

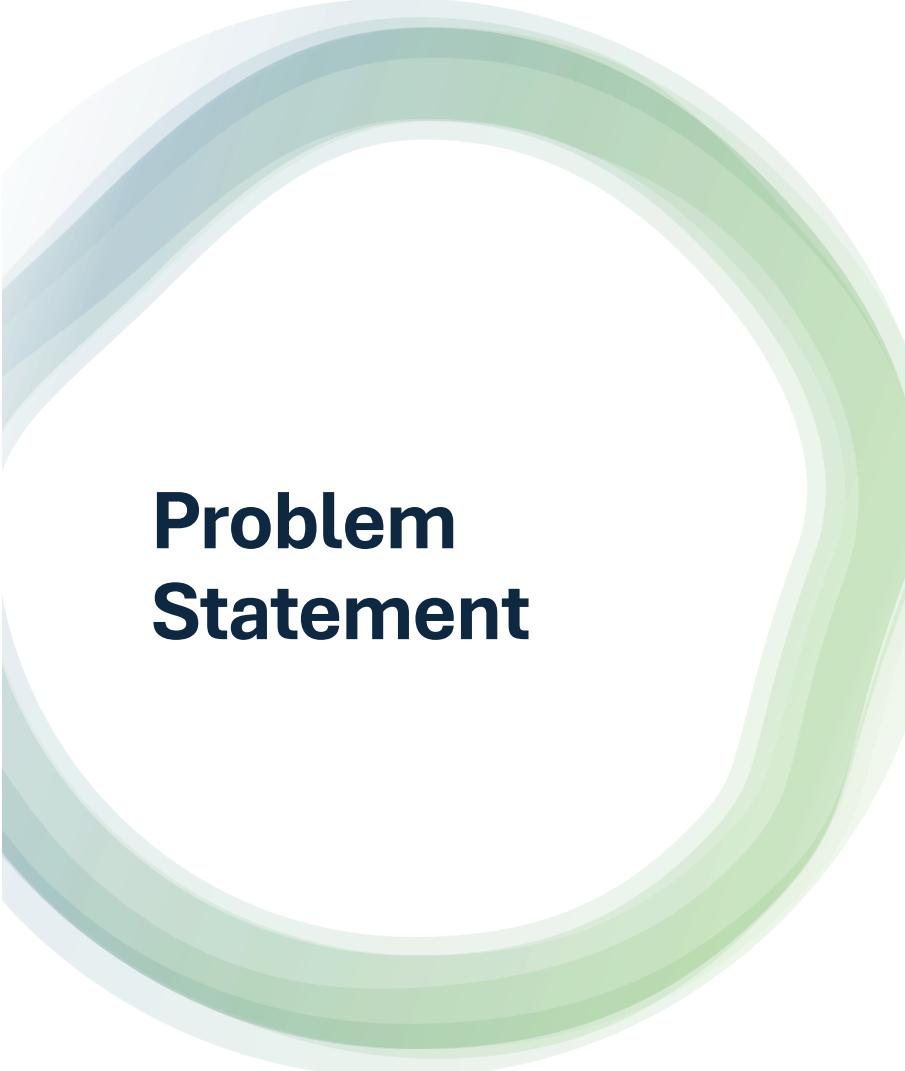
Lending Club Case Study

EDA

Naresh Padiyar



- **Problem Statement**
- **Data Summary**
- **Data Cleaning**
- **Data conversions and Derived Columns**
- **Dropping / Imputing the Rows / Outliers**
- **Univariate Analysis**
- **Bivariate Analysis**
- **Multivariate Analysis**
- **Correlations**
- **Summary**



Problem Statement

Problem

You work for a consumer finance company which specialises in lending various types of loans to urban customers. When the company receives a loan application, the company has to make a decision for loan approval based on the applicant's profile. Two types of risks are associated with the bank's decision:

- If the applicant is likely to repay the loan, then not approving the loan results in a loss of business to the company
- If the applicant is not likely to repay the loan, i.e. he/she is likely to default, then approving the loan may lead to a financial loss for the company

Objective

The aim is to identify patterns which indicate if a person is likely to default, which may be used for taking actions such as denying the loan, reducing the amount of loan, lending (to risky applicants) at a higher interest rate, etc.

Constraints

When a person applies for a loan, there are two types of decisions that could be taken by the company:

- ❖ **Loan accepted:** If the company approves the loan, there are 3 possible scenarios described below:
 - **Fully paid:** Applicant has fully paid the loan (the principal and the interest rate)
 - **Current:** Applicant is in the process of paying the instalments, i.e. the tenure of the loan is not yet completed. These candidates are not labelled as 'defaulted'.
 - **Charged-off:** Applicant has not paid the instalments in due time for a long period of time, i.e. he/she has defaulted on the loan
- ❖ **Loan rejected:** The company had rejected the loan (because the candidate does not meet their requirements etc.).



Dataset Summary

- ✓ The **Loan.csv** file contains **111 columns and 39717 Rows**.
- ✓ The dataset contains **Loan Attributes and Customer Attributes**.
- ✓ The **Data_Dictionary.xls** file provided along with the dataset, has a description of each column of the **Loan.csv** file and was used to make a decision if the column would contribute to the EDA analysis.

Data Cleaning

- ✓ There are **54** columns which have all rows containing NA/Null values. These columns were removed/dropped.
- ✓ Columns **id, member_id and url** have unique value in every row so these 3 columns were also dropped.
- ✓ Columns **desc, title and emp_title** contain descriptions so these columns were dropped as they won't contribute to the EDA analysis.
- ✓ Columns **pymnt_plan, initial_list_status, collections_12_mths_ex_med, policy_code, application_type, acc_now_delinq, chargeoff_within_12_mths, delinq_amnt and tax_liens** have only 1 unique value and hence were dropped as they won't contribute to the EDA analysis.
- ✓ Columns **mths_since_last_delinq, mths_since_last_record, and next_pymnt_d** were dropped as they were having more than **50%** NA/Null Rows.
- ✓ Columns **delinq_2yrs, earliest_cr_line, last_pymnt_amnt, inq_last_6mths, open_acc, pub_rec, revol_bal, revol_util, total_acc, 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_d, last_credit_pull_d** were dropped as these columns either represent behavioural data or post loan approval data, hence these columns were dropped as they won't contribute to the EDA analysis.

Data Conversion and Derived Columns

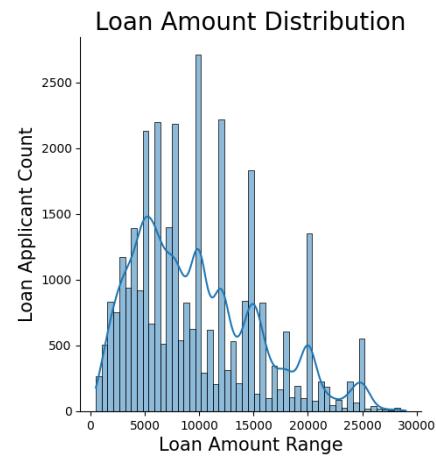
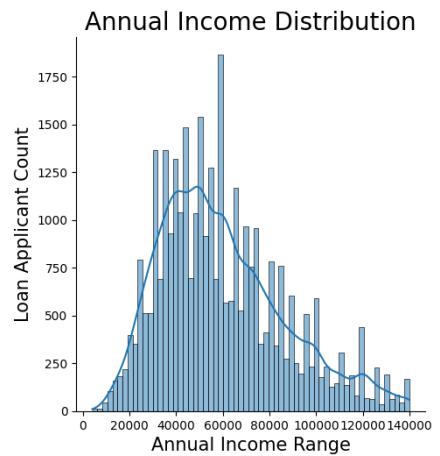
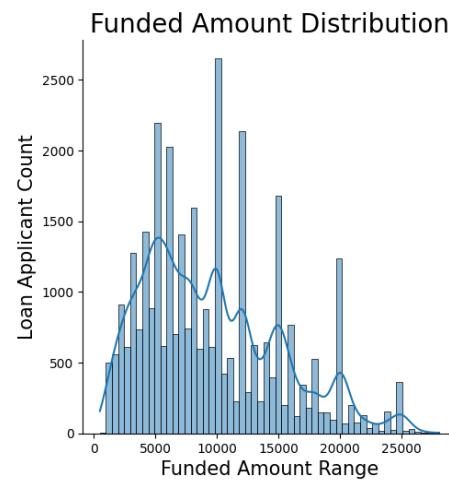
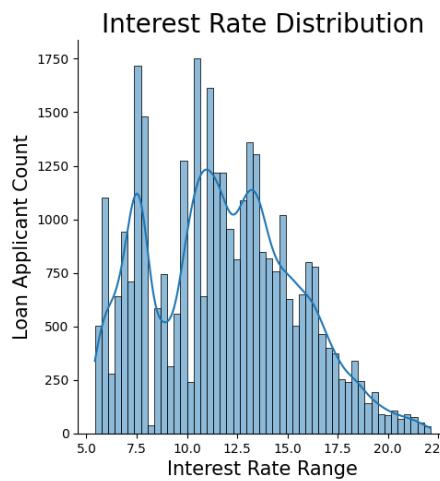
- ✓ Additional string value in **term** column has been trimmed and column has been converted to int data type.
- ✓ Additional string value in **int_rate** column has been trimmed and column has been converted to int data type.
- ✓ Columns **loan_amnt** and **funded_amnt** have been converted to int data type.
- ✓ Column **issue_d** converted to datetime data type.
- ✓ Created derived columns **issue_month** and **issue_year** from **issue_d** column.
- ✓ Created derived columns **loan_amnt_bucket**, **annual_income_bucket**, **int_rate_bucket**, **dti_bucket**, **installment_bucket** and **funded_amnt_inv_bucket** from columns **loan_amnt**, **annual_income**, **int_rate**, **dti**, **installment** and **funded_amnt_inv** respectively for **bucketing** purpose.

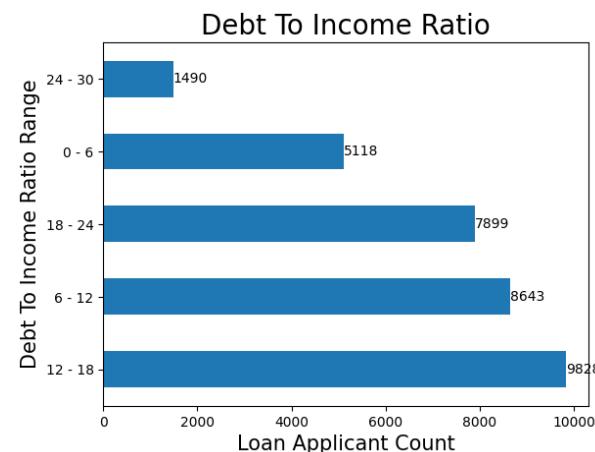
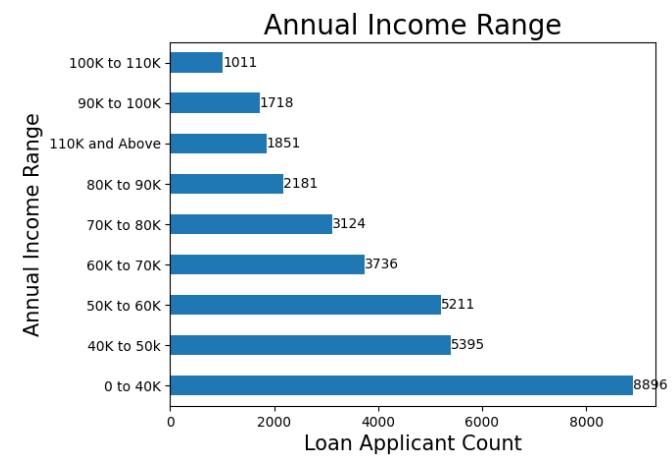
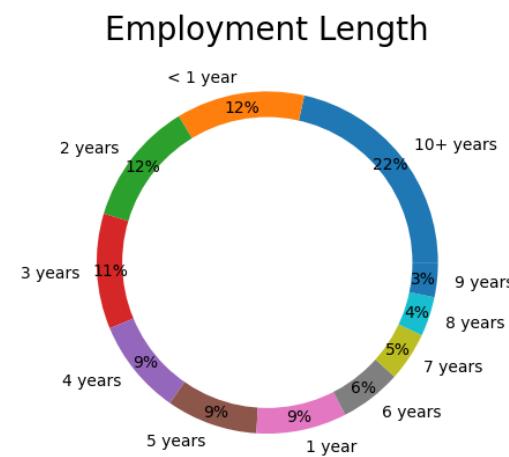
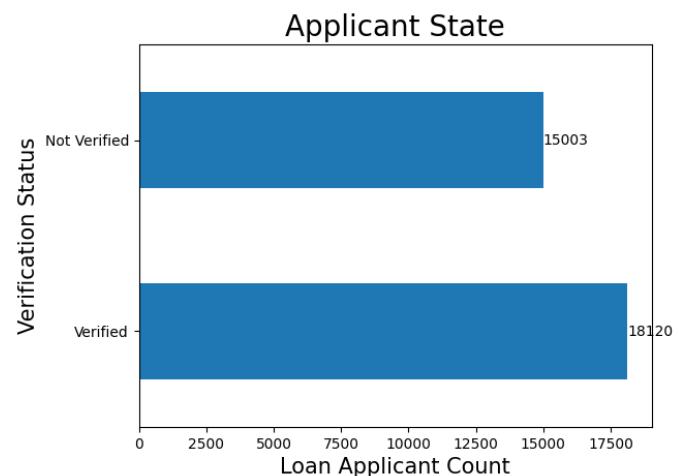
Dropping / Imputing Rows / Outliers

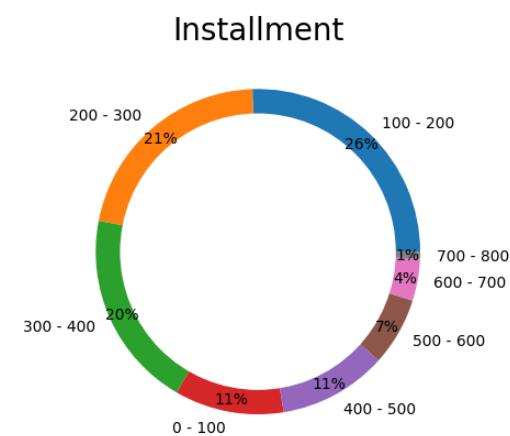
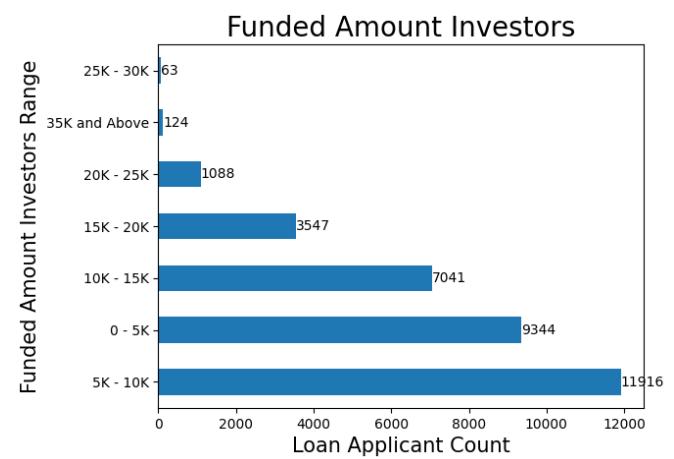
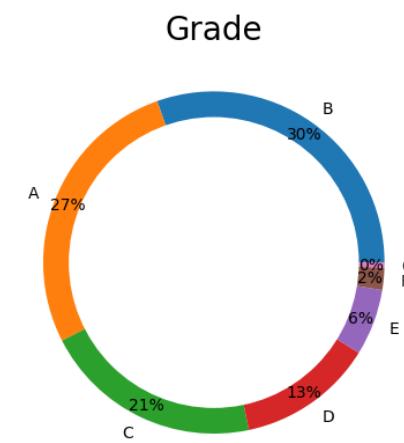
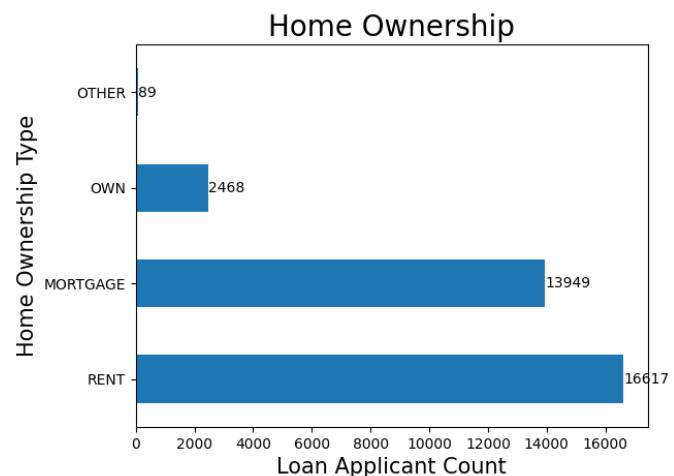
- ✓ Since the goal of the EDA is to identify who is likely to default and this can only be said in case of either fully paid or charged off loans, hence we dropped all rows where **loan_status** is equal to "Current"
- ✓ **verification_status** column has 3 categories **Not Verified**, **Verified** and **Source Verified**. **Verified** and **Source Verified** both mean the same thing, hence we have replaced the value **Source Verified** to **Verified**. This leaves us with only 2 categories in the **verification_status** column, i.e. **Not Verified** and **Verified**
- ✓ Columns **emp_length** and **pub_rec_bankruptcies** both contains 2.67% and 1.80% of rows as null respectively. We have dropped the rows where the value for these columns was NA/Null.
- ✓ Outliers exists for numeric data **loan_amnt**, **funded_amnt**, **funded_amnt_inv**, **int_rate**, **installment** and **annual_inc**. We have treated the outliers in these columns using the using quantile mechanism.

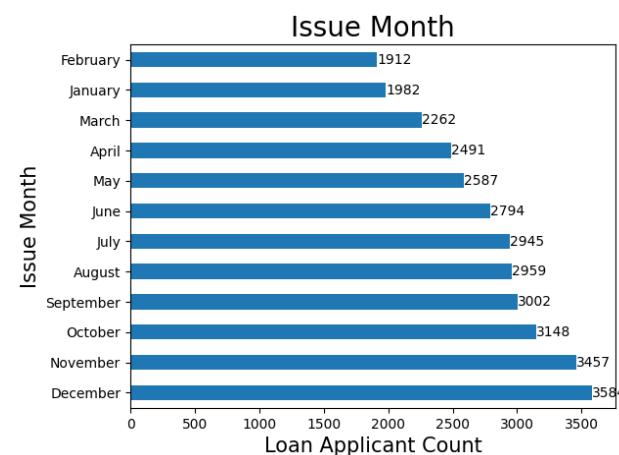
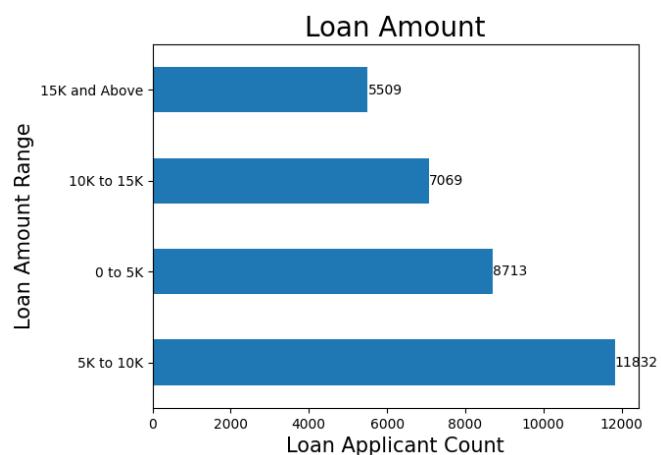
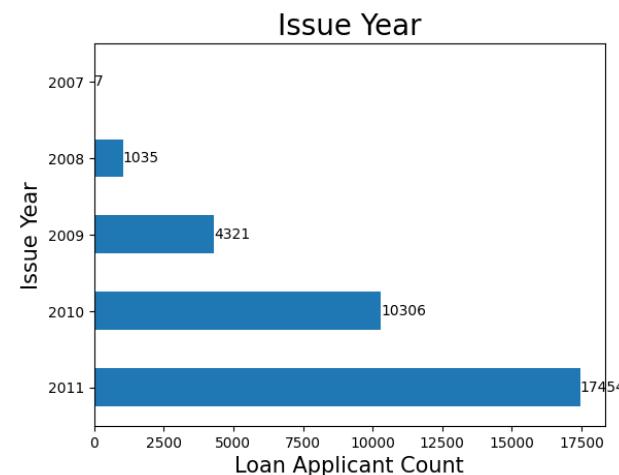
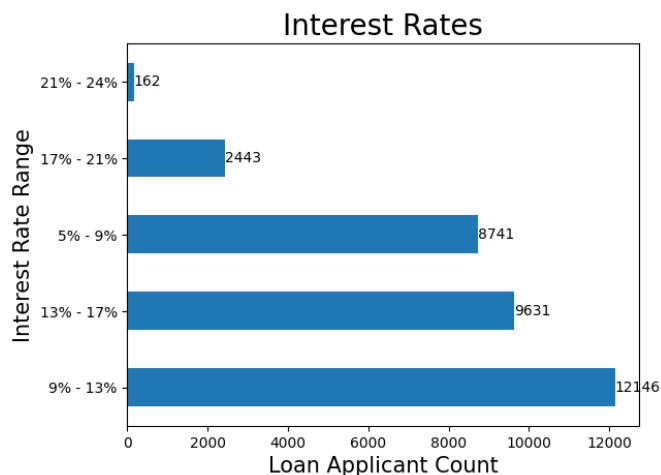


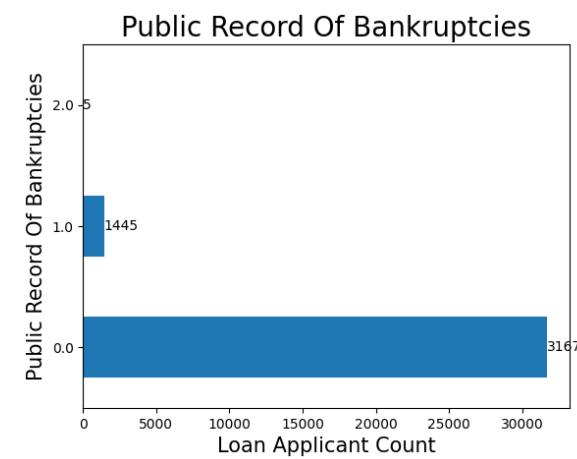
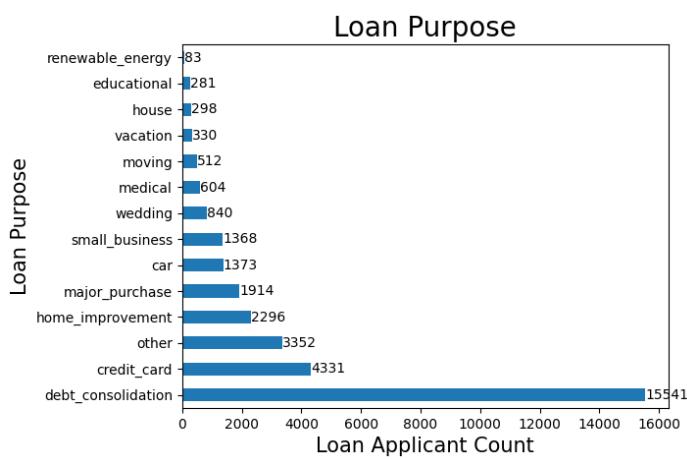
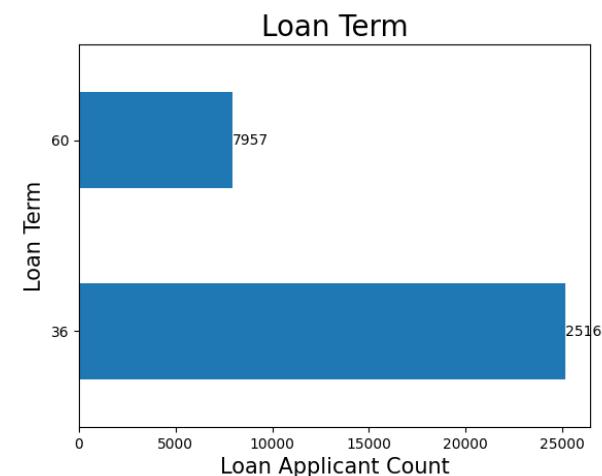
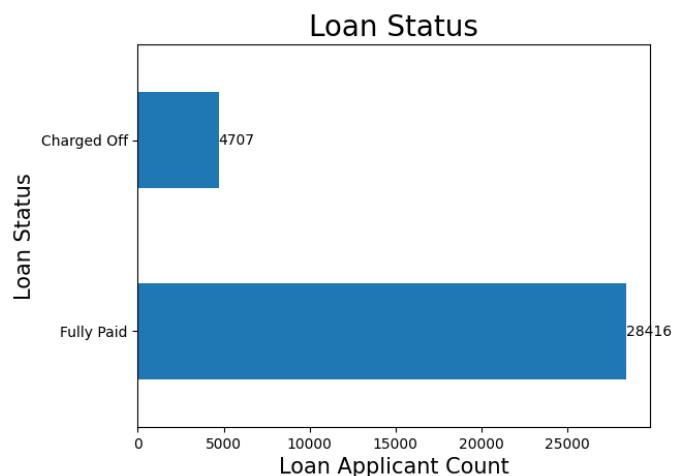
Univariate Analysis









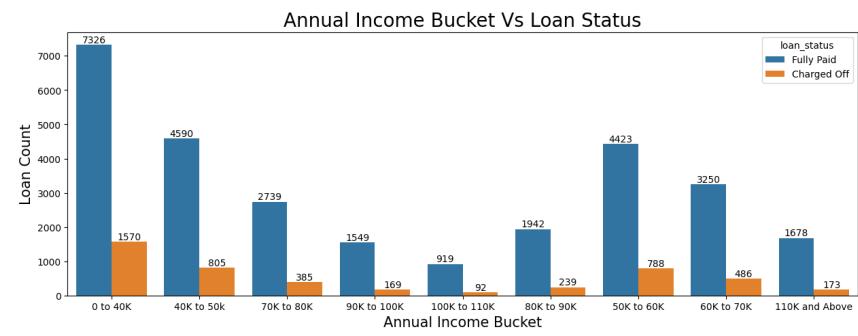
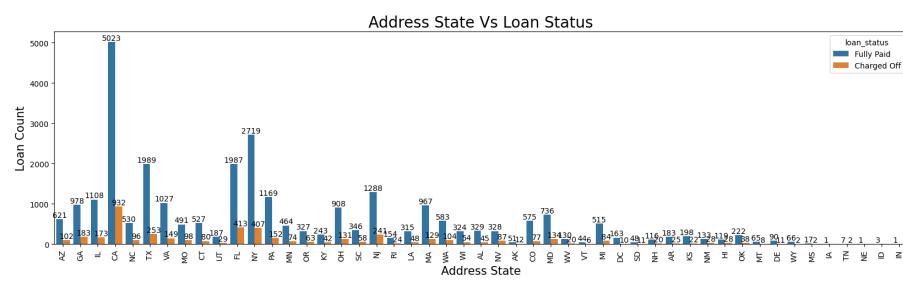
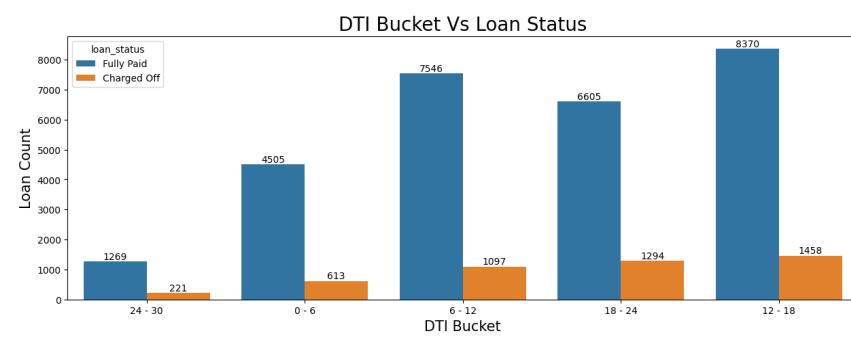
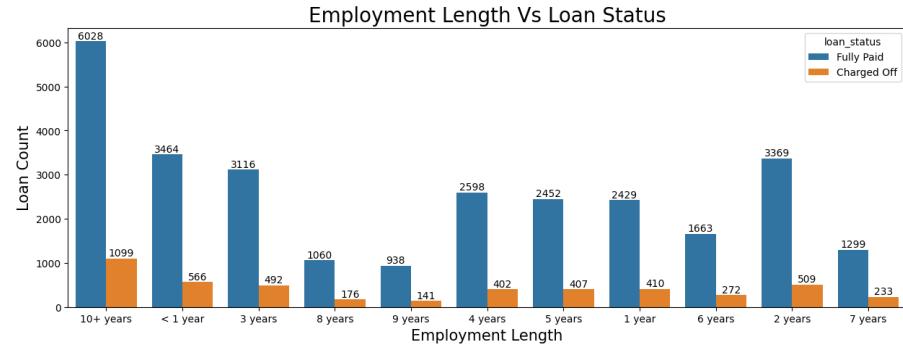


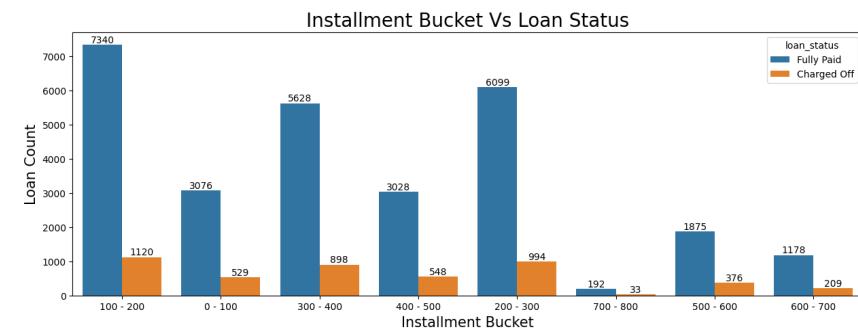
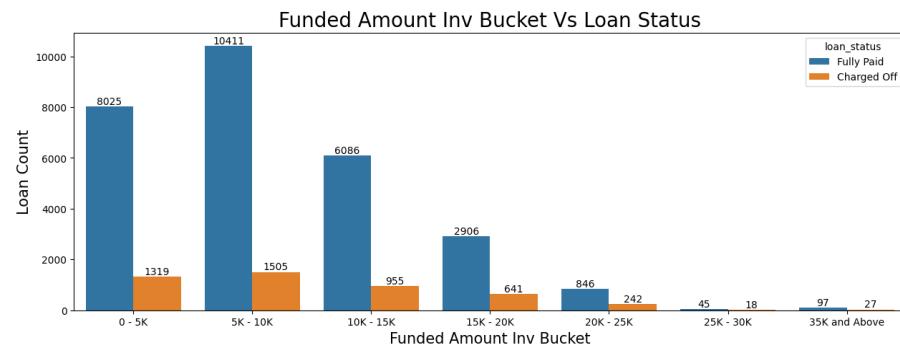
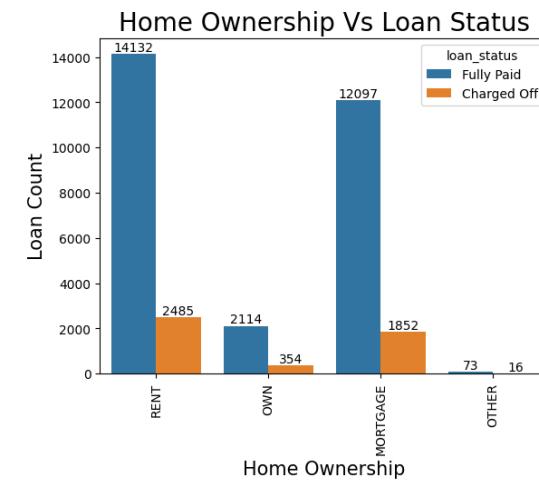
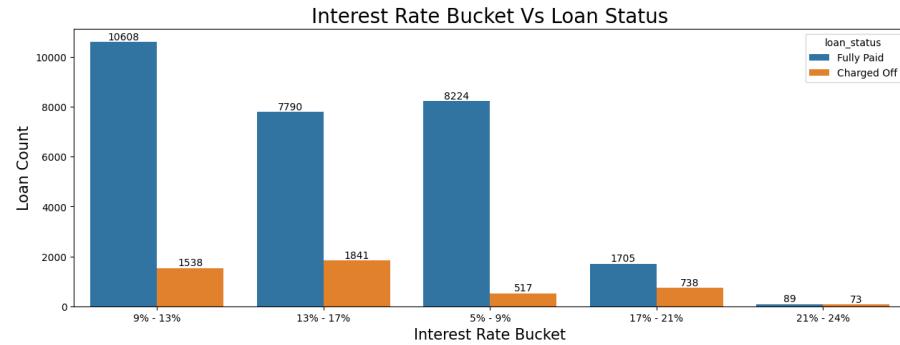
Conclusion from Univariate Analysis

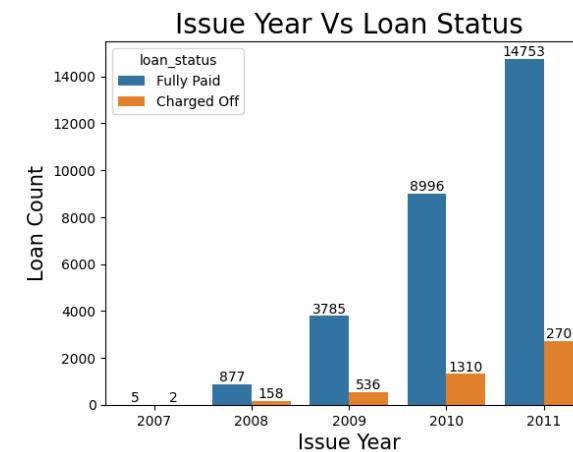
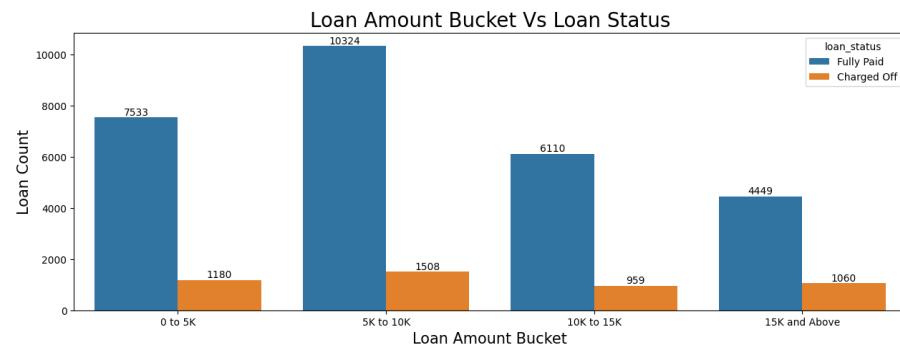
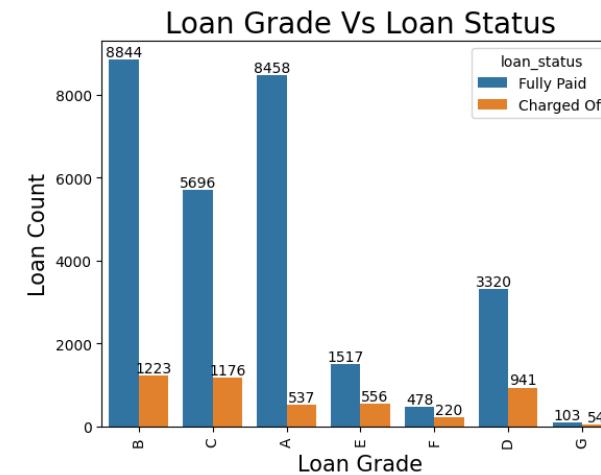
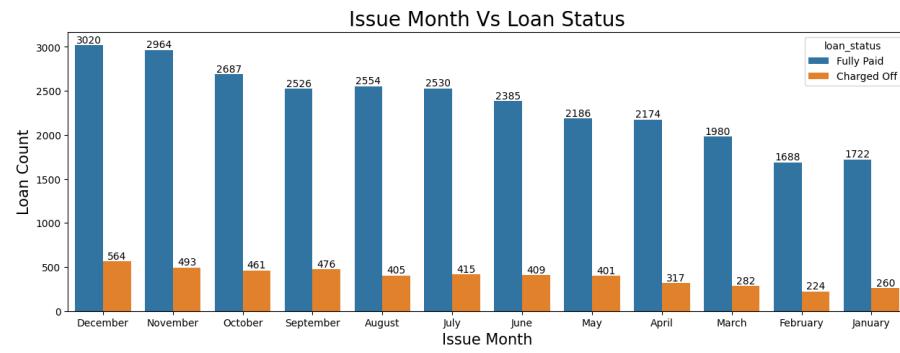
- ✓ 1. The annual income of most of the loan applicants is between 40k-75k and the average annual income is 59840.
- ✓ 2. Most of the loan amount applied was in the range of 5K to 14K and maximum loan amount applied was ~27K.
- ✓ 3. The funded amount for most of the loan applicants is between the range of 5K to 13K and maximum funded amount is ~25K.
- ✓ 4. Most of the loan applicants have interest rate between range of 8% - 14%. The average rate of interest is 11.7%.
- ✓ 5. Majority of loan applicants are living either on Rent or on Mortgage.
- ✓ 6. The purpose of most of the loan applicants is debt_consolidation.
- ✓ 7. Majority of the loan applicants have status as Fully Paid.
- ✓ 8. Highest number of loan applicants are from CA state.
- ✓ 9. Majority of the loan applicants have 10+ years of experience..
- ✓ 10. Majority of the loan applicants belong to Grade B.
- ✓ 11. Majority of the loan applicants will payment installment between 100-200.
- ✓ 12. Majority of the investors have committed for an amount between 5K to 10K.
- ✓ 13. Majority of the loan applicants belong to Verified status.
- ✓ 14. Majority of the loans were funded in the month of December, followed by November and then October.
- ✓ 15. Majority of loans funded were issued in the year 2011.
- ✓ 16. Majority of the Loans have a term of 36 months.
- ✓ 17. Annual Income of majority of the loan applicants is below 40K.
- ✓ 18. Public Record of Bankruptcies for majority of the loan applicants is 0.
- ✓ 19. Majority of the loan applicants have taken loan amount between 5K to 10k.
- ✓ 20. Majority of the loan applicants have debt to income ratio between the range of 12 – 18.

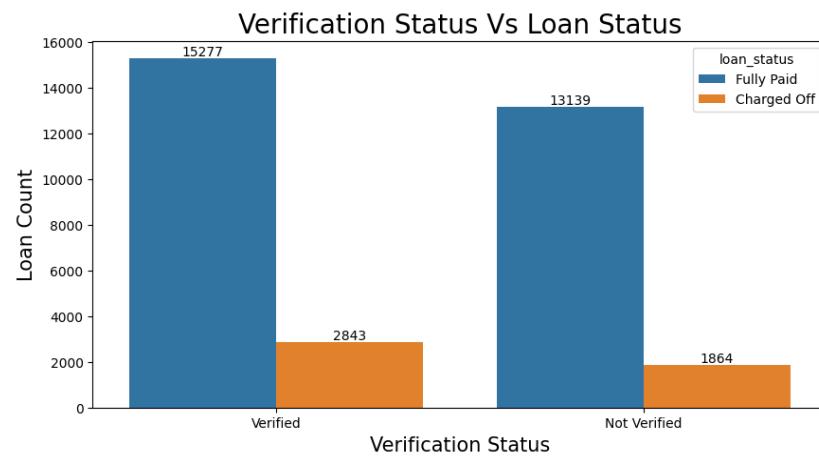
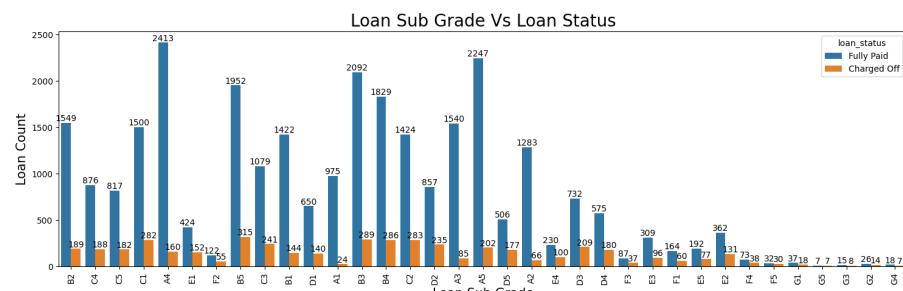
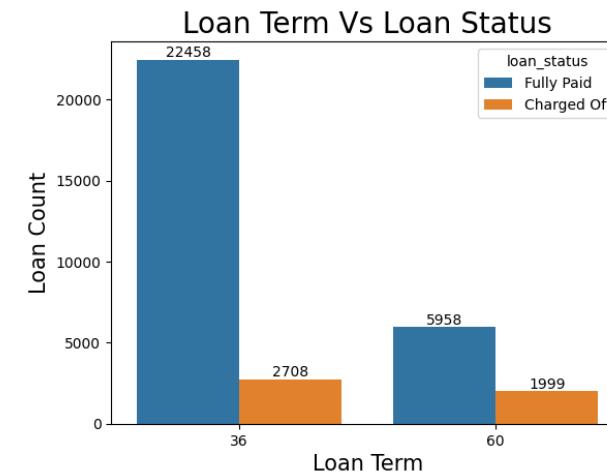
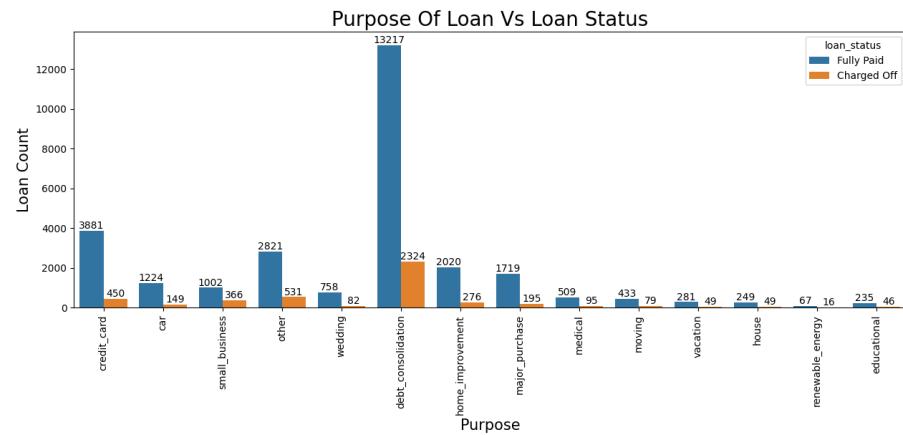


Bivariate Analysis







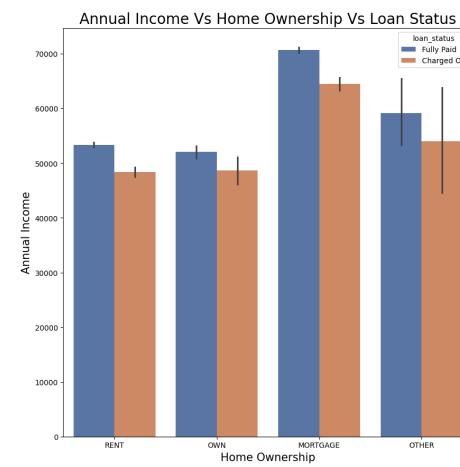
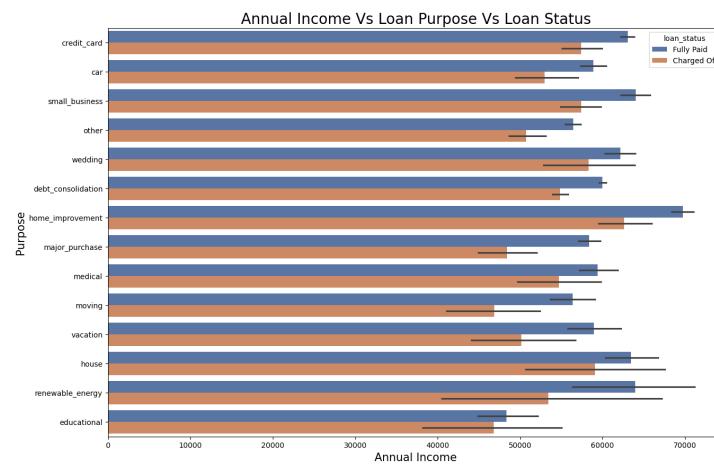
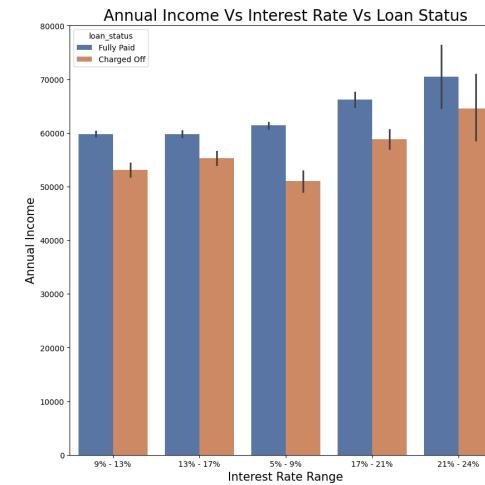
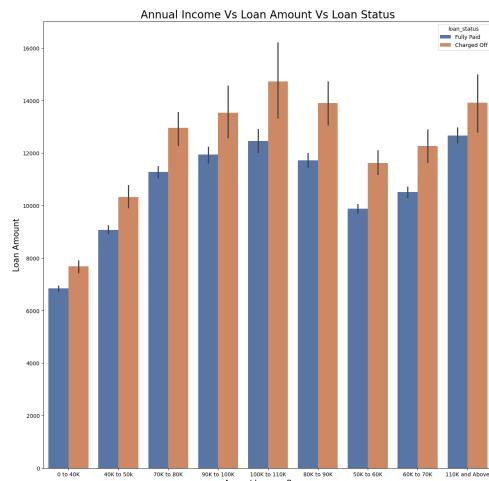


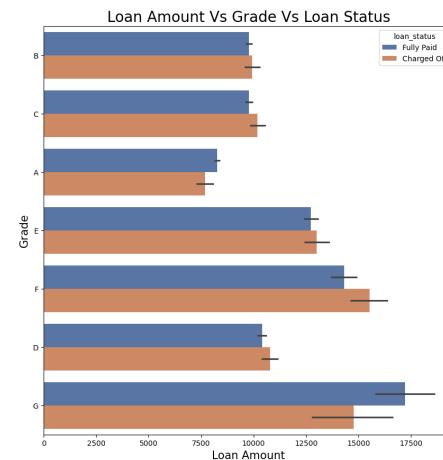
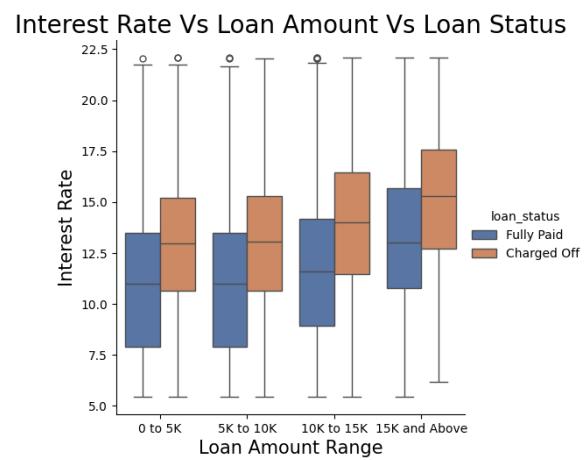
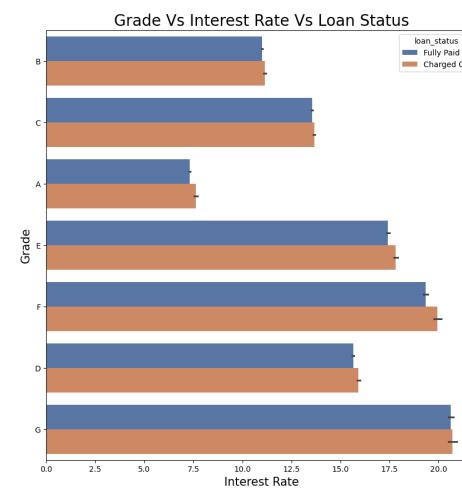
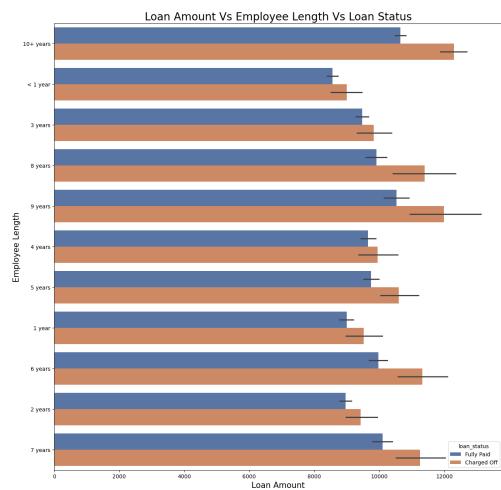
Conclusion from Bivariate Analysis

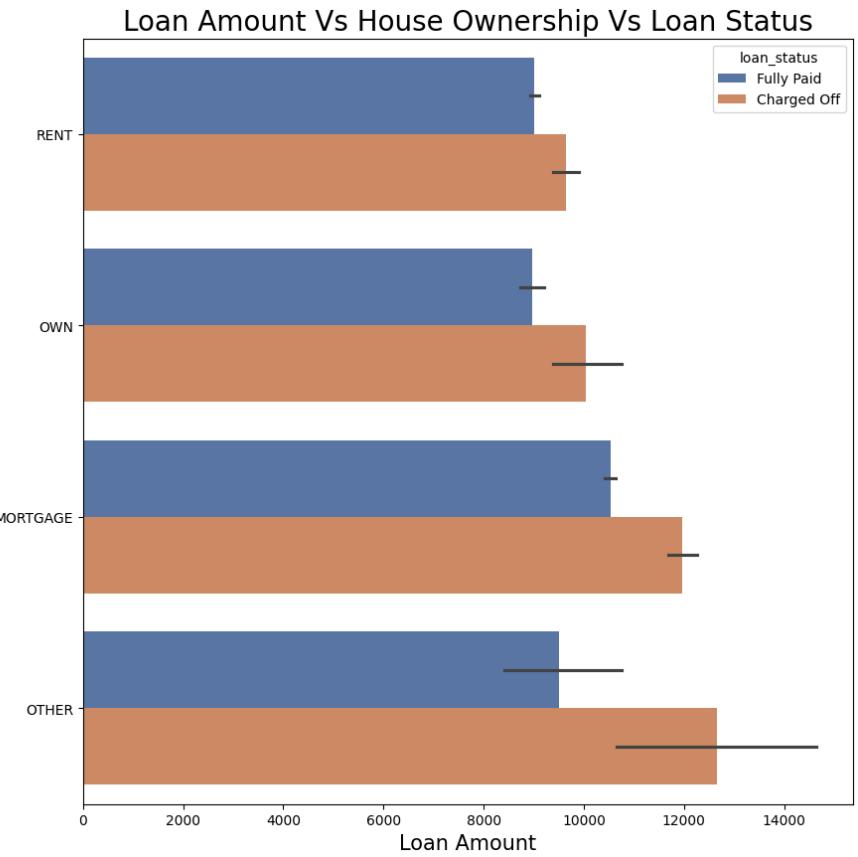
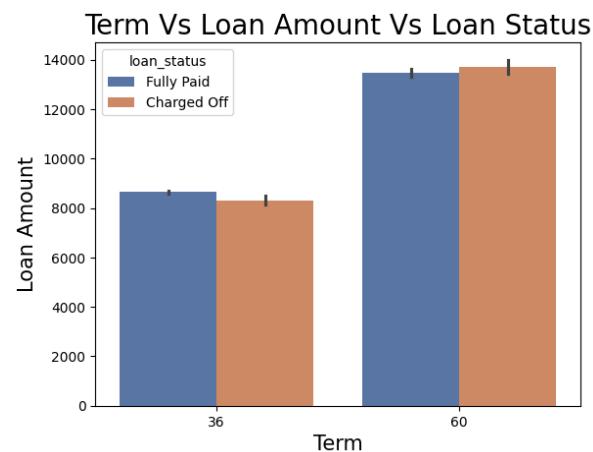
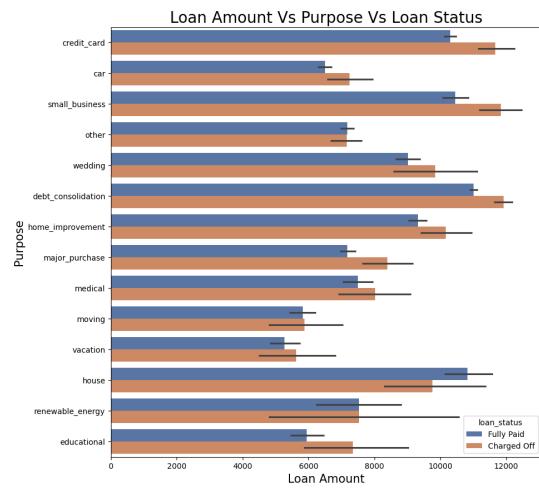
- ✓ 1. Debt Consolidation is the category where the maximum number of loans are 36-month and people have defaulted in the same category.
- ✓ 2. Loan applicants with home ownership as RENT and MORTGAGE are more likely to default.
- ✓ 3. Verified loan applicants are defaulting more than who are Not Verified.
- ✓ 4. Loan applicants from states of California (CA), Florida (FL), and New York (NY) are most likely to default.
- ✓ 5. Loan applicants belonging to Grades B, C and D contribute to the most of the Charged Off loans.
- ✓ 6. Loan applicants belonging to Sub Grades B3, B4 and B5 contribute to the most of the Charged Off loans.
- ✓ 7. Loan applicants applying for 36-month term are more likely to default than those taking loans for 60 months.
- ✓ 8. Loan applicants with more than 10 years employment length contribute the most to the charge off loans.
- ✓ 9. The number of loan applicants have steadily increased from 2007 to 2011, indicating a positive trend in upcoming years, 2011 year has the highest number of loan applicants and defaulters.
- ✓ 10. December is the most preferred month for taking loans.



Multivariate Analysis



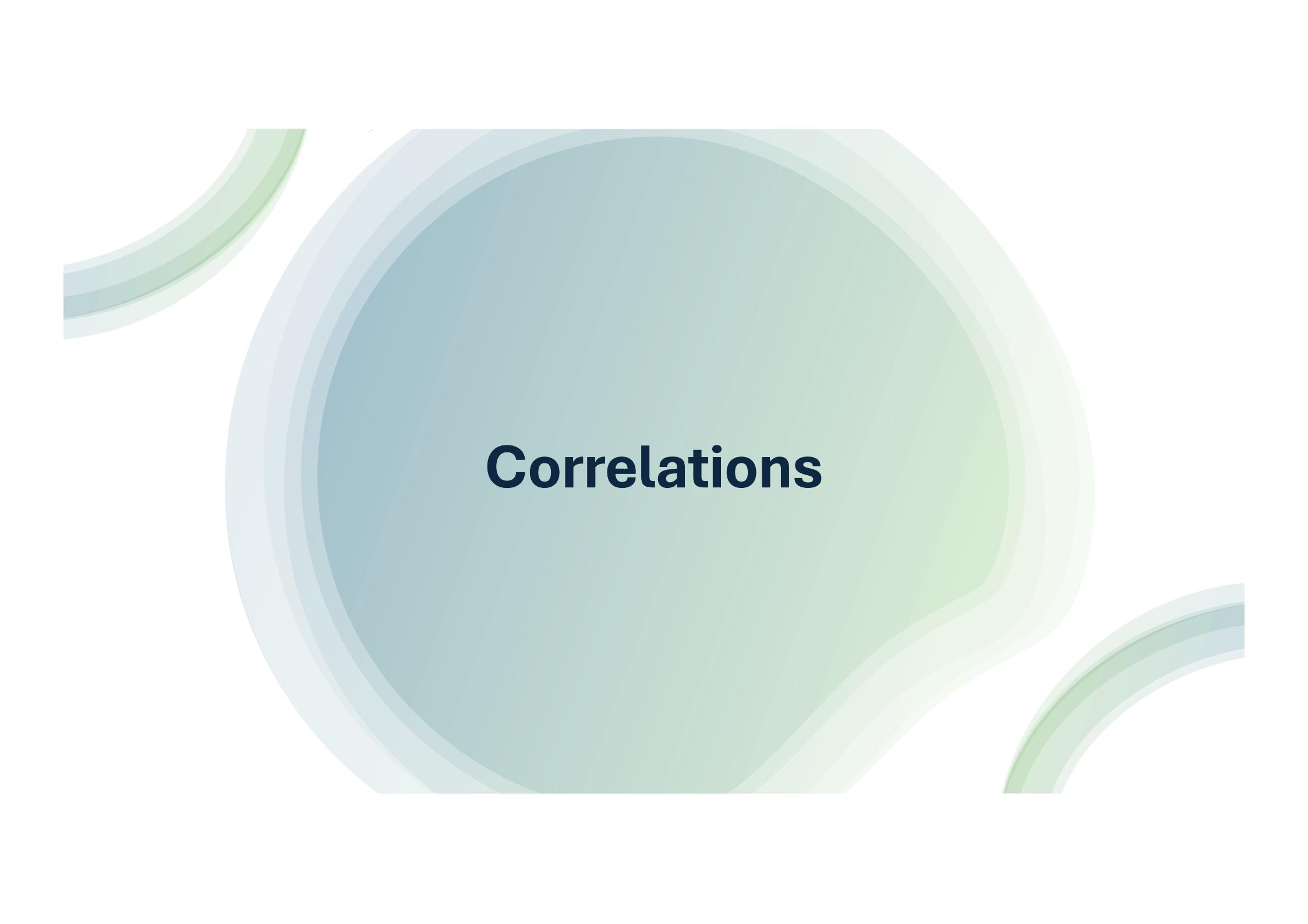




Conclusions from Multivariate Analysis

The below analysis is with respect to the charged off loans. There is a more probability of defaulting when :

- ✓ 1. Though the number of loans applied and defaulted are the highest in number for "debt_consolidation". The annual income of those who applied under the "debt_consolidation" purpose isn't the highest. Applicants with higher annual income applied loans for home_improvement, house, renewable_energy and small_businesses.
- ✓ 2. The highest number of loan applicants who have defaulted have a house ownership of 'MORTGAGE' and have the highest Annual Income.
- ✓ 3. Across all income groups, the number of defaulters and the loan amount is higher in higher income groups.
- ✓ 4. The interest rate for charged off loans at a rate of 21-24% and have an income of 60-70k.
- ✓ 5. Loan applicants with loan amount of 15K and above have an interest rate between 14 - 16%.
- ✓ 6. Loan applicants with purpose of 'small_business' have taken highest loan amount of 12k.
- ✓ 7. Loan applicants whose ownership is OTHER have taken the highest loan amount between 12-14k while the applicants whose ownership is 'MORTGAGE' have taken loan amount between 12K to 14K.
- ✓ 8. For grade of F the loan amount is between 15K to 17.5K.
- ✓ 9. For applicants with employment length above 10 years, the loan amount is greater than 12K.
- ✓ 10. Applicants with Grade G have interest rate more than 20%.
- ✓ 11. As the Loan Amount increases, the Interest Rate also increases. The interest rate for charged off loans is higher than that of fully paid loans in all the loan_amount groups.
- ✓ 12. Loan applicants who have applied for long term of 60 Months have applied for higher loan amounts.

The background features a large, central circle composed of several concentric rings in shades of light blue and green. Two curved, thin lines in the same color palette extend from the top left and bottom right towards the center of the circle.

Correlations

Correlation Analysis

Strong Correlation

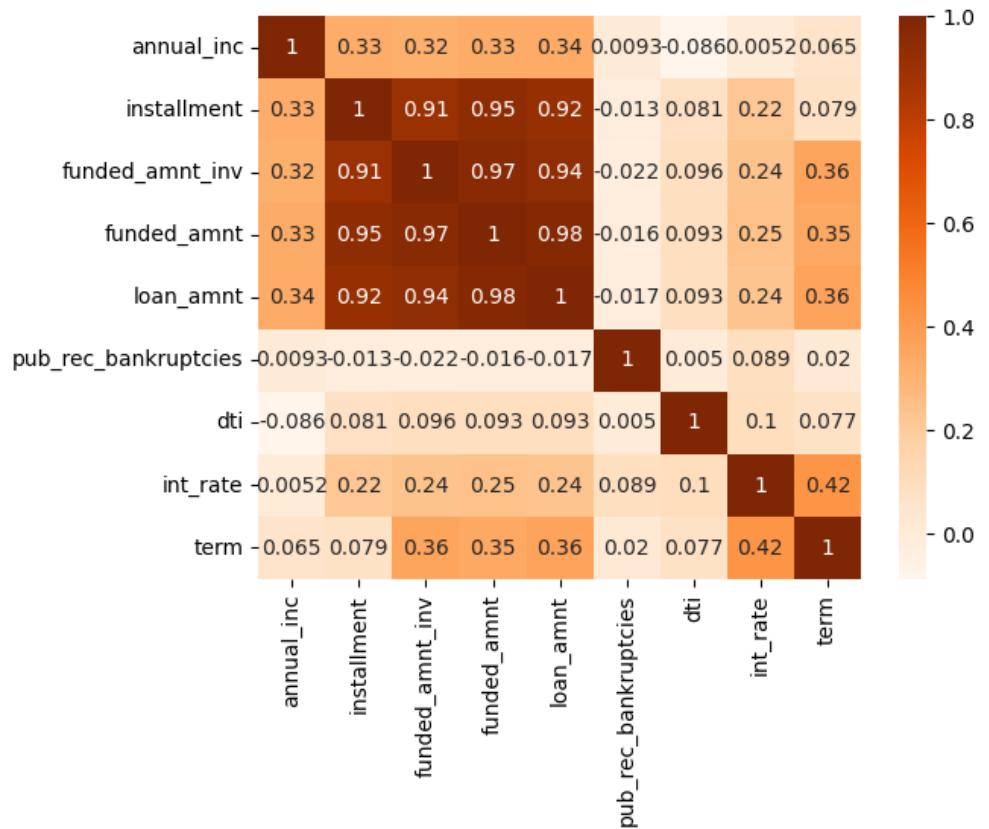
- ✓ **installment** has strong correlation with **funded_amnt**, **loan_amnt** and **funded_amnt_inv**
- ✓ **term** has a strong correlation with **int_rate**
- ✓ **annual_inc** has a strong correlation with **loan_amnt**

Weak Correlation

- ✓ **dti** has weak correlation with most of the fields.

Negative Correlation

- ✓ **annual_inc** has negative correlation with **dti**
- ✓ **pub_rec_bankruptcies** has a negative correlation with almost every field



Summary

- ✓ 1. There has been a positive trend in loan applications from 2007 to 2011. The company should capitalize on the market's growth trend by meeting the increased demand.
- ✓ 2. Highest number of loan applications are seen in the month of December. Possibly due to holiday season. Company can use this trend and create new marketing strategy.
- ✓ 3. Grades B, C and D have a high possibility of defaulting and hence needs stricter risk assessment.
- ✓ 4. Sub Grades B3, B4 and B5 have a high possibility of defaulting and hence needs stricter risk assessment.
- ✓ 5. Company should carefully evaluate loan applications applied for debt consolidation purpose.
- ✓ 6. Company should evaluate the home ownership status of the loan applicant to access housing stability, and its impact on the applicant's ability to repay the loan.
- ✓ 7. It has been seen that verified applicants were the highest in loan defaulters. Company needs to consider improving their verification process.
- ✓ 8. It has been seen that the highest number of defaulter belong to CA, NY and FL states. Company should consider implementing stricter risk assessment for these states.
- ✓ 9. Loan applicants with short term, 36 months term have been seen defaulting more compared to long term, 60 months tenure. Company can think of increasing the interest rates for short term loans.