

# EDA-Lending Club Case Study

Settu Murugan

Sid Sinari

# Problem Statement

- 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
- The data in "loan.csv" contains information about past loan applicants and whether they 'defaulted' or not. 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,

# Objectives:

- **Using Exploratory Data Analysis** we need to understand how consumer attributes and loan attributes influence the tendency of default.
- When a person applies for a loan, there are **two types of decisions** that could be taken by the company:
  1. **Loan accepted:** If the company approves the loan, there are 3 possible scenarios described below:
    1. **Fully paid:** Applicant has fully paid the loan (the principal and the interest rate)
    2. **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'.
    3. **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
  2. **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)

# Loan Data:

- Shared loan.csv file contains 39717 Rows and 111 columns.
- Data Dictionary provided to explain the details of 111 columns
- This file contains several details about the loan like, amount, interest rate, distribution date, terms and etc..
- This file contains several details about customers like, ID, Address, grade, annual income, employment length, monthly debt paying ratio and etc..

# 1. Data Cleaning

- The first part of EDA is Data cleaning. Looking at the head and tail of data frame Looks like there are no header and footer present. Also no summary columns available in the data frame.
- After checking for duplicate, could not find any duplicate rows.
- Found 54 null columns.  
['mths\_since\_last\_major\_derog', 'annual\_inc\_joint', 'dti\_joint', 'verification\_status\_joint', 'tot\_coll\_amt', 'tot\_cur\_bal', 'open\_acc\_6m', 'open\_il\_6m', 'open\_il\_12m', 'open\_il\_24m', 'mths\_since\_rcnt\_il', 'total\_bal\_il', 'il\_util', 'open\_rv\_12m', 'open\_rv\_24m', 'max\_bal\_bc', 'all\_util', 'total\_rev\_hi\_lim', 'inq\_fi', 'total\_cu\_tl', 'inq\_last\_12m', 'acc\_open\_past\_24mths', 'avg\_cur\_bal', 'bc\_open\_to\_buy', 'bc\_util', 'mo\_sin\_old\_il\_acct', 'mo\_sin\_old\_rev\_tl\_op', 'mo\_sin\_rcnt\_rev\_tl\_op', 'mo\_sin\_rcnt\_tl', 'mort\_acc', 'mths\_since\_recent\_bc', 'mths\_since\_recent\_bc\_dlq', 'mths\_since\_recent\_inq', 'mths\_since\_recent\_revol\_delinq', 'num\_accts\_ever\_120\_pd', 'num\_actv\_bc\_tl', 'num\_actv\_rev\_tl', 'num\_bc\_sats', 'num\_bc\_tl', 'num\_il\_tl', 'num\_op\_rev\_tl', 'num\_rev\_accts', 'num\_rev\_tl\_bal\_gt\_0', 'num\_sats', 'num\_tl\_120dpd\_2m', 'num\_tl\_30dpd', 'num\_tl\_90g\_dpd\_24m', 'num\_tl\_op\_past\_12m', 'pct\_tl\_nvr\_dlq', 'percent\_bc\_gt\_75', 'tot\_hi\_cred\_lim', 'total\_bal\_ex\_mort', 'total\_bc\_limit', 'total\_il\_high\_credit\_limit']
- Deleted all the above columns having Null values as it is not useful for analysis.
- Deleted the columns having unique values since it may not be useful for analysis ( *id, member\_id, url, chargeoff\_within\_12\_mths*)
- Now Delete all the columns which are having all the values as same.. also may not be usefull for further analysis(*pymnt\_plan', 'initial\_list\_status', 'collections\_12\_mths\_ex\_med*).
- Some columns having same values across all loan Ids, hence may not be useful in the analysis. Hence removing them (*sub\_grade', 'total\_rec\_late\_fee', 'last\_pymnt\_amnt', 'last\_pymnt\_d', 'next\_pymnt\_d', 'last\_credit\_pull\_d', 'policy\_code', 'application\_type', 'acc\_now\_delinq', 'delinq\_amnt', 'tax\_liens', 'earliest\_cr\_line'*)
- *#Now Lets drop all columns having more than 40% values are null (mths\_since\_last\_delinq', 'mths\_since\_last\_record)*
- *Some columns values are grouped together for better analysis example., Purpose column has home\_improvement , house both are having the same meaning in the purpose. So treating them as same value.*
- *Now we do not need a title or description of the loan as it is already categorized properly with "purpose" (desc', 'title', 'emp\_title). Hence delete it.*

## 2. Data Standardization

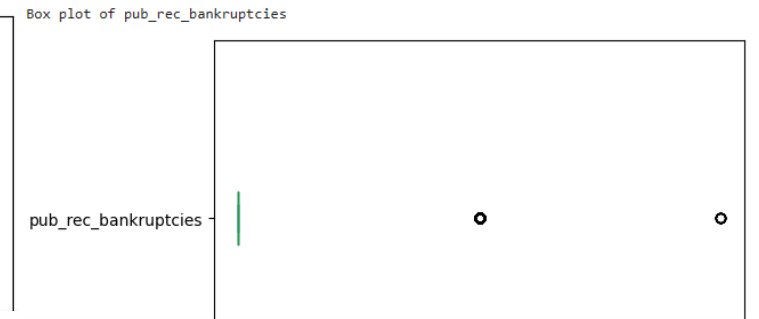
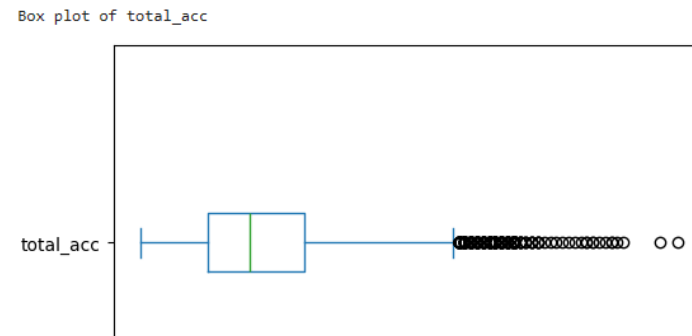
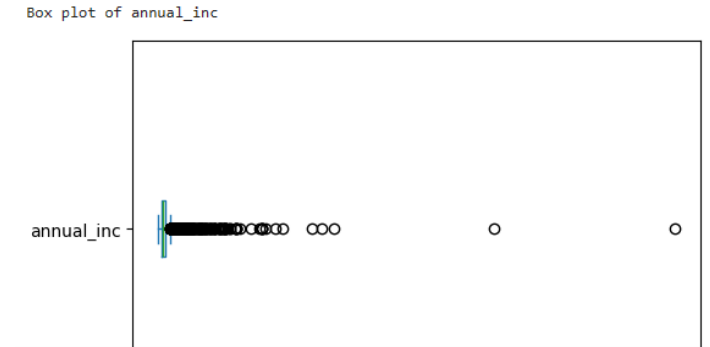
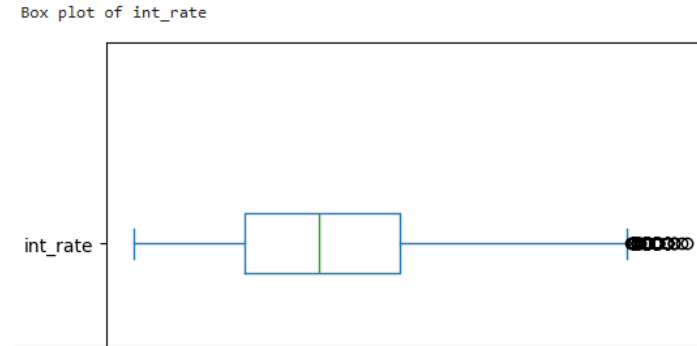
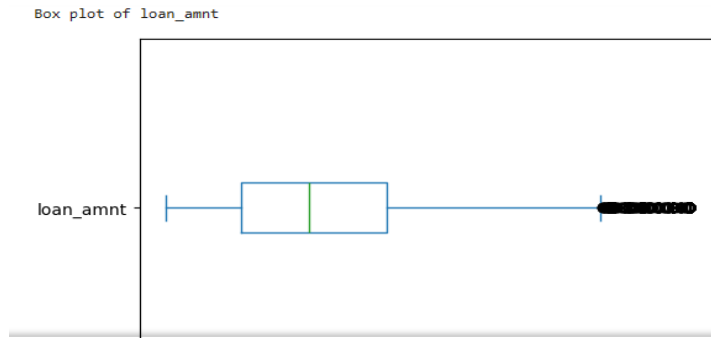
- We need to trim some characters or strings in the values for proper analysis ex, remove months in term, % in *int\_rate* and *revol\_util*.
- ***emp\_length*** has 2.706650% values as empty. Covert the column dtype to float, find the median(4.0 here) and fill the missing values with it. Similarly do it for other colums (***revol\_util***, ***pub\_rec\_bankruptcies***).
- converting Columns to proper data type. *loan\_amnt*, *funded\_amnt*, *int\_rate*, *revol\_util* to float , *term* to int and *issue\_d* to datetime.

# Segmentation

- Segmentation done based on numerical and categorical variables
- numerical\_column=  
`['loan_amnt','funded_amnt','funded_amnt_inv','term','int_rate','installment','emp_length','annual_inc','dti','delinq_2yrs','inq_last_6mths','open_acc','pub_rec','total_acc','pub_rec_bankruptcies']`
- categorical\_column =  
`['grade','home_ownership','verification_status','issue_d','loan_status','purpose','addr_state']`
- *Followingly we do univariate analysis. We draw the box plot to find out the outliers and remove the outliers.*

# 3.Univariate Analysis

- Univariate Analysis on Numerical Columns: Draw Box plot for numerical columns to understand the outliers values .



- Outliers removed using following formula**

$Q1 = df[col].quantile(0.25)$

$Q3 = df[col].quantile(0.75)$

$IQR = Q3 - Q1$

$Upper\_Boundary = Q3 + threshold * IQR$

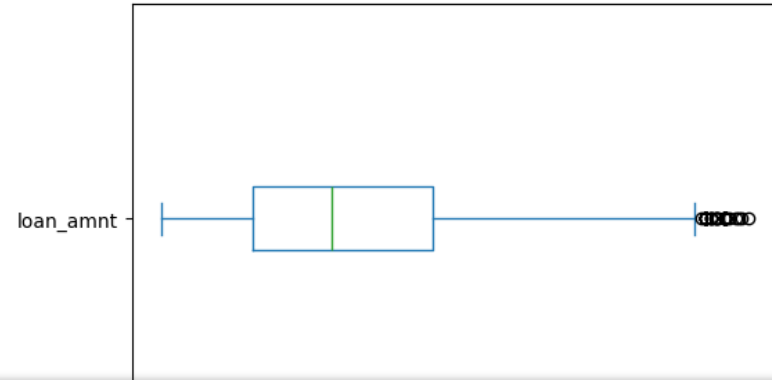
$Lower\_Boundary = Q1 - threshold * IQR$

$df = df [(df[col] \geq Lower\_Boundary) \& (df[col] \leq Upper\_Boundary)]$

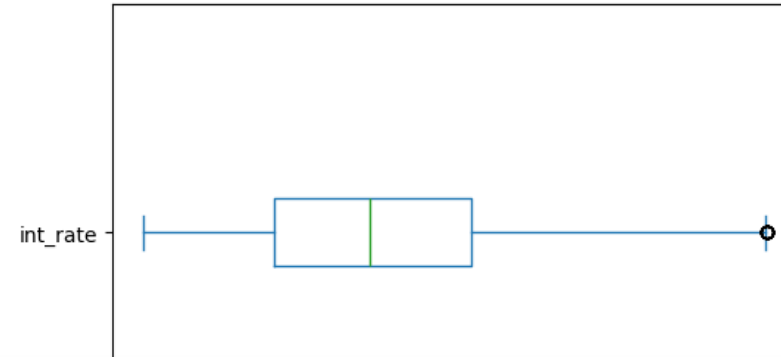


# After removing outliers...

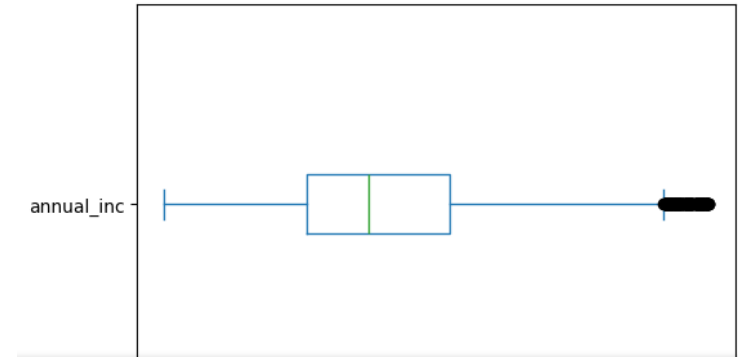
Box plot of loan\_amnt



Box plot of int\_rate

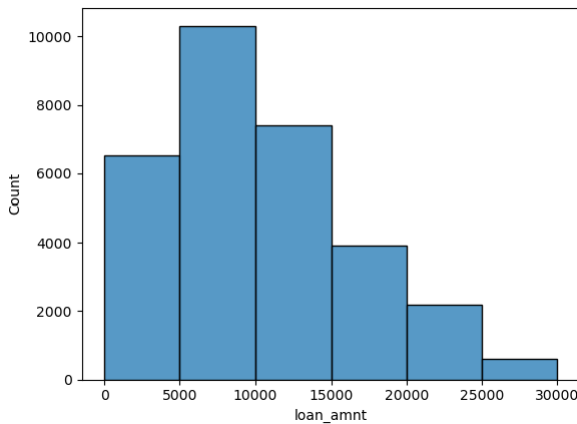


Box plot of annual\_inc



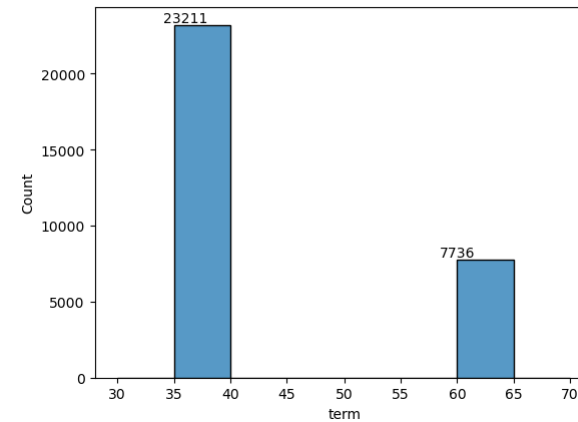
## Distribution of loan amount

Hist plot of loan\_amnt



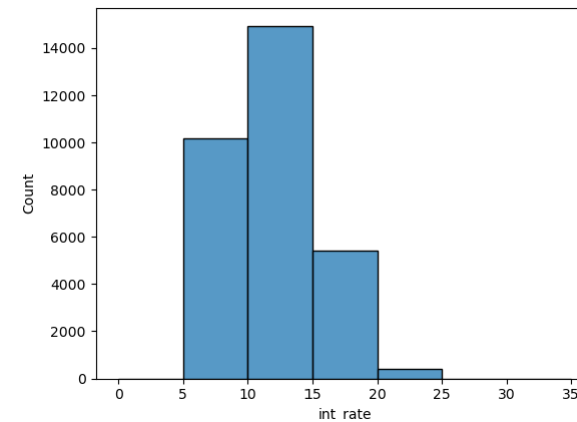
## Loan Term

Hist plot of Term  
value 36 23211  
value 60 7736



## Interest rate

Hist plot of Interest Rate



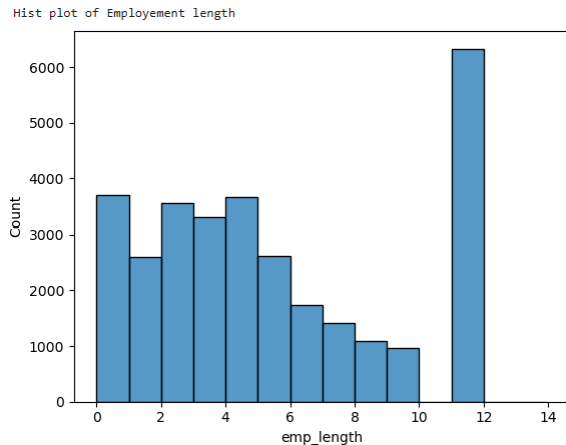
- Max number customer got loan from 5000 to 15000.
- Min is 500 and Max is 29100.

Only two terms of loans

- 36 months and 60 months

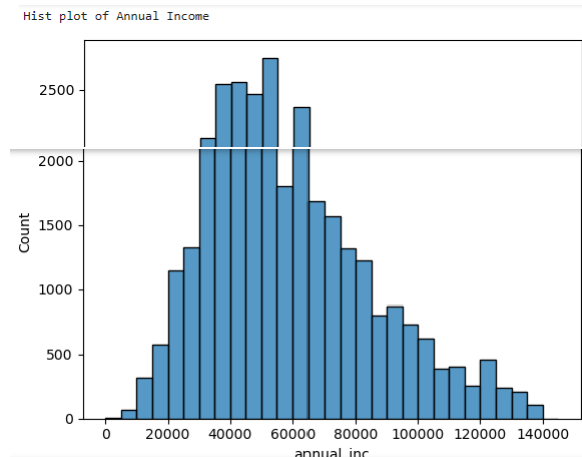
Most of the loan distributed with interest rate ranging from 10 to 15

## Employment length



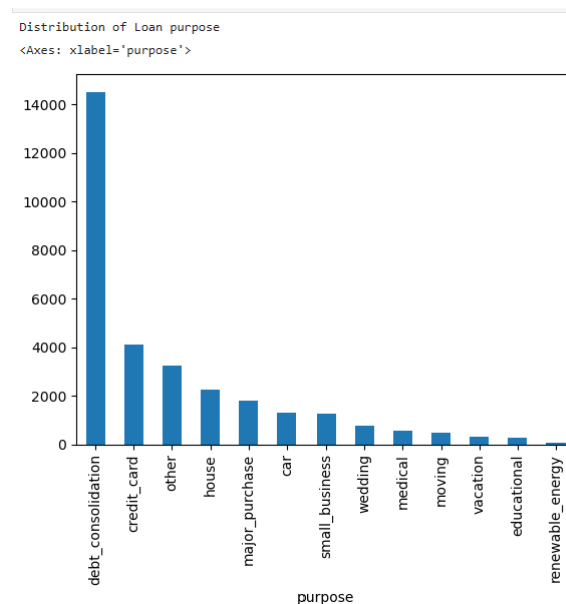
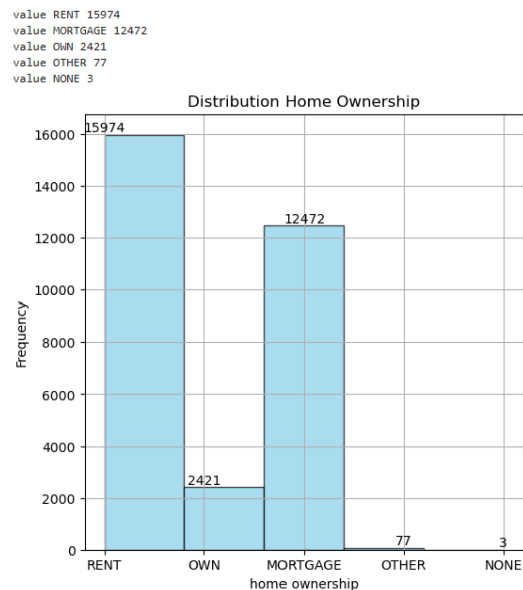
Max, Customer's employment length  
Varies from 1 year to 6 years

## Annual income distribution

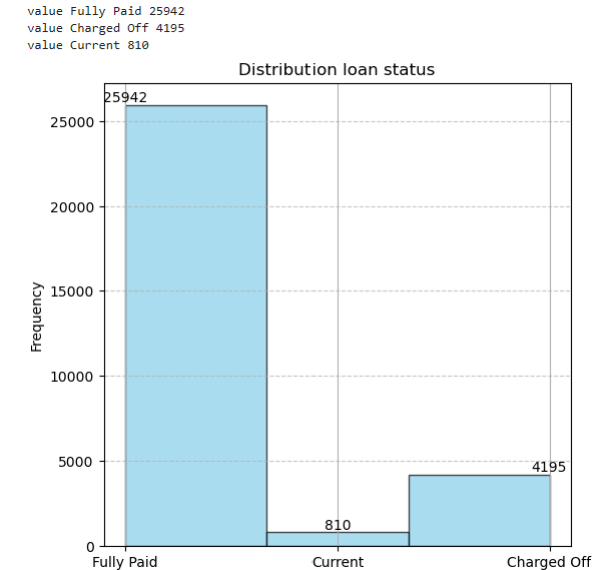


- Min income 4k, Max 137k.
- Most of the them in range of 20k to 85k

## Un ordered Categorical – Univariate Analysis



## Un ordered Categorical – Univariate Analysis



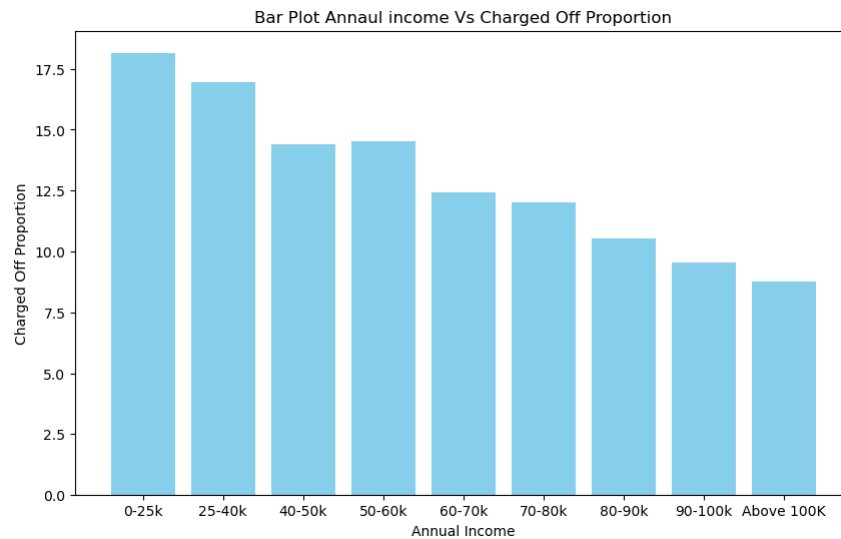
## Observations:

- Most of the customers fully paid their loan.
- Less than 5% of customers had charged off.
- Most of the customers living in Rent and Mortgaged houses
- Maximum loan is taken for the purpose of debt consolidation.

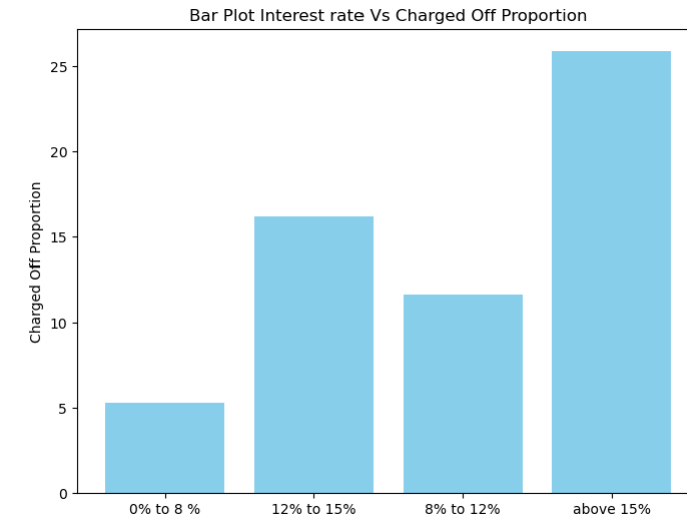
## 4. Bivariate Analysis – Relationship between two variables

- Annual income, loan amount, interest rate are continuous variables. So for better analysis let's bucketize these values and create different data sets to get better insights.

**Annual Income vs Charged Off Proportion**

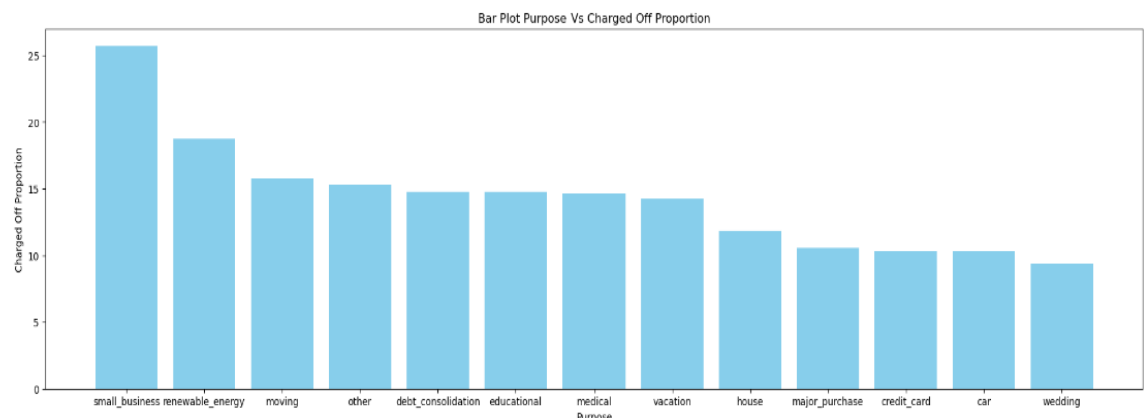


**Interest rate vs Charged Off Proportion**



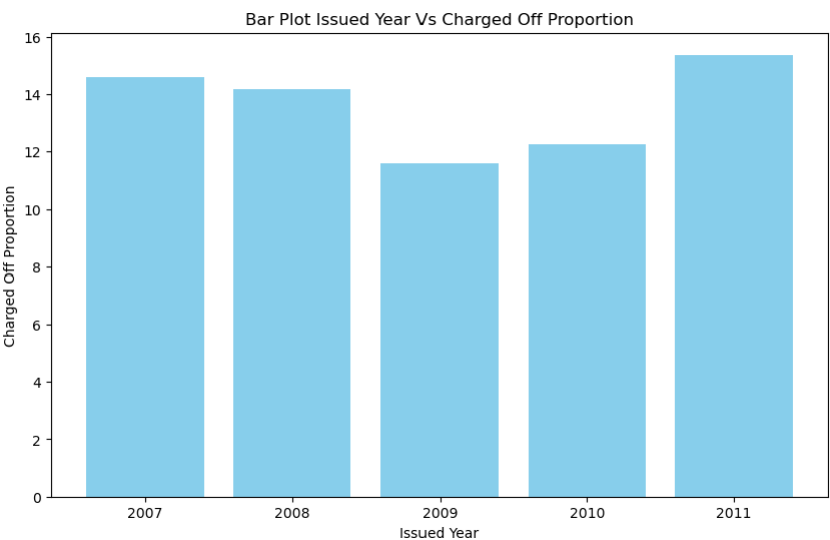
- Customers who get less salary likely to default the loan.
- Higher interest rate loans - high percentage of charged off
- 35% of customers who defaulted having salary less than 40k
- Int rate greater 12% having higher percentage of charged off

## Purpose vs Charged off proportion



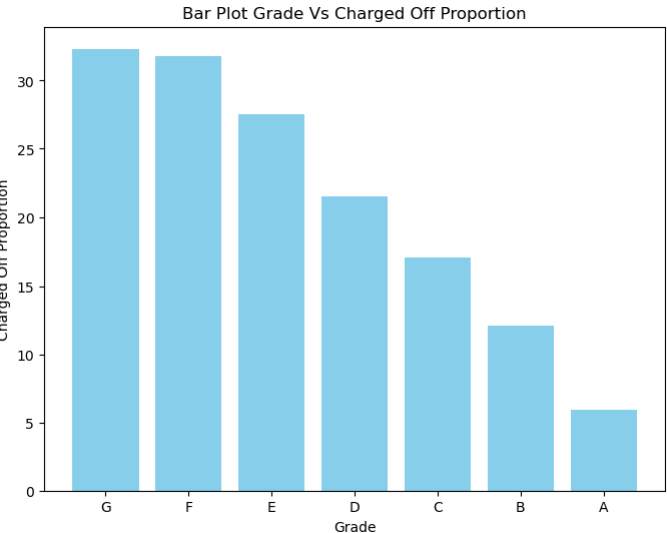
Loan given for Small business and renewable energy having highest percentage of charged off

## Loan issued year Vs Charged off

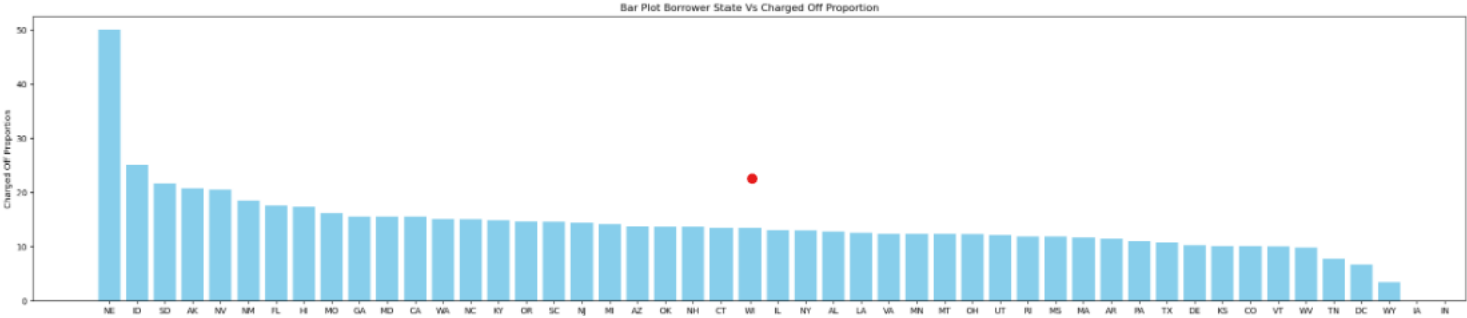


Loan issued in 2011 having higher charged off rate.

## Customers grade vs Charged off proportion

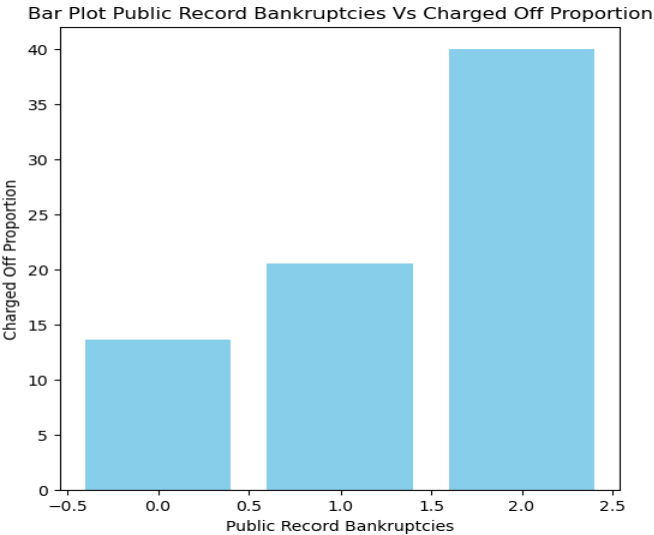


## Customers State vs Charged off proportion



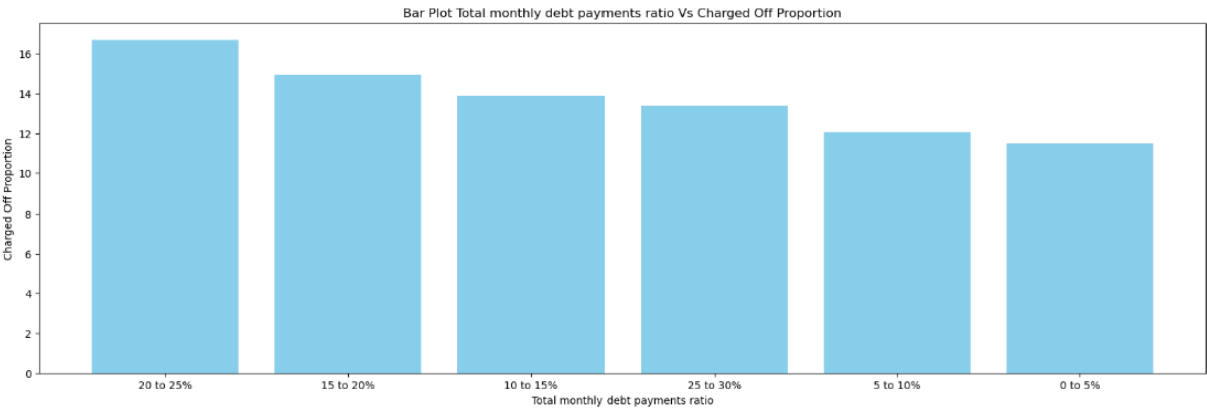
- G and F grade employees having higher charged off proportion
- Customer who belong to NE state having higher charged off proportion

# Public Record Bankruptcies VS Charged off proportion



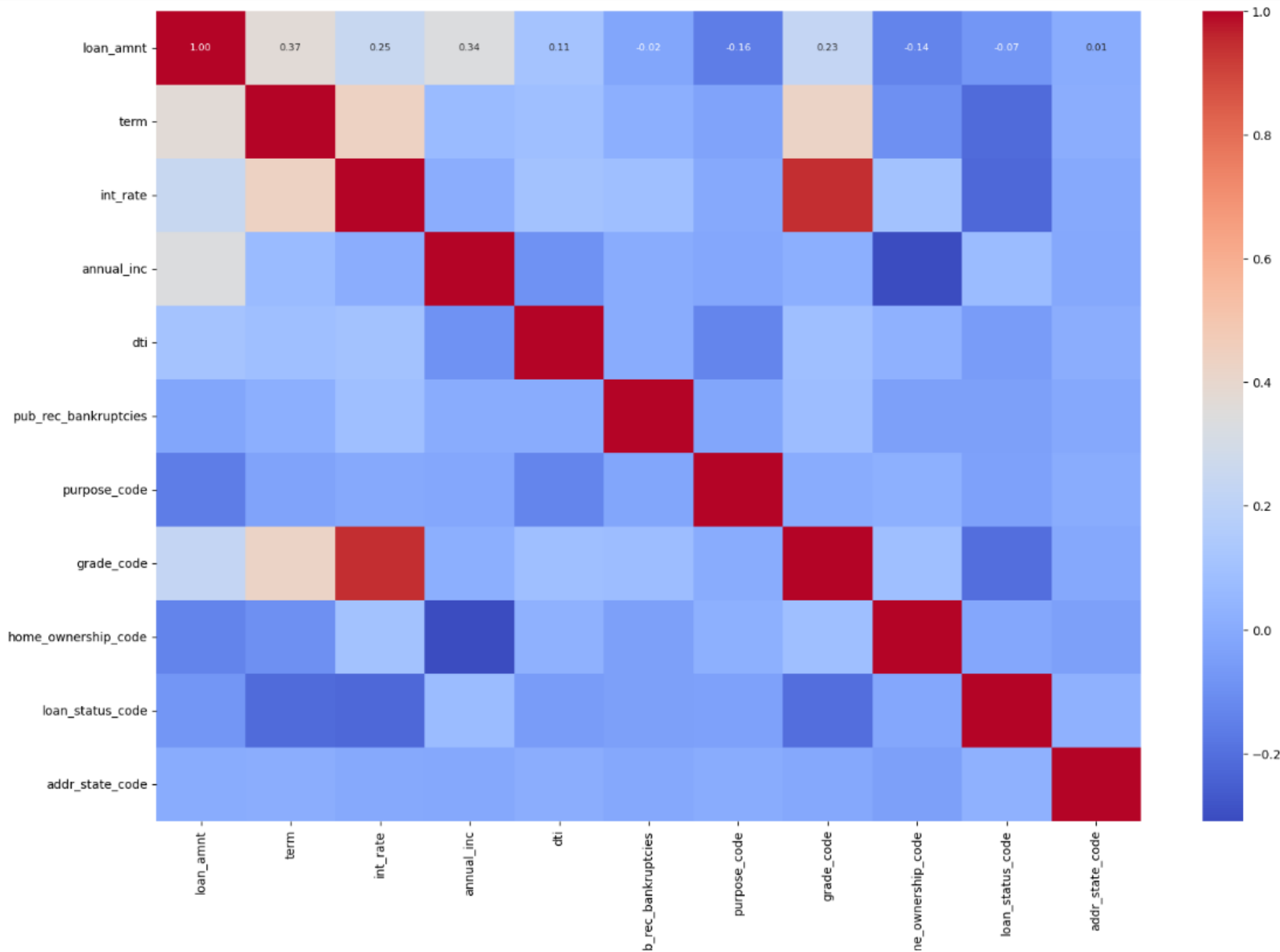
Customers having more public record bankruptcies  
having higher charged off proportion.

# Monthly debt payments ratio VS Charged off proportion



Customers having monthly debt payment ratio  
Greater than 10% having higher charged off proportion

# Correlation:



## ***Positive Correlation***

- Loan amount and term having positive correlation
- Annual income and loan amount having positive correlation
- Interest rate and term having positive correlation

## **Negative Correlation:**

- Loan amount and loan purpose having negative correlation
- Annual income and home ownership having strong negative correlation.

# Summary:

- Lower the income group having higher percentage of charged off.
- Interest rate above 12% is having higher percentage of Charged off
- Small business and renewable energy loans having higher percentage of charged off.
- The loan distributed in 2011 and Dec month having higher percentage of charged off.
- G and F grade customers having higher percentage of charged off.
- Customers who are having higher number of public record bankruptcies having higher percentage of charged off.
- Customers having above 15% total montly monthly debts payment having higher percentage of charged off.