

## **Case Study: Lending Club Case Study**

**Submitted By: Utkarsh Chaturvedi**

## **Problem Statement:**

### **Business Understanding - Lending Club Case Study**

We are working for a consumer finance company which specializes in lending various types of loans to urban customers. When we receive a loan application, we have to make a decision for loan approval based on the applicant's profile. Two types of risks are associated with our 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

### **Business Objective**

The objective is 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.

We need to research about risk analytics by understanding the types of variables and their significance should be enough.

### **Analysis Approach:**

The following will be the approach that will be followed to analyze the data and then provide a solution:

1. Understand the data properly
2. Clean and manipulate data as needed
3. Data analysis
4. Provide result of analysis in form of visualization

### **Analysis Results:**

The first exercise is to analyze all the data that has been provided and only consider the data that is relevant. The objective is to find the applicants which have strong probability of defaulting and also to identify the applicants which can repay their loan.

The objective is that we want to know which loan applications are risky. So, the fields that are created after a loan application is approved will not be part of our analysis. Also, I have considered some customer related fields that are relevant and will provide us with intelligence related to demographics.

Based upon that I have identified that the following data is relevant for our analysis:

- loan\_amnt
- funded\_amnt
- funded\_amnt\_inv
- term
- int\_rate
- installment
- grade
- sub\_grade

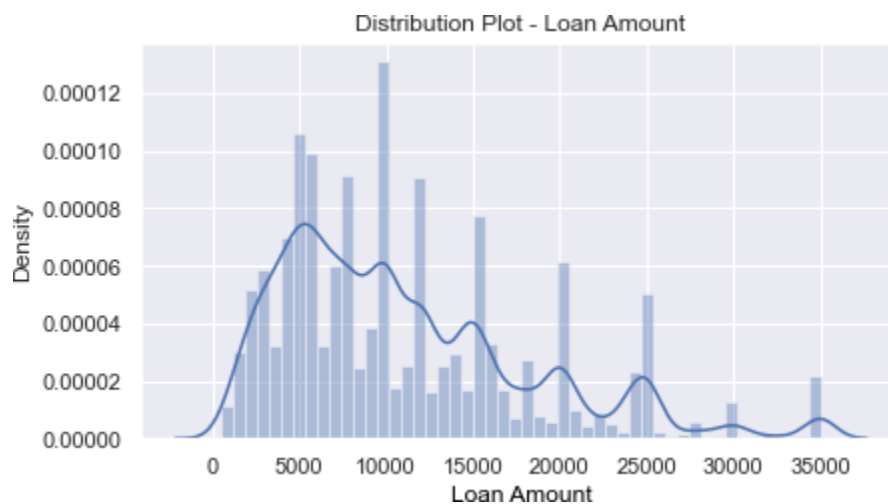
- emp\_title
- emp\_length
- home\_ownership
- annual\_inc
- verification\_status
- issue\_d
- loan\_status
- url
- desc
- purpose
- title
- zip\_code
- addr\_state
- dti
- addr\_state
- purpose
- loan\_amnt
- funded\_amnt
- funded\_amnt\_inv
- dti
- term

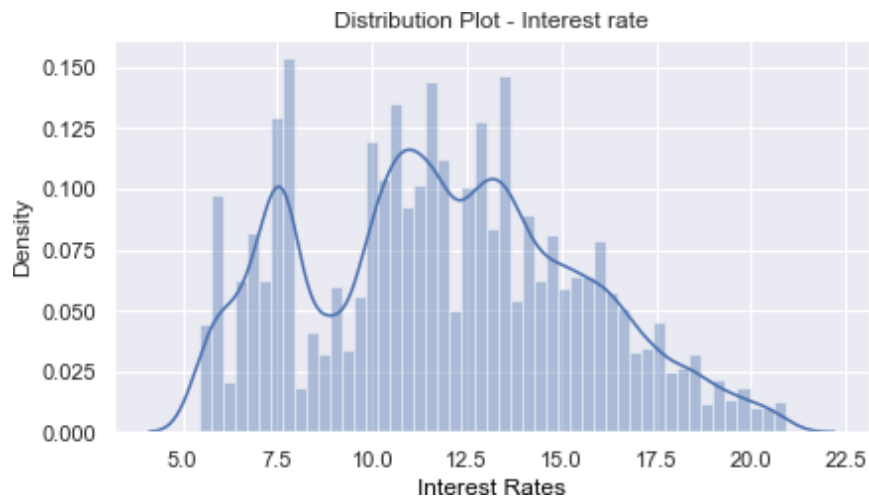
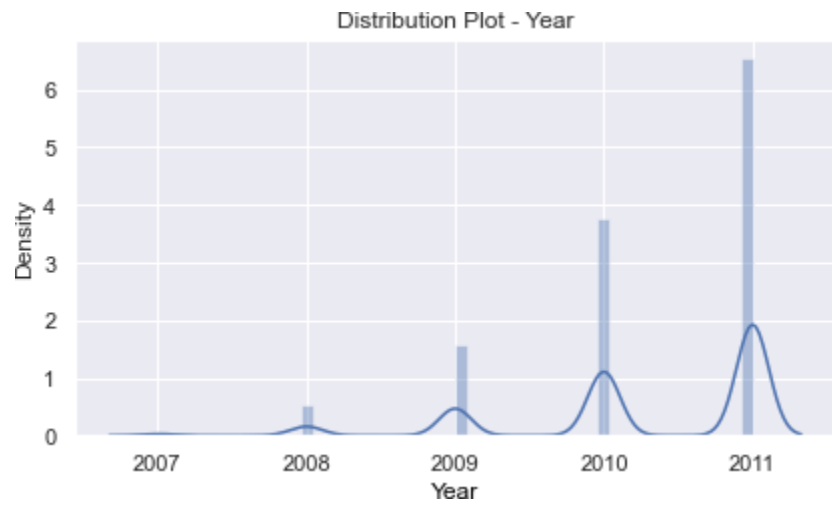
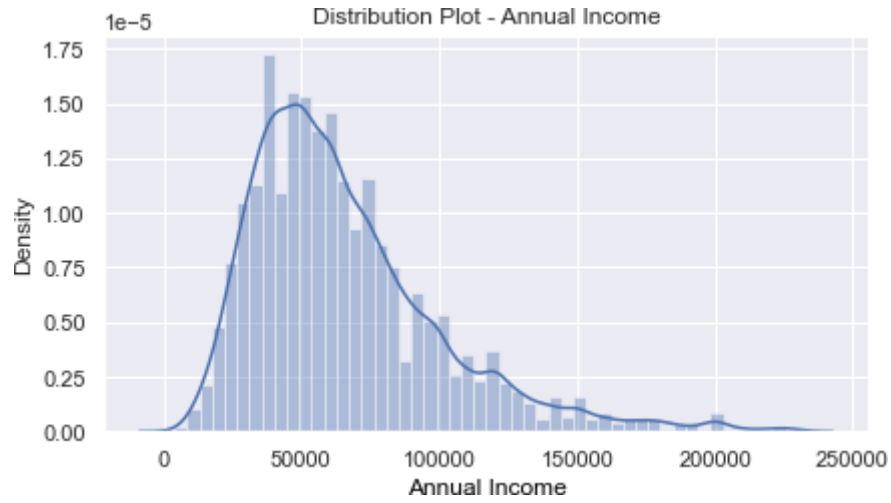
#### Data was cleaned:

1. Removed blank entries
2. Did some manipulation on emp\_length, int\_rate, term and zip code fields
3. Created bins to be further used during bivariate analysis.

#### Univariate Analysis:

1. Outlier analysis was done using box plot on int\_rate and annual\_inc
2. Univariate Analysis on Loan Amount, Annual Income, Loan Issue Year (Derived metrics), Interest Rate, Loan Status, year(derived)





Based upon the above distribution plots the following are the conclusions based upon **Univariate analysis**:

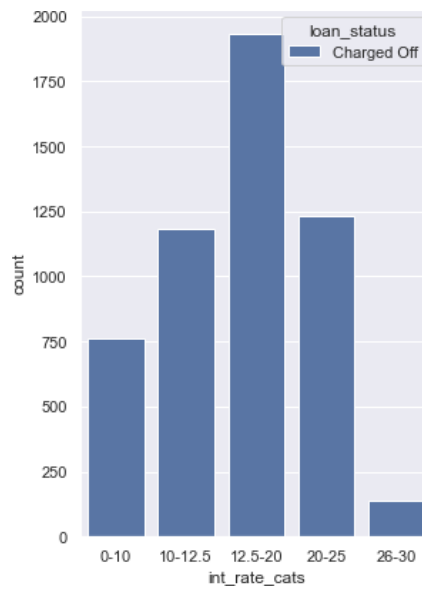
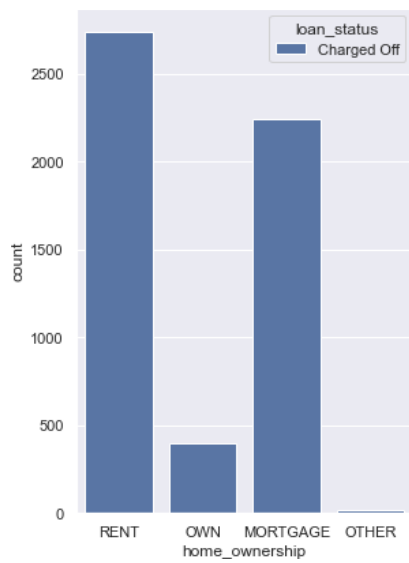
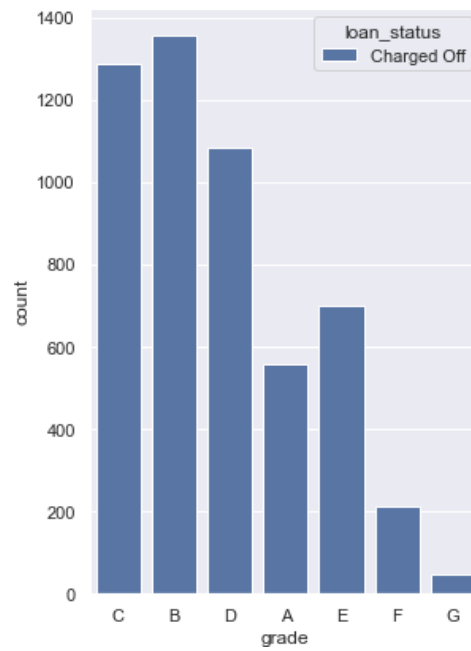
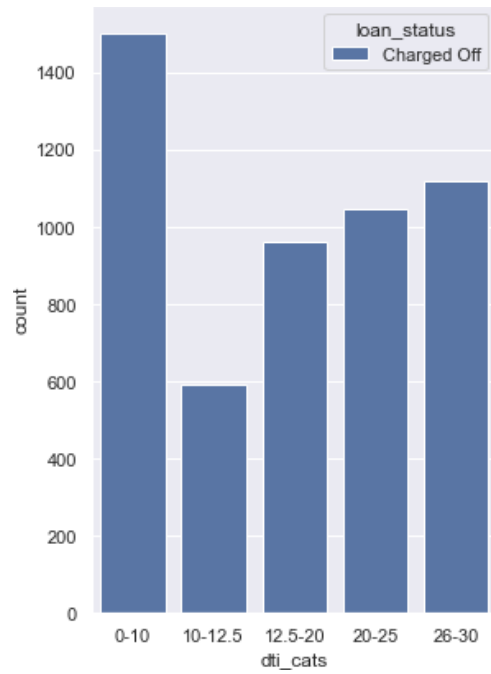
1. Most loans that are being taken are between \$5000 - \$10000
2. Annual income of customers that are taking loan is around \$50000 range
3. Numbers of loans distributed are increasing by each year. Year 2011 had the greatest number of loans disbursed.
4. Interest rates are mostly between 10% - 15%

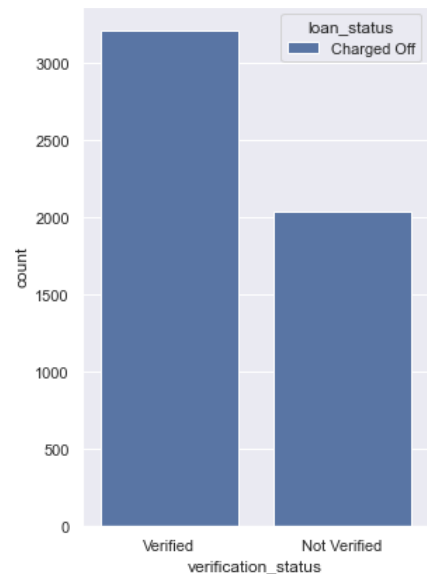
#### **Bivariate Analysis:**

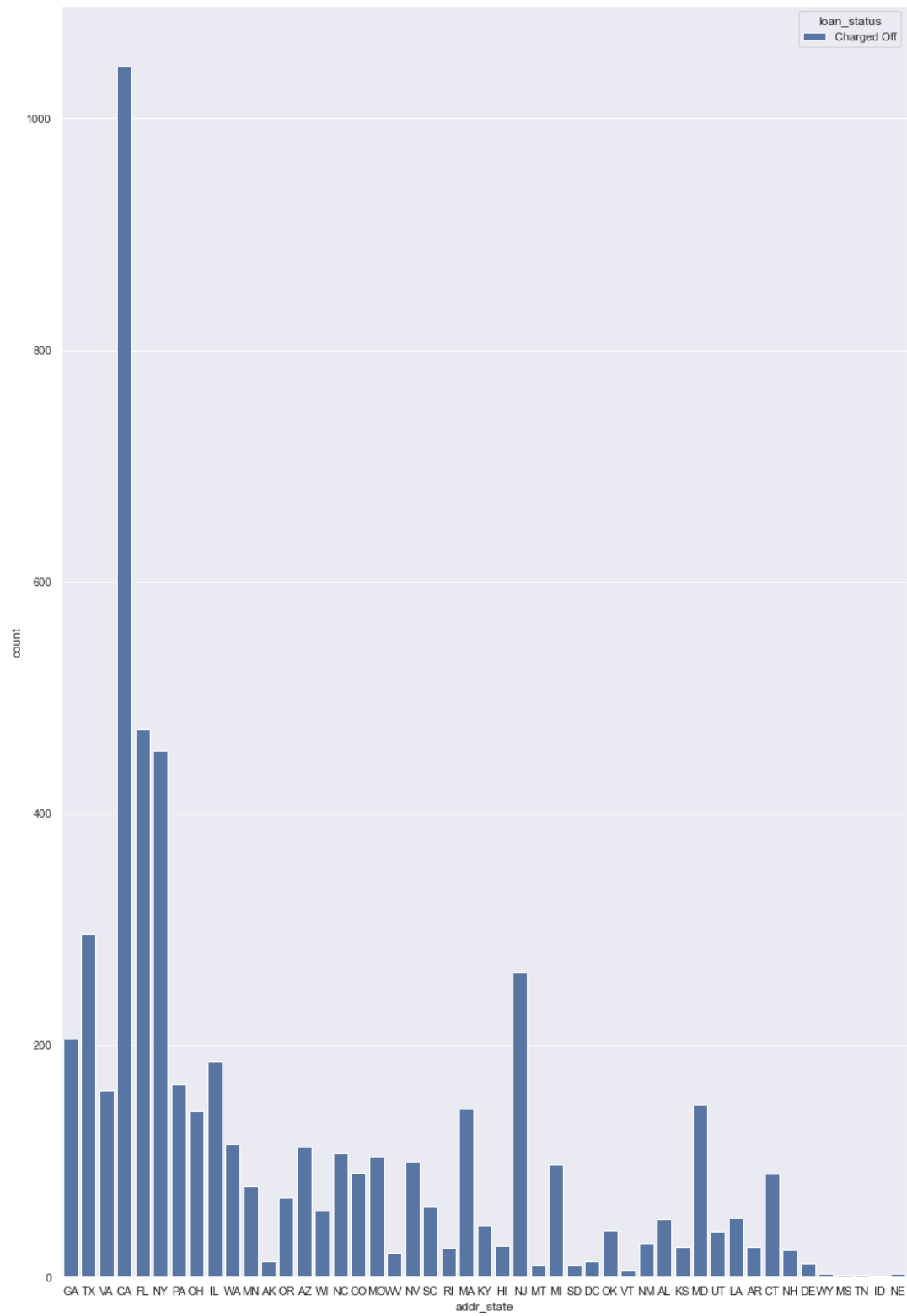
For Bivariate Analysis we will focus on the loan status of "Charged Off" so that we can analyze which variables have more impact. Based upon this analysis bank can make intelligent decision during loan approvals.

Bivariate analysis was done on the following variables:

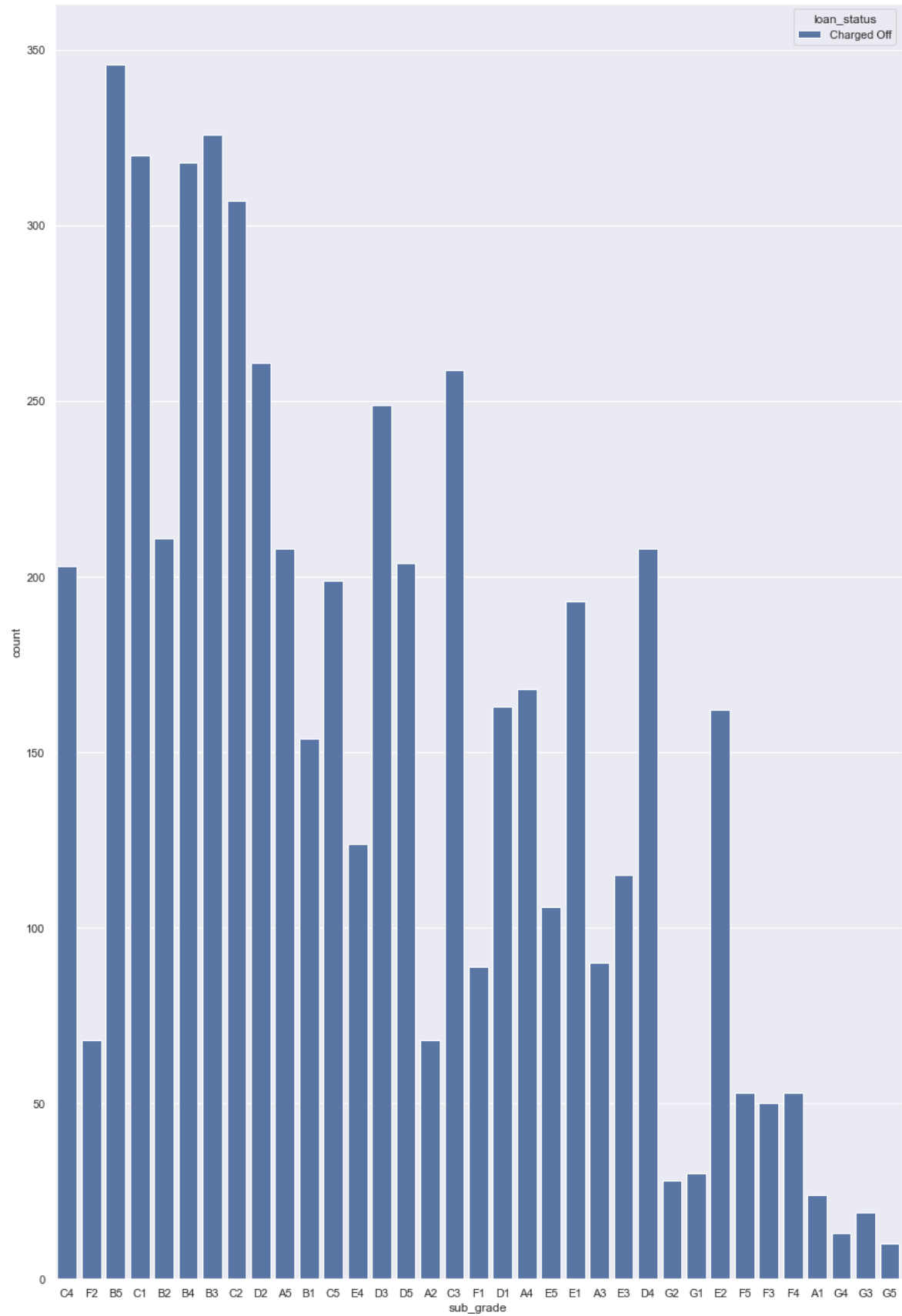
1. Loan Status (Charged off) with Grade
2. Loan Status (Charged off) with Verification Status
3. Loan Status (Charged off) with Sub Grade
4. Loan Status (Charged off) with Interest rate categories
5. Loan Status (Charged off) with dti categories
6. Loan Status (Charged off) with home ownership categories
7. Loan Status (Charged off) with state categories









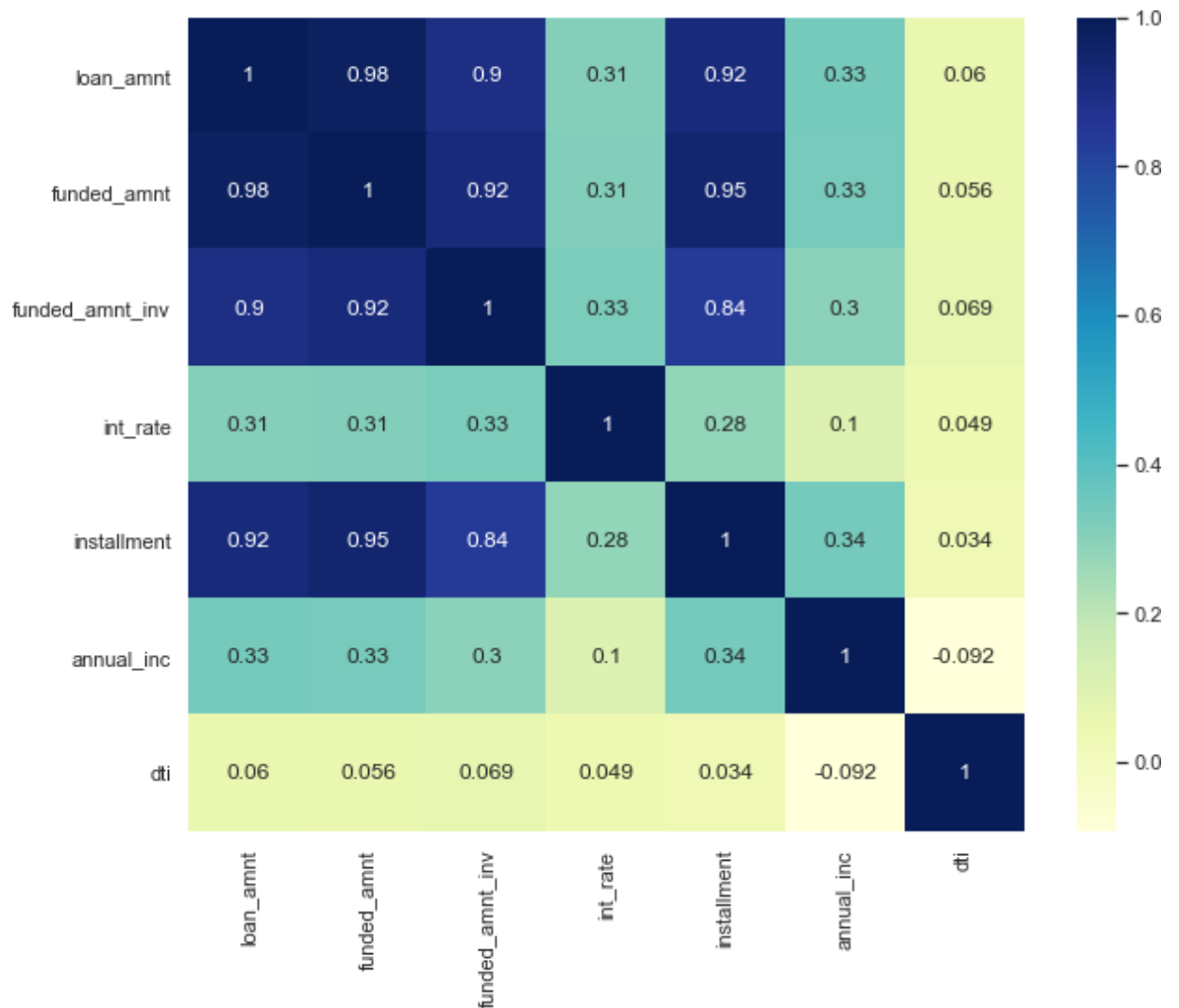


Based upon the above count plots the following are the conclusions based upon **Bivariant analysis**:

1. Customers that were verified have more chances of loan being charged off
2. Customers with sub grade B3, B4 and B5 have very more chances of loan being charged off
3. Customers with Grade G have less chances of loan being charged off
4. Customers with interest rates between 12.5% - 16% have high chances of loan being charged off
5. Customers with dti between less than 10% have high chances of loan being charged off
6. Customers who have rented accommodation have more chances of loan being charged off
7. Customers that are living in California have high chances of loan being charged off

### Multivariant Analysis:

Multivariant analysis has been done using heatmap.



Based upon the above heat map the following are the conclusions based upon **Multivariant analysis**:

1. Loan amounts have high correlation with Funded amount, Funded amount committed by investor and installment
2. Funded amounts have high correlation with Loan amount, Funded amount committed by investor and installment
3. Funded amount committed by investor have high correlation with Loan amount, Funded amount and installment