

Lending Club Case Study: Exploratory Data Analysis

Submitted by

Himanshu Agrawal

Introduction

- The dataset contains 39,717 rows and 111 columns.
- Goal: Identify factors leading to customers being charged off and gain insights into loan defaults.
- Target Variable: 'loan_status' (Fully Paid, Charged Off, Current).

Loan Status	Count	Percentage
Fully Paid	32950	82.96
Charged Off	5627	14.17
Current	1140	2.87

Steps Followed

1. Null Values Treatment
2. Data Type Treatment
3. Univariate, Bivariate, and Multivariate Analysis
4. Insights

Null Values Treatment

- **Dropped Columns with High Null Values:**
Removed columns with more than 30% missing values.
- **Imputed Remaining Null Values:**
Used median for numerical columns and mode for categorical columns.
- **Removed Low Variance Columns:**
Dropped columns containing only a single unique value.

Data Type treatment

- **Converted Date Columns:**
Transformed object columns representing dates into datetime objects.
- **Separated Numerical and Categorical Columns:**
Categorized columns as numerical or categorical for further analysis.
- **Validated Data Types:**
Ensured numerical columns are correctly typed.
Rechecked if categorical columns are appropriately labeled.
- **Removed High-Cardinality Columns:**
Dropped categorical columns with more than 50% unique values.
- **Created a Clean Dataset:**
Combined refined numerical and categorical columns to form the final clean dataset.

- **Univariate Analysis:**

Explored individual columns for distribution and trends.

Plotted histograms, boxplots, and pie charts to identify patterns and outliers.

- **Bivariate Analysis:**

Analyzed relationships between two variables.

Used scatter plots, bar charts, and pivot tables to identify correlations.

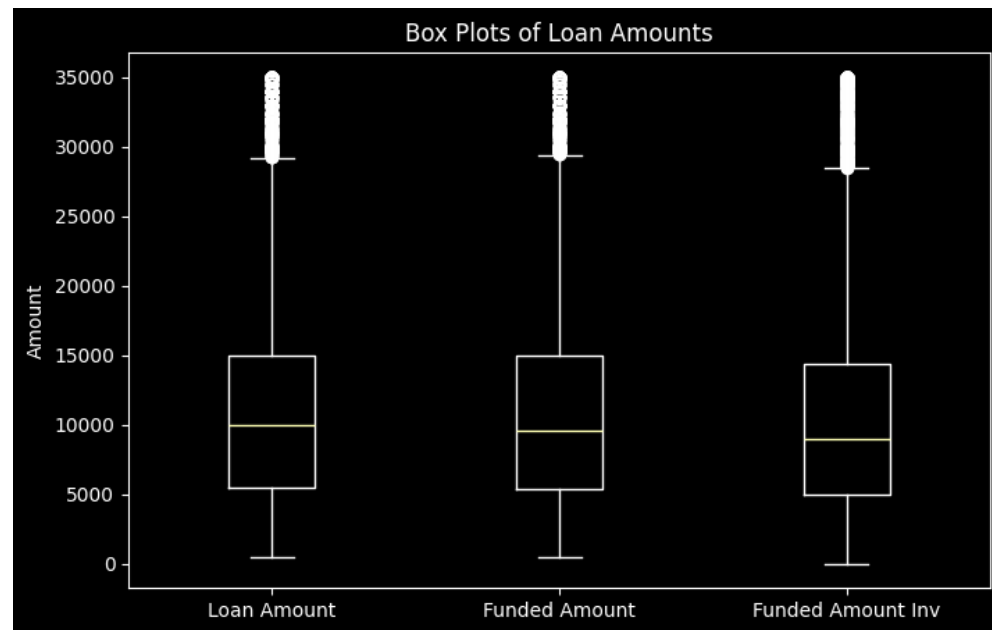
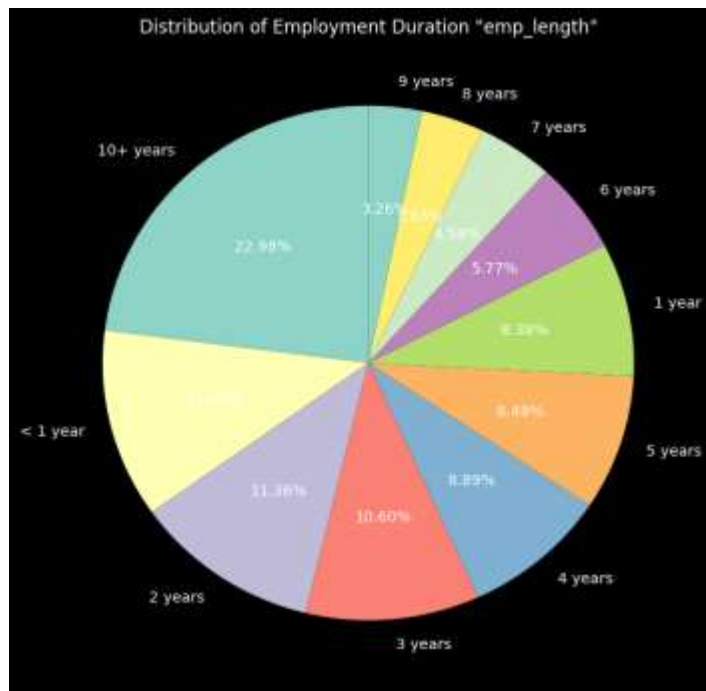
Compared "Charged Off %" across different categories of key columns.

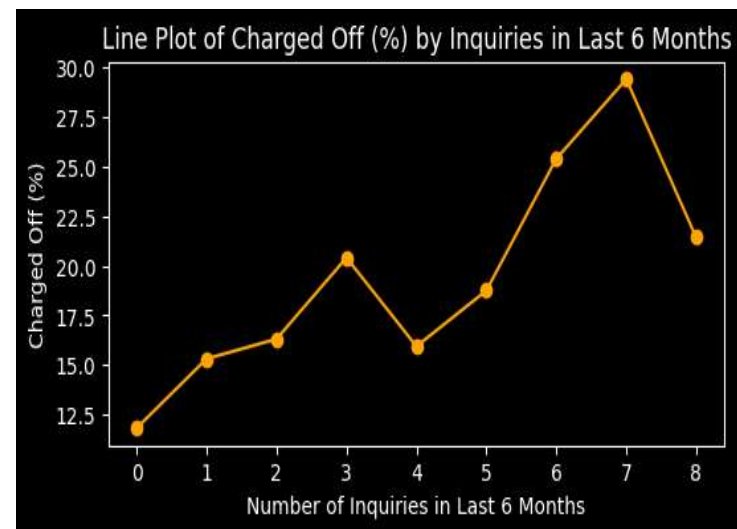
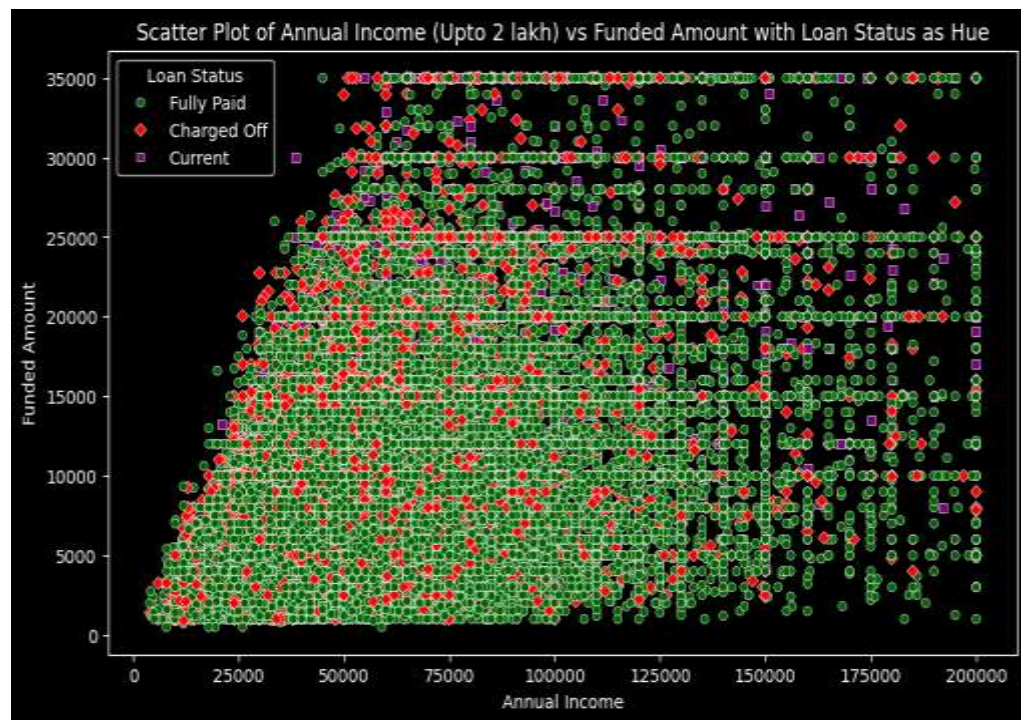
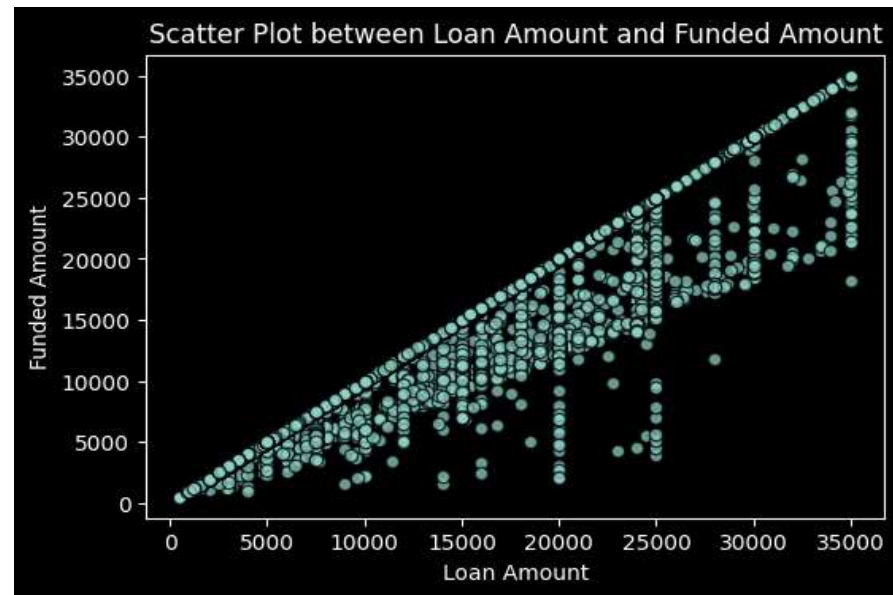
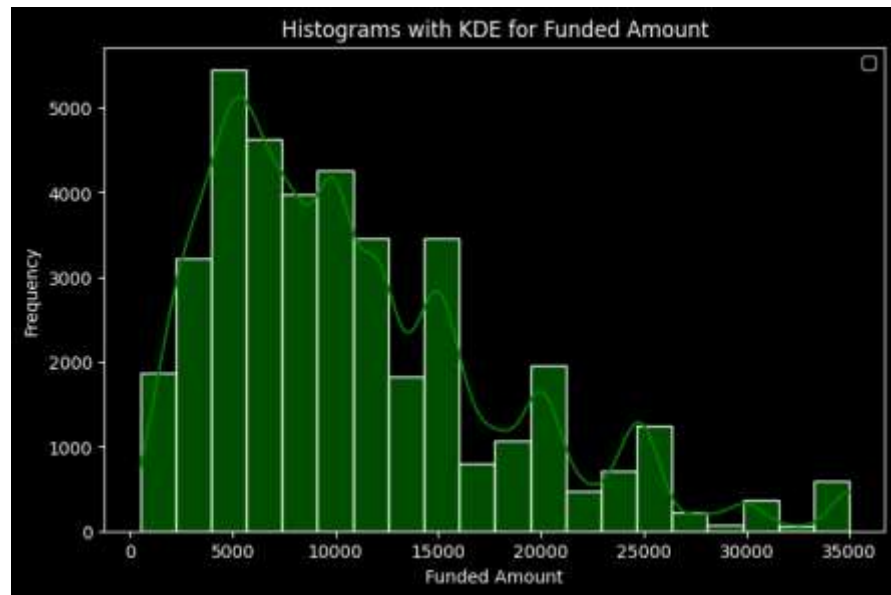
- **Multivariate Analysis:**

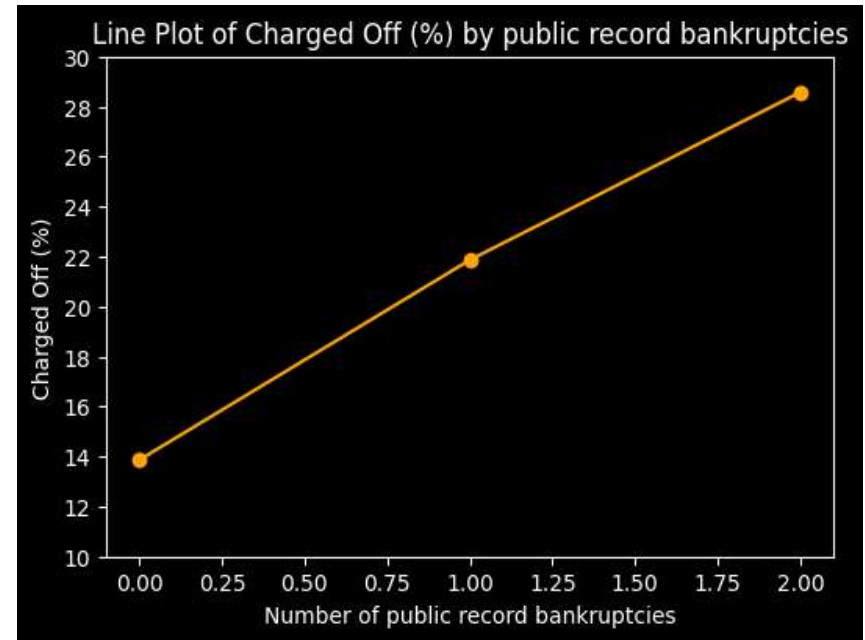
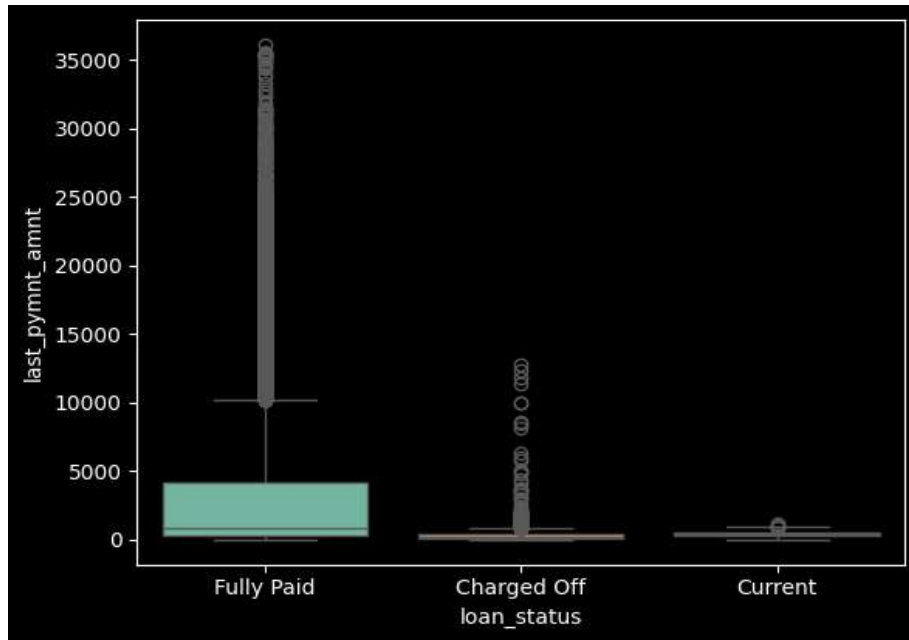
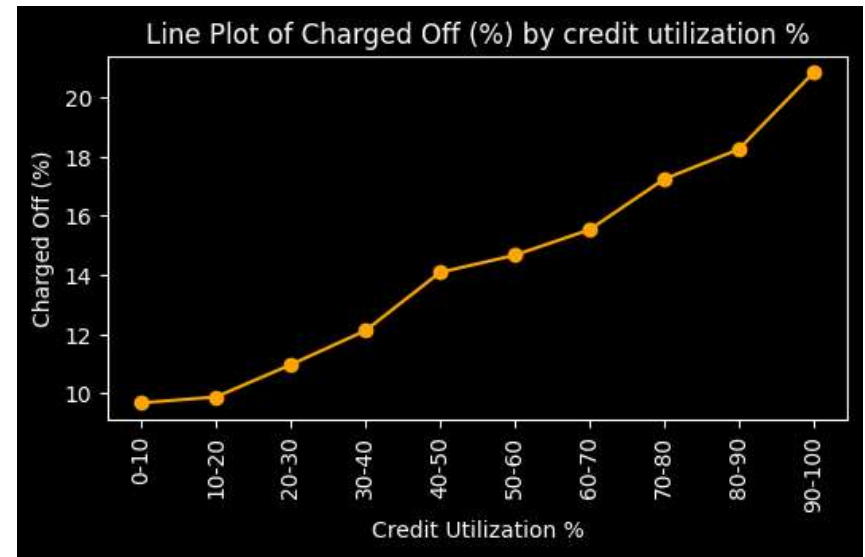
Studied interactions among multiple variables.

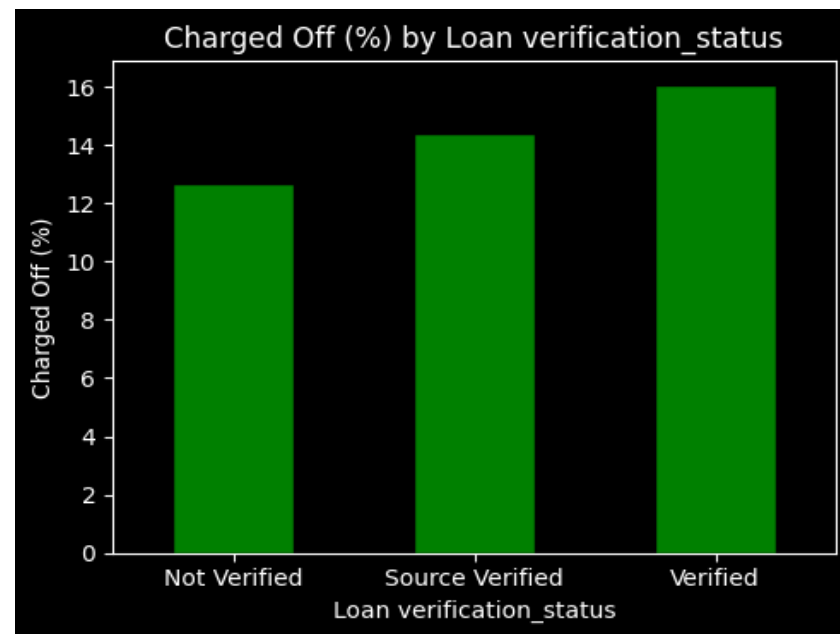
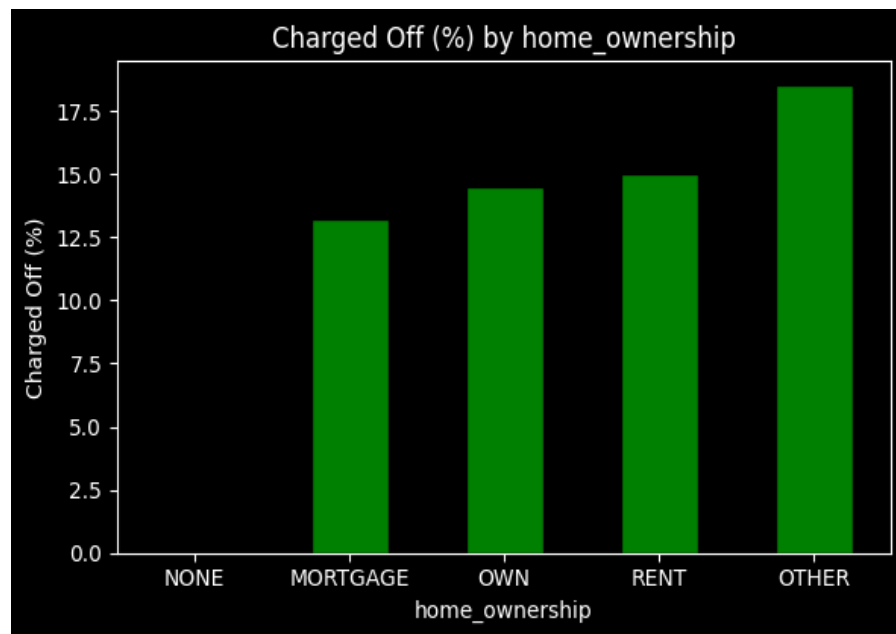
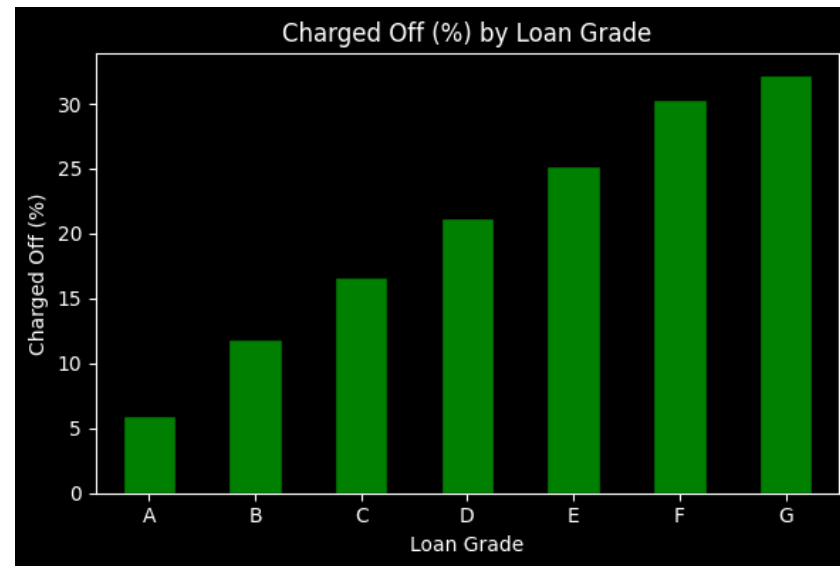
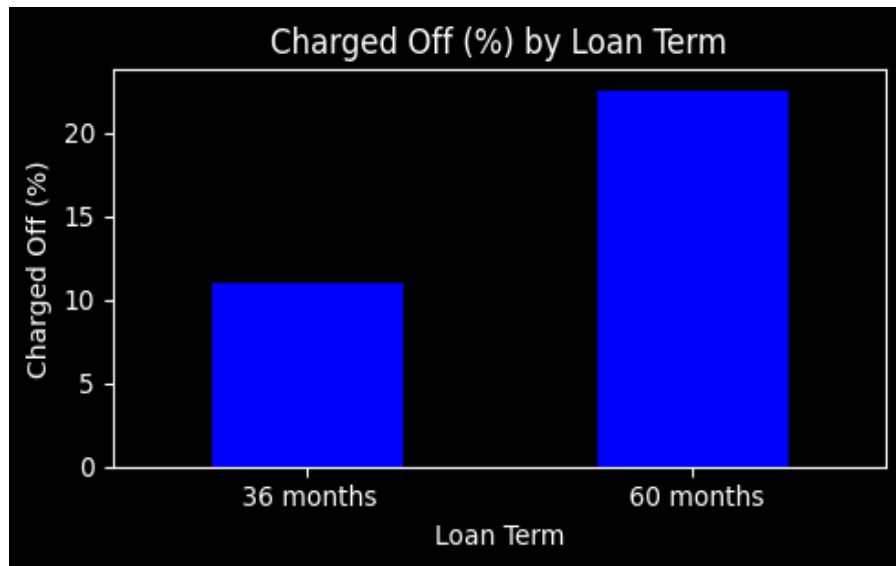
Combined visualizations like scatter plots for deeper insights.

Identified key drivers influencing the "Charged Off" status.









Key Insights

1. As Annual income is rising, installment amount is also rising till $\text{annual_inc} < 1$ Lakh. After that it is random.
2. As annual income is increasing, funded amount has also increased but this is limited to annual income < 1 Lakh, after that it is random
3. With increasing number of inquiries in last 6 months, there are more chances of being charged off
4. Higher the credit utilization %, higher the chances of being charged off.
5. Those who have paid (total payment) equal to or less than 100% of the funded amount, they are charged off
6. Those who have paid almost 80% or lesser of the principal amount, are 'Charged Off'
7. Last payment amounts are way higher in case of 'Fully Paid' than the other two categories - Charged Off, Current

8. With increasing number of public record bankruptcies, there are more chances of being charged off.
9. Loans with a 60-month term have twice the likelihood of being charged off compared to those with a 36-month term.
10. From A to G in Loan Grades, likelihood of being charged off is also increasing.
11. As we go from A1 to G5 in Loan subgrade, the % of being charged off is increasing.
12. There are slight variations of being charged off in case of home_ownership which is OTHER > RENT > OWN > MORTGAGE. OTHER has the highest chances of being charged off.
13. Surprisingly, the loans which are verified have higher chances of being charged off than those which are not verified.
14. Loans for small business purpose is at the most risk and chances are high of being charged off than other categories.

Conclusion

- 1. Longer loan terms (60 months) carry higher default risk.
- 2. High credit utilization and low last payment amounts signal risk.
- 3. Loan verification processes need reevaluation.

Thank You