# Lending Club Case Study

Data Understanding

- The aim is to find the loan applicants that defaulted or not based on the past history . In this case we need to consider "loan_status"

- Dataset has around 40K entries with 111 features .We need to remove the columns which are not impacting loan status .

- Check the data types for better understanding .

- Check the data dictionary to understand the columns .

- Loan status has 3 values "Current","Fully paid" and "Charged off", since we are not sure on the current status if they will pay up or not , we will ignore them and consider only full paid and charged off status.

# Data cleaning

- Drop columns containing all null value, max null values ,unique value etc .

- Filter rows with missing values

- Change the data types for the columns whereever required eg : float64,int64 etc

- Convert certain column values to numeric variables for better analysis such as remove % or year etc .

# Data Analysis

- Derived Columns
- Univariate Analysis
- Segmented Univariate Analysis
- Bivariate Analysis
- Multivariate Analysis

# Derived Metrics

- Derive columns for Month and Year from column "issue_d"
- Derive column for "loan_amnt" to "annual_inc" ratio as " ratio_loan_to_income"
- Create new column loan_status_code with 0 and 1 values based on loan_status column where 0="Fully Paid" and 1="Charged off"
- Creating different groups for interest rate
- Combine "Source Verified" and "Verified" into "Verified"
- Create Group annual_income
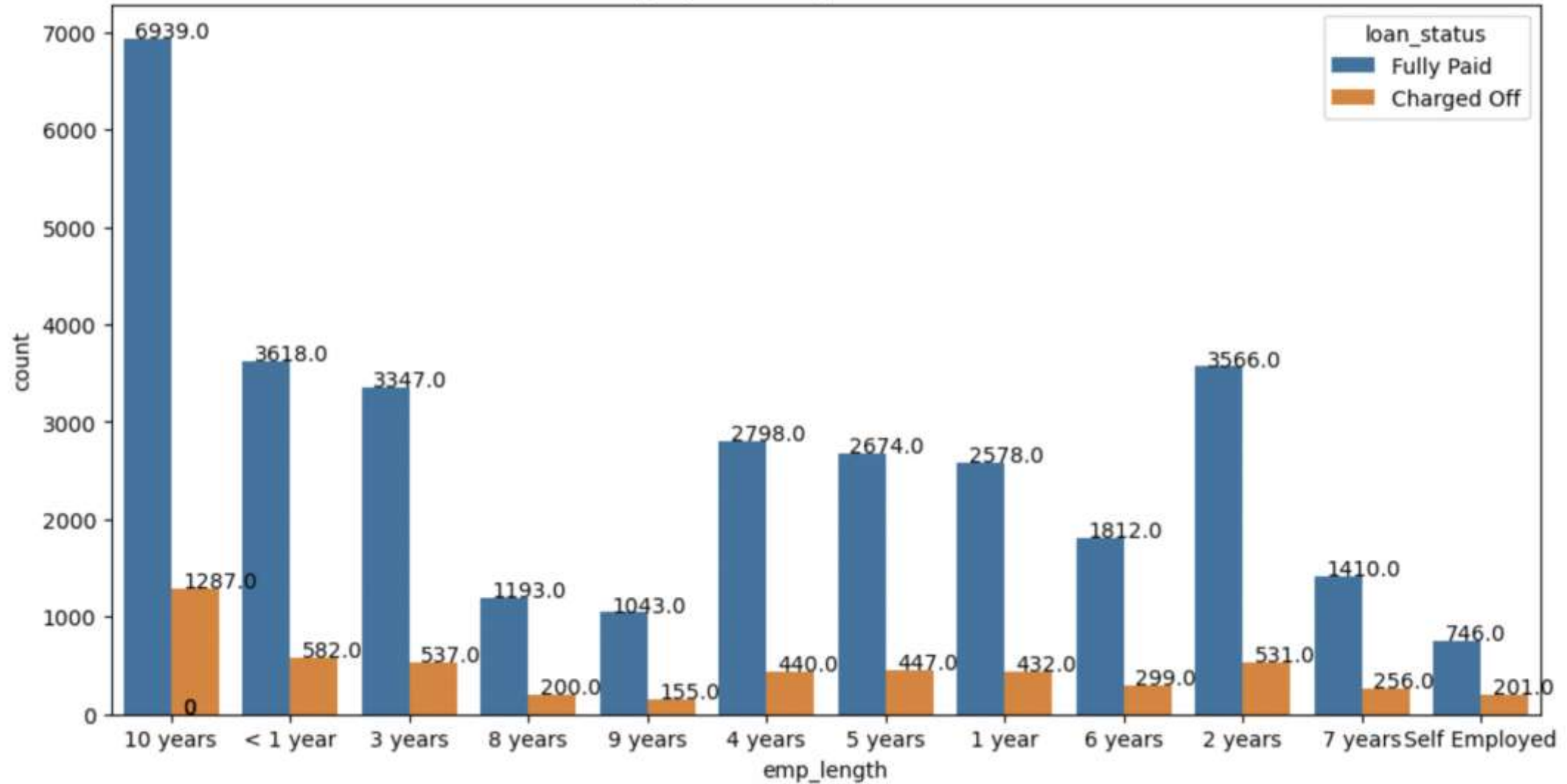
# Univariate Analysis

- Categorical variables
  - Ordered categorical data
    - Grade
    - Sub grade
    - Term (36 / 60 months)
    - Employment length
    - Loan issue year
    - Loan issue month
  - Un-ordered categorical data
    - State
    - Loan purpose
    - Home Ownership
    - Loan status
- Quantitative
  - Interest rate group
  - Annual income group
  - Loan amount
  - Funded amount
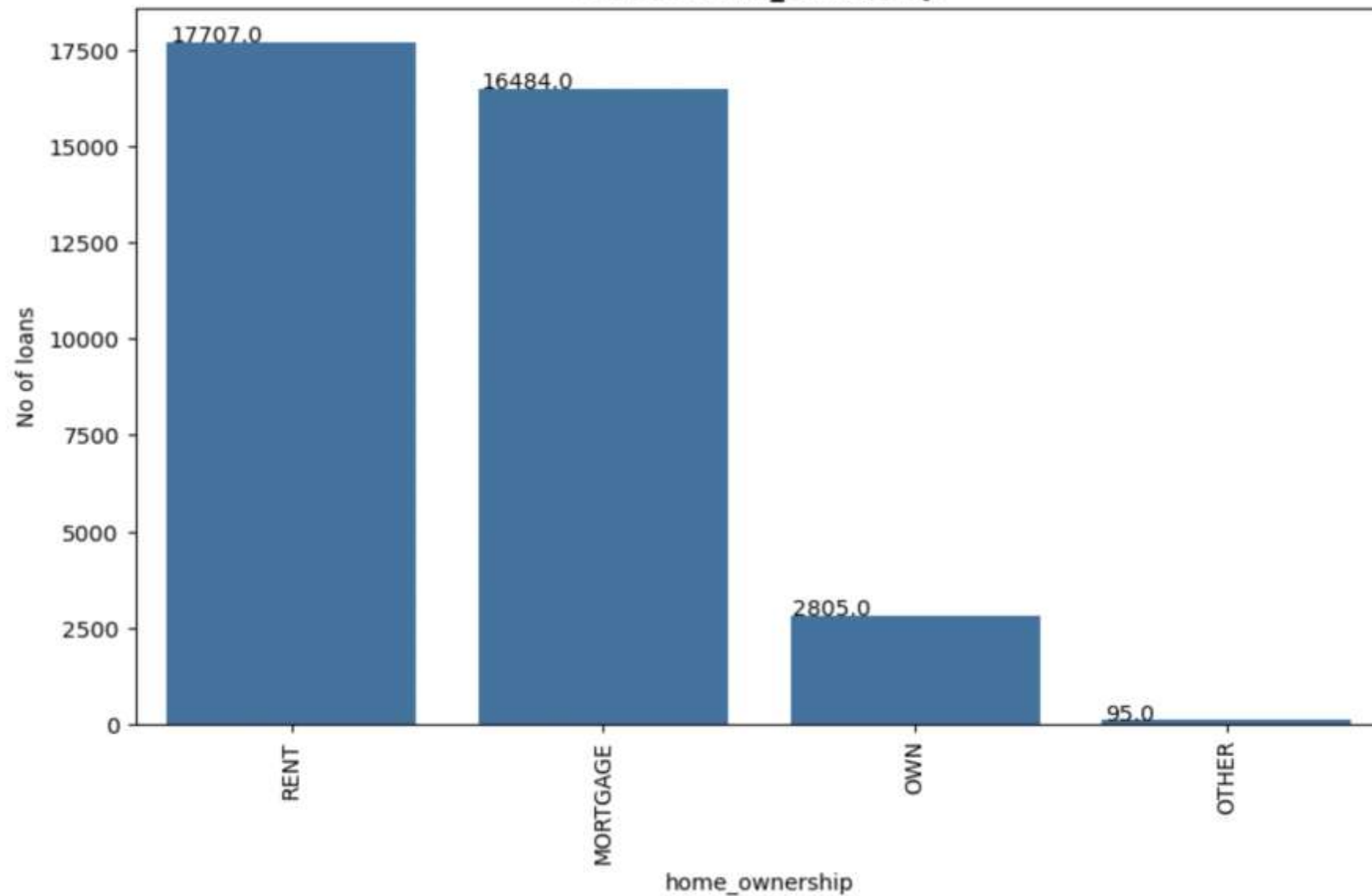  - Loan amount to annual income ratio

# Univariate and Segmented Univariate analysis

- Grade A and B are given more loans compared to other grades
- Sub Grade A4, B3, A5, B5, B4 are given more loans compared to other grades
- term shows 36 months loans are issued more compared to 60 months loan
- employees with 10 years and above are given loan comapred with lesser experience
- year shows maximum loans were taken in the year 2011 and is in increasing trend since 2007
- Maximum loans were given in the month of Oct, Nov, Dec
- 14% of the total loans are charged off and 86% are fully paid
- States CA, NY, FL ,TX ,NJ are the top 5 states where maximum loans been issued .
- Maximum loans are given for debt consolidation .
- People who are in Rented house or Mortgate have availed maximum loans
- Funded amount is ranging from 500 to 35000 USD
- Installment amount is ranging from ~15 to 1300 USD
- The loan amount to income ratio mean is around 0.18
- Interest rate range 10 to 15 is the range where maximum loans have been issued
- 20 - 25% is the range where minimum loans have been issued
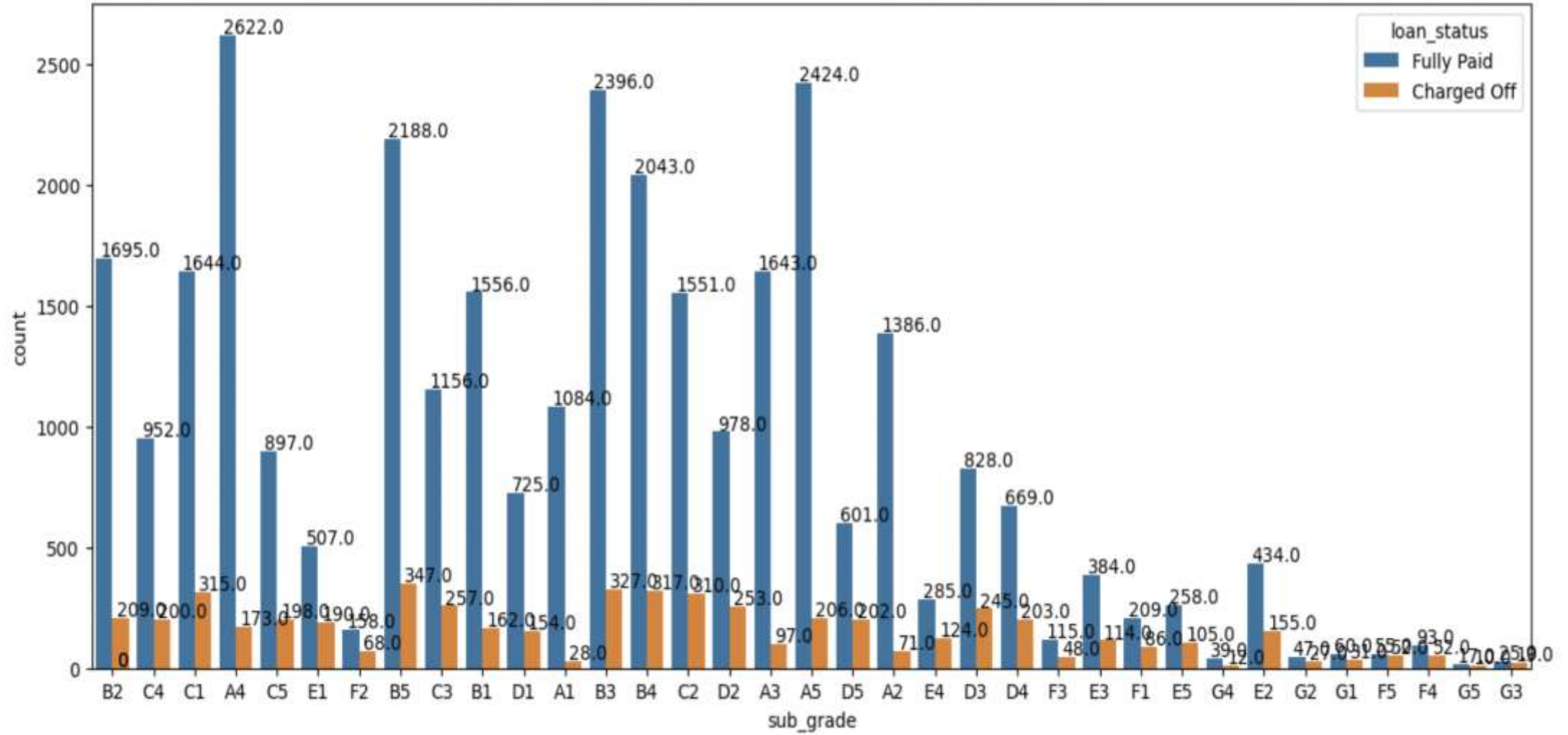
Employment length vs Loan Status
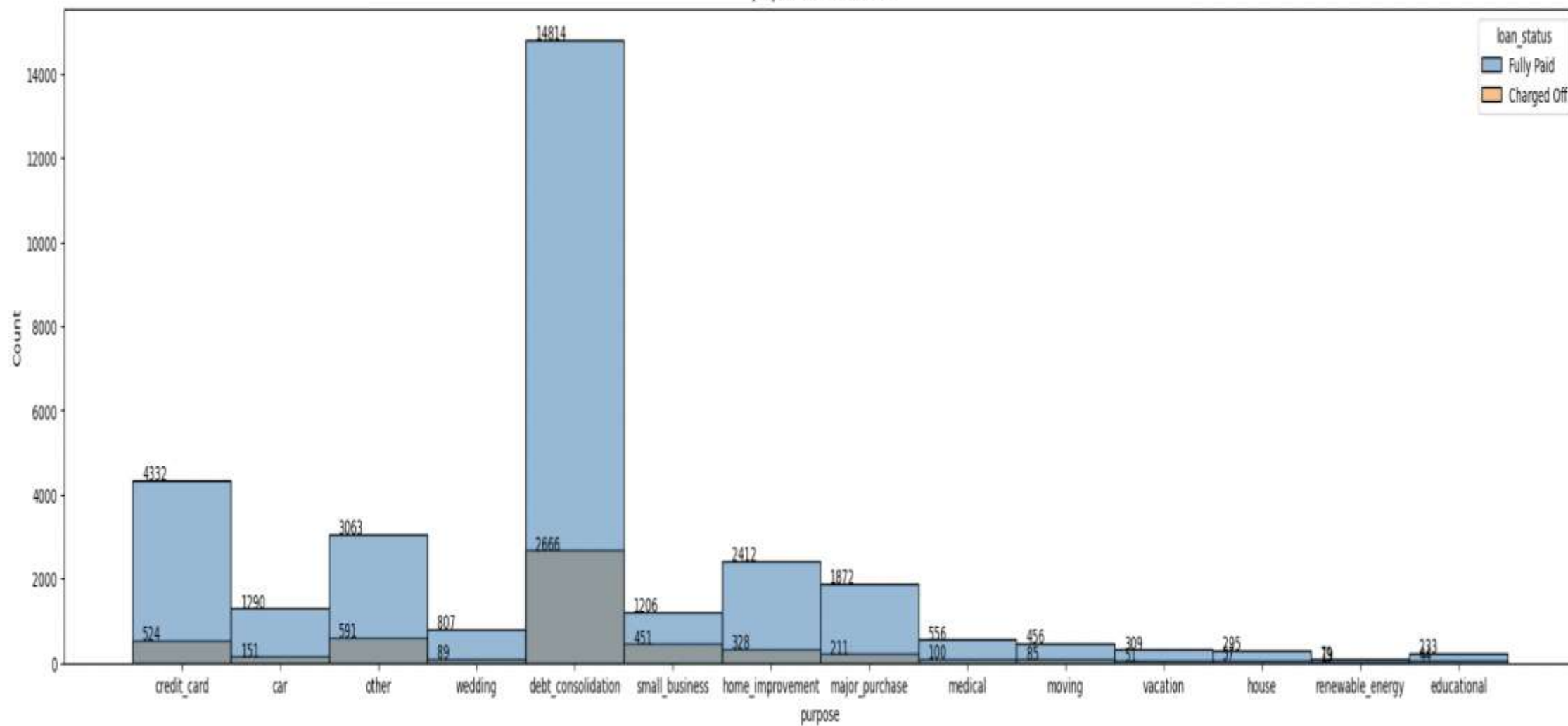
# Plot of home_ownership

Bar chart titled "Plot of home_ownership" with y-axis "No of loans" and x-axis "home_ownership":
- RENT: 17707.0
- MORTGAGE: 16484.0
- OWN: 2805.0
- OTHER: 95.0

Grade vs Loan Status

purpose vs Loan Status

# Summary of Bivariate Analysis

- Based on the counts, Grade B, C and D are top three in Charged Off
- A,B have have also very good numbers of paying it off fully
- Based on the counts sub grades C1,B5,C3,B3,B4,C2 and D2 have higher charged off count
- Maximum loans were issued to 10 or above category and the charged off is also highest in this category
- Defaulters are more in 36 months category
- In 2011 the defaulters are more
- Don't find any conclusive answer looking at the numbers for months.
- The defaulters are more when the purpose is debt_consolidation .
- People of RENT or MORTGAGE failed to pay
- The verified loans are charged off more than Not verified
- Employment above 10 years have got more loans and more defaulters too
- States from CA,FL and NY have more defaulters .
- The median of 10 years and more is highest
- 6,7,8 and 9 having almost same median value .
- The median is more in case of defaulter .
- The IQR is more in case of charged off

# Summary of Multivariate Analysis

- Higher the interest rate higher charged off ratio
- Higher the annual income higher the loan amount slightly.
- increase in number of charged off with increase in year.
- interest rate is increasing with loan amount increase