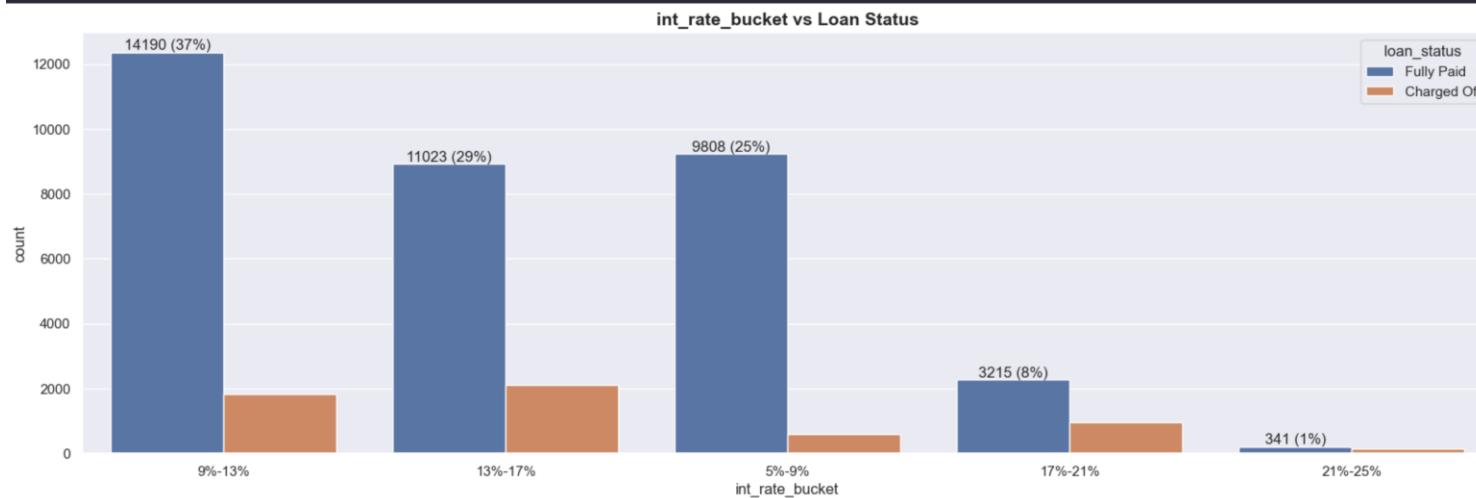
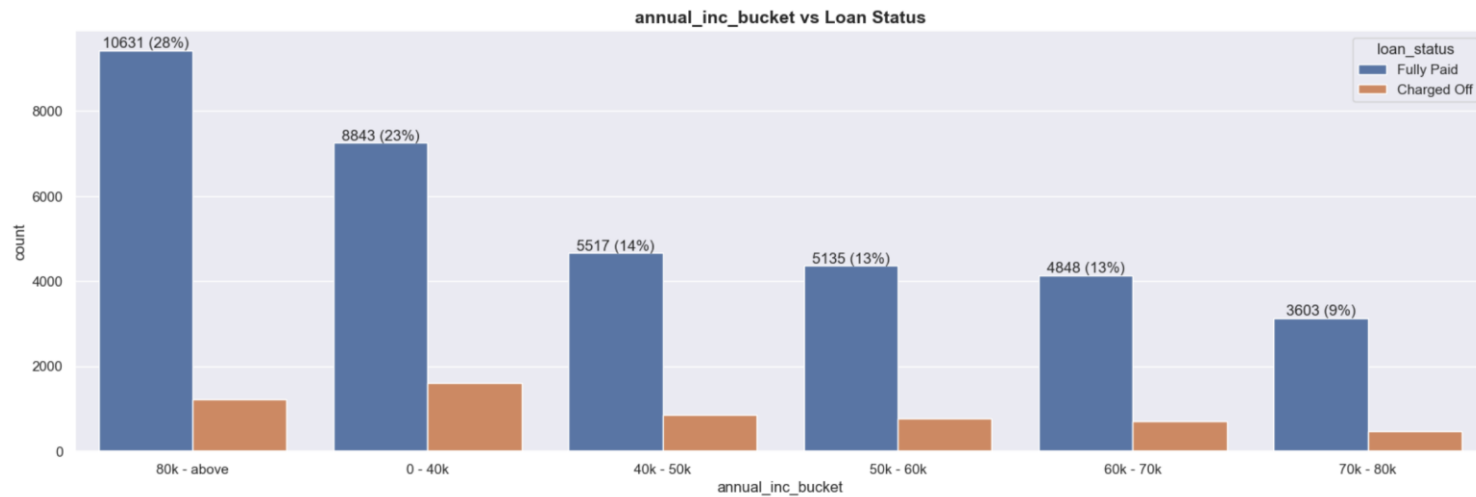
Dti Distribution by Loan Status



- There is a higher tendency to default when Interest rates are higher as box plot suggest Lower income levels have higher tendency to default.
- Debt-to-Income Ratio does not provide any indication of default.

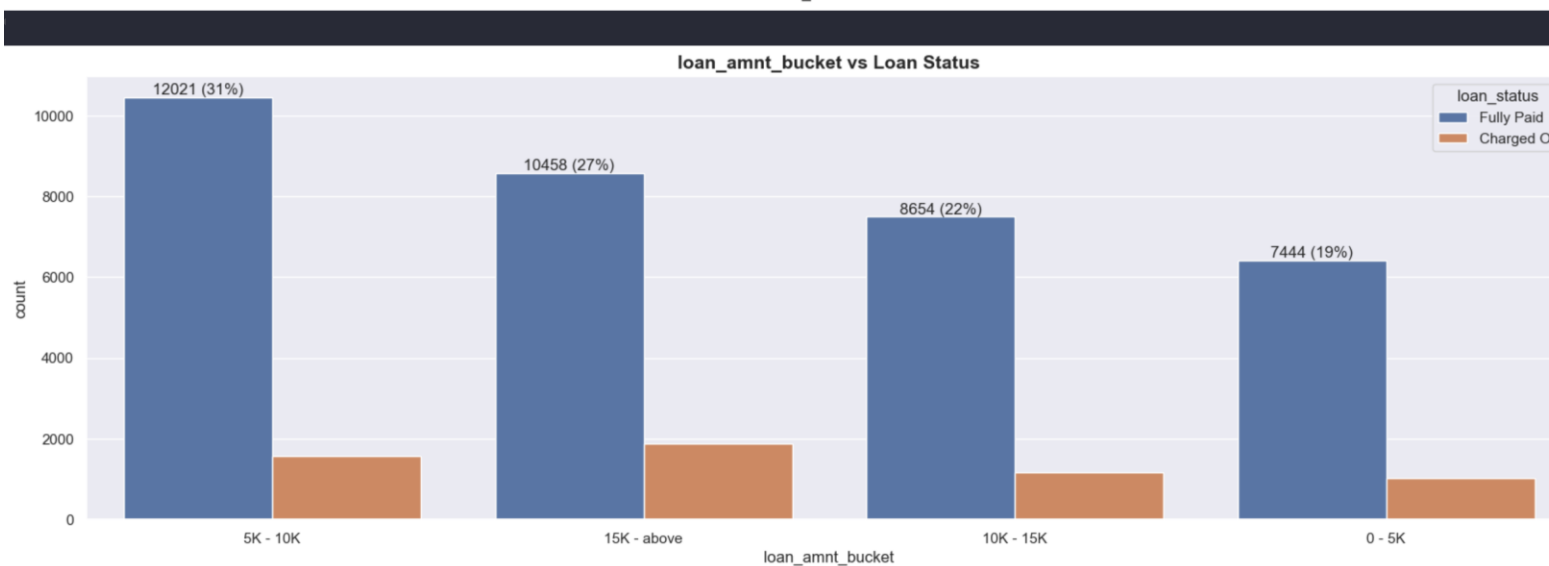
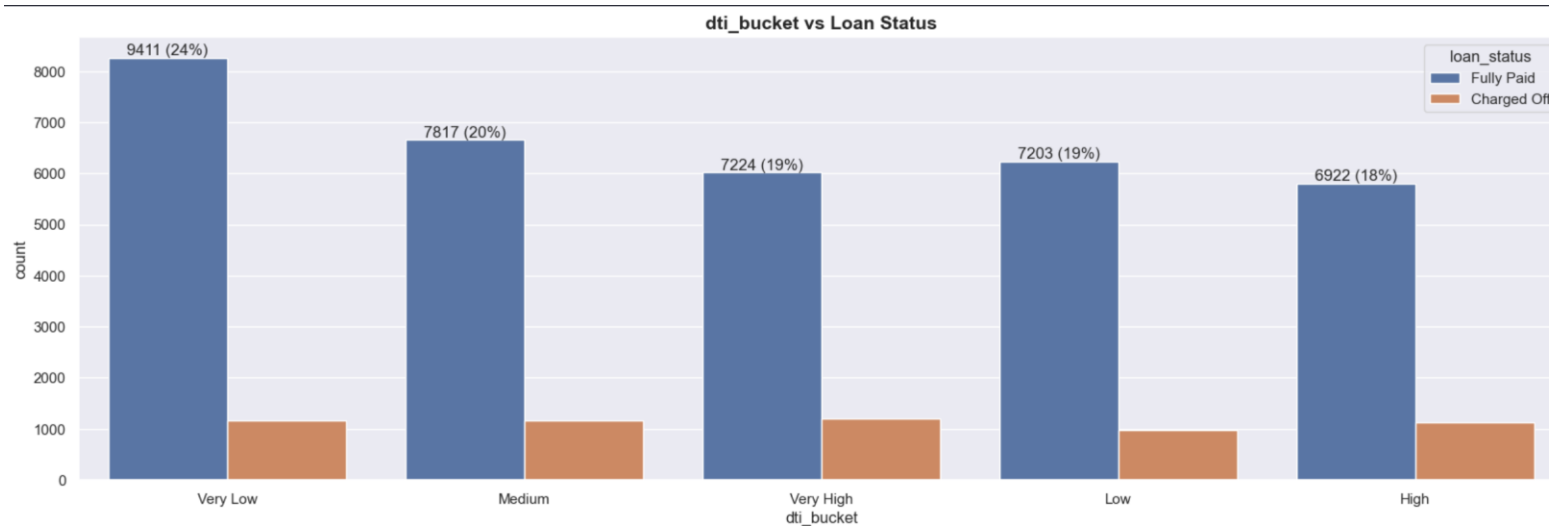


- Installment does not provide any indication of default.



- A significant portion of loan applicants who defaulted received loans with interest rates falling within the range of 13% to 17%.

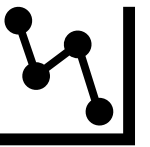
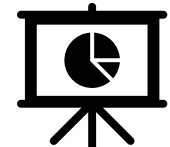
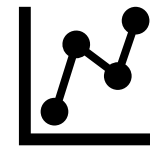
- A majority of the loan applicants who charged off reported an annual income of less than \$40,000.



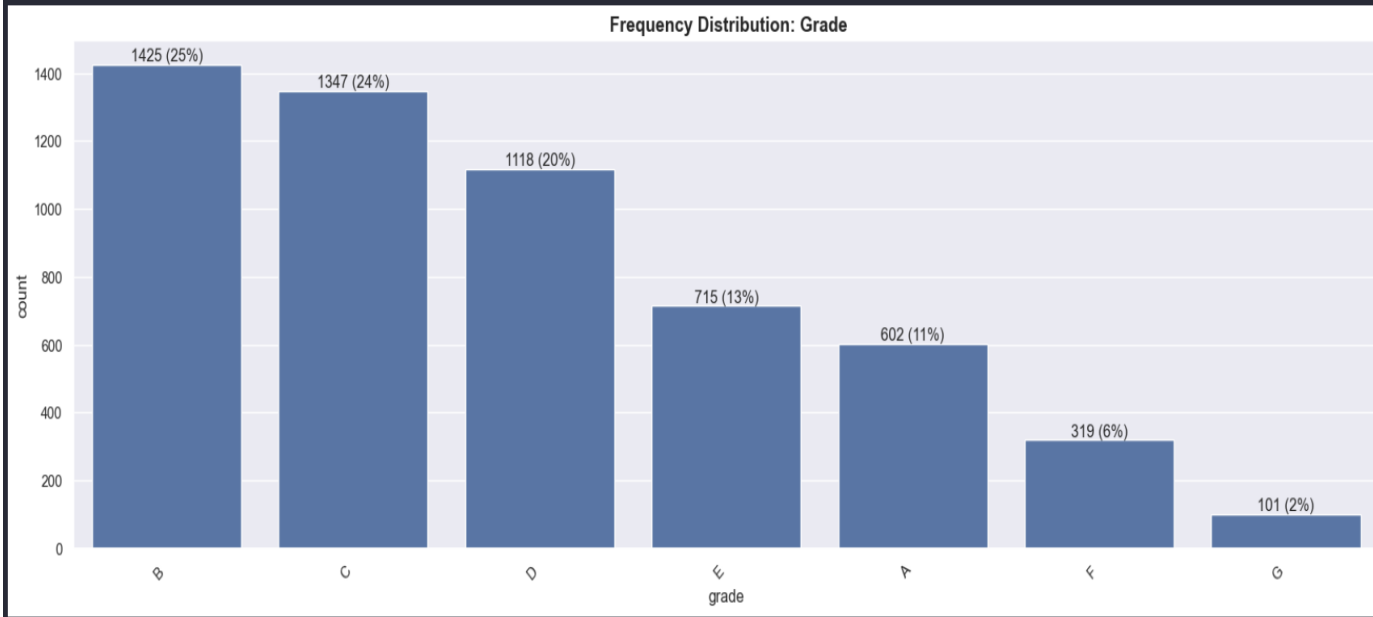
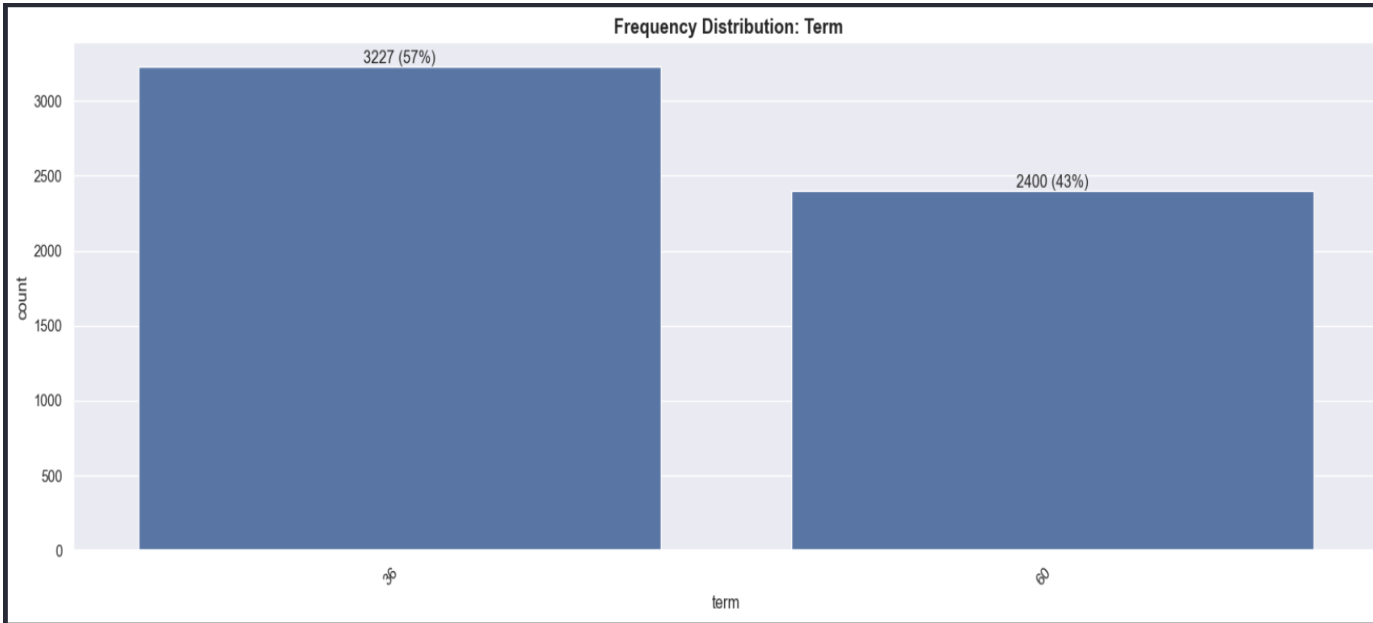
- A majority of the loan applicants who defaulted received loan amounts of \$15,000 or higher.

- The majority of loan applicants who charged off had significantly high Debt-to-Income (DTI) ratios.

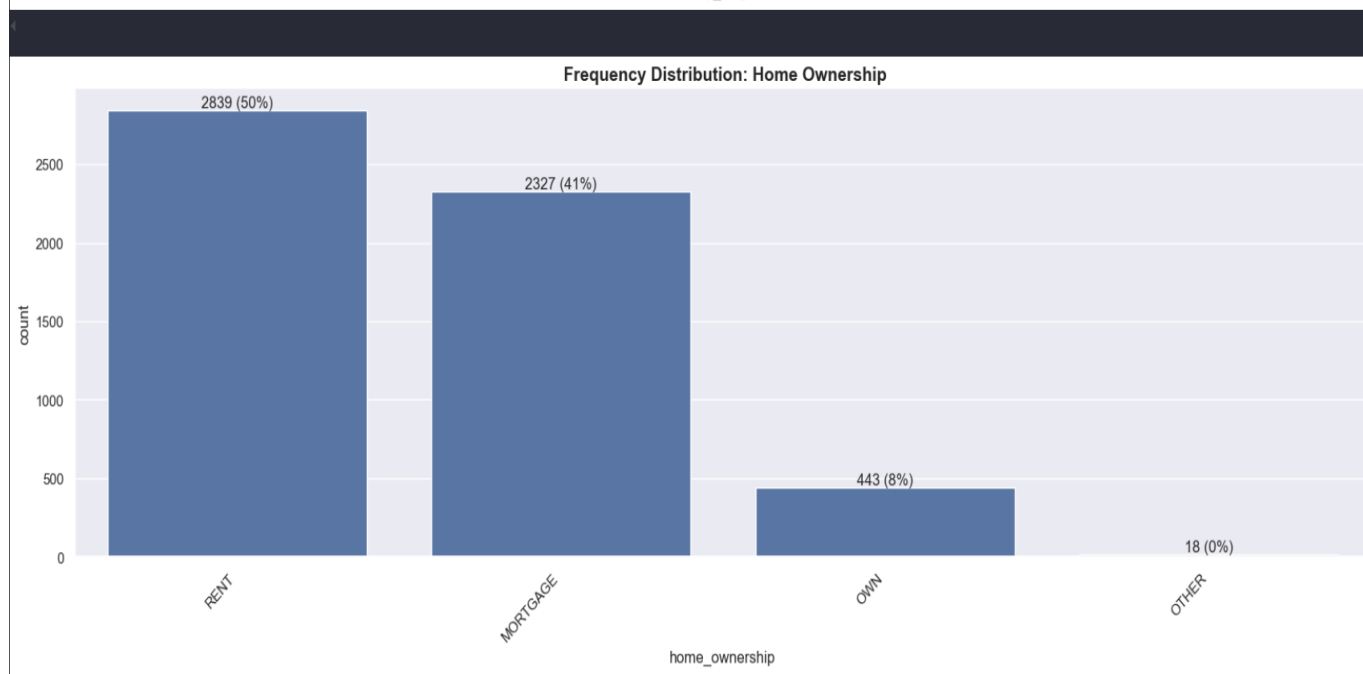
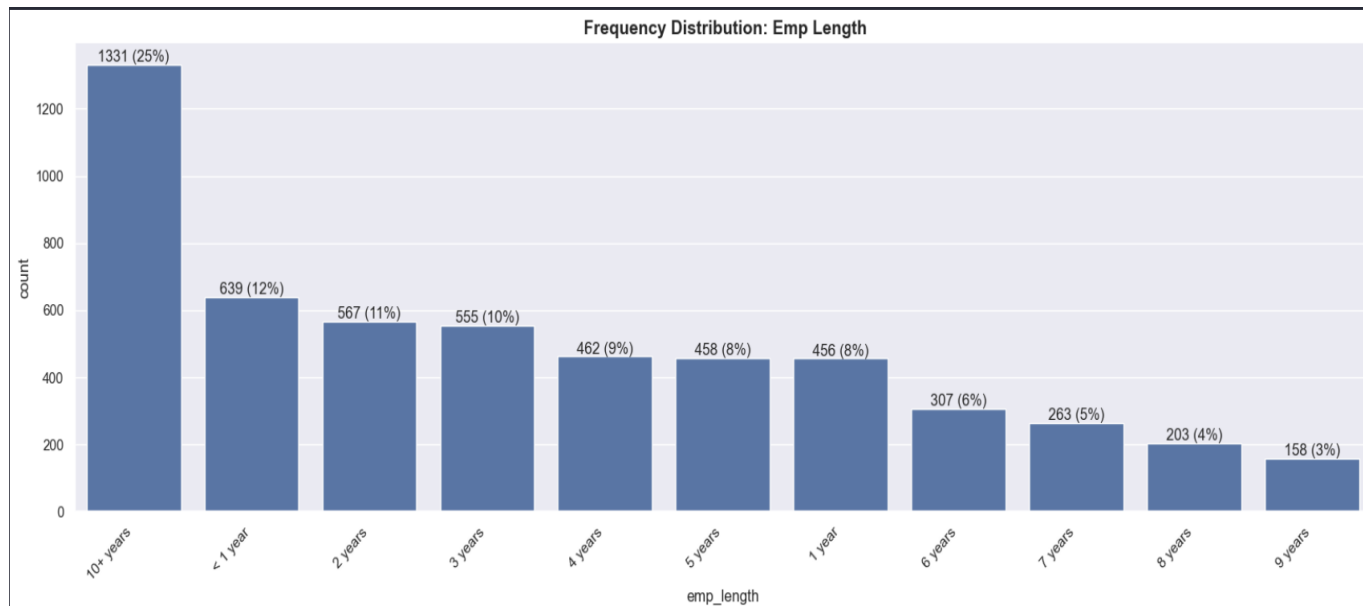
# Univariate Analysis (Only for Charged Off Customers) of Term, Grade, Emp Length, Home Ownership and Verification Status



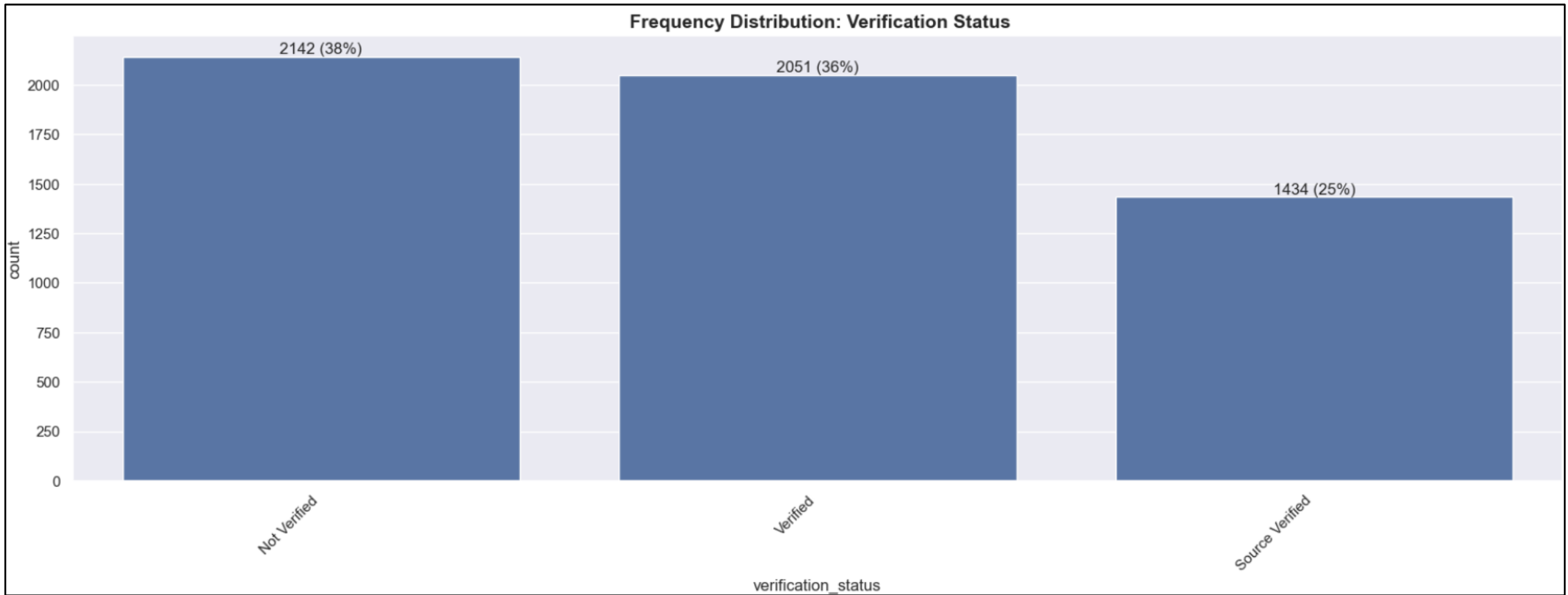




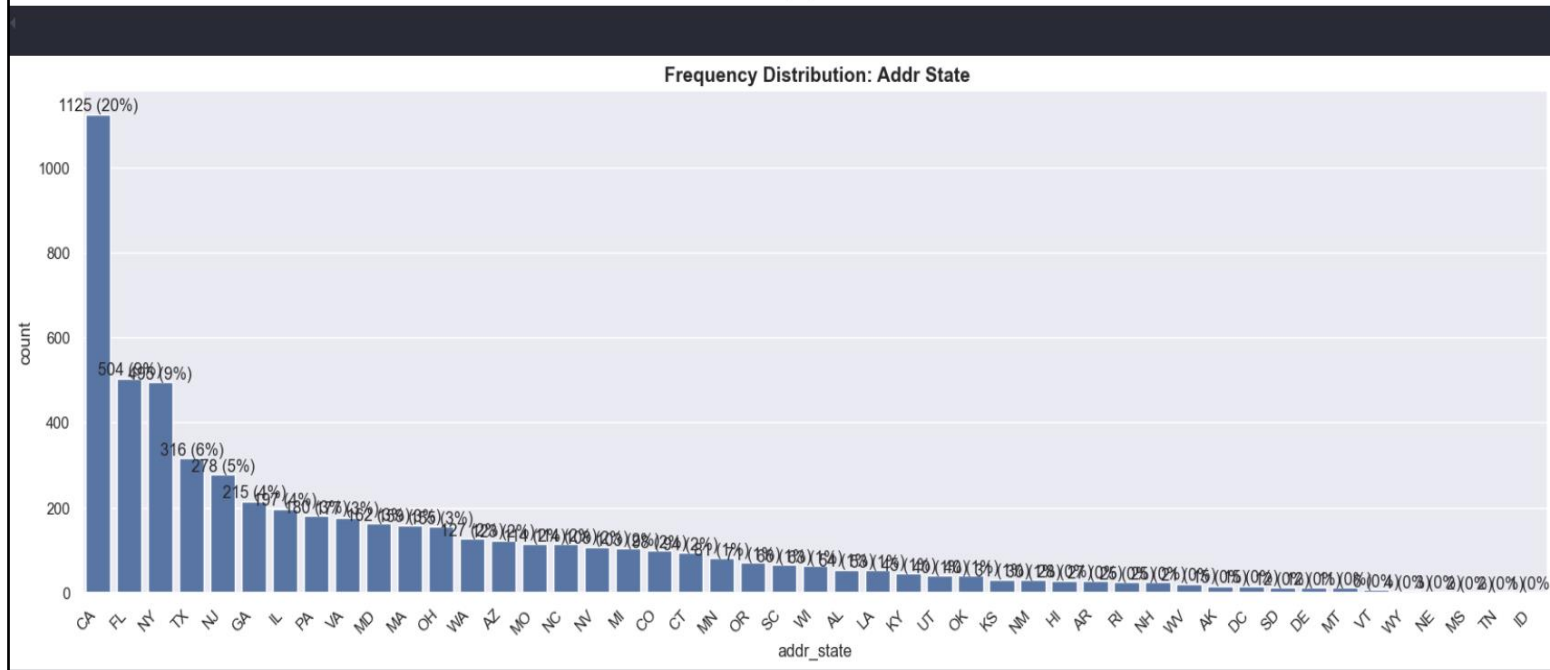
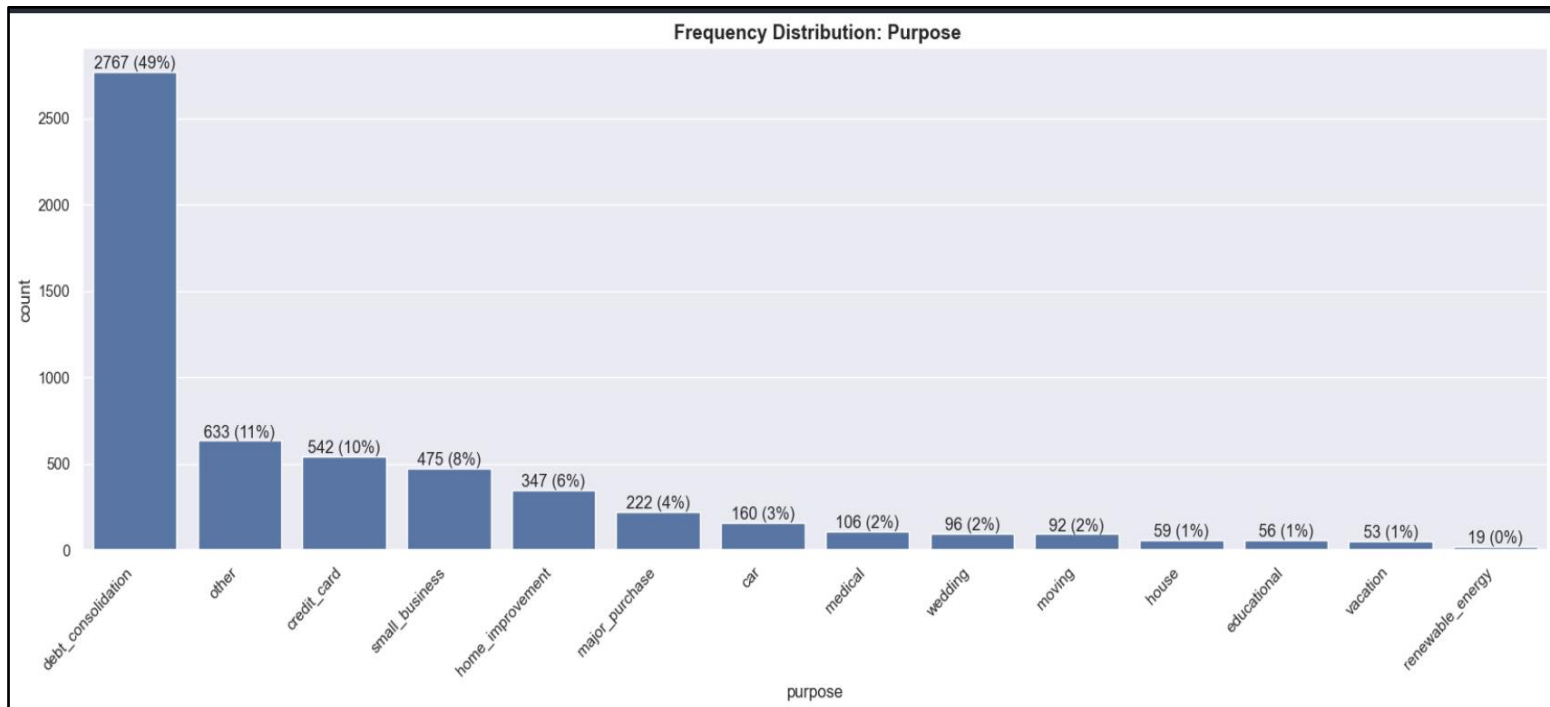
- There is a higher tendency to default when term of loan is for 36 months.
- There is a higher tendency to default when grade of customer is either B or C



- There is a higher tendency to default when employment length is 10+ years.
- There is a higher tendency to default when home ownership is either on Rent or Mortgage

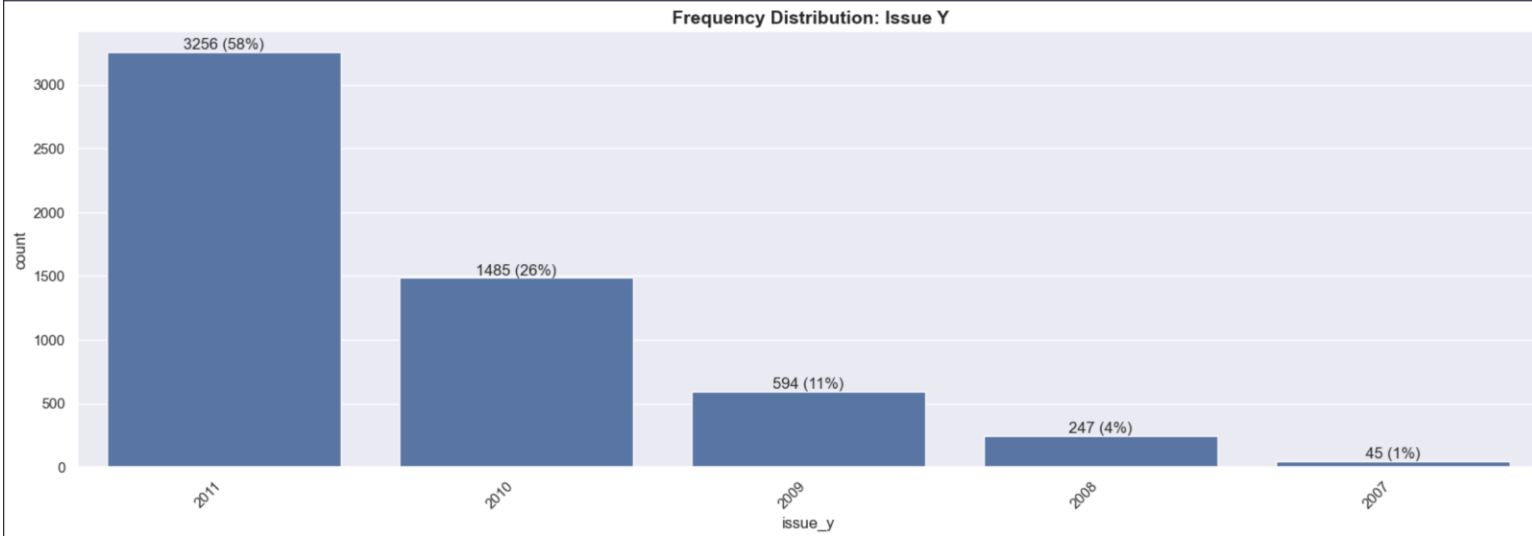
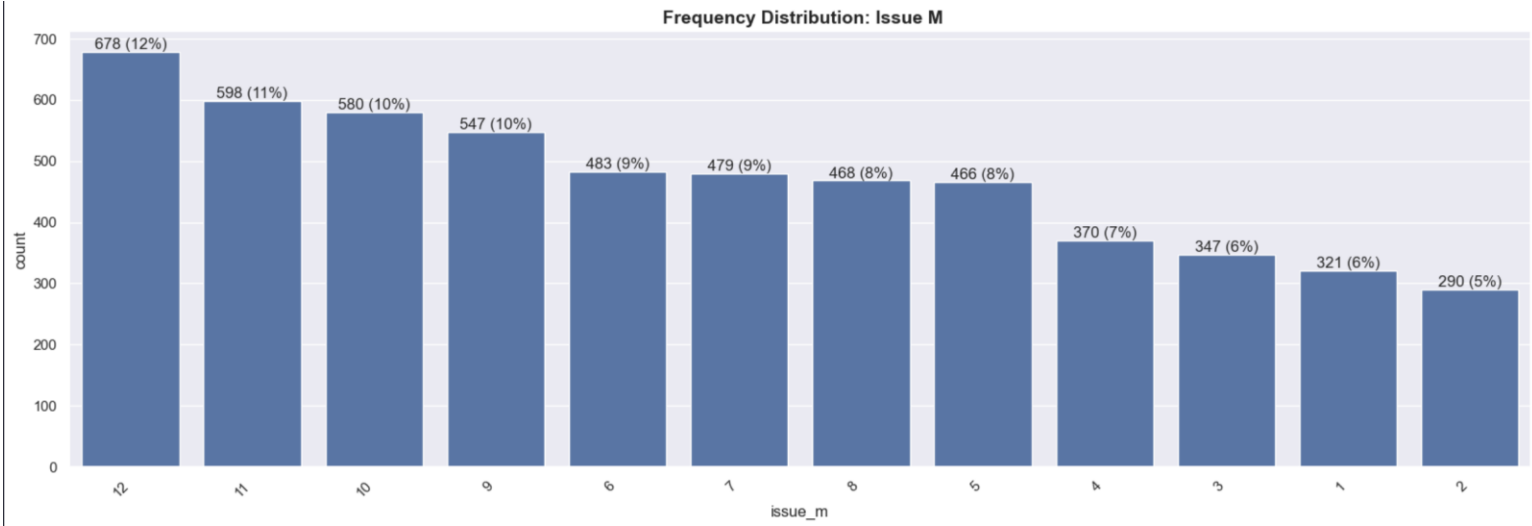


- There is a higher tendency to default when verification status is not 'Source Verified'

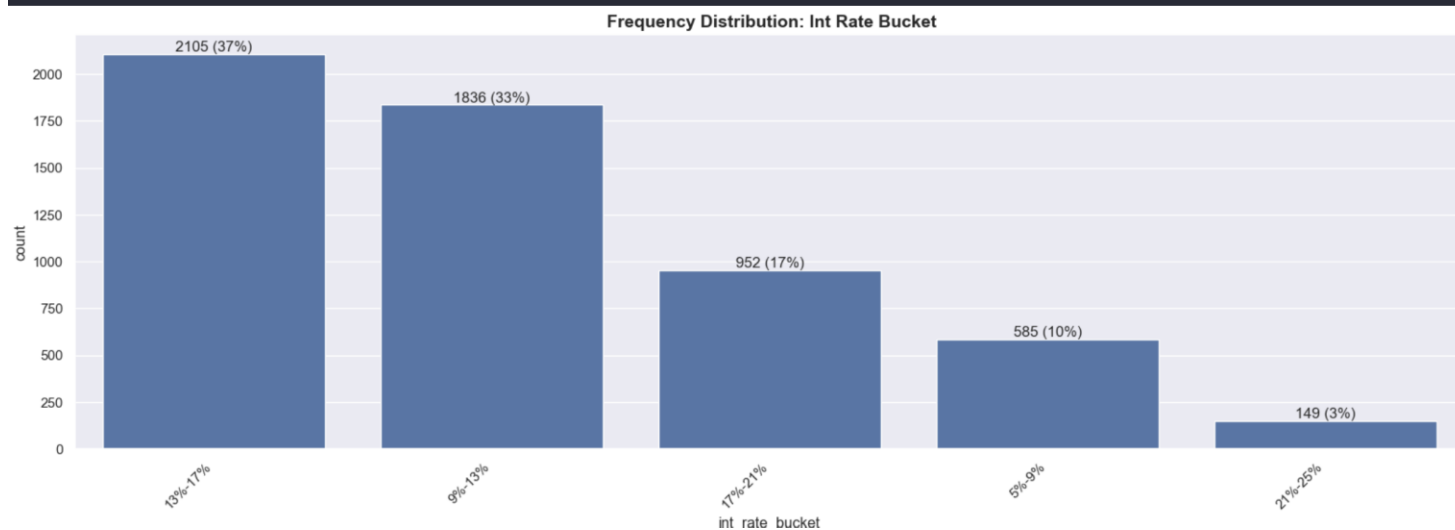
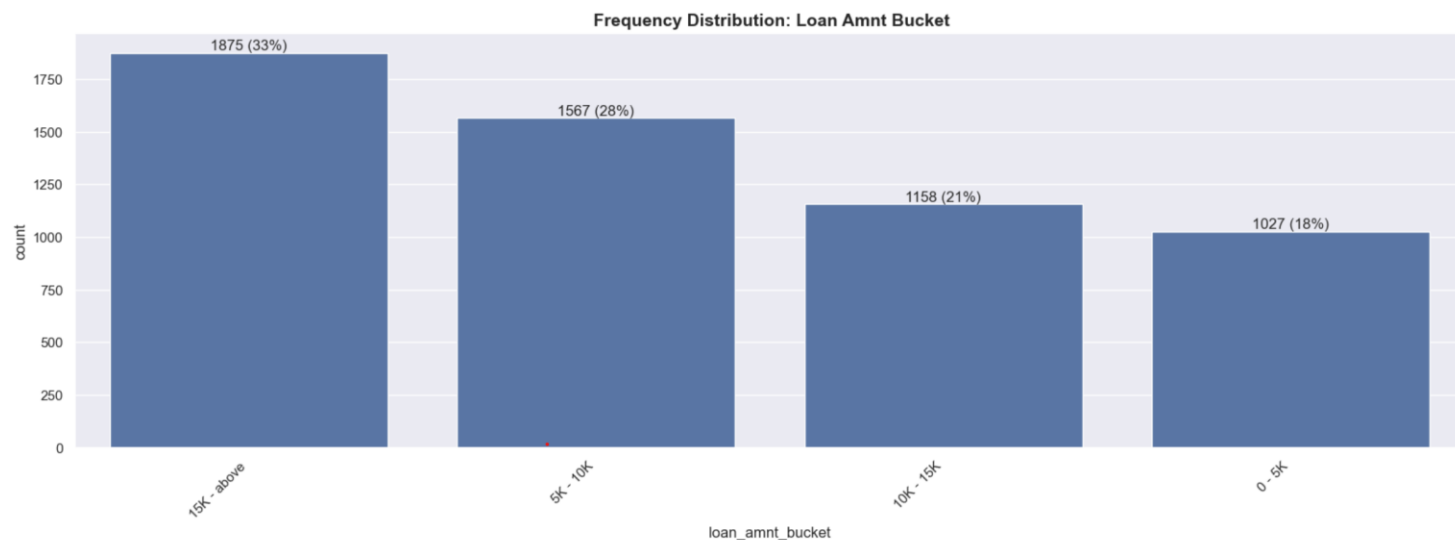


- There is a higher tendency to default when purpose is of loan is Debt Consolidation
- There is a higher tendency to default when Customer is from CA, FL or NY

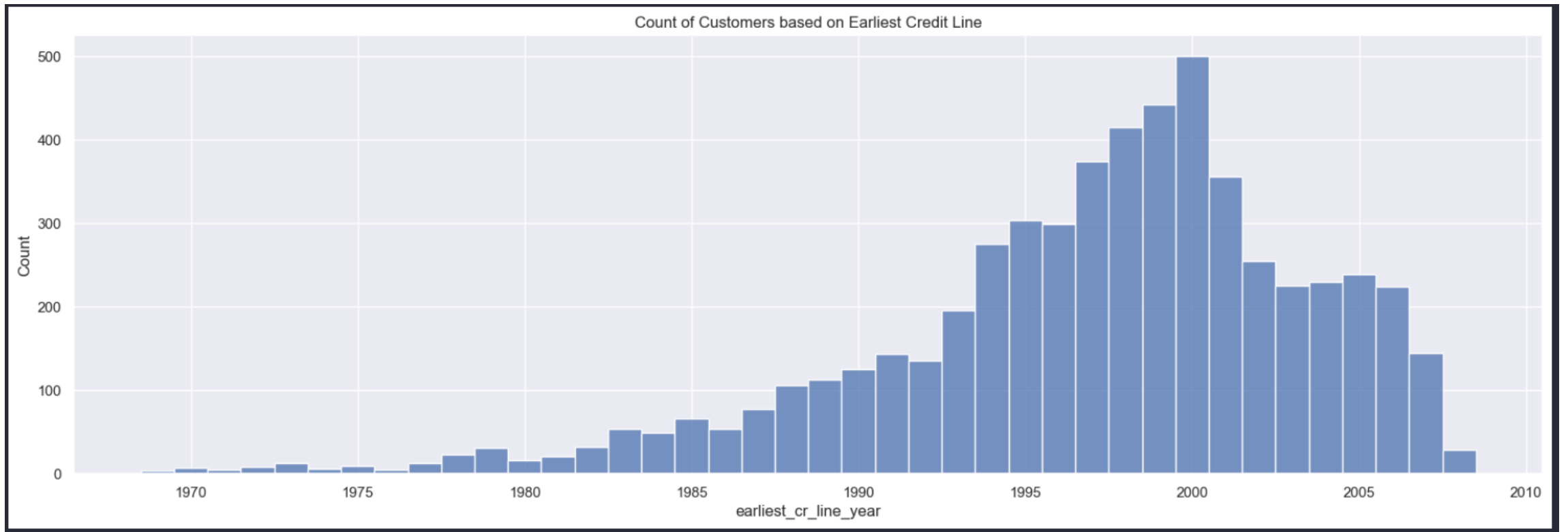
Note: CA: California , FL: Florid, NY: New York



- Higher tendency to default was observed for year 2011 and for later months of the year

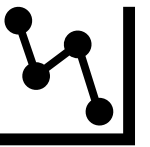
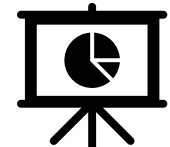
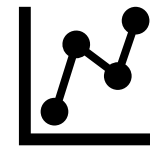


- Higher tendency to default is observed for customers with higher loan amount i.e. 15k – above.
- Higher tendency to default is observed for int bucket 13% to 17%

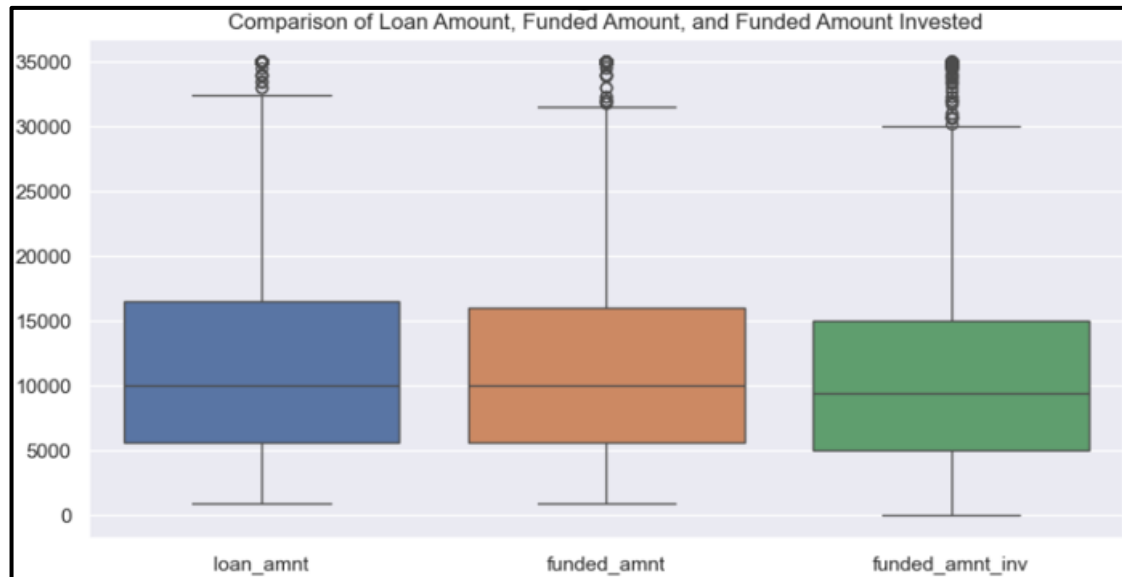


The above frequency distribution chart clearly shows most of the default customers were the ones which have there first credit line year between 1997 - 2000

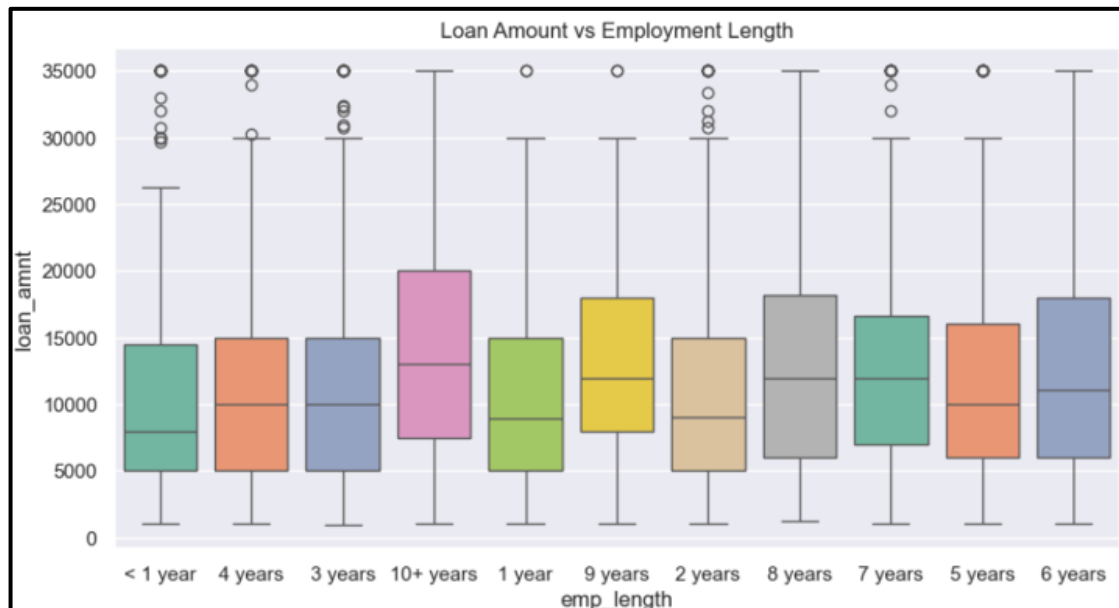
# Bivariate Analysis (for Charged Off Customers) of Employment Length, State with respect to Loan Amount



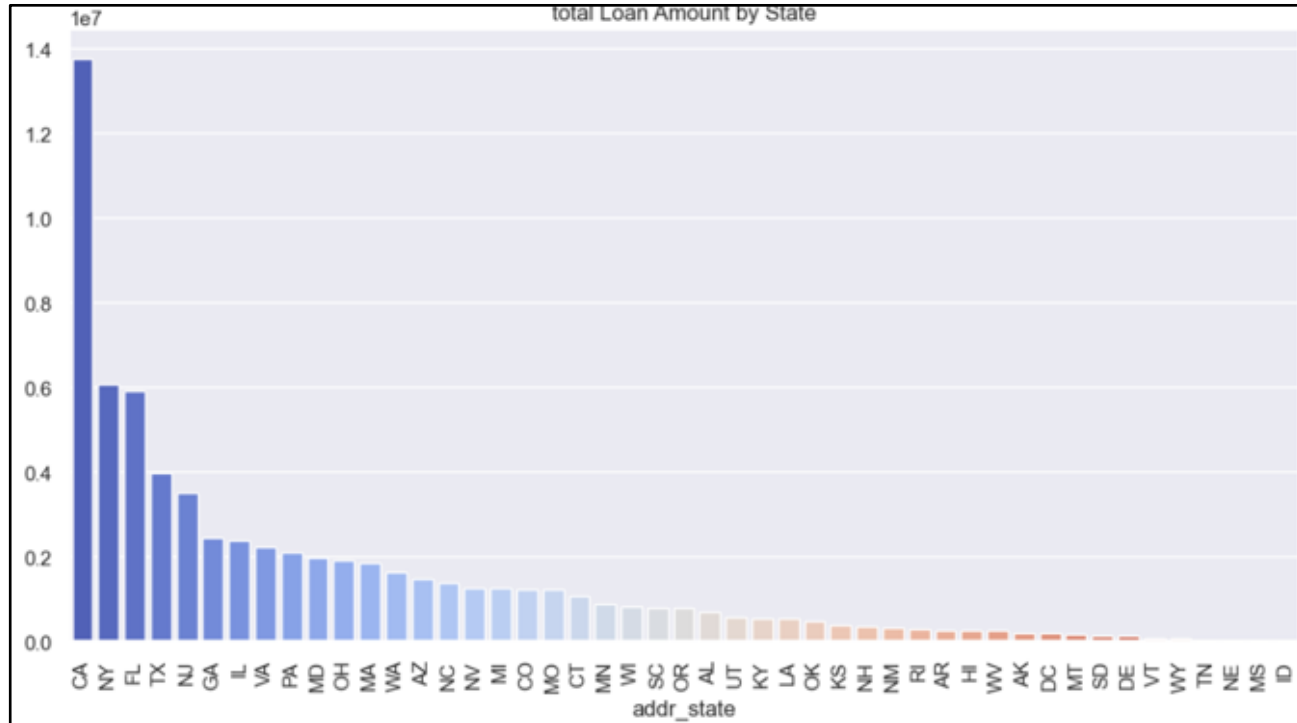




- There is no much difference in loan amount, funded amount and funded amount investor this behavior is in sync with the overall analysis as well so we can infer for defaulters the behavior is same there is no much difference between these amounts



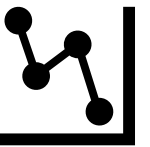
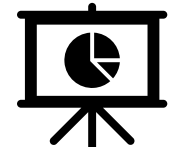
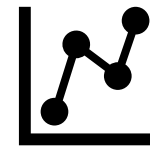
- The Loan amount Average value is higher for customers with employment length 10+, 9 and 8 years. So, this infer the customer that have higher employment length intend to take more loan

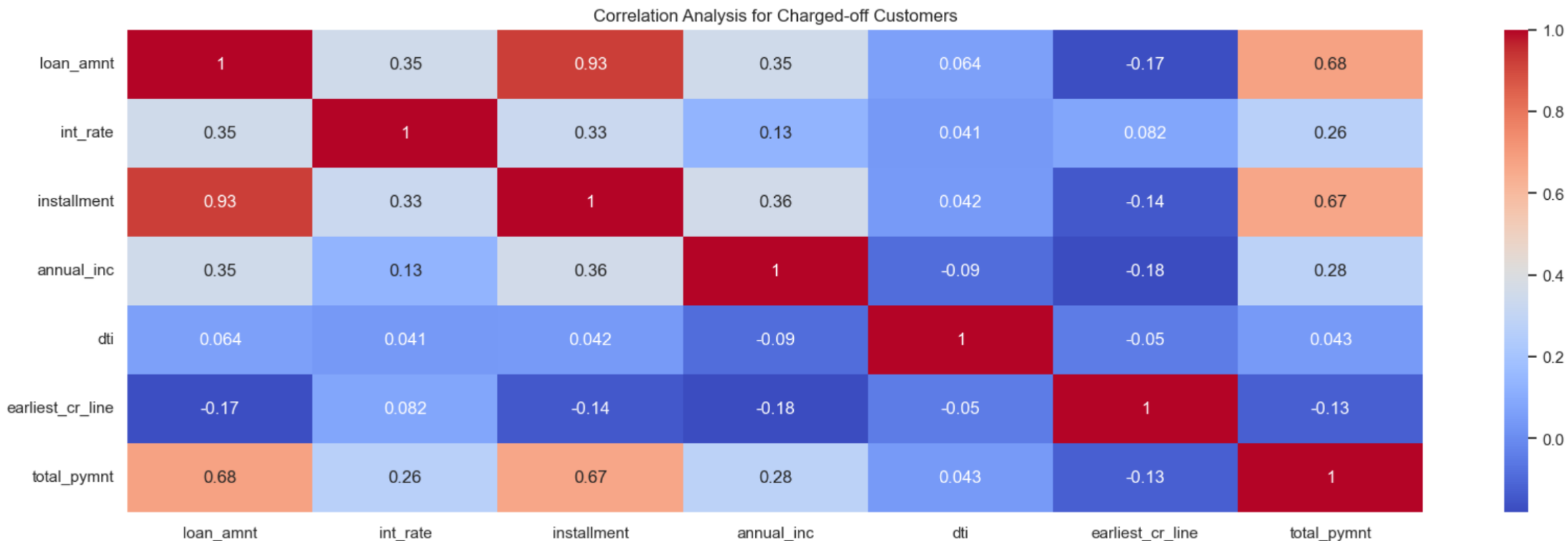


- Customers from State : CA, FL or NY have tendency to take more loan amount and this is in sync with total count of customers as well.

Note: CA: California , FL: Florid, NY: New York

Multivariate Analysis: Analysing multiple variables and checking correlation between them on a heat map





The above heat map shows there is not much correlation between the numerical variables except a high correlation between loan amount and instalment i.e. 93% but the correlation between loan amount and total payment is only 68% which clearly infer that charged off customers failed to pay the full payment on time.