# LENDING CLUB CASE STUDY
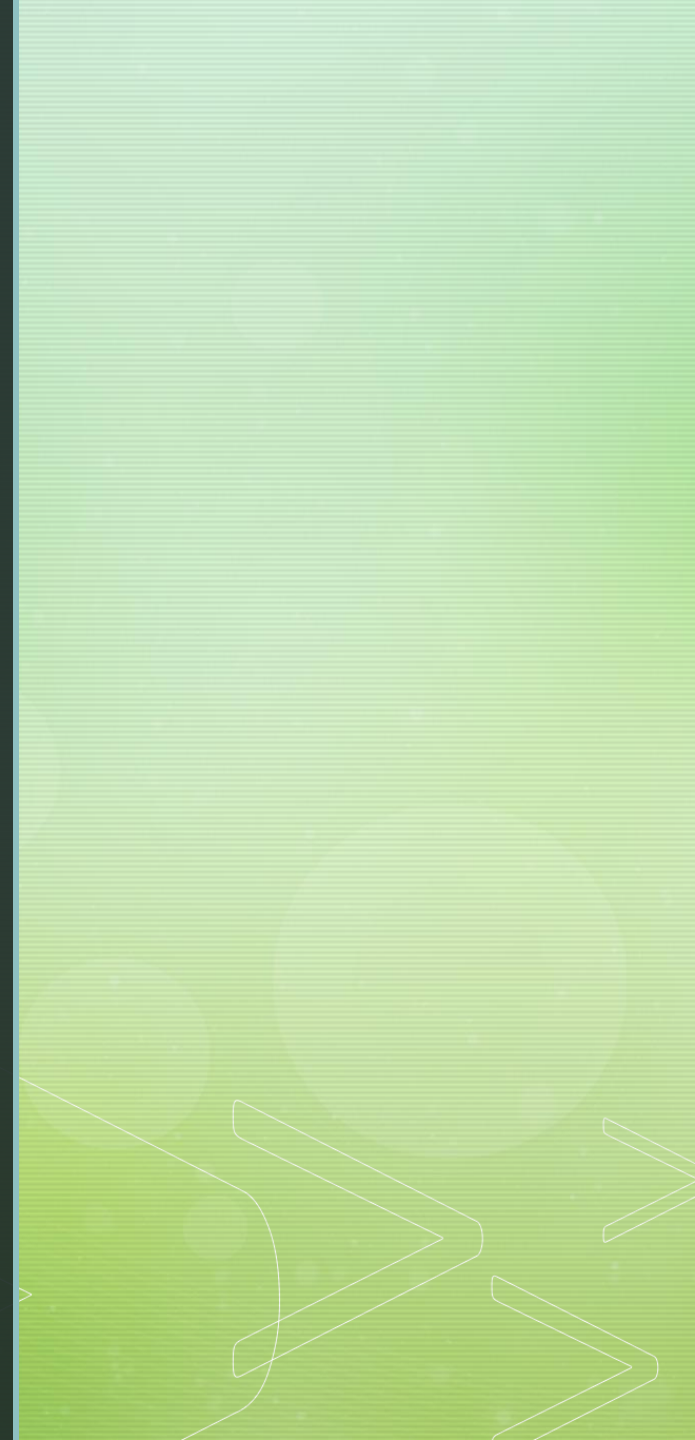
Submitted by:

HariKrishna Kotha

Hemant R

# Contents

- Problem Statement
- Data Description
- Data Understanding
- Data Cleaning & Pre-processing
- Univariate Analysis
- Bivariate Analysis
- Multivariate Analysis
- Correlation Analysis
- Suggestions

# Problem Statement

**Lending Club**, a Consumer Finance marketplace specializing in offering a variety of loans to urban customers, faces a critical challenge in managing its loan approval process. When evaluating loan applications, the company must make sound decisions to minimize financial losses, primarily stemming from loans extended to applicants who are considered **"Risky"**.

➢ These financial losses, referred to as **Credit Losses**, occur when borrowers fail to repay their loans or default. In simpler terms, borrowers labelled as **"Charged-Off"** are the ones responsible for the most significant losses to the company.

➢ The primary objective of this exercise is to assist Lending Club in mitigating credit losses. This challenge arises from two potential scenarios:

   1. **Identifying applicants likely to repay their loans is crucial, as they can generate profits for the company through interest payments. Rejecting such applicants would result in a loss of potential business.**

   2. **On the other hand, approving loans for applicants not likely to repay and at risk of default can lead to substantial financial losses for the company.**

➢ The objective is to pinpoint applicants at risk of defaulting on loans, enabling a reduction in credit losses. This case study aims to achieve this goal through Exploratory Data Analysis (EDA) using the provided dataset.

➢ In essence, the company wants to understand the driving factors (or driver variables) behind loan default, i.e. the variables which are strong indicators of default. The company can utilize this knowledge for its portfolio and risk assessment.

# DataDescription

- Lending Club provided us with customer's historical data. This dataset contained information pertaining to the borrower's past credit history and ending Club loan information. The total dataset consisted of over 39717 records and 111 columns, which was sufficient for our team to conduct analysis. Variables present within the dataset provided an ample amount of information which we could use to identify relationships and gauge their effect upon the success or failure of a borrower fulfilling the terms of their loan agreement.
- There are lot of columns with no data available, so only considered the columns which have the data for analysis

# Data Understanding

**Dataset Facts:**

**Primary Fact:**

**Loan Status: The Principal Attribute of Interest (loan_status). This column consists of three distinct values:**

- ✓ **Fully-Paid:** Signifies customers who have successfully repaid their loans.

- ✓ **Charged-Off:** Indicates customers who have been labeled as "Charged-Off" or have defaulted on their loans.

- ✓ **Current:** Represents customers whose loans are presently in progress and, thus, cannot provide conclusive evidence regarding future defaults.

For the purposes of this case study, rows with a "Current" status will be excluded from the analysis.

**Decision Matrix:**

**Loan Acceptance Outcome** - There are three potential scenarios:

**Fully Paid** - This category represents applicants who have successfully repaid both the principal and the interest rate of the loan.

**Current** - Applicants in this group are actively in the process of making loan installments; hence, the loan tenure has not yet concluded. These individuals are not categorized as 'defaulted.'

**Charged-off** - This classification pertains to applicants who have failed to make timely installments for an extended period, resulting in a 'default' on the loan.

**Loan Rejection** - In cases where the company has declined the loan application (usually due to the candidate not meeting their requirements), there is no transactional history available for these applicants. Consequently, this data is unavailable to the company and is not included in this dataset.

# Data Understanding

Key Dimension of table analyzed and segregated by demography

| Customer | Loan |
|---|---|
| Annual income | Loan Amount |
| Home ownership | Grade |
| Employment Length | Term |
| Debt To income | Loan Date |
| State | Purpose Of loan |
| | Verification Status |
| | Interest Rate |
| | Installment |

# Data Cleaning & Pre-processing

1. Loading data from loan CSV
2. Check for null values in the dataset
3. Check for unique values
4. Check for duplicated rows in data
5. Dropping Unnecessary Records & Columns
6. Common Functions
7. Data Conversion
8. Outlier Treatment
9. Imputing values in Columns

# Data Cleaning & Pre-processing

1. **Loading data from loan CSV: After loading data into notebook identified some data types are not what expected**. So changed the data types of the columns

2. **Checking for null values in the dataset:** There're many columns with null values. Dropped them as they can't contribute any kind of analysis. Most of the columns in given data has null values.

3. **Checking for unique values:** If the column has only a single unique value, it does not make any sense to include it as part of our data analysis. We need to find out those columns and drop them from the dataset. 9 columns had such unique values and they were removed.

4. **Checking for duplicated rows in data:** No duplicate rows were found.

5. **Dropping Records and Columns:**

    1. Dropped records where **loan_status="Current"**. Current loan status applications will not be of interest to us as they are still in repayment mode.

    2. Dropping columns where missing data is **>=65%** as these columns will skew our data analysis and they need to be removed.

    3. Dropping extra columns containing text like **collection_recovery_fee, delinq_2yrs, desc, earliest_cr_line, emp_title, id, inq_last_6mths, last_credit_pull_d, last_pymnt_amnt, last_pymnt_d, member_id, open_acc, out_prncp, out_prncp_inv, pub_rec, recoveries, revol_bal, revol_util, title, total_acc, total_pymnt, total_pymnt_inv, total_rec_int, total_rec_late_fee, total_rec_prncp, url, zip_code** as these will not contribute to loan pass or fail.

# Data Cleaning & Pre-processing

6. **Common Functions:** Common functions were created for repeating common operations like plotting bar graphs, box plots, histograms, countplots, binning etc.

7. **Data Conversion**: Converted columns like **debt to income (dti), funded amount (funded_amnt), funded amount investor (funded_amnt_inv) and loan amount (loan_amnt) to float** to match the data. Also converted **loan date (issue_d)** to **DateTime (format: yyyy-mm-dd).**

8. **Outlier Treatment:** Calculated the **Inter-Quartile Range (IQR)** and filtering out the outliers outside of lower and upper bound. During Outlier analysis the following observations were made

   ✓ The annual income of most of the loan applicants is between 40K - 75K USD

   ✓ The loan amount of most of the loan applicants is between 5K - 15K

   ✓ The funded amount of most of the loan applicants is between 5K - 14K USD

   ✓ The funded amount by investor for most of the loan applicants is between 5K - 14K USD

   ✓ The interest rate on the loan is between 9% - 14%

   ✓ The monthly installment amount on the loan is between 160 - 440

   ✓ The debt to income ration is between 8 - 18

# Data Cleaning & Pre-processing

9. **Imputing values in Columns:**

   ✓ **Replaced missing values of annual_inc with the corresponding mode value of annual_inc of the emp_length annual_inc field:** They Employment length has **1015** missing values, which means either they are **not employed or self-employed (business owners).** Considering they have a decent average annual income, we have assumed that these are business owners and we have added their employment duration with the mode value of **emp_length** which is **10+ years**.

   ✓ Mapped employment length with the respective number of years in int.

   ✓ Imputed **NONE** values as **OTHER** for **home_ownership.**

   ✓ Replaced the **'Source Verified'** values as **'Verified'** since both values mean the same thing i.e. the loan applicant has some source of income which is verified.

   ✓ There are **660 null values** for **pub_rec_bankruptcies**. Dropped those rows as they cannot be imputed.

Post Data cleaning and Pre-processing of dataset, we were left with **36094** rows × **18** columns.

# Univariate Analysis

✓ **Univariate analysis** is a statistical method used to analyze and summarize data sets consisting of **one variable**. It deals with the analysis of a single variable, rather than multiple variables, to understand its distribution, central tendency and dispersion.

✓ It was carried out for both **Categorical** and **Quantitative** Variables

A. **Categorical Variables:**

| Ordered | Unordered |
|---|---|
| ✓ Grade (grade)<br>✓ Sub grade (sub_grade)<br>✓ Term (36 / 60 months) (term)<br>✓ Employment length (emp_length)<br>✓ Issue year (issue_y)<br>✓ Issue month (issue_m)<br>✓ Issue quarter (issue_q) | ✓ Address State (addr_state)<br>✓ Loan purpose (purpose)<br>✓ Home Ownership (home_ownership)<br>✓ Loan status (loan_status)<br>✓ Loan paid (loan_paid) |

B. **Quantitative Variables:**

✓ **Interest rate bucket (int_rate_bucket)**

✓ **Annual income bucket (annual_inc_bucket)**

✓ **Loan amount bucket (loan_amnt_bucket)**

✓ **Funded amount bucket (funded_amnt_bucket)**

✓ **Debt to Income Ratio (DTI) bucket (dti_bucket)**

✓ **Monthly Installment (installment)**

# Univariate Analysis (Unordered Categorical)

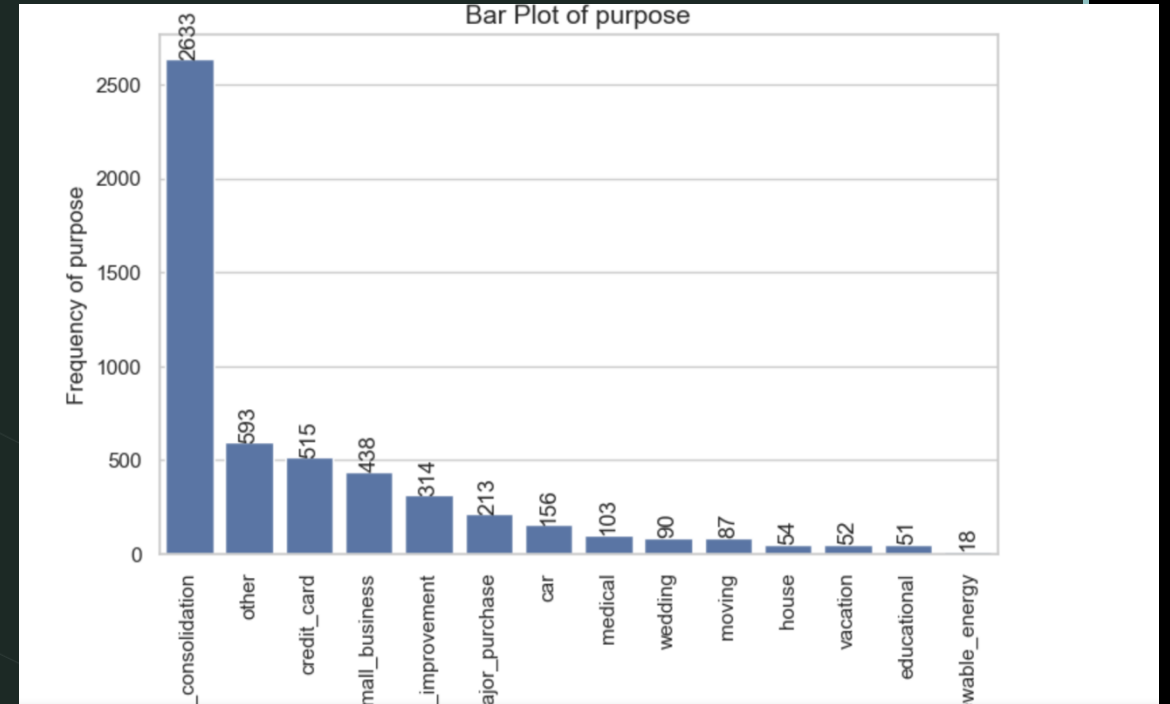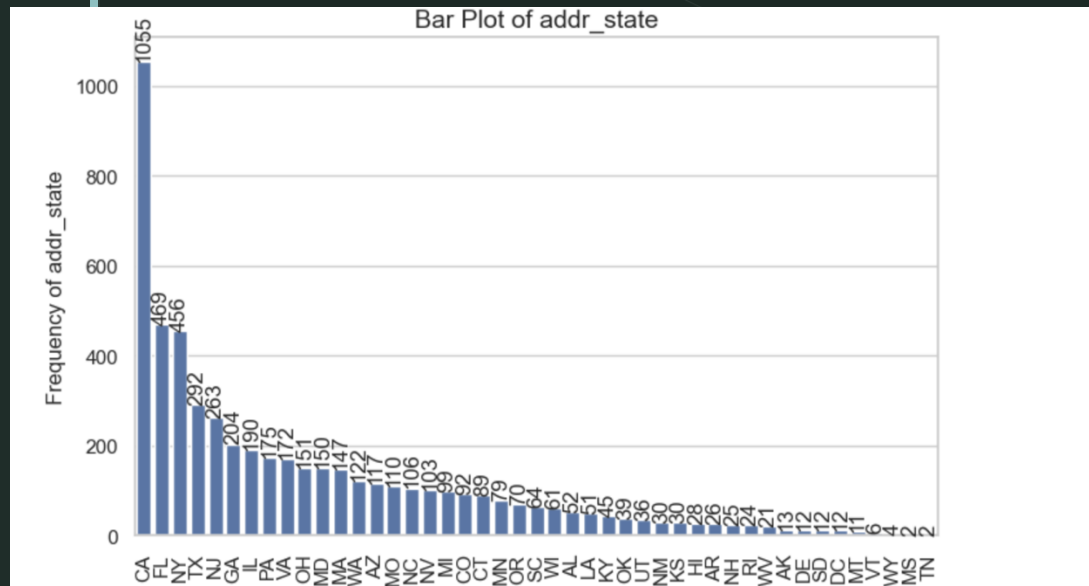Bar Plot of sub_grade

Bar Plot of home_ownership


Bar Plot of purpose


Bar Plot of addr_state

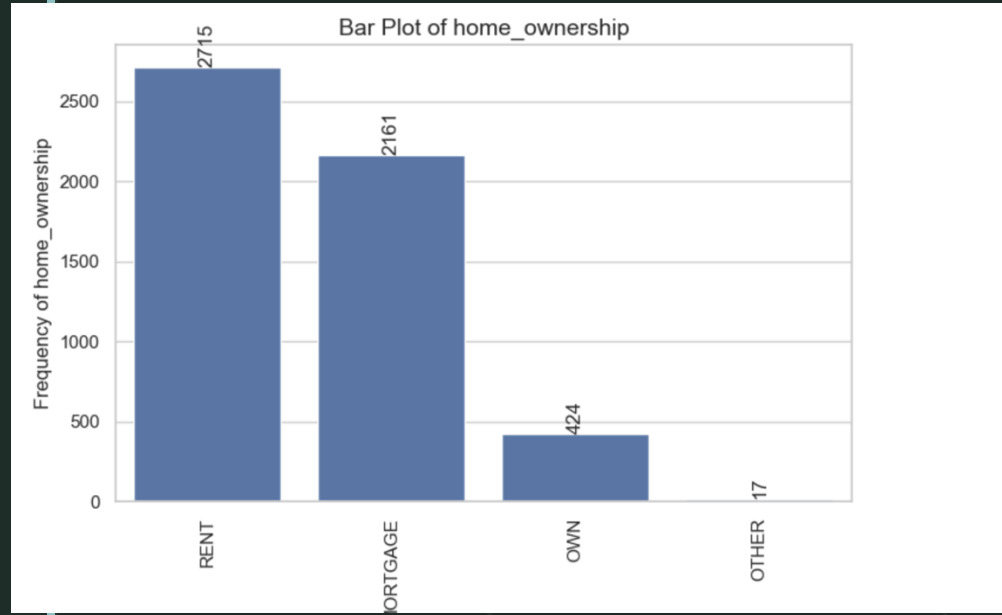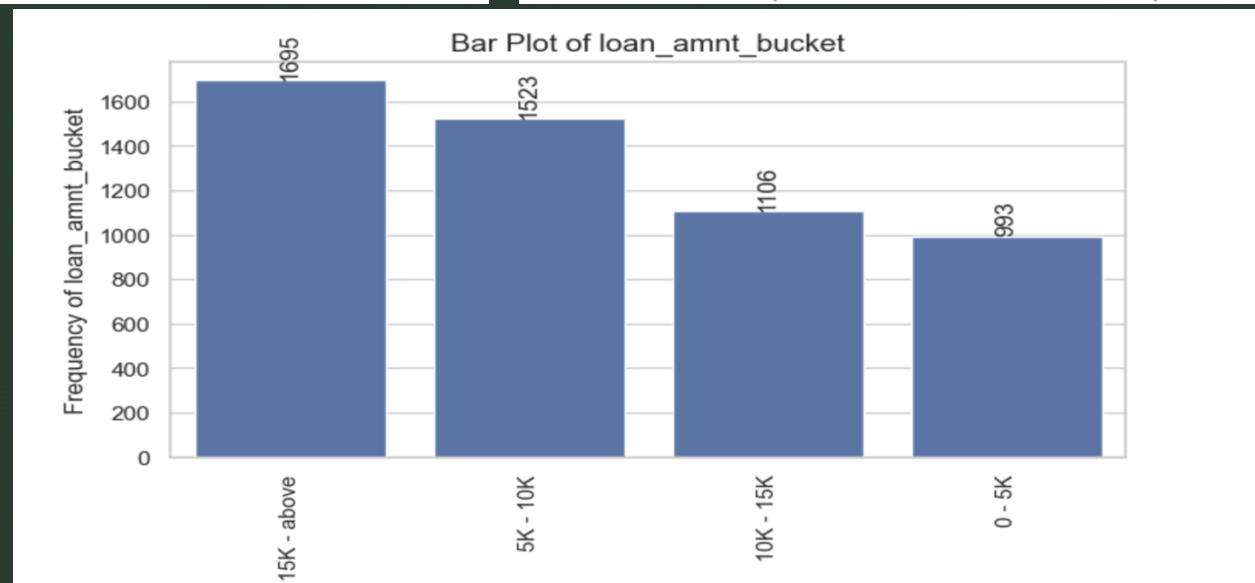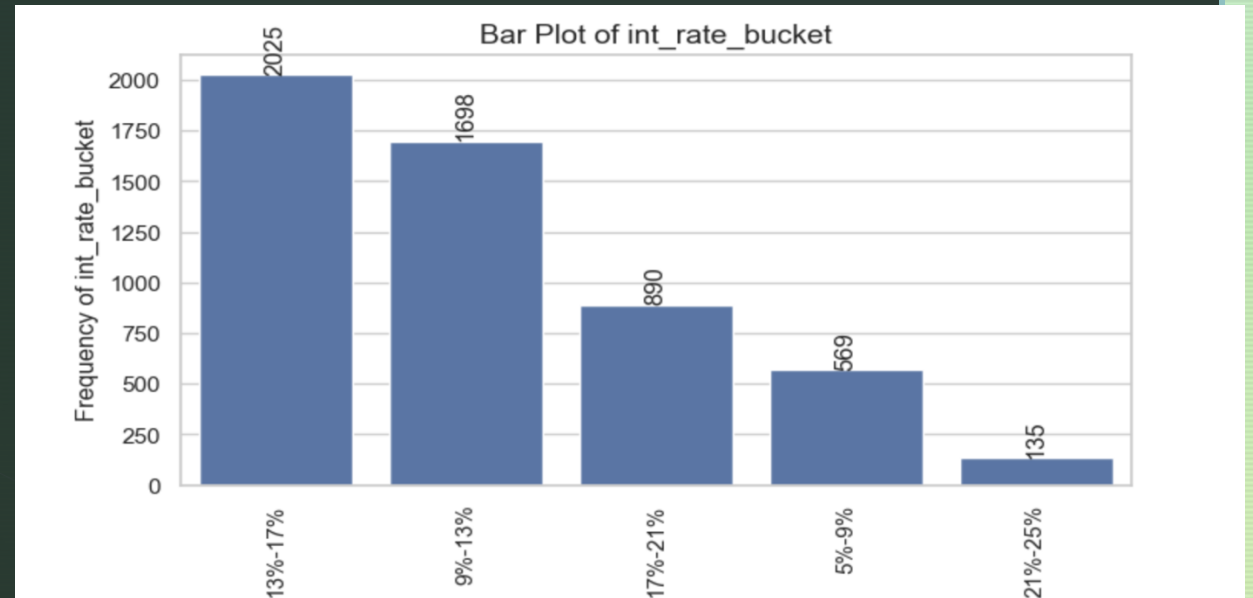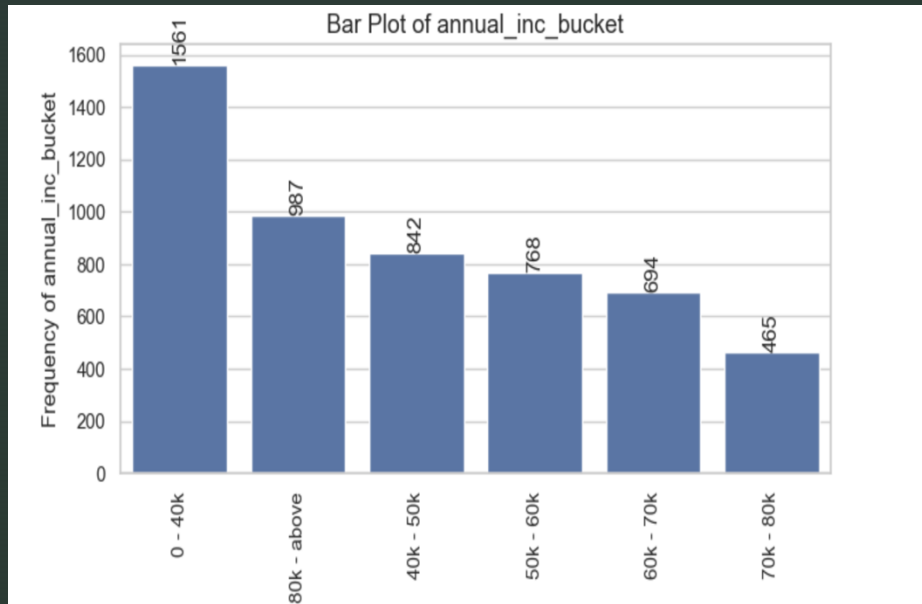# Observations & Inferences:

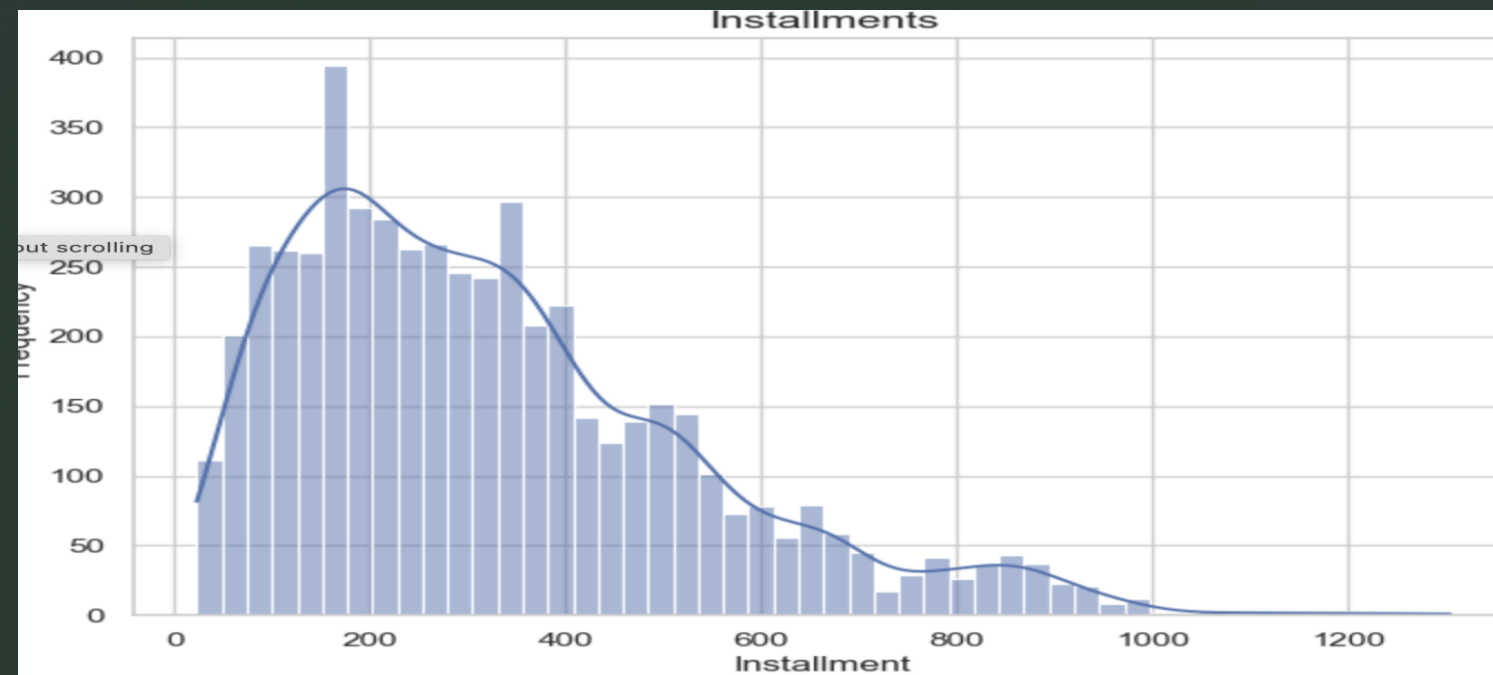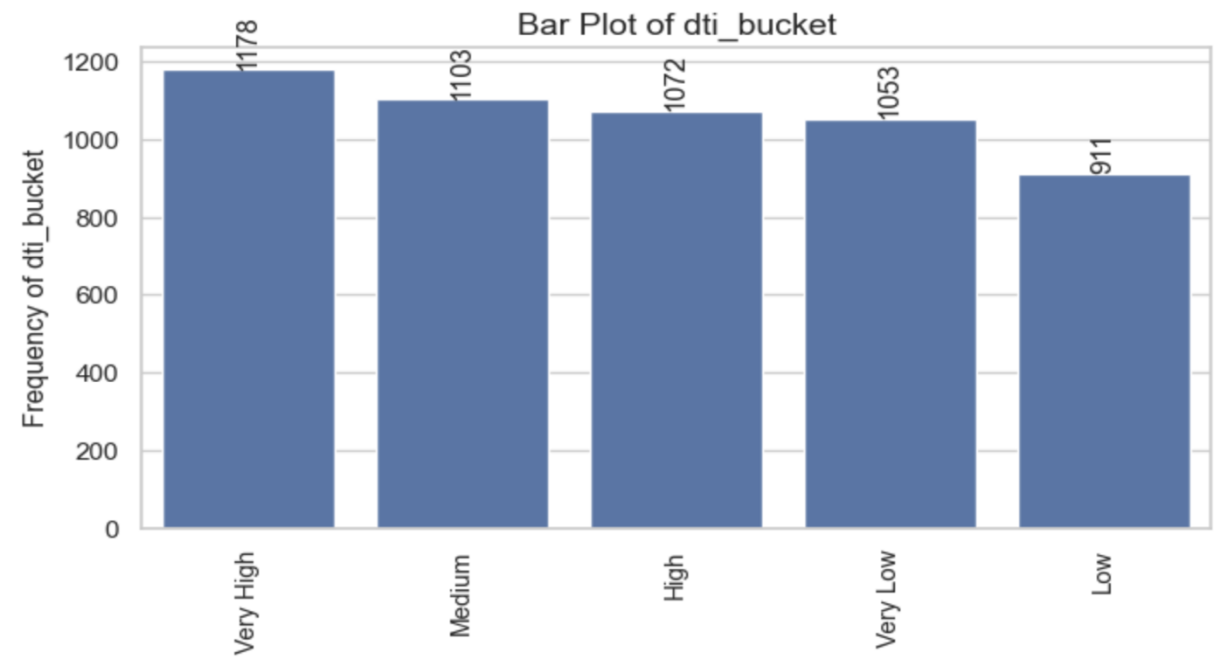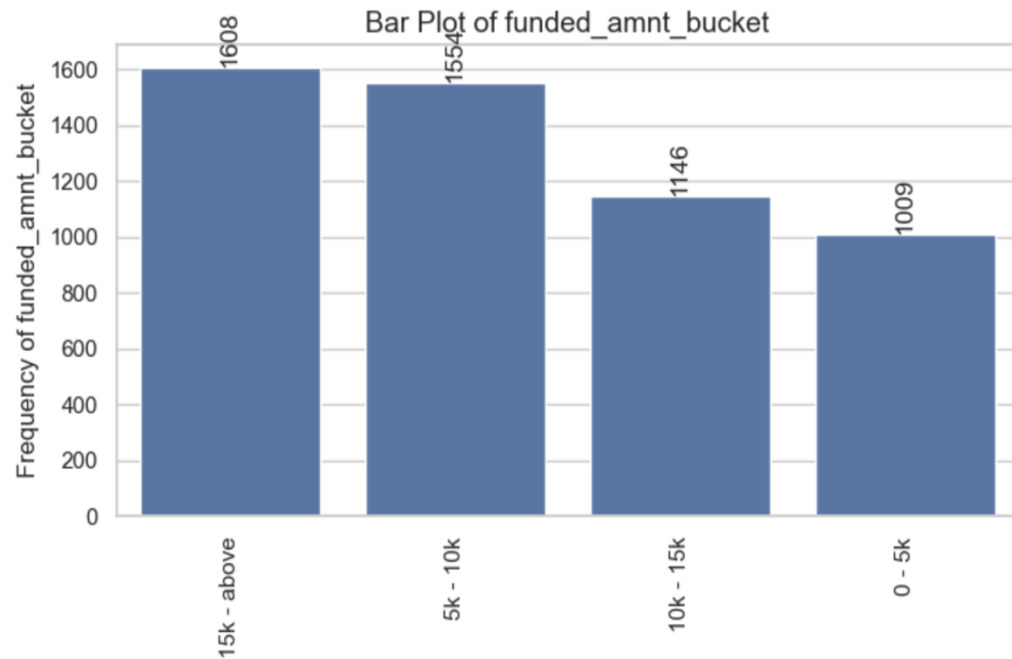### A. Ordered Categorical Variables:

- Grade B had the highest number of "Charged off" loan applicants, with a total of 1,352 applicants, indicating that applicants with this credit grade faced challenges in repaying their loans.
- Short-term loans with a duration of 36 months were the most "Charged off" applicants, with 3,006 applications. This suggests that a significant portion of applicants who experienced loan default chose shorter repayment terms.
- Applicants who had been employed for more than 10 years accounted for the highest number of "Charged off" loans, totaling 1,474. This indicates that long-term employment history did not necessarily guarantee successful loan repayment.
- The year 2011 recorded the highest number of "Charged off" loan applications, totaling 3,152, signaling a positive trend in the number of applicants facing loan defaults over the years. This could be indicative of economic or financial challenges during that year.
- "Charged off" loans were predominantly taken during the 4th quarter, with 2,284 applications, primarily in December. This peak in loan applications during the holiday season might suggest that financial pressures during the holidays contributed to loan defaults.

### B. Unordered Categorical Variables:

- California had the highest number of "Charged off" loan applicants, with 1,055 applicants. For such applicants, the lending company needs to implement stricter eligibility criteria or credit assessments due to a higher number of "Charged off" applicants from this state.
- Debt consolidation was the primary loan purpose for most "Charged off" loan applicants, with 2,633 applicants selecting this option. The lending company needs to exercise caution when approving loans for debt consolidation purposes, as it was the primary loan purpose for many "Charged off" applicants.
- The majority of "Charged off" loan participants, totaling 2,715 individuals, lived in rented houses. The lending company must assess the financial stability of applicants living in rented houses, as they may be more susceptible to economic fluctuations.
- A significant number of loan participants, specifically 5,317 individuals, were loan defaulters, unable to clear their loans. The lending company should enhance risk assessment practices, including stricter credit checks and lower loan-to-value ratios, for applicants with a history of loan defaults. They should offer financial education and support services to help borrowers manage their finances and improve loan repayment outcomes.
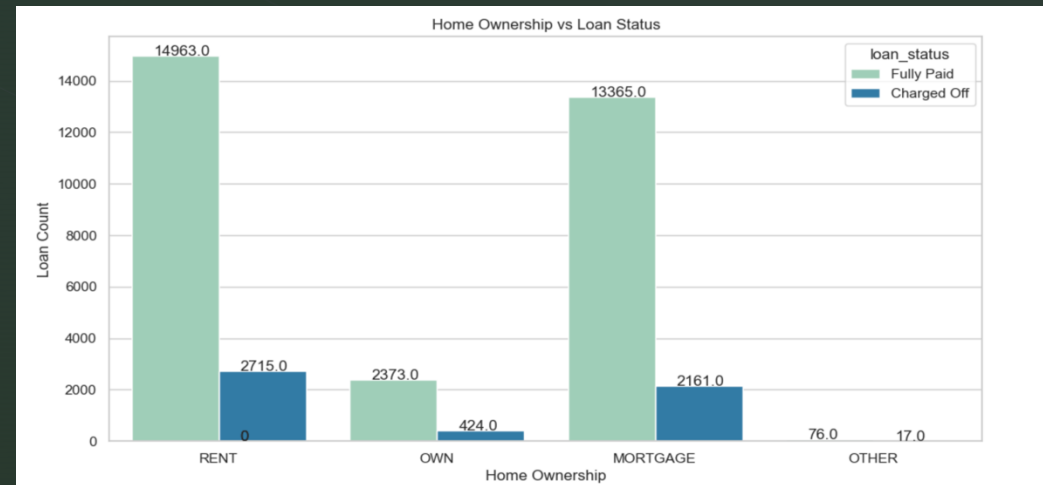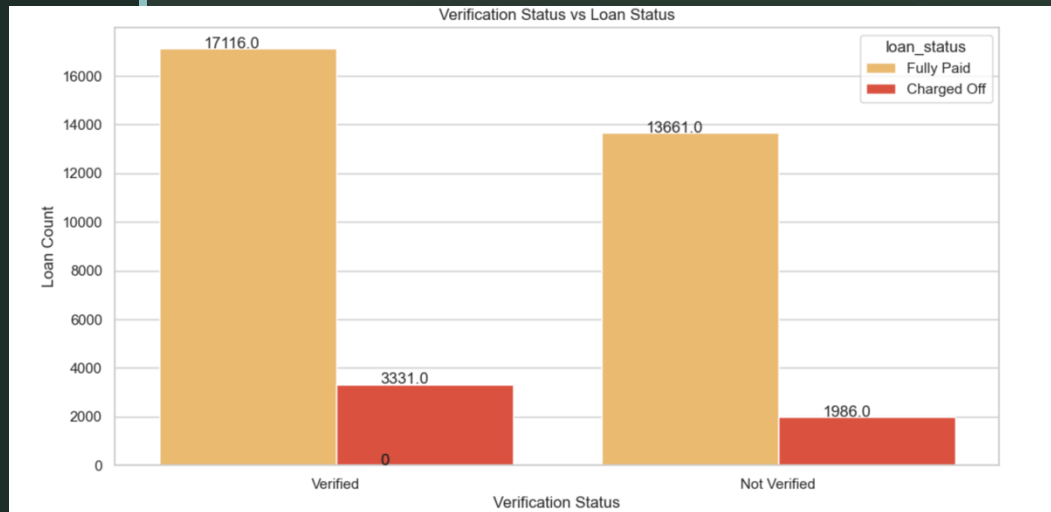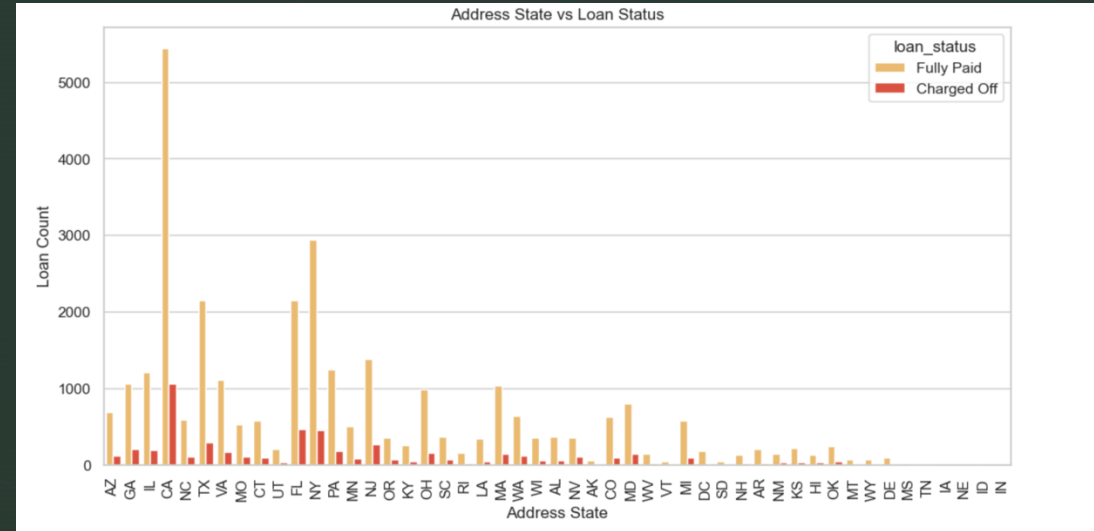
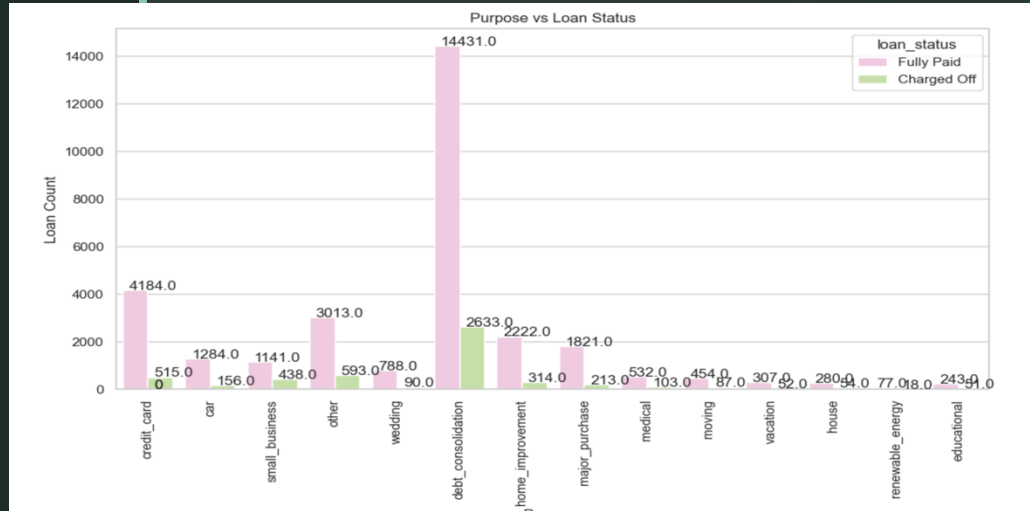# Univariate Analysis (Quantitative Variables)

# Observations

1. 1,561 loan applicants who charged off had annual salaries less than 40,000 USD. The lending company should exercise caution when lending to individuals with low annual salaries. They should implement rigorous income verification and assess repayment capacity more thoroughly for applicants in this income bracket.

2. Among loan participants who charged off (2,025), a considerable portion belonged to the interest rate bucket of 13%-17%.To reduce the risk of default, the lending company should consider offering loans at lower interest rates when possible.

3. 1,695 loan participants who charged off received loan amounts of 15,000 USD and above. The lending company should evaluate applicants seeking higher loan amounts carefully. They should ensure the applicants must have a strong credit history and repayment capability to handle larger loans.

4. 1,608 loan participants who charged off received funded amounts of 15,000 USD and above. The lending company should ensure that the funded amounts align with the borrower's financial capacity. They should conduct thorough credit assessments for larger loan requests.

5. Among loan participants who charged off, 1,178 loan applicants had very high debt-to-income ratios. The lending company

6. should implement strict debt-to-income ratio requirements to prevent lending to individuals with unsustainable levels of debt relative to their income.

7. Among loan participants who charged off, it's observed that the majority of them had monthly installment amounts falling within the range of 160-440 USD. The lending company should closely monitor and assess applicants with similar installment amounts to mitigate the risk of loan defaults.
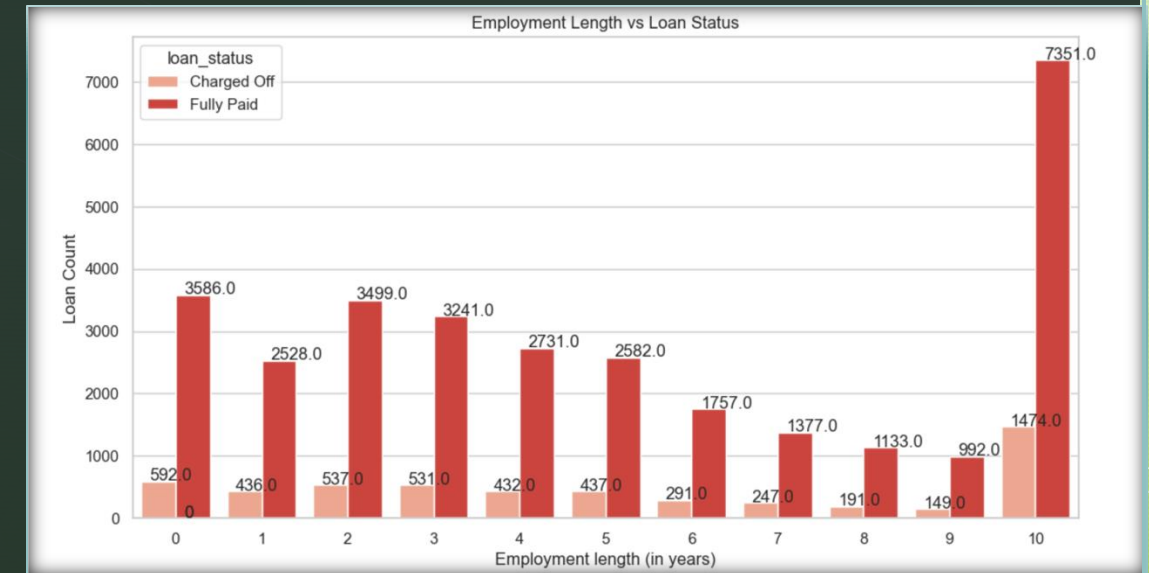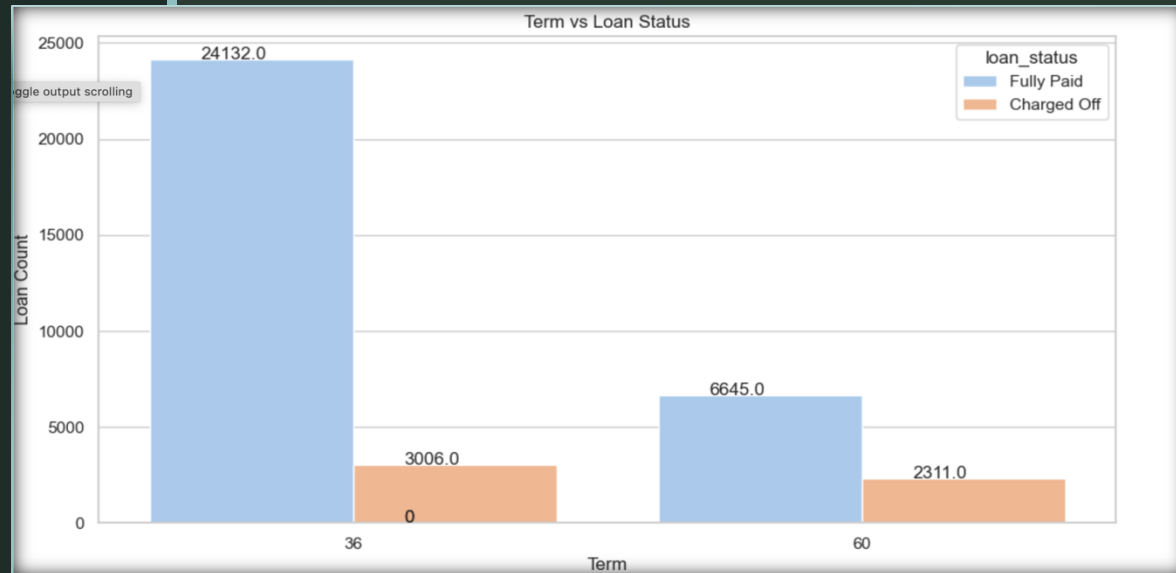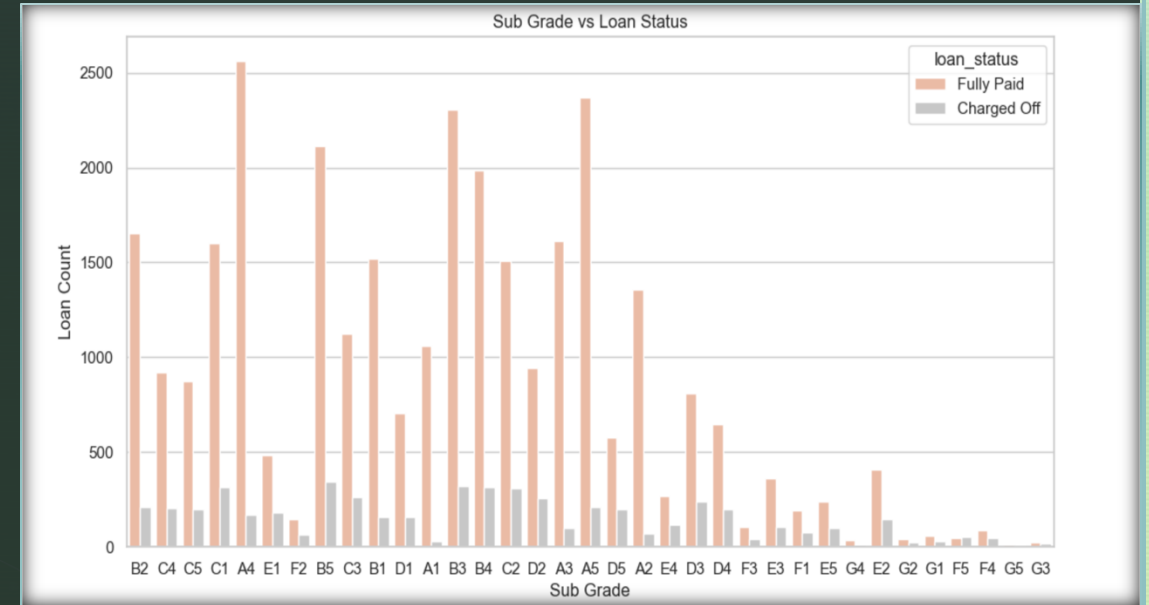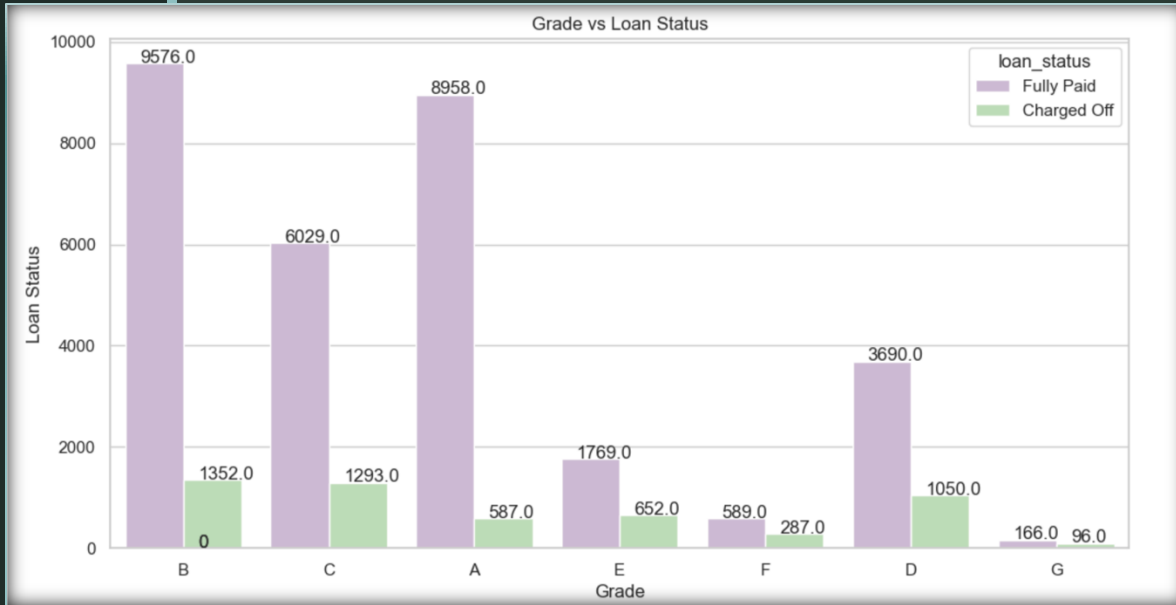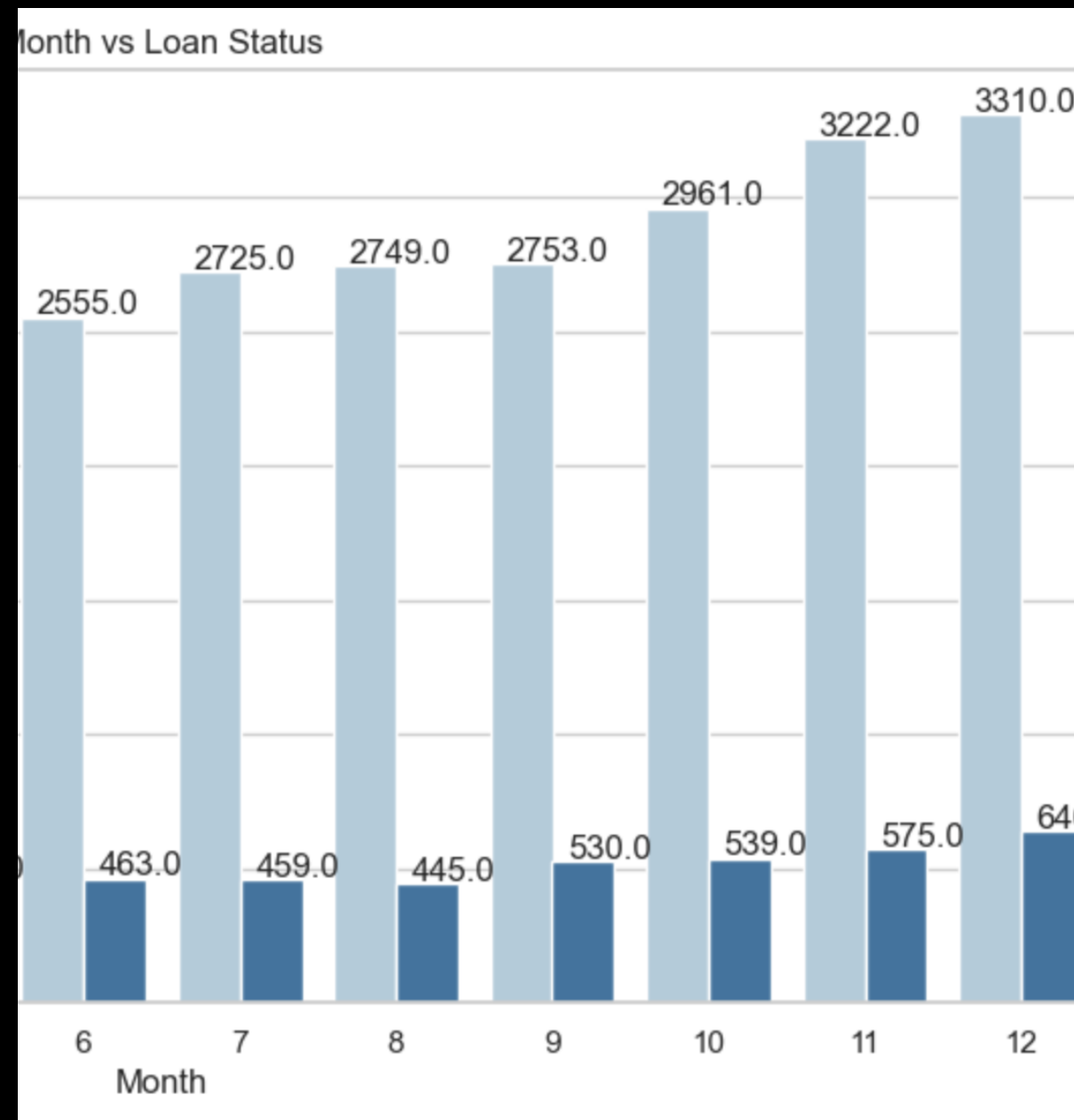
# Bivariate Analysis

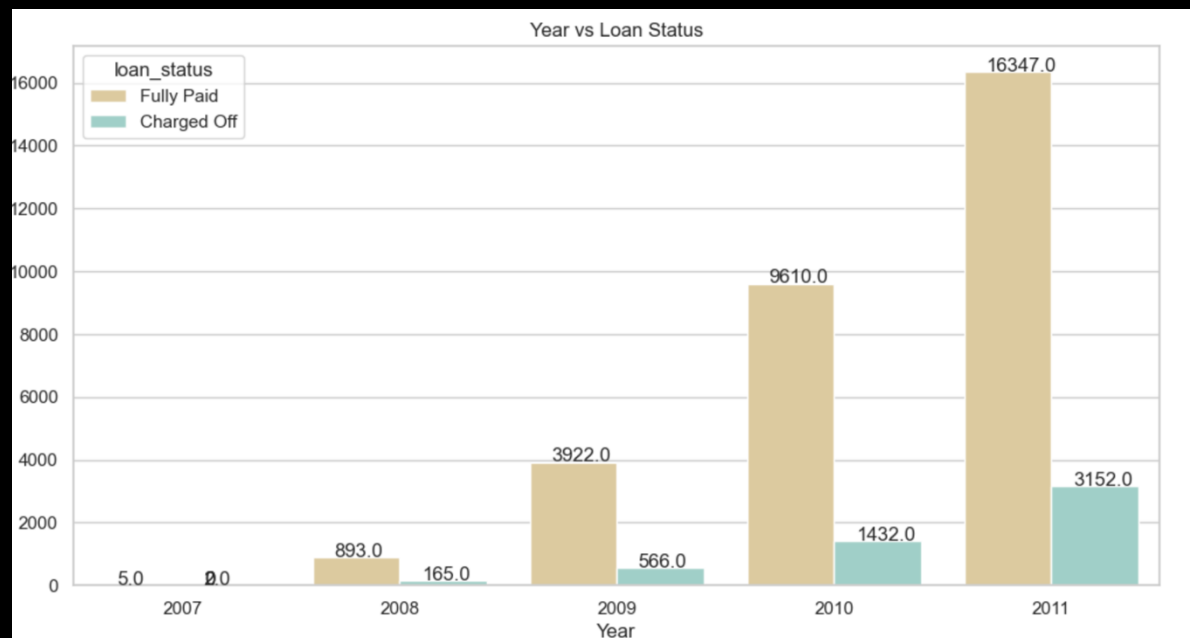- Bivariate analysis is a statistical method that involves the simultaneous analysis of two variables (factors). It aims to determine the empirical relationship between them. The analysis can be used to test hypotheses, identify patterns, or explore relationships between the variables.
- It was carried out for both Categorical and Quantitative Variables

# Bivariate Analysis (Unordered Categorical)

# Bivariate Analysis (Ordered Categorical)

**Year vs Loan Status**

loan_status
- Fully Paid
- Charged Off

| Year | Fully Paid | Charged Off |
|------|-----------|-------------|
| 2007 | 5.0 | 0 |
| 2008 | 893.0 | 165.0 |
| 2009 | 3922.0 | 566.0 |
| 2010 | 9610.0 | 1432.0 |
| 2011 | 16347.0 | 3152.0 |

**Quarter vs Loan Status**

loan_status
- Fully Paid
- Charged Off

| Quarter | Fully Paid | Charged Off |
|---------|-----------|-------------|
| Q4 | 12246.0 | 2284.0 / 0 |
| Q3 | 5474.0 | 904.0 |
| Q2 | 7251.0 | 1278.0 |
| Q1 | 5806.0 | 851.0 |

**Month vs Loan Status**

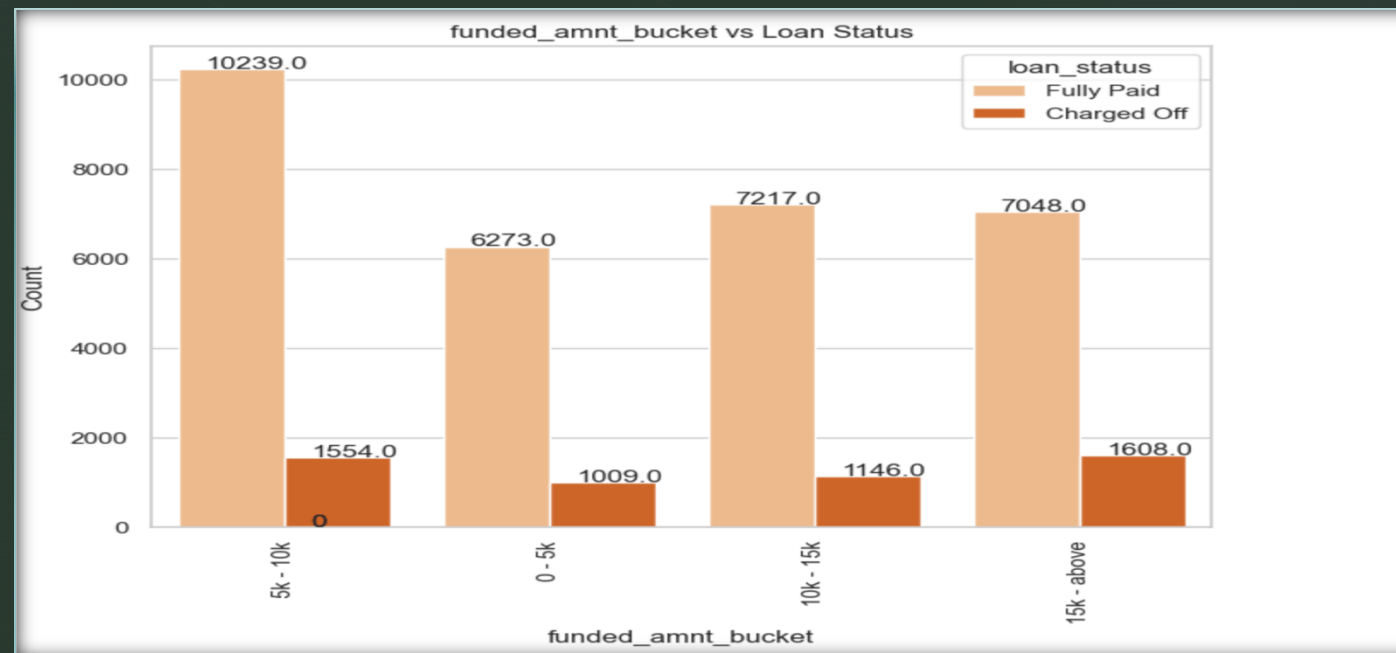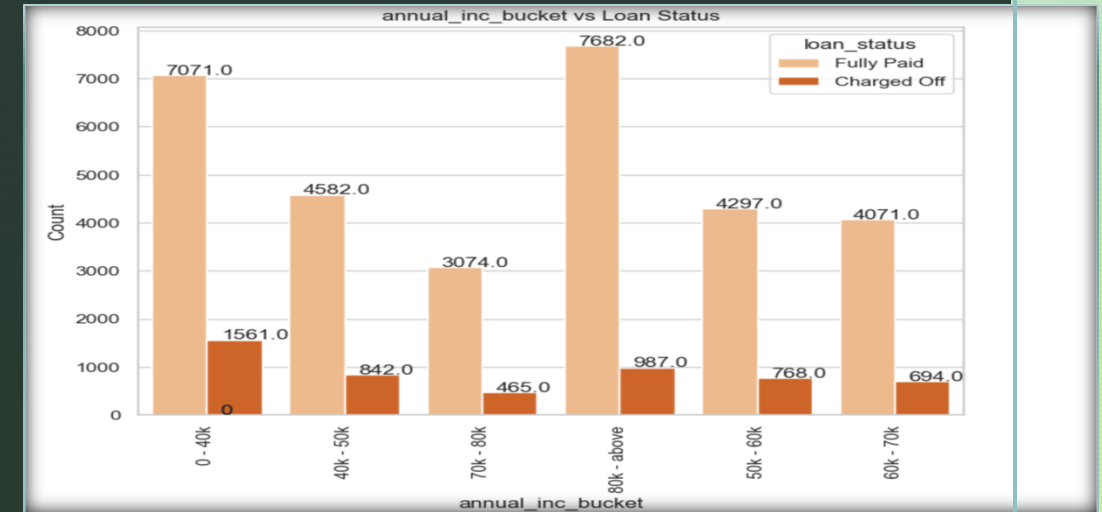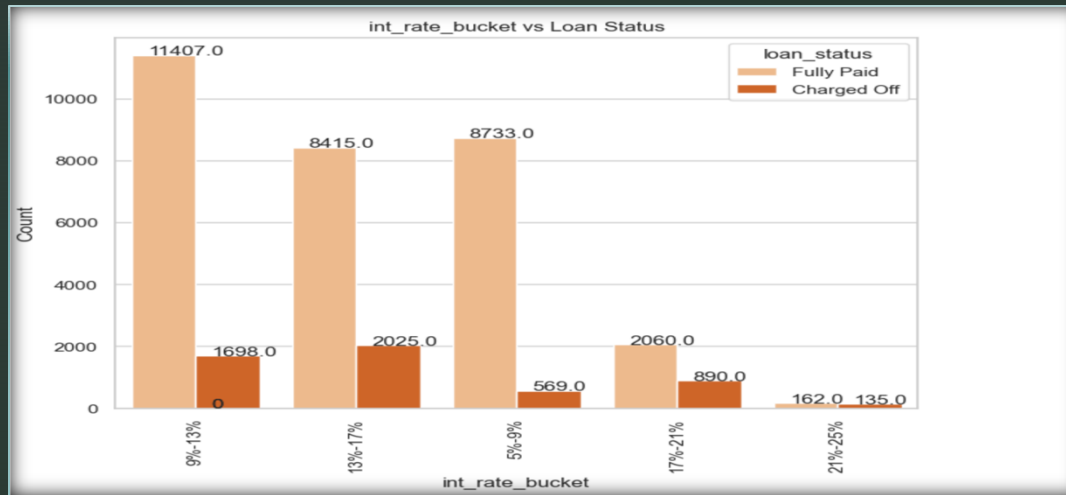| Month | | |
|-------|------|-------|
| 6 | 2555.0 | |
| 7 | 2725.0 | 463.0 |
| 8 | 2749.0 | 459.0 |
| 9 | 2753.0 | 445.0 |
| 10 | 2961.0 | 530.0 |
| 11 | 3222.0 | 539.0 |
| 12 | 3310.0 | 575.0 / 64... |

# Inferences

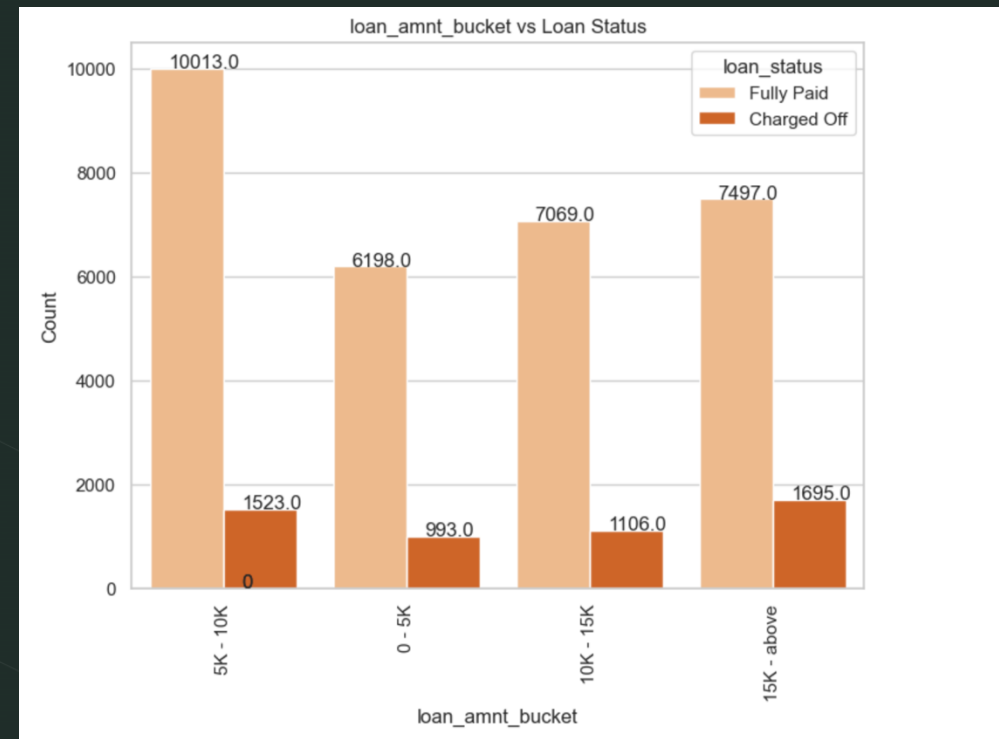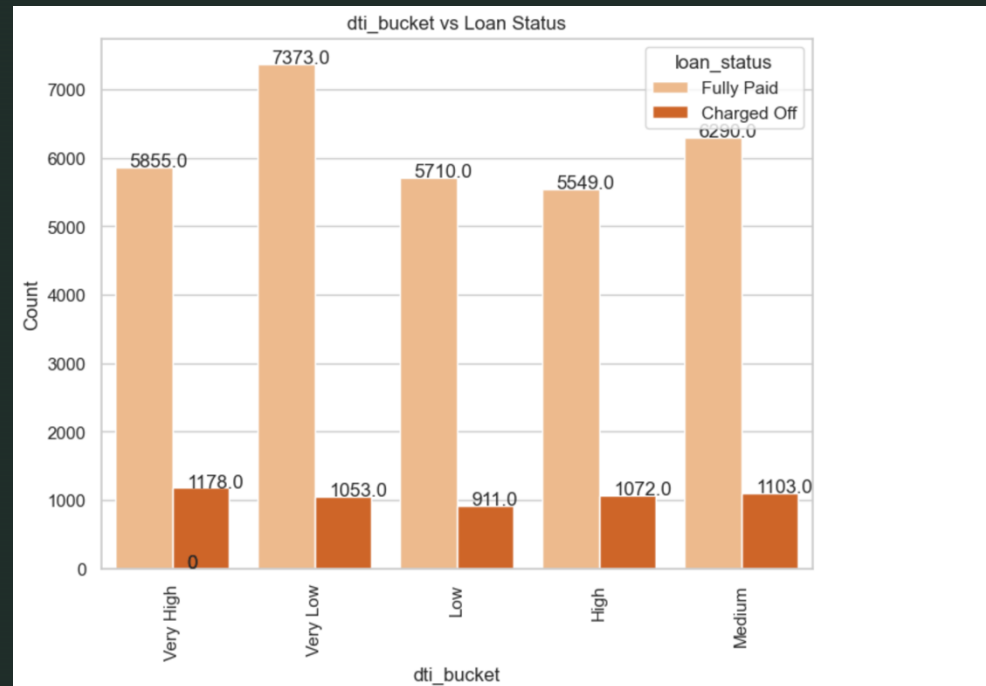A. **Ordered Categorical Variables:**

- Since loan applicants from Grades B, C, and D contribute to most of the "Charged Off" loans, the company should consider implementing stricter risk assessment and underwriting criteria for applicants falling into these grades.
- Pay special attention to applicants with Subgrades B3, B4, and B5, as they are more likely to charge off. Implementing additional risk mitigation measures or offering them lower loan amounts could be considered.
- Given that applicants opting for 60-month loans are more likely to default, the company should consider evaluating the risk associated with longer-term loans and potentially limiting the maximum term or adjusting interest rates accordingly.
- Loan applicants with ten or more years of experience are more likely to default. This suggests that experience alone may not be a reliable indicator of creditworthiness. The company should use a more comprehensive credit scoring system that factors in other risk-related attributes.
- The steady increase in the number of loan applicants from 2007 to 2011 indicates growth in the market. The company can capitalize on this trend by maintaining a competitive edge in the industry while keeping risk management practices robust.
- December and Q4 are peak periods for loan applications, likely due to the holiday season. The company should anticipate increased demand during these periods and ensure efficient processing to meet customer needs

B. **Unordered Categorical Variables:**

- Since debt consolidation is the category with the maximum number of loans and high default rates, the company should carefully evaluate applicants seeking debt consolidation loans and potentially adjust interest rates or offer financial counseling services.
- Applicants living in rented or mortgaged houses are more likely to default. This information can be considered in the underwriting process to assess housing stability and its impact on repayment ability.
- Verified loan applicants are defaulting more than those who are not verified. The company should review its verification process to ensure it effectively assesses applicant creditworthiness and consider improvements or adjustments.
- Loan applicants from states like California (CA), Florida (FL), and New York (NY) are more likely to default. The company should monitor regional risk trends and adjust lending strategies or rates accordingly in these areas.

# Bivariate Analysis (Quantitative Variables)

dti_bucket vs Loan Status
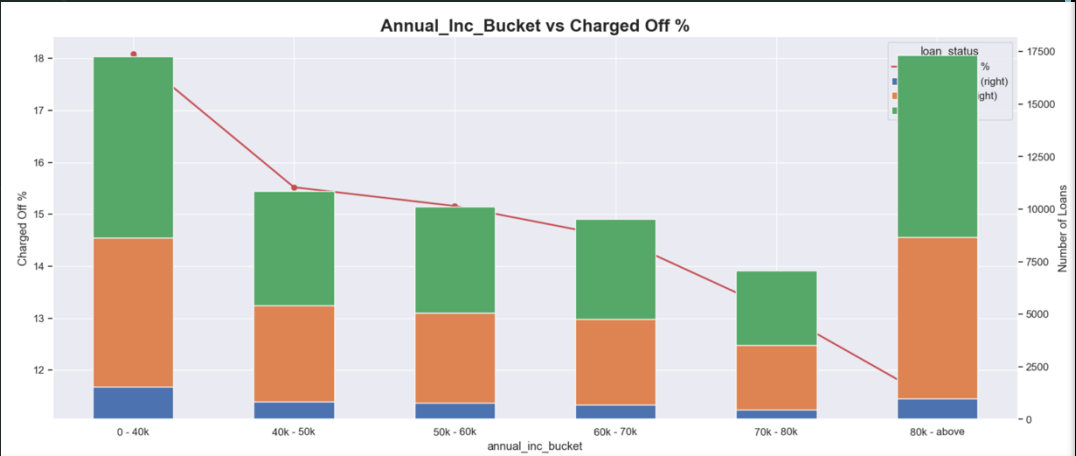
loan_amnt_bucket vs Loan Status

## Observations:
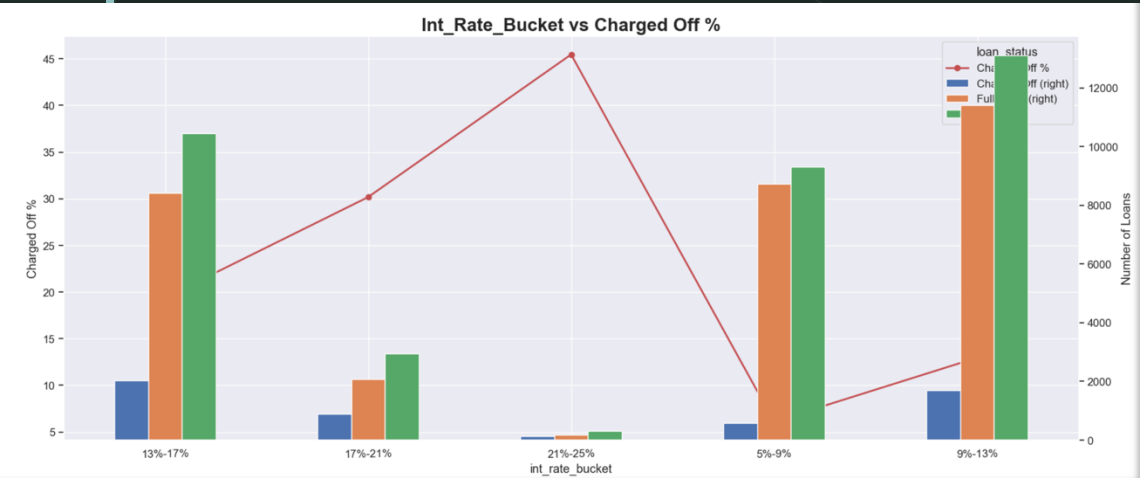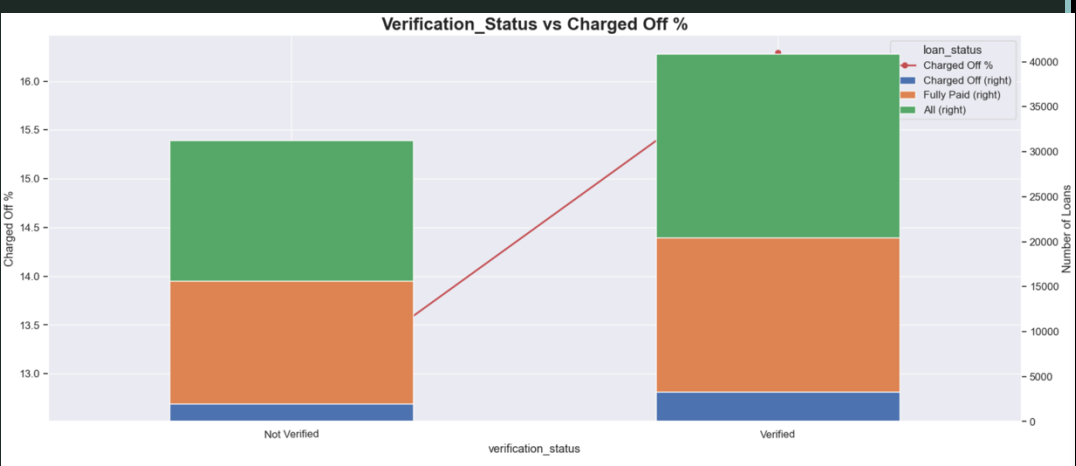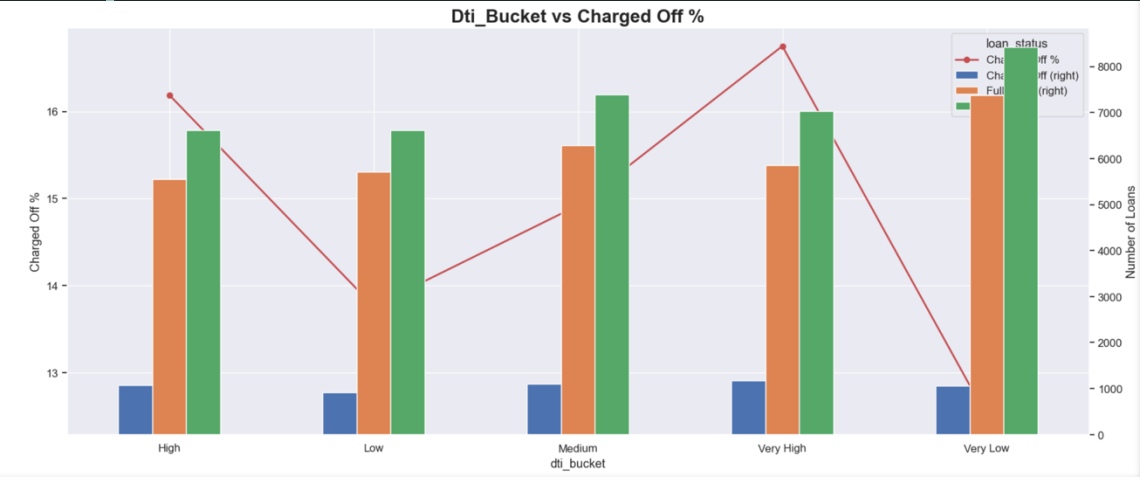
✓ A majority of the loan applicants who defaulted received loan amounts of $15,000 or higher.

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

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

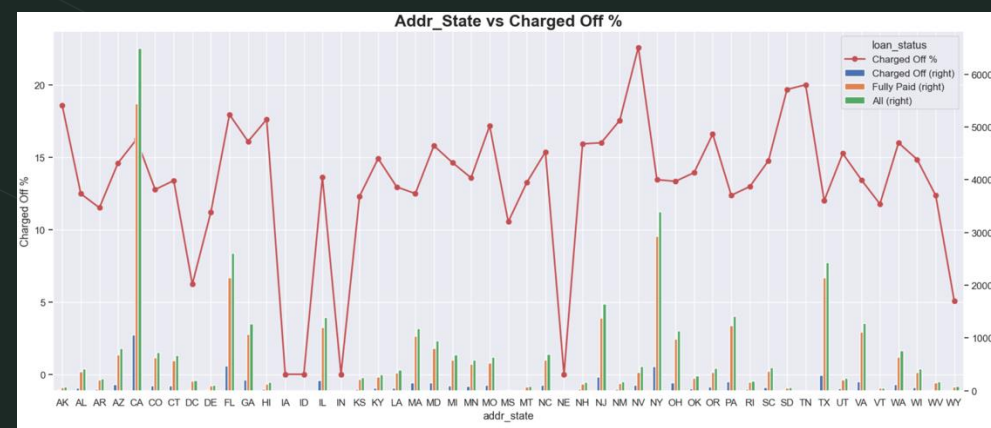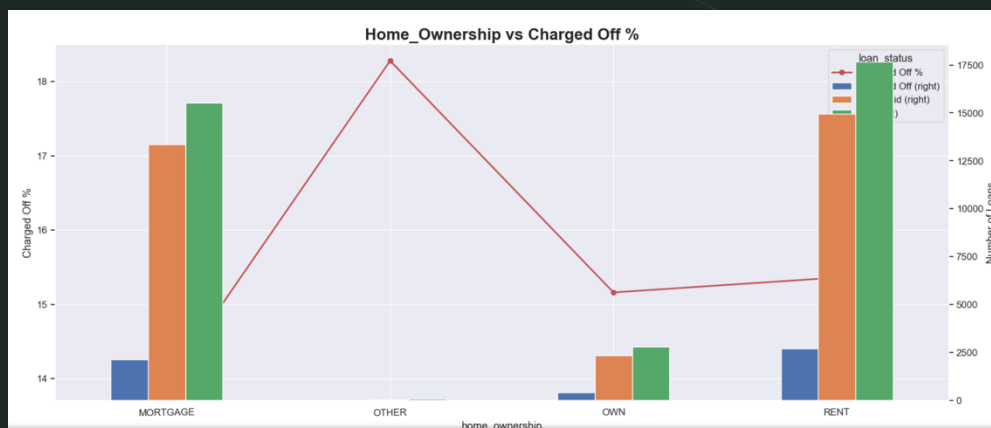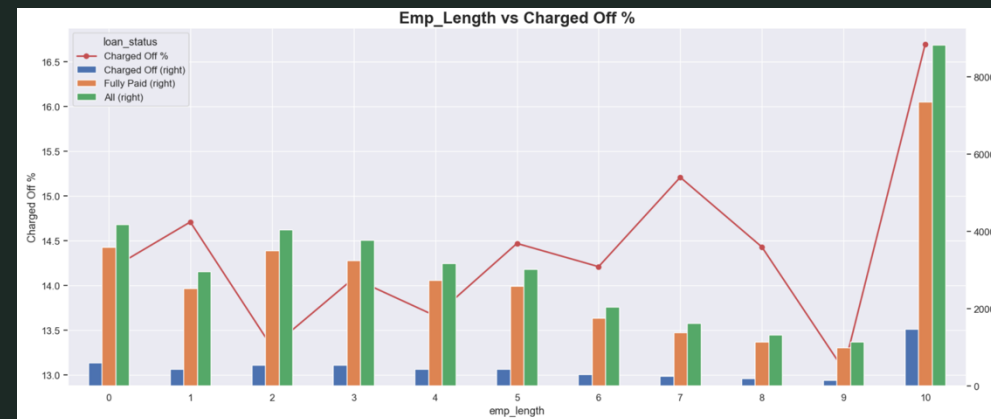✓ A majority of the loan applicants who charged off reported an annual income of less than $40,000.

## Inferences:

✓ Applicants receiving loan amounts of $15,000 or higher are more likely to default. The company can mitigate this risk by conducting more thorough assessments for larger loan requests and potentially capping loan amounts for higher-risk applicants.

✓ High Debt-to-Income (DTI) ratios and interest rates in the 13%-17% range are associated with defaults. The company should review its interest rate determination process and consider adjusting rates based on DTI ratios to better align with the borrower's ability to repay.

✓ Applicants with annual incomes less than $40,000 have a higher likelihood of defaulting. The company should consider offering financial education resources or setting maximum loan amounts based on income levels to ensure affordability for borrowers.

# Multivariate Analysis

Multivariate analysis is a statistical technique used to analyze data that involves more than two variables.

**Dti_Bucket vs Charged Off %**

**Verification_Status vs Charged Off %**

**Int_Rate_Bucket vs Charged Off %**

**Annual_Inc_Bucket vs Charged Off %**

# Correlation Analysis

# SUGGESTIONS

1. Grade B,C and D customers are more prone to loan default. Implement stricter guidelines when disbursing loan to mentioned grades
2. Sub grades B3, B4, and B5 are also more prone to loan default. Implement stricter guidelines when disbursing loan to mentioned grades
3. Longer term loans are more prone to loan default. Move with caution when dealing with them or charge high interest so that loss can be minimum.
4. Customers with more experience are prone to loan default. Which suggest that experience alone cannot trusted when disbursing loan consider other indicators also.
5. Q4 are peak periods for loan applications, likely due to the holiday season. The company should anticipate increased demand during these periods and ensure efficient processing to meet customer needs.
6. debt consolidation loan has more prone to loan default. Maintain caution when dealing with such type of loans
7. Customers living in rented or mortgaged are prone to loan default.
8. As expected Verified customers are less likely to default. So bring in more robust verification process and maximum try to verify all customers
9. Applicants receiving loan amounts of $15,000 or higher are more likely to default.
10. High Debt-to-Income (DTI) ratios and interest rates in the 13%-17% range are associated with defaults.
11. Applicants with annual incomes less than $40,000 have a higher likelihood of defaulting. The company should consider offering financial education resources or setting maximum loan amounts based on income levels to ensure affordability for borrowers.
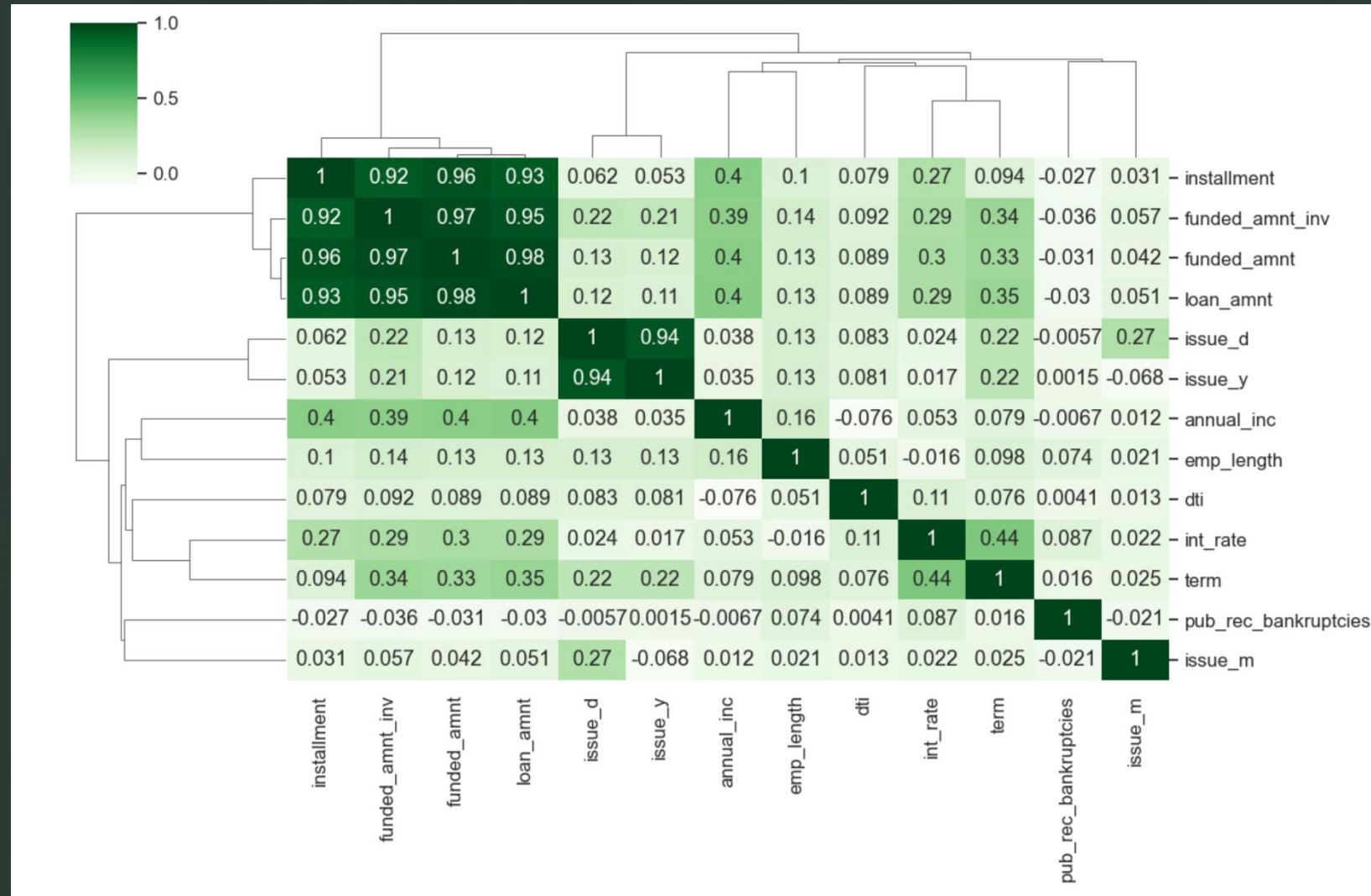12. The steady increase in the number of loan applicants from 2007 to 2011 indicates growth in the market. The company can capitalize on this trend by maintaining a competitive edge in the industry while keeping risk                                   management                                   practices                                   robust.

# Thank You