

# Lending Club Case Study

**Manjunatha P**  
**Pallavi D S**

- Abstract
- Data Insight
- Data Cleaning
- Derived Columns&Dropping the Rows
- Outliers
- Univariate Analysis
- Bivariate Analysis
- Conclusions

# Abstract

## 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:

- Use EDA to understand how consumer attributes and loan attributes influence the tendency of default

## 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.). Since the
    - loan was rejected, there is no transactional history of those applicants with the company and so this data is not available with the company (and thus in this dataset)

# Data Insight

- Loan.csv file contains total **39717** rows and **111** columns.
- There are two types of data related to Loan Attribute and Customer attributes.

# Data Cleaning

- There were no duplicates rows found.
- There were 1140 rows present of loan\_status equal to 'current' which has been deleted as it does n't participate in analysis. There were
- 55 columns which is having all the rows values as null/blank and doesn't participate in analysis has been removed.
- 'URL', 'desc' and 'title' text/description values and doesn't participate has been dropped from analysis.
- Limiting our analysis till 'Group' level only so sub group has been dropped.
- Using domain knowledge, some data columns are related to post the loan approval and doesn't participate in analysis. 19 data columns has deleted.
- 8 columns whose values were 1, and has only one value for all records has been dropped from analysis.
- There were two columns which is having more that 50% of data as null has been removed.
- After all the Data cleaning process we are left with 38577 rows and 19 columns.

# Derived Columns and Dropping Rows

- Additional string value has been trimmed from 'term' column and has been converted to int data types.
- 'int\_rate' has been converted from string to int. Additional '%' has been trimmed.
- Column 'loan\_funded\_amnt' and 'funded\_amnt' converted to float.
- 'loan\_amnt', 'funded\_amnt', 'funded\_amnt\_inv', 'int\_rate', 'dti' columns valued rounded off to two decimal points.
- 'issue\_d' has been converted to datatype.
- Creating a derived columns for 'issue\_year' and 'issue\_month' from 'issue\_d' which will be using for further analysis.
- 'loan\_amnt\_b', 'annual\_inc\_b', 'int\_rate\_b', and 'dti\_b' derived columns has been created for better analysis.

# Outliers

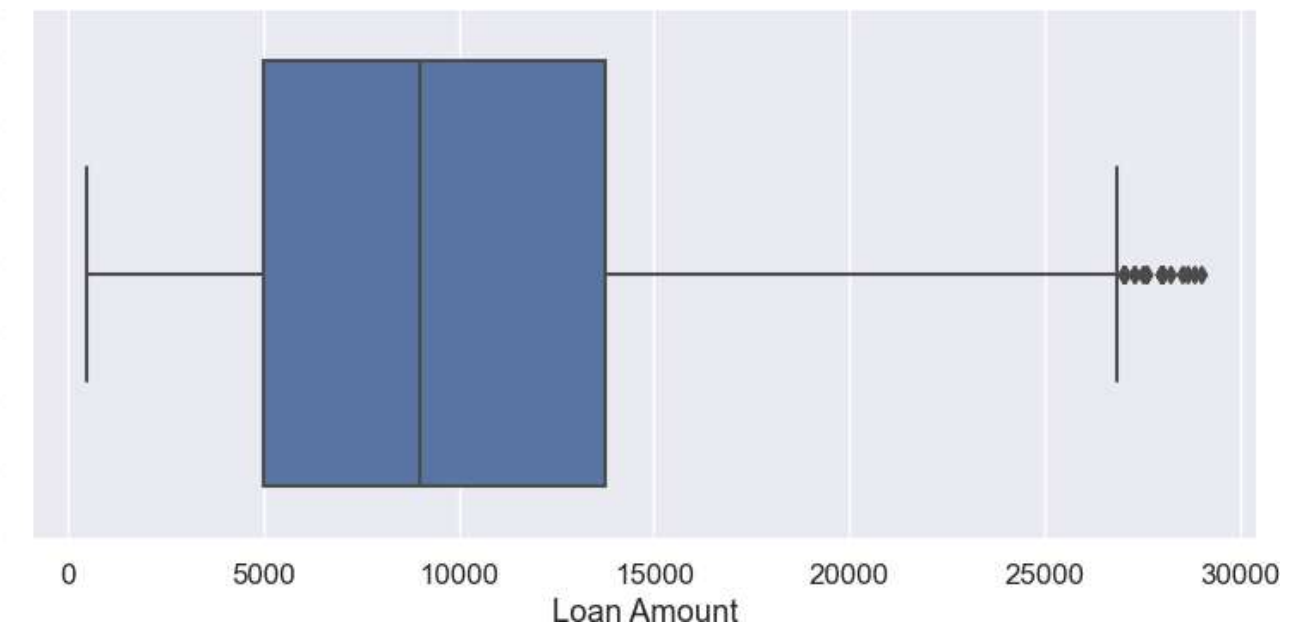
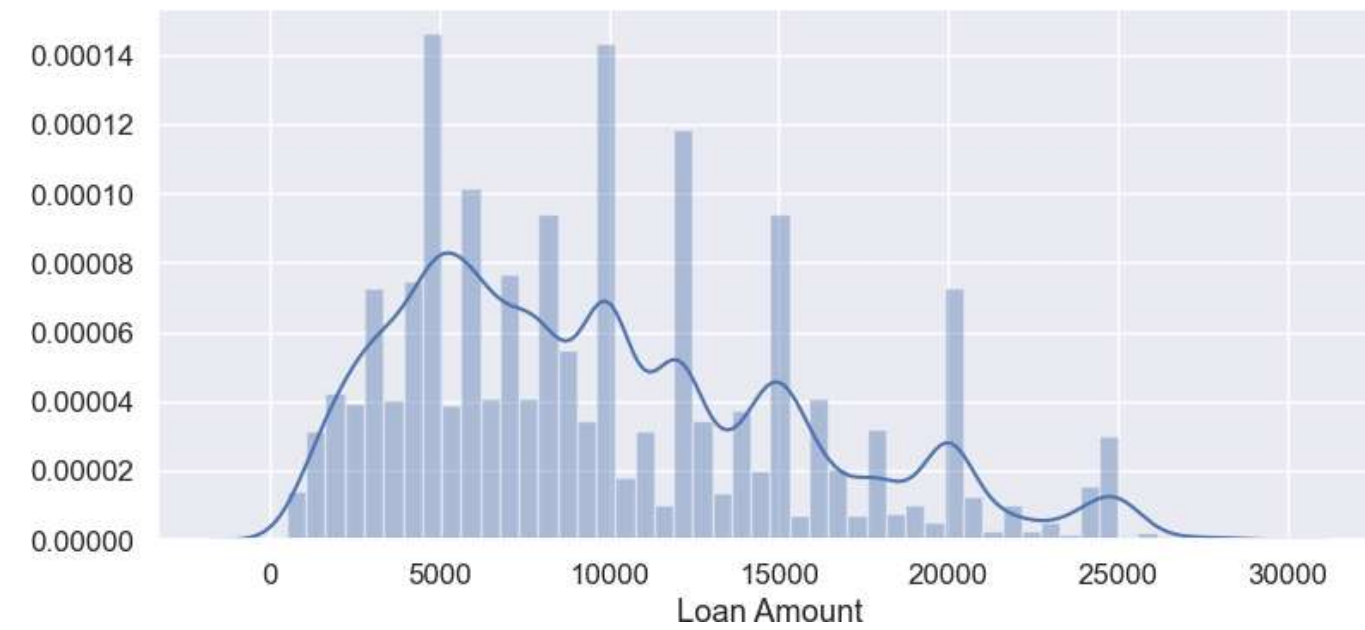
- 'emp\_length' and pub\_rec\_bankruptcies contains 2.67% and 1.80% of rows as null, which is very small percetnage of data which we can drop it.
- Outliers exists for numeric data 'loan\_amnt', 'funded\_amnt', 'funded\_amnt\_inv', 'int\_rate', 'installment', 'annual\_inc'.
- Outliers treatment has been done for above fields using box plot .

# Univariate Analysis



# Loan Amount

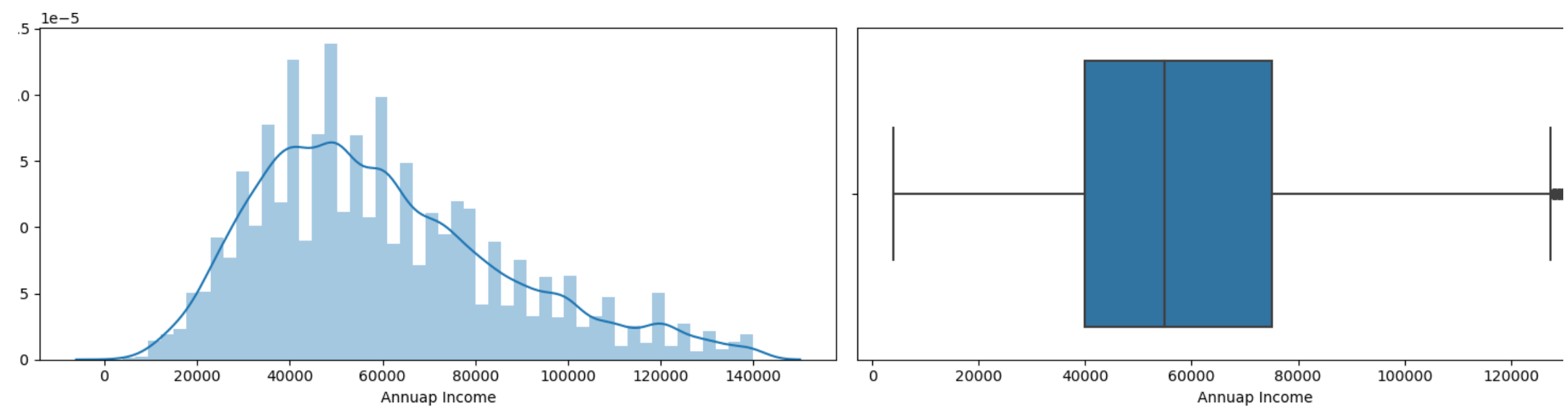
- **Observations:**
  - Most of the loan amount applied was in the range of 5k-14k.
  - Max Loan amount applied was 29k.



# Annual Income

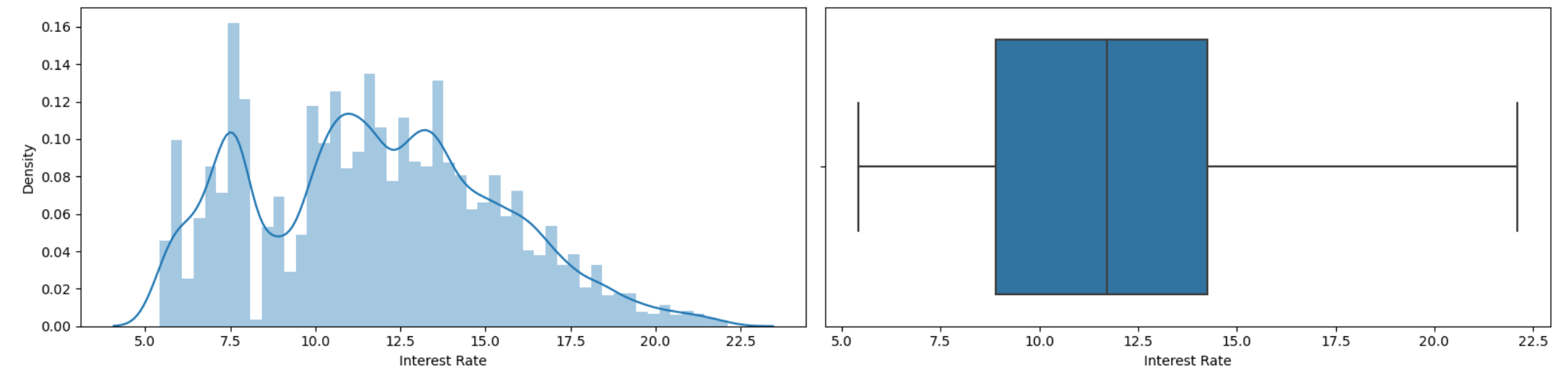
- **Observations:**

- The Annual income of most if applicants lies between 40k-75k.
- Average Annual Income is : 59883



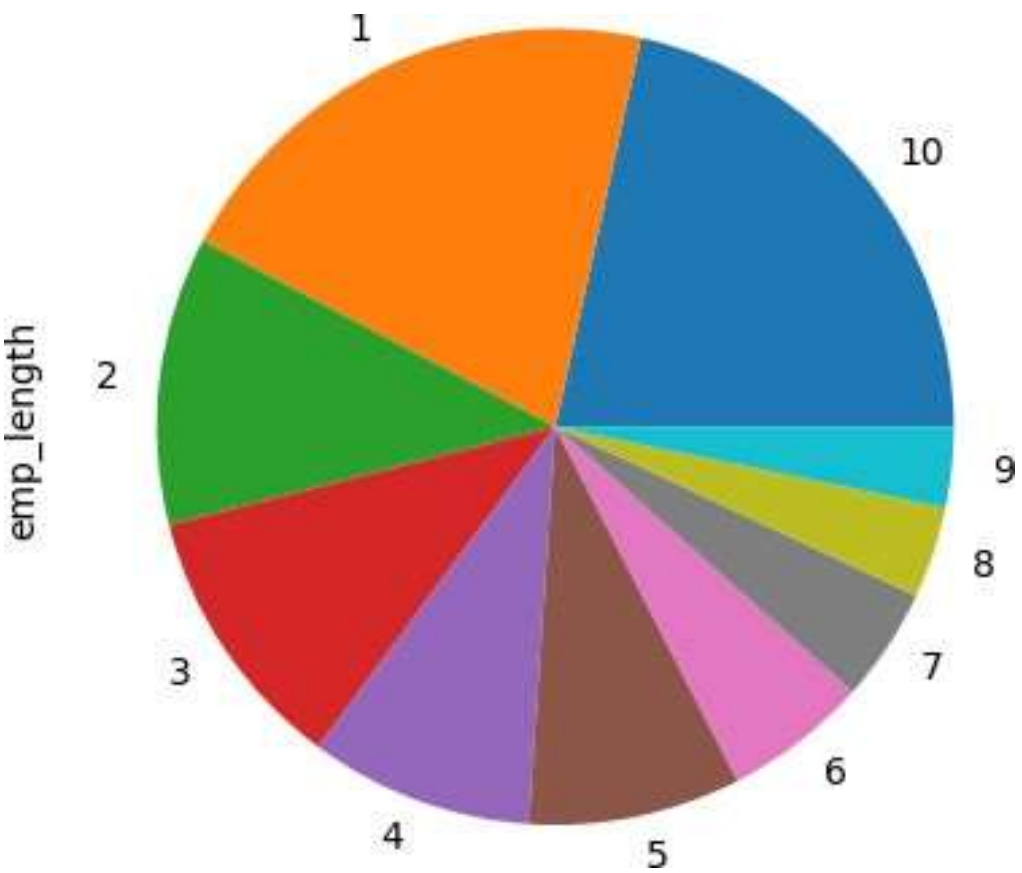
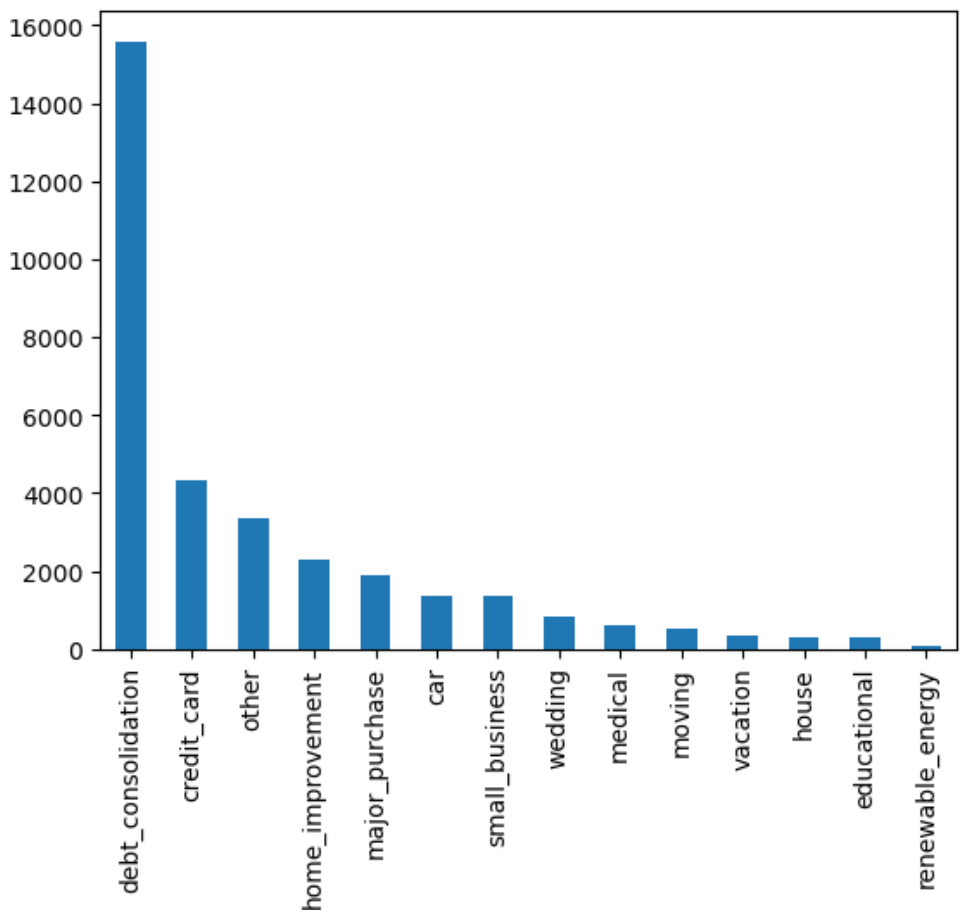
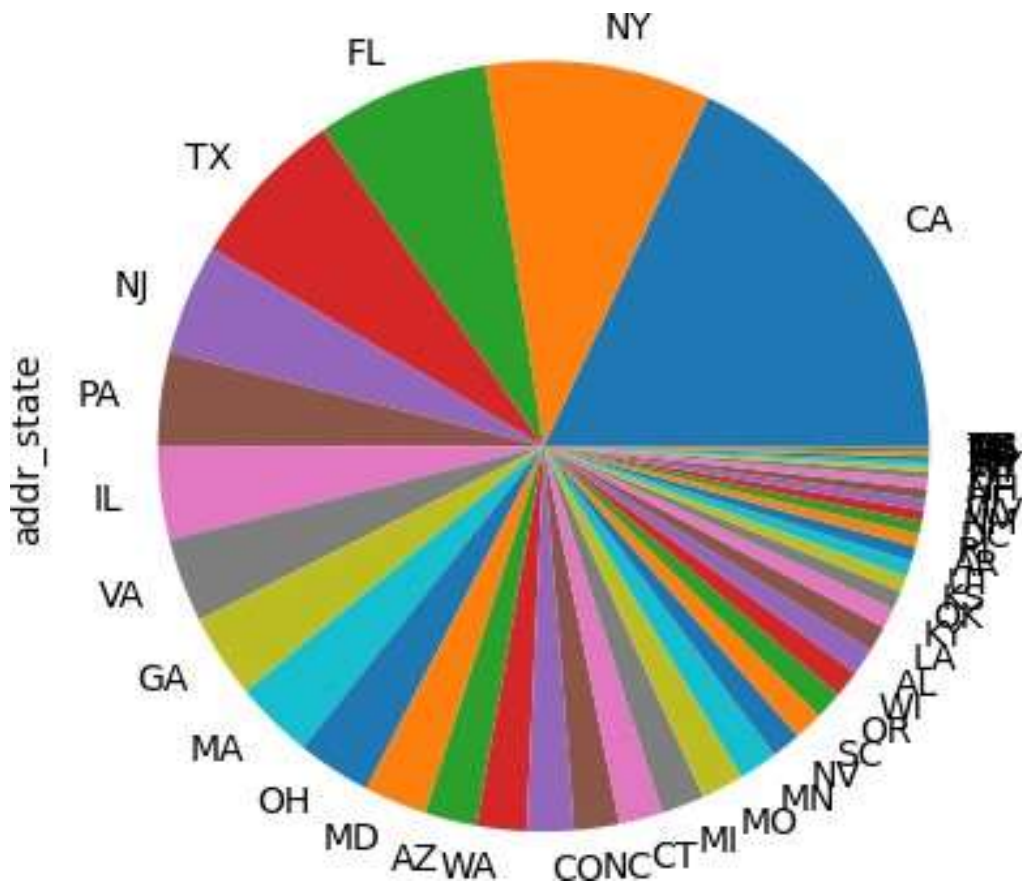
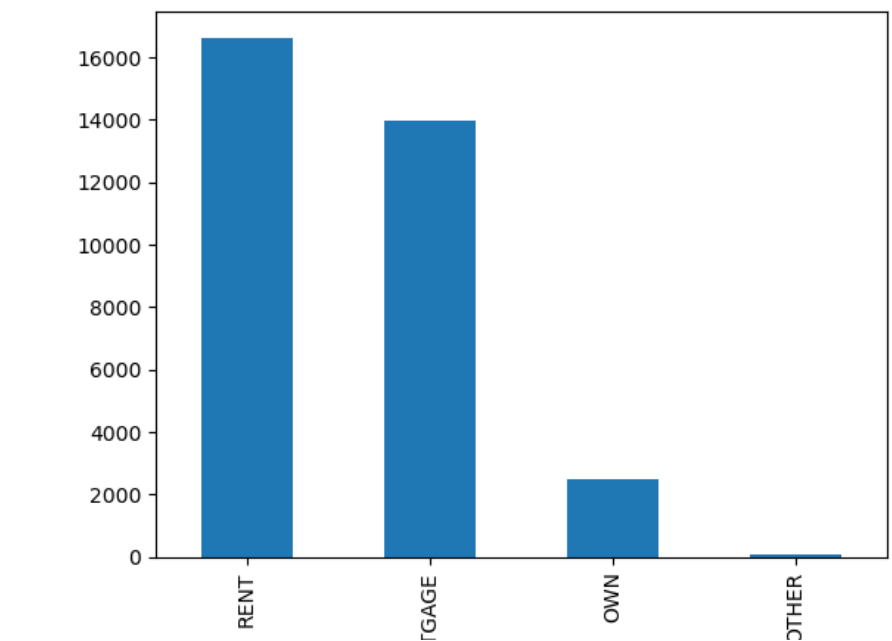
# Interest Rate

- **Observations:**
  - Most of the applicant's rate of interest is between in the range of 8%-14%.
  - Average Rate of interest of rate is 11.7 %



# Univarient Analysis

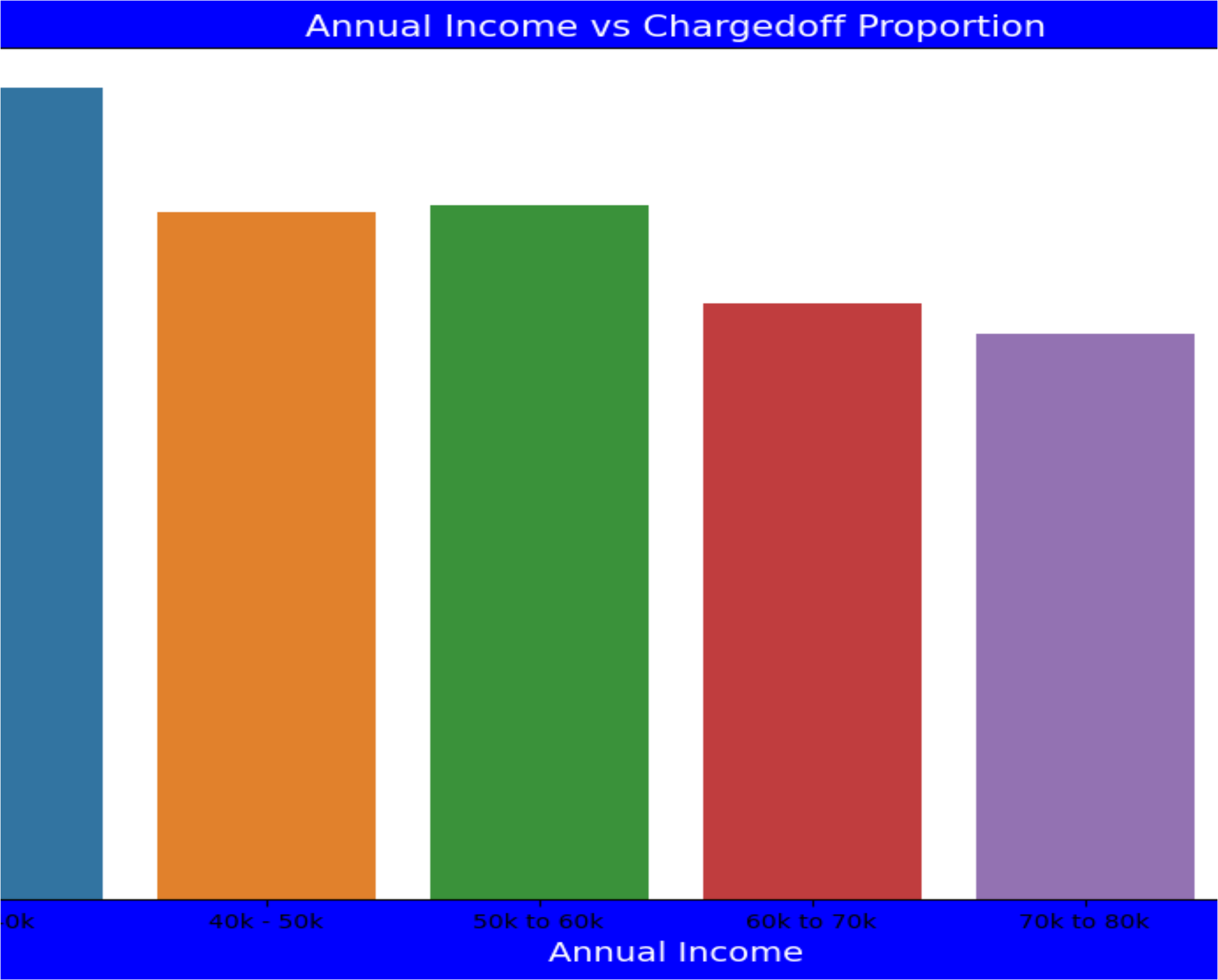
- **Observations:**
  - Majority of loan applicants are either living on Rent or on Mortgage
  - Most of the loan applicants are for debt\_consolidations
  - Most of the Loan applicants are from CA(State).
  - Most of the applications are having 10+ yrs of Exp.



# Bivariate Analysis

# Annual income vs Charged Off

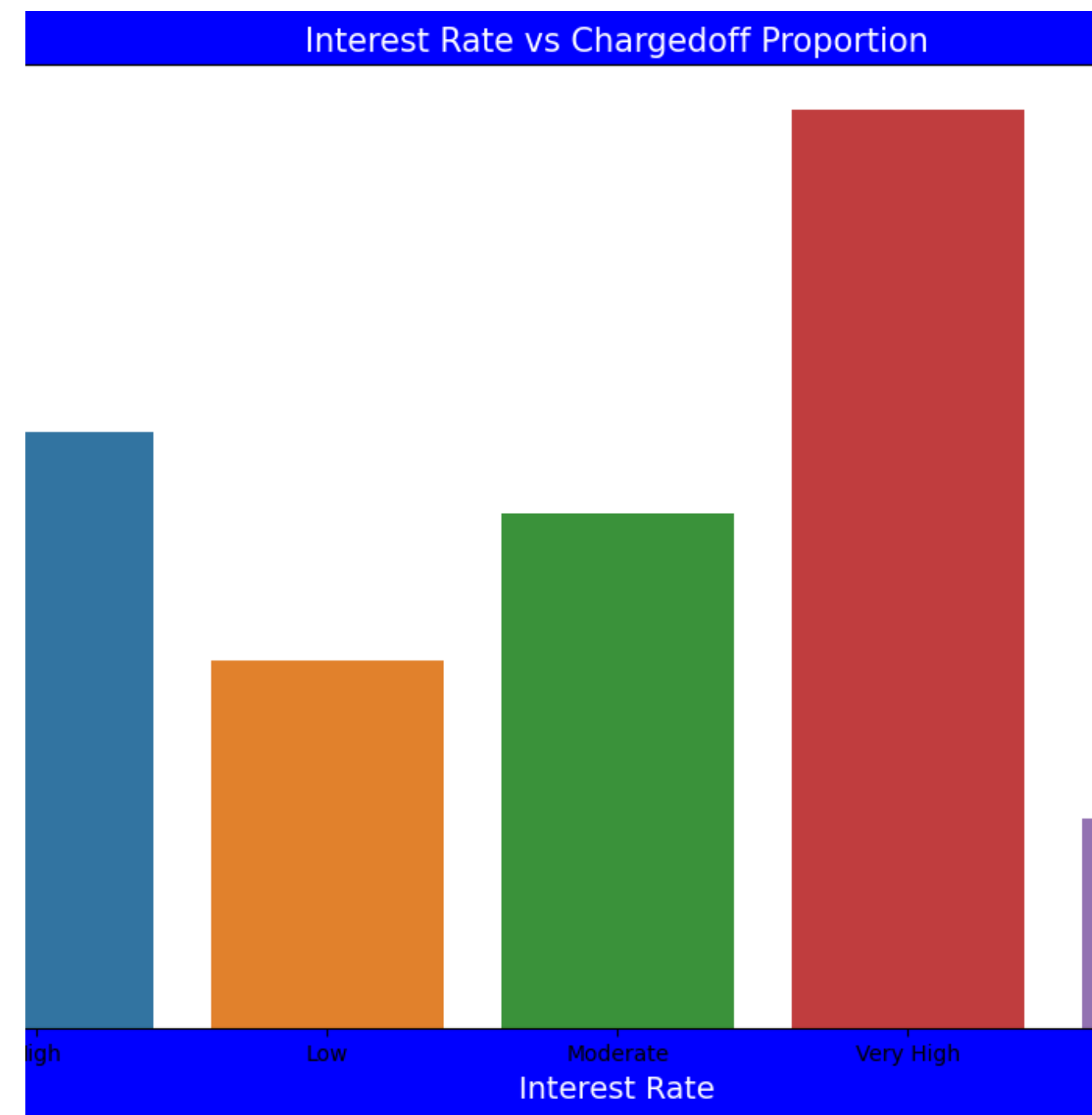
- **Observations:**
  - Income range 80k+ has less chances of charged off.
  - Income range 0-20k has high chances of charged off.
  - Notice that with increase in annual income charged off proportion got decreased.



# Interest Rate vs Charged off

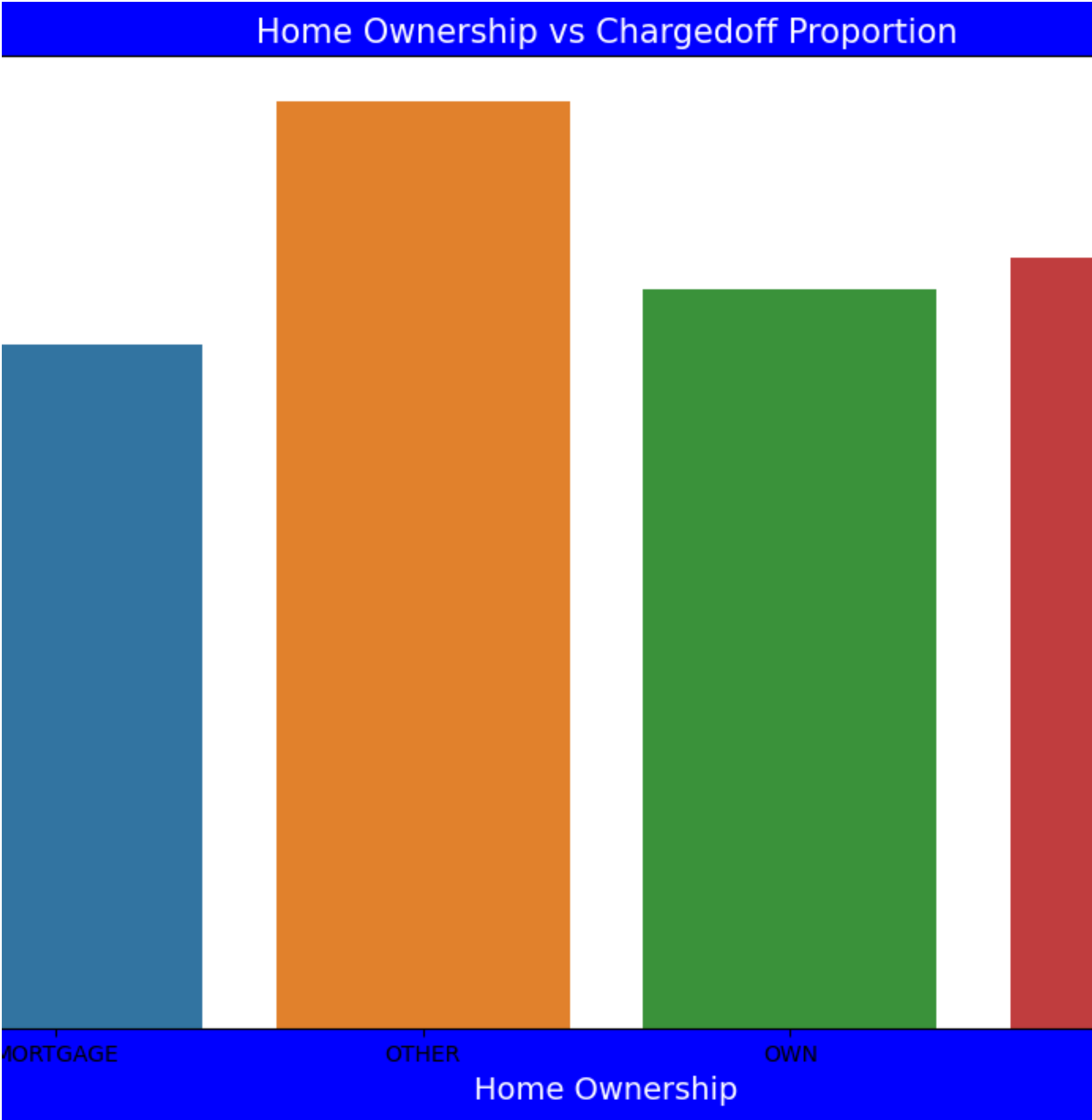
- **Observations:**

- Interest rate less than 10% or very low has very less chances of charged off. Interest rates are starting from minimum 5 %.
- Interest rate more than 16% or very high has good chances of charged off as compared to other category interest rates.
- Charged off proportion is increasing with higher interest rates.



# Home Ownership vs Charged off

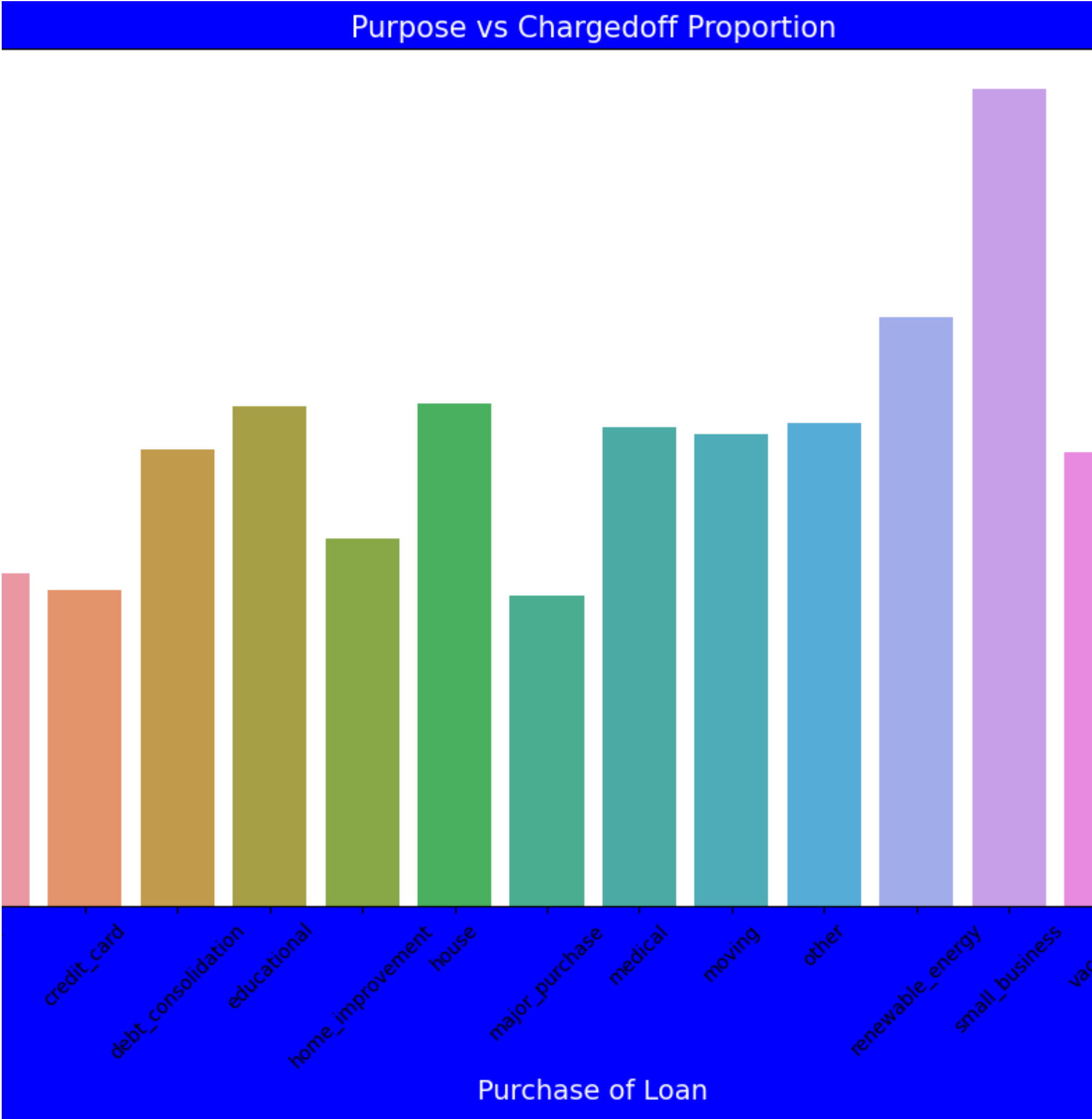
- **Observations:**
  - Those who are not owning the home is having high chances of loan defaulter.
  - From the graph even shows high chances of charged off. Proportions, but data available is very limited compared to other points





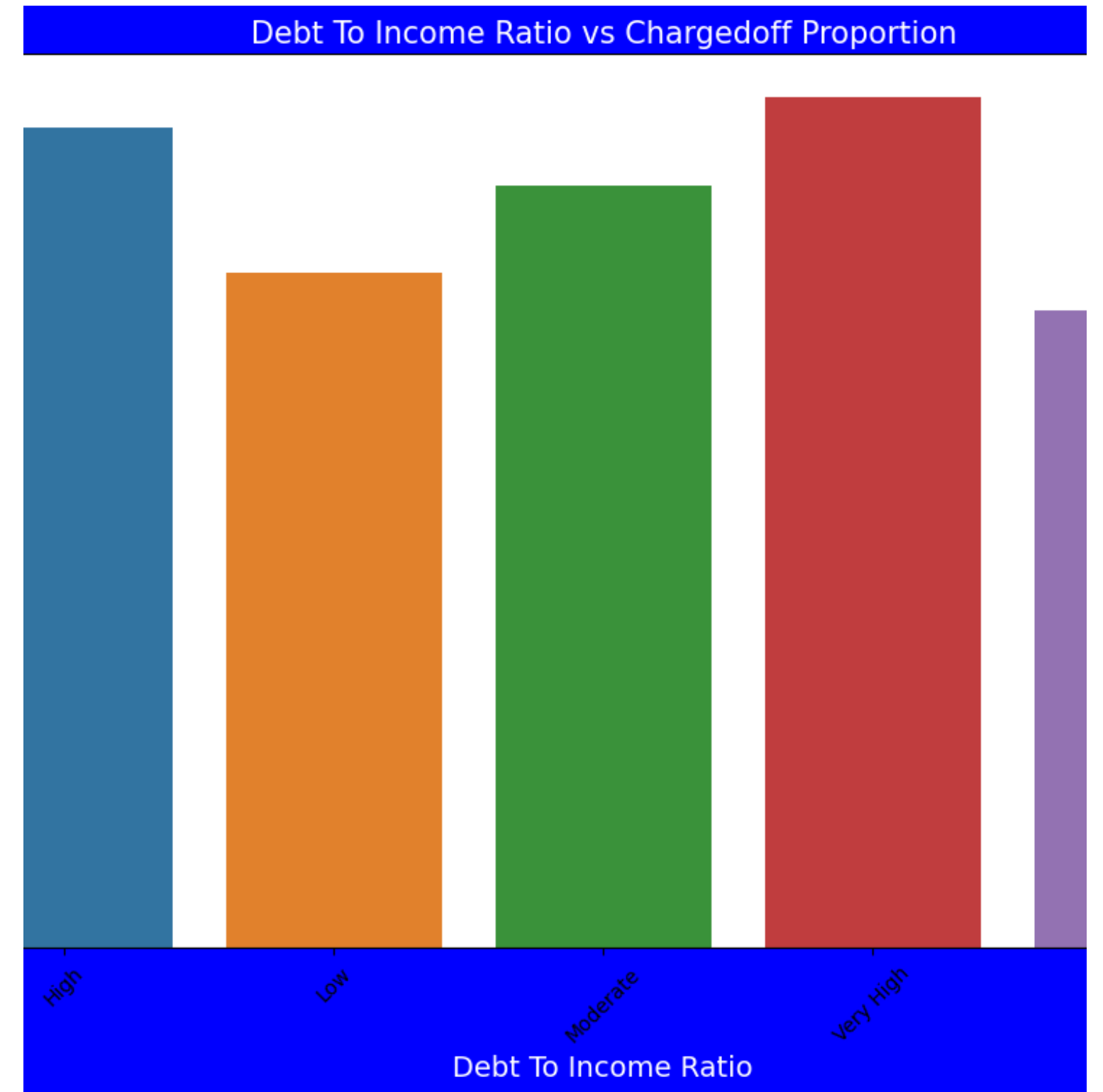
# Purpose vs Charged Off

- **Observations:**
  - Those applicants who is having home loan is having low chances of loan defaults.
  - Those applicants having loan for small business is having high chances for loan defaults.



# DTI Vs Charged off

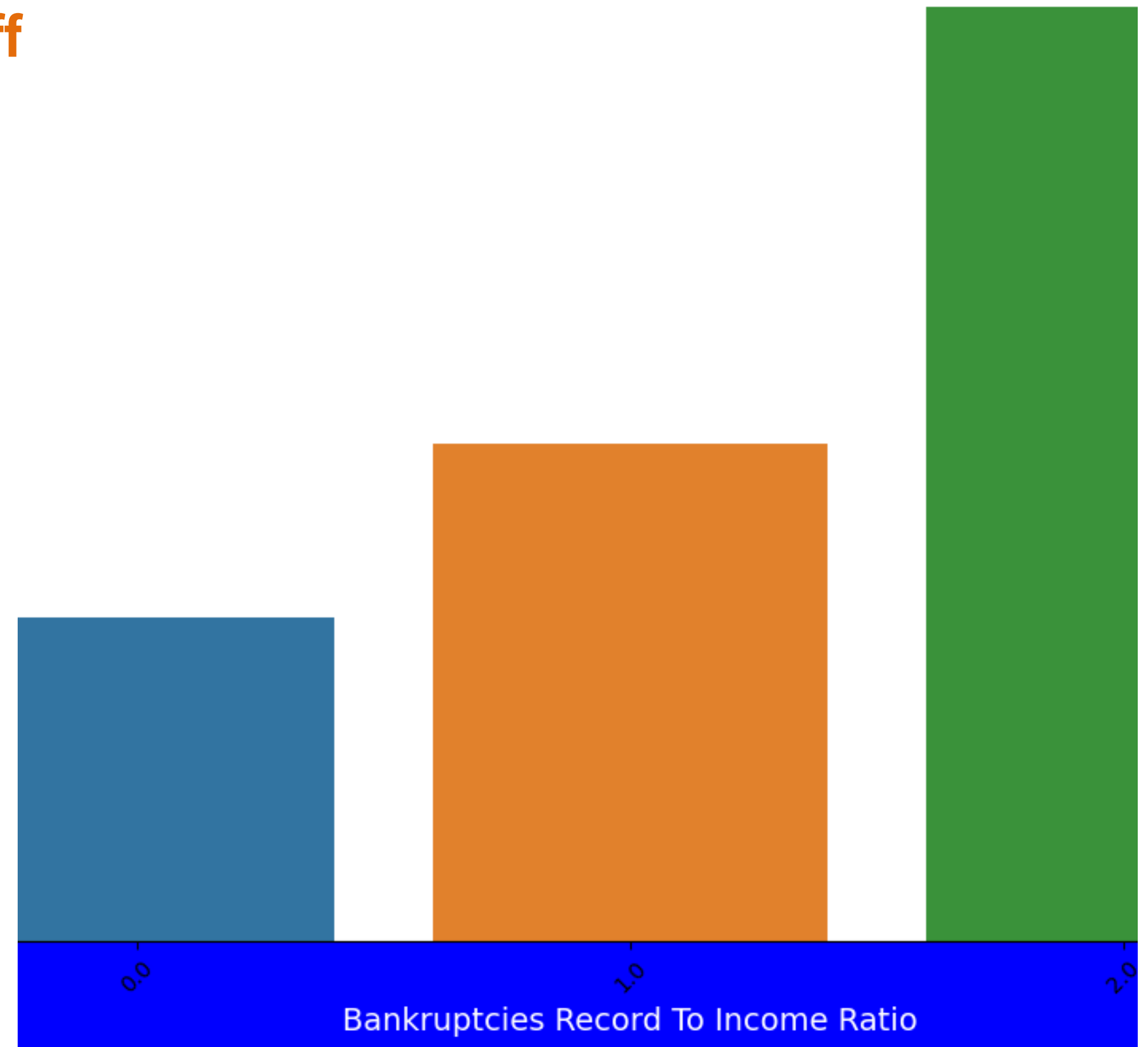
- **Observation:**
  - High DTI value having high risk of defaults
  - Lower the DTO having low chances loan defaults.



# Bankruptcies Record vs Charged off

- Observations:**

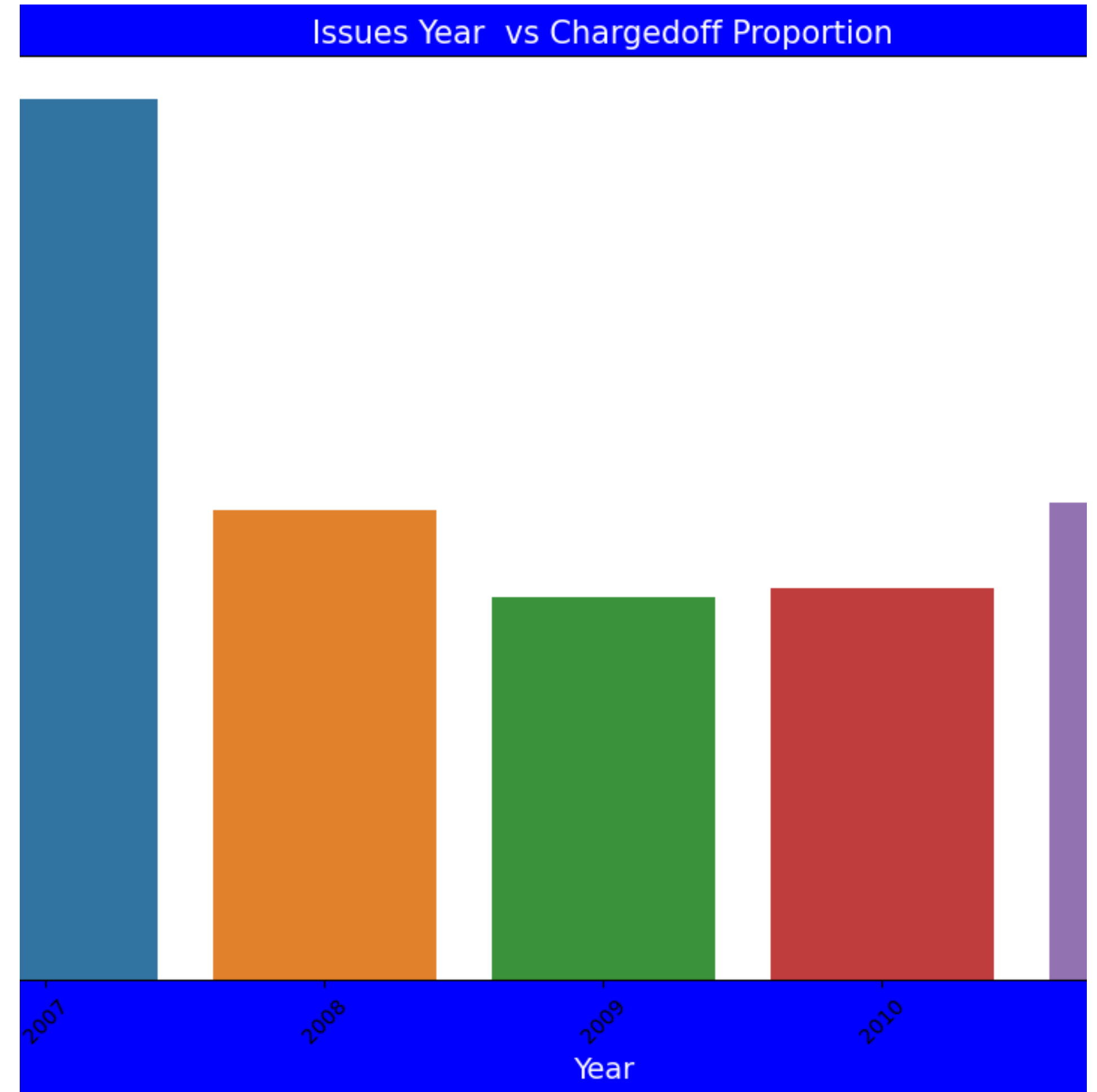
- Bankruptcies Record with 2 is having high impact on loan defaults
- Bankruptcies Record with 0 is low impact on loan defaults
- Lower the Bankruptcies lower the risk.



# Issue Year vs Charged off

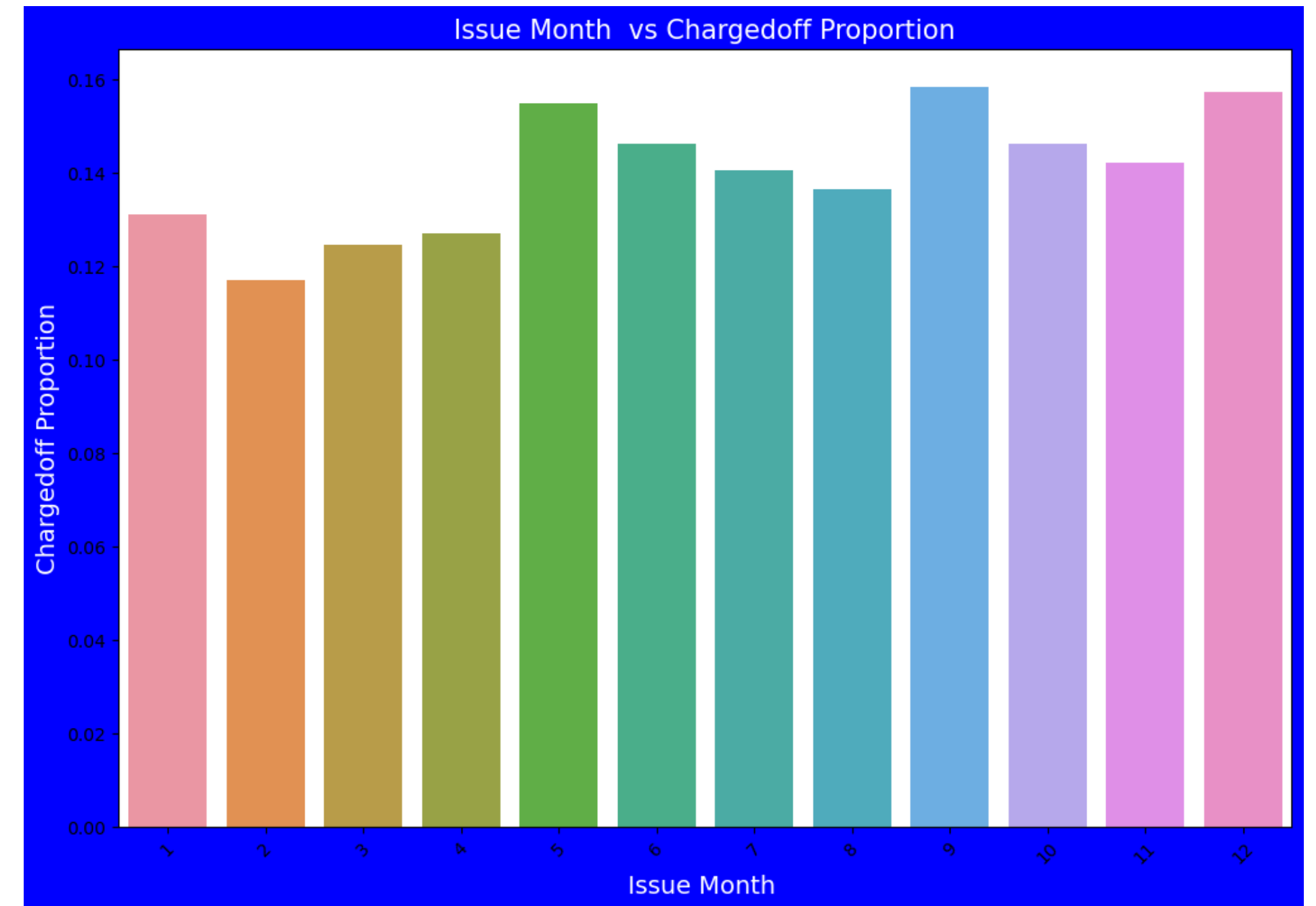
- **Observations:**

- Year 2007 is highest loan defaults.
- 2009 is having lowest loan defaults.



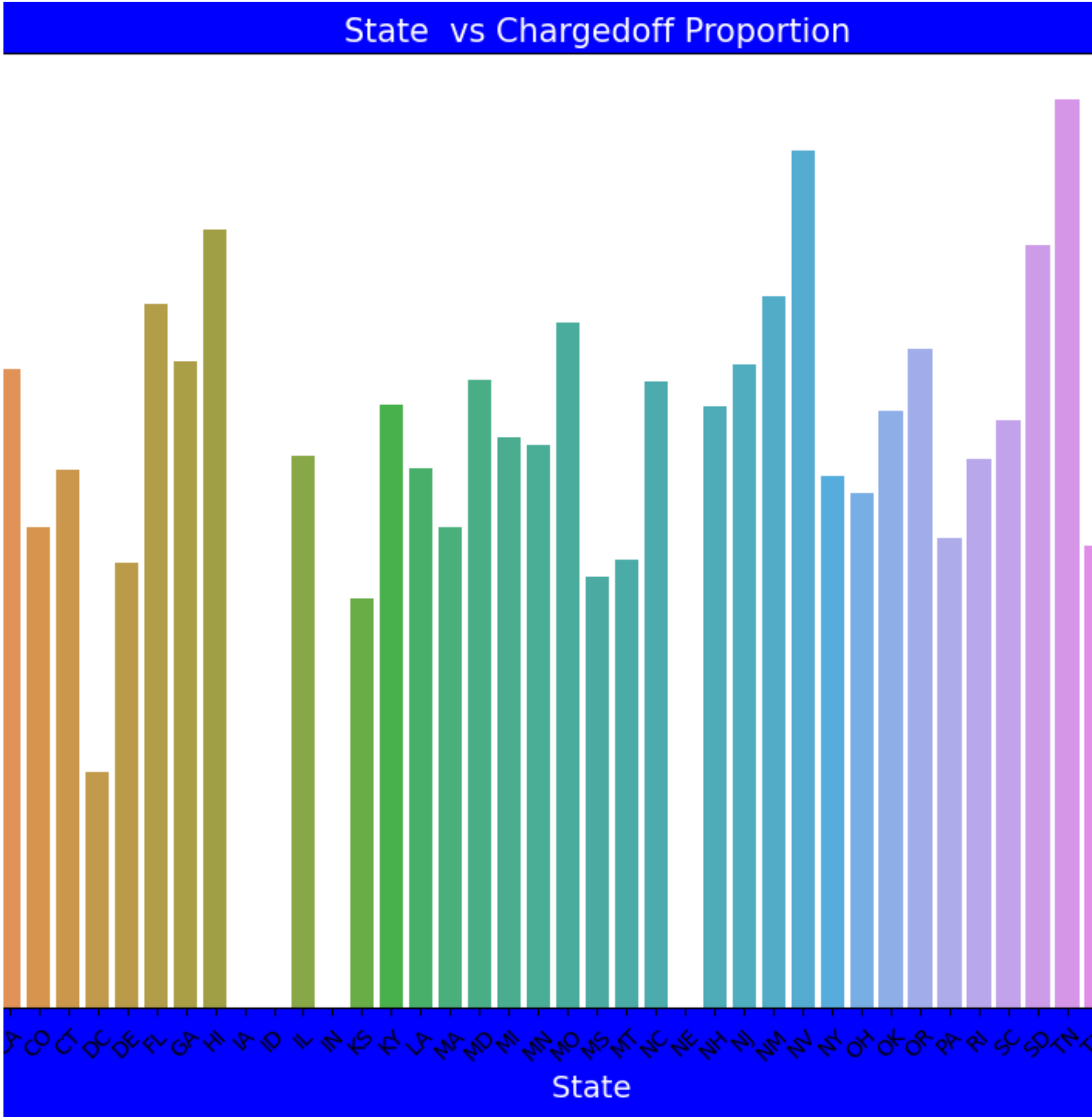
# Issue Month Vs Charged off

- **Observations:**
  - Those loan has been issued in May, September and December is having high number of loan defaults
  - Those loan has been issued in month of February is having high number of loan defaults
  - Majority of loan defaults coming from applicants whose loan has been approved from September-to December



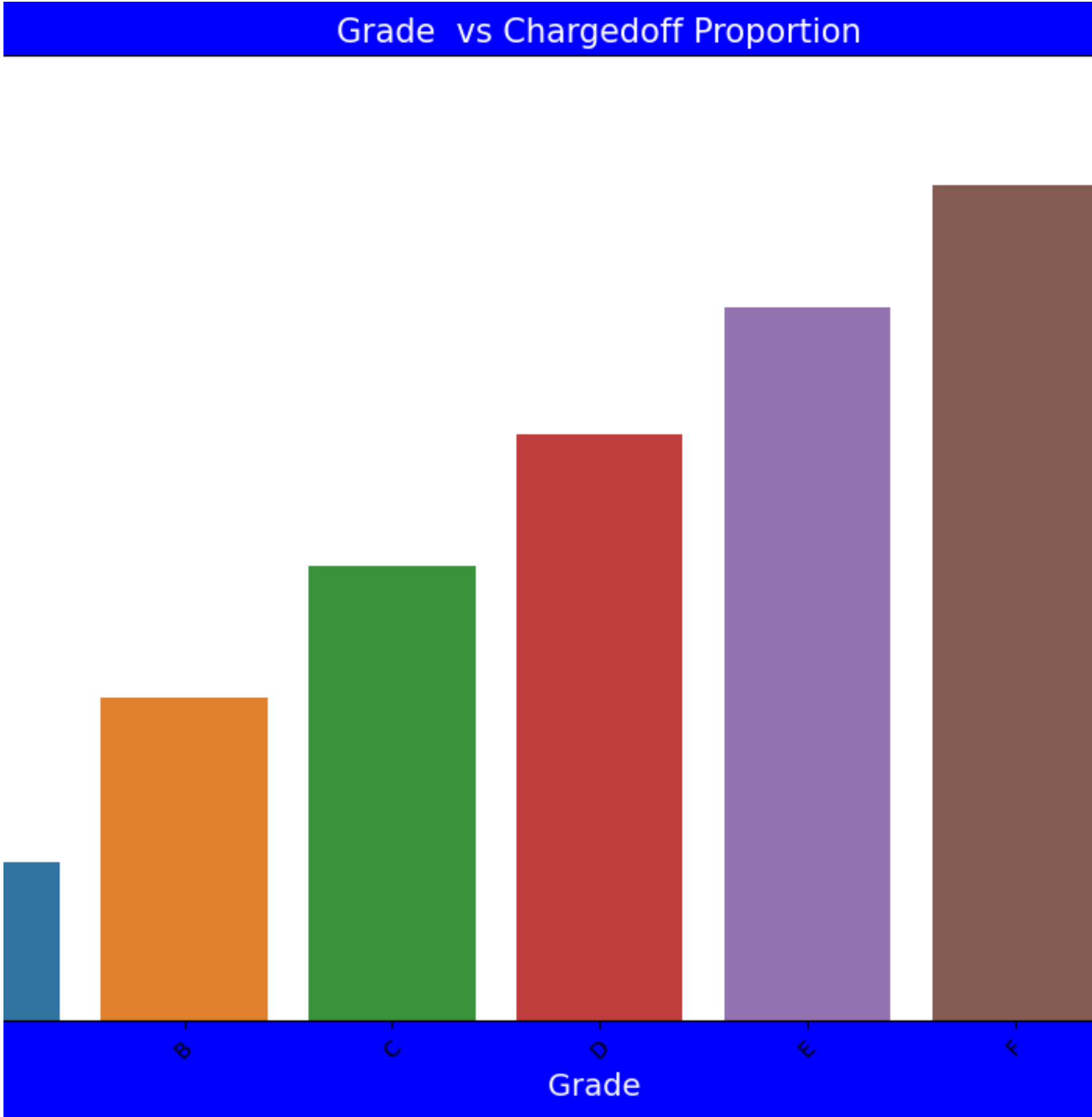
# State vs Charged off

- **Observations:**
  - DE States is holding highest number of loan defaults.
  - CA is having low number of loan defaults



# Grade vs ChargedOff

- **Observations:**
  - The Loan applicants with loan Grade G is having highest Loan Defaults.
  - The Loan applicants with loan A is having lowest Loan Defaults.



# Conclusions

- Income range between 0-20k has high chances of charged off.
- Interest rate more than 16% has good chances of charged off as compared to other category interest rates.
- Those who are not owning the home is having high chances of loan defaulter.
- Those applicants having loan for small business is having high chances for loan defaults.
- High DTI value having high risk of defaults.
- Higher the Bankruptcies record higher the chance of loan defaults.
- DE States is holding highest number of loan defaults.
- The Loan applicants with loan Grade G is having highest Loan Defaults.