# Lending Club Case Study

## 1. Reading & Understanding the data

- **Importing the input files**
- **Inspect data**

## 2. Data Cleaning and Manipulation

- **Remove columns where NA values are more than or equal to 30%**

- **Remove irrelevant columns.**

  Now let's look at each column from business perspective if that is required or not for our analysis such as Unique ID's, URL. As last 2 digits of zip code is masked 'xx', we can remove that as well.

- **Remove irrelevant records**

  Purpose of loan: Drop records where values are less than 0.75% We will analyse only those categories which contain more than 0.75% of records. Also, we are not aware what comes under 'Other' we will remove this category as well.

## 3. Derived Metrics

We will now derive some new columns based on our business understanding that will be helpful in our analysis.

1. **Loan amount to Annual Income ratio**
2. **Extract Year & Month from Issue date**
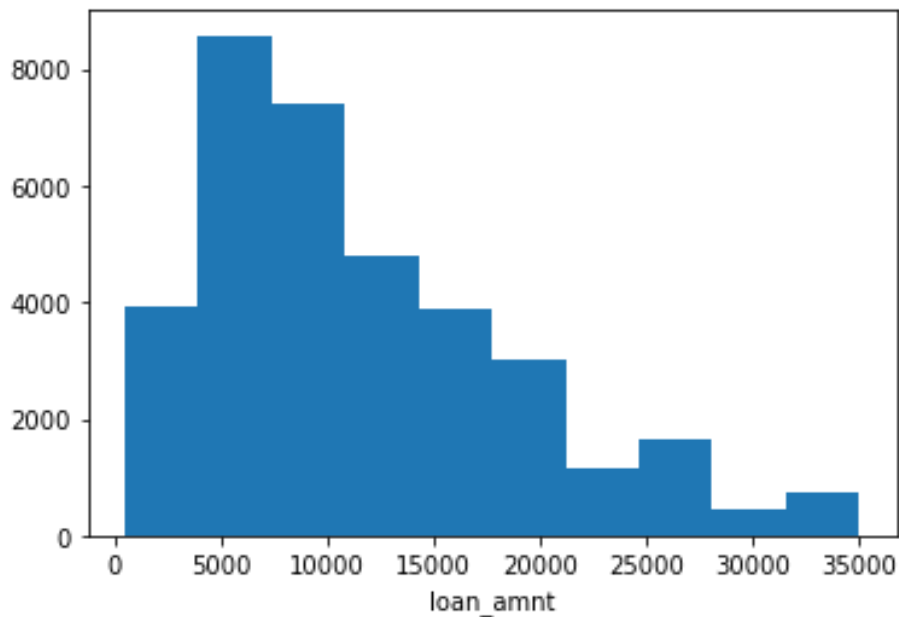3. **Remove '%' from int_rate column**

## 4. Univariate Analysis

- Continuous Variables In case of continuous variables, we need to understand the central tendency and spread of the variable.These are measured using various statistical metrics visualization methods such as Boxplot,Histogram/Distribution Plot etc.

- Categorical Variables For categorical variables, we'll use frequency table to understand distribution of each category. It can be be measured using two metrics, Count and Count% against each category. Countplot or Bar chart can be used as visualization.
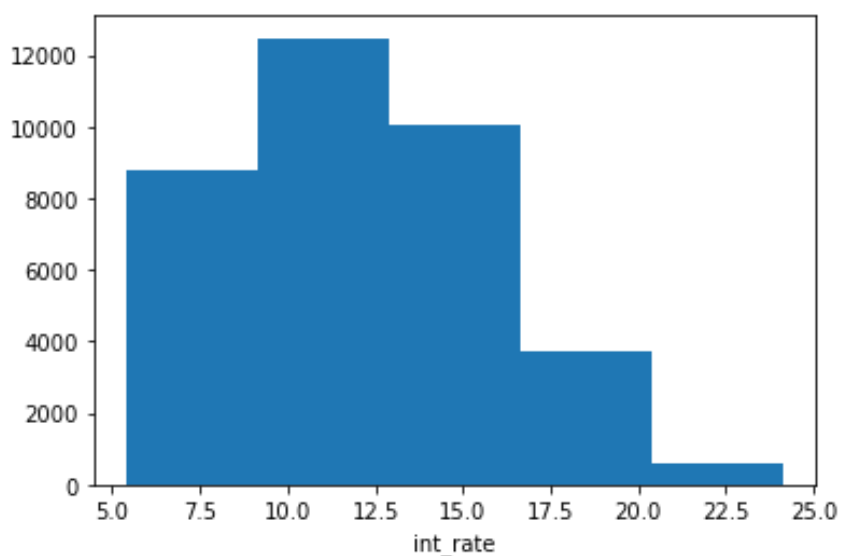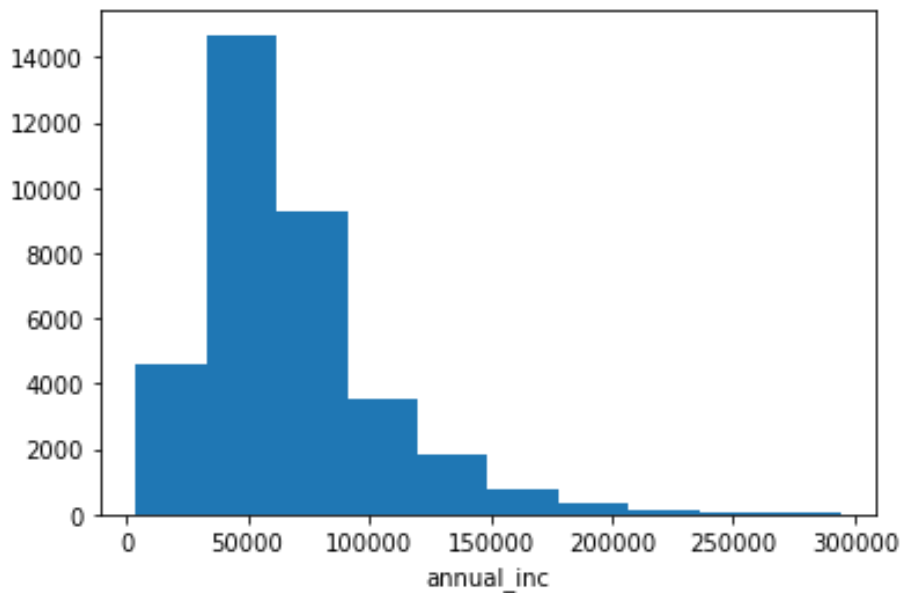
# Continuous Variables

## 1. Loan Amount



**Insights: Most of the loan amounts are distributed between 5000 to 10000 USD.**

## 2. Interest Rate

### 3. Annual Income

## Categorical Variables

### 4. Loan Status

## 5. Purpose of loan

## 6. Home Ownership wise Loan

## 7. Year wise Loan
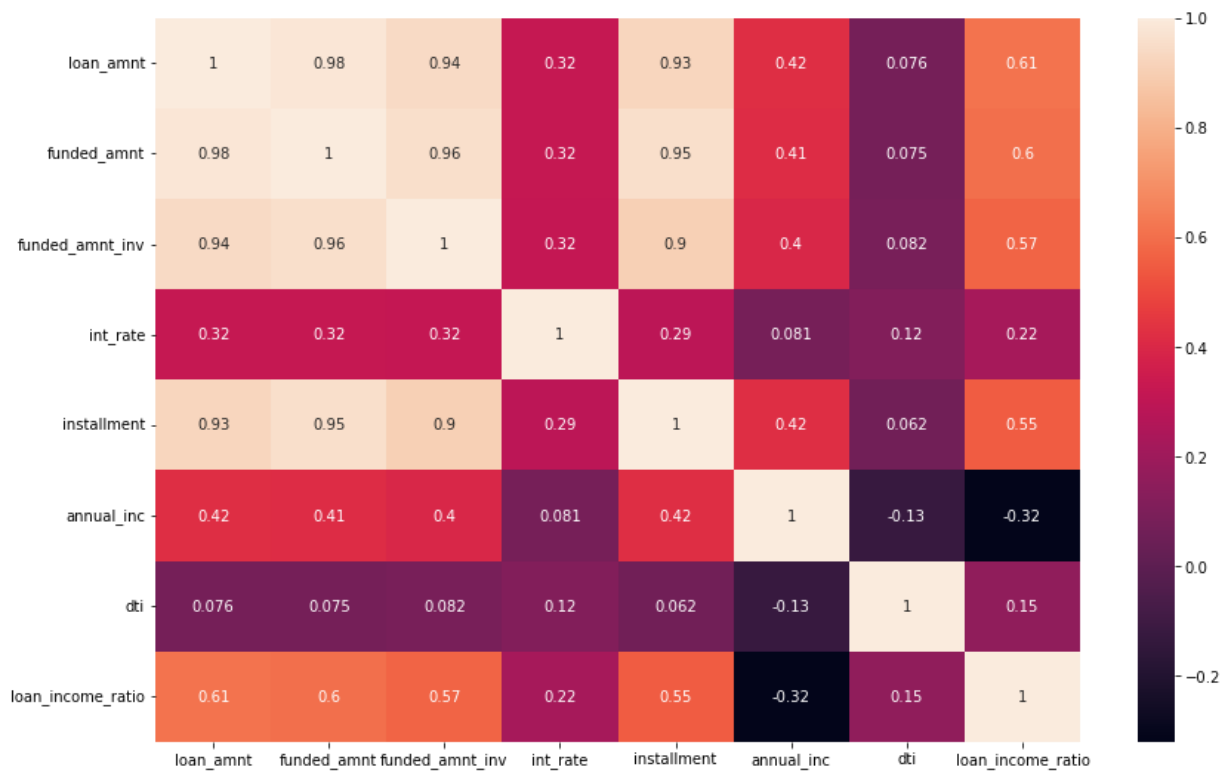
### 8. Loan Term



**Insights: 70% of applicants applied loan for 36 months term period**
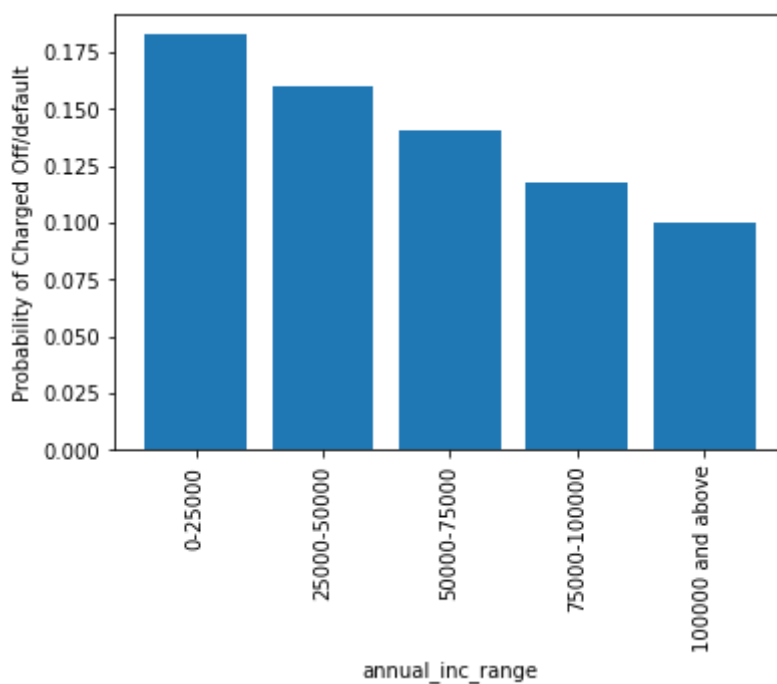
## 5. Bivariate/Multivariate Analysis

Bivariate/Multivariate Analysis finds out the relationship between two/two or more variables. We can perform Bivariate/Multivariate analysis for any combination of categorical and continuous variables. The combination can be: Categorical & Categorical, Categorical & Continuous and Continuous & Continuous.

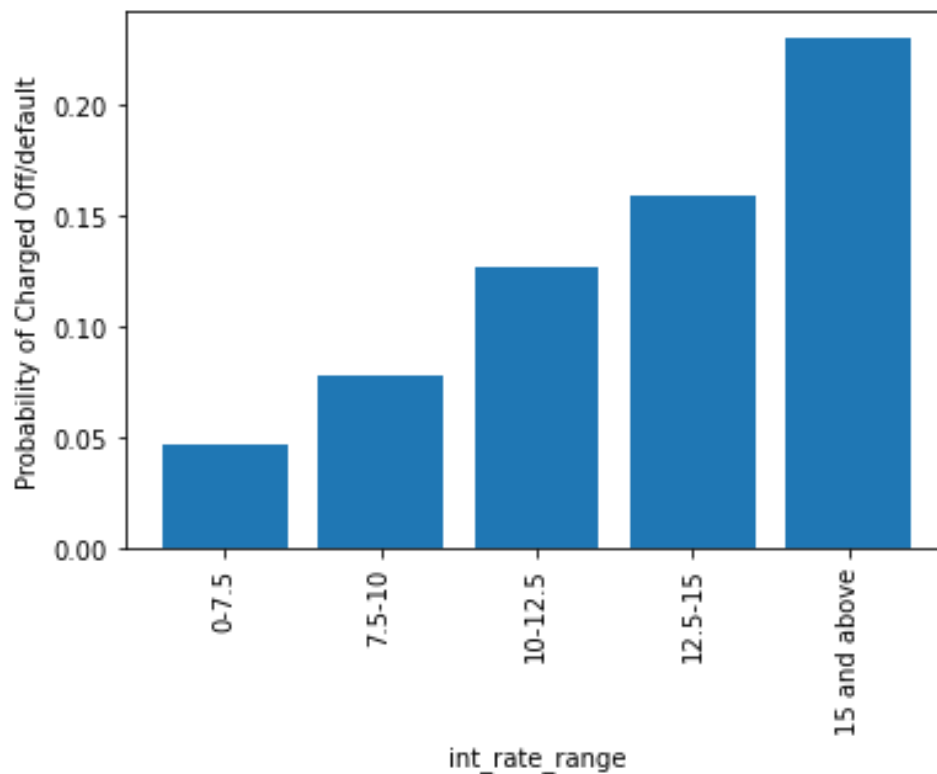### 1. Correlation Matrix: All Continuous (Numeric) Variables

**Insights: It is clear from the Heatmap that how 'loan_amnt','funded_amnt' & 'funded_amnt_inv' are closely interrelated.So we can take any one column out of them for our analysis**

### 3. Annual Income Range vs Probability Charge Off

**Insights: As the annual income is decreasing the probability that person will default is increasing with highest at (0 to 25000) salary bracket.**

4. Interest rate Vs Probability Charge Off



**Insights: As the interest rate is increasing the probability that person will default is increasing with highest of 9% at 15% & above bracket**