

LENDING CLUB ASSIGNMENT SUBMISSION

Name:

VINYAS SHUKLA

SRIKANTH PADMANABHUNI

Abstract :

- A consumer finance company provides loans to different kinds of people in urban. But company feels difficult in finding whether a person pays/defaults the loan
- If person doesn't pay the loan, then the investor faces the loss. If a person who pays complete loan doesn't get approval for loan then also investor and bank gets loss
- We as a Data scientist need to find the factors that are responsible for finding the person can pay loan or default based on the data provide by the company

Approach for Analysis:

- Load the data from the csv to start the analysis.
- **Cleaning Data:**
 - Find the percentage of null data available in the dataset. We can see that a large of data is null.
 - There are some columns which are having 100% null data in them. So, keeping them in data set is not much useful. So, lets drop those columns
 - There are some more columns which are having >30% of null data. So, lets drop those columns instead of imputing and making data biased towards the median
 - After removing those null data still, we can find data with <10% null for following columns '**emp_title, emp_length, title, revol_util, last_payment_d, laste_credit_pull_d, collections_12_mths_ex_med, chargeoff_within_12_mths, pub_rec_bankruptcies, tax_lines**'.
 - **pub_rec_bankruptcies** is a kind of categorical data and most of the columns are with data '0'. So, instead of imputing data and create bias over the 0 lets drop rows with null data
 - Most of the other columns are kind of categorical/non numeric data. So, instead of imputing them lets drop the rows and make data clean.
 - After cleaning all the data we can find total columns available are 53.
 - Among them, we can further drop the columns based on domain knowledge and the data dictionary provided. We can understand which are not required for our analysis.
 - After dropping columns which are not required for our analysis, we have the total remaining columns are 23.
 - Since we need to work on the data related to Charged Off and Fully Paid loan data, so lets remove the rows which are having loan status '**Current**' [Since, theya re currently paying and not marked as Default].

- There are columns like emp_length which are having data like <1 and >10+, so, let's clean them by removing <,>,+ symbols and let's make it as numerical categorical data
- There is one more column called term, which is having data like 36months, 60months. Let's remove months from it and make it as numerical categorical data

Deriving New Variables:

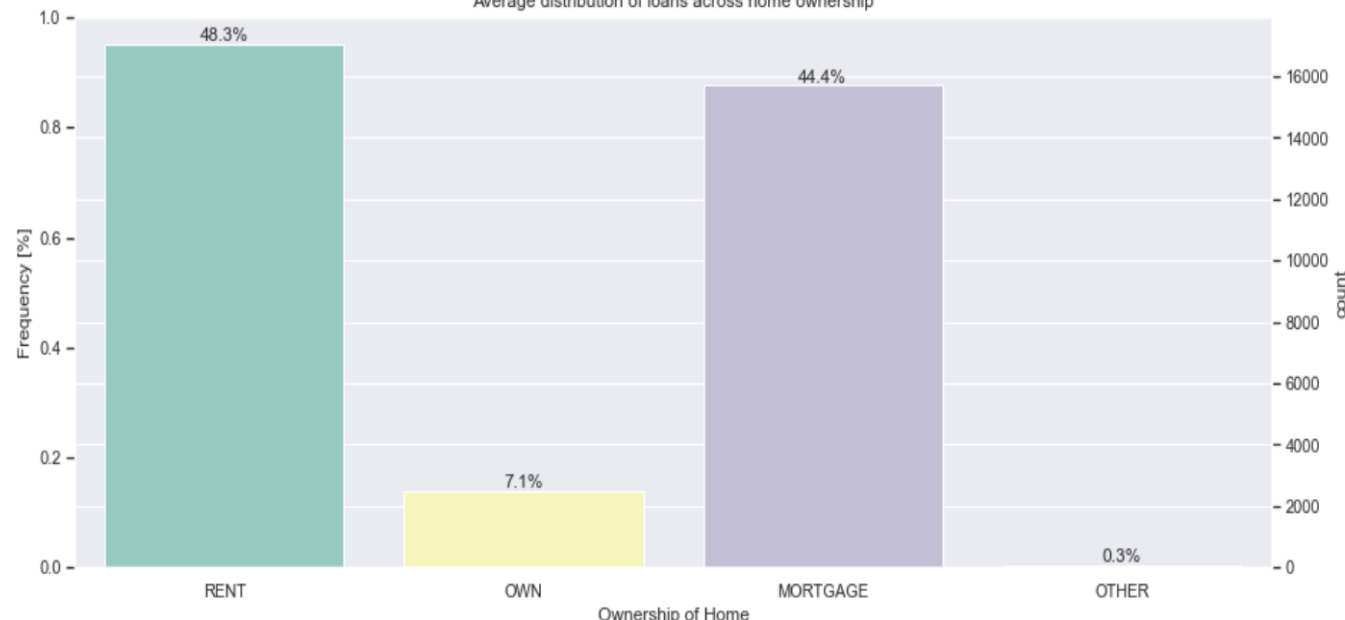
- Now, from the existing we can derive new column, which can help further in our analysis part.
- Deriving '**year**' from the 'issue_d' column, which helps in knowing in which year that loan is approved instead of a date.
- Deriving '**l_t_ai**' from the data of 'annual_inc' and 'loan_amnt'. This column says the ratio of the loan amount taken and the annual income, which helps in deriving further more analysis.
- Deriving column called '**income_bin**' from income. It is providing the category like 1,2,.. Based on the ranges of annual income. This helps in making the analysis on the persons based on the range of annual income instead of one particular annual income
- Deriving the columns '**dti_range**', '**l_t_ai_range**', '**exp_level**' from **dti**, **l_t_ai**, **emp_length** columns respectively. It provides the categories to those columns data and helps in making analysis based on that range instead of a particular value of the data.
- Now, based on the above derived columns and the existing columns we can start working on the Univariate and Bivariate analysis and Plot graphs from them to get more insights.



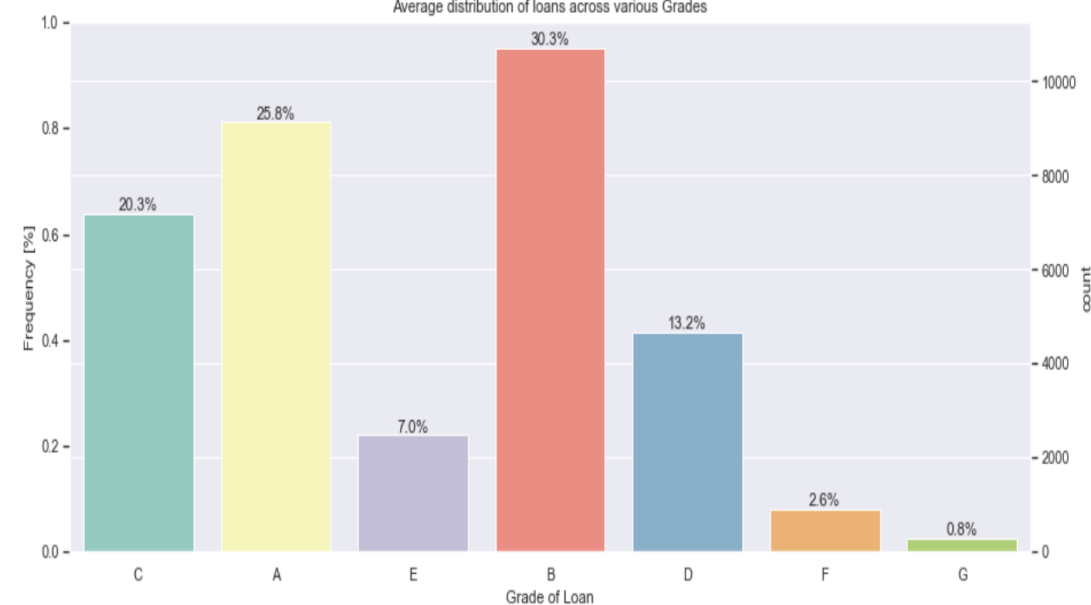
Univariate Analysis:

- Plot the graph for different columns to find its frequency among the data. Like, how such kind of data is spread among the dataset.
- **Loan Status:** Among the dataset 14% of people are having with default loan status, where as 86% of people paid loan Fully
- **Issue Year:** From the analysis we can find that people started taking loans from the year 2008, but it gradually increases while the year increase. From 2010 to 2011 people who are taking loans got increased enormously to 54%.
- **Purpose of Loan:** Most of the people are taking loan with the purpose of 'debt_consolidation' (Means paying outside debts with loan) and 'credit_card' (loans on credit card).
- **Grades and Sub grades:** Finance companies provided some grades and subgrades to people based on some conditions. Among them people with grades B,A,C and with subgrades A4, B3, A2 are the people taking/getting approvals for loans more frequently
- **Years of Experience/Experience levels:** People having the 10(≥ 10) and 1(≤ 1) years of experience are taking/getting approvals for loans more frequently. When it comes to derived column of experience level Juniors are getting more approvals for loans frequently
- **Home ownership:** People who are rented/mortgage home are having more approvals for loans than own house people.
- **State:** People from states 'CA', 'NY', 'FL' are having more approvals for loans than for other state people
- **Debt to loan and loan to income ratio:** People with medium DTI range and Medium, High loan to income ratio people are having more approvals for loans than other people.
- From the above analysis we can see that
 - people having less years of employment experience
 - People having rented/mortgaged house
 - People with the grades B,A,C and subgrades A4, B3, A2.
 - People from states CA, NY, FL.
 - People with medium DTI range and Medium, High loan to income ratio are having more approvals for loans.
- Now, let's see the above data graphically.

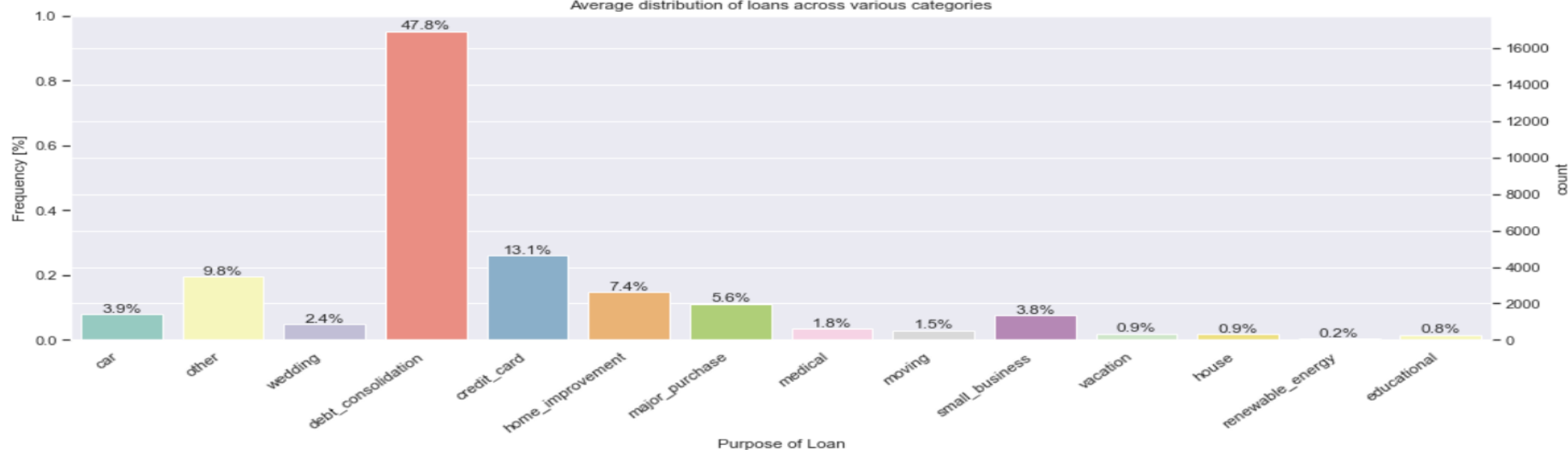
Average distribution of loans across home ownership

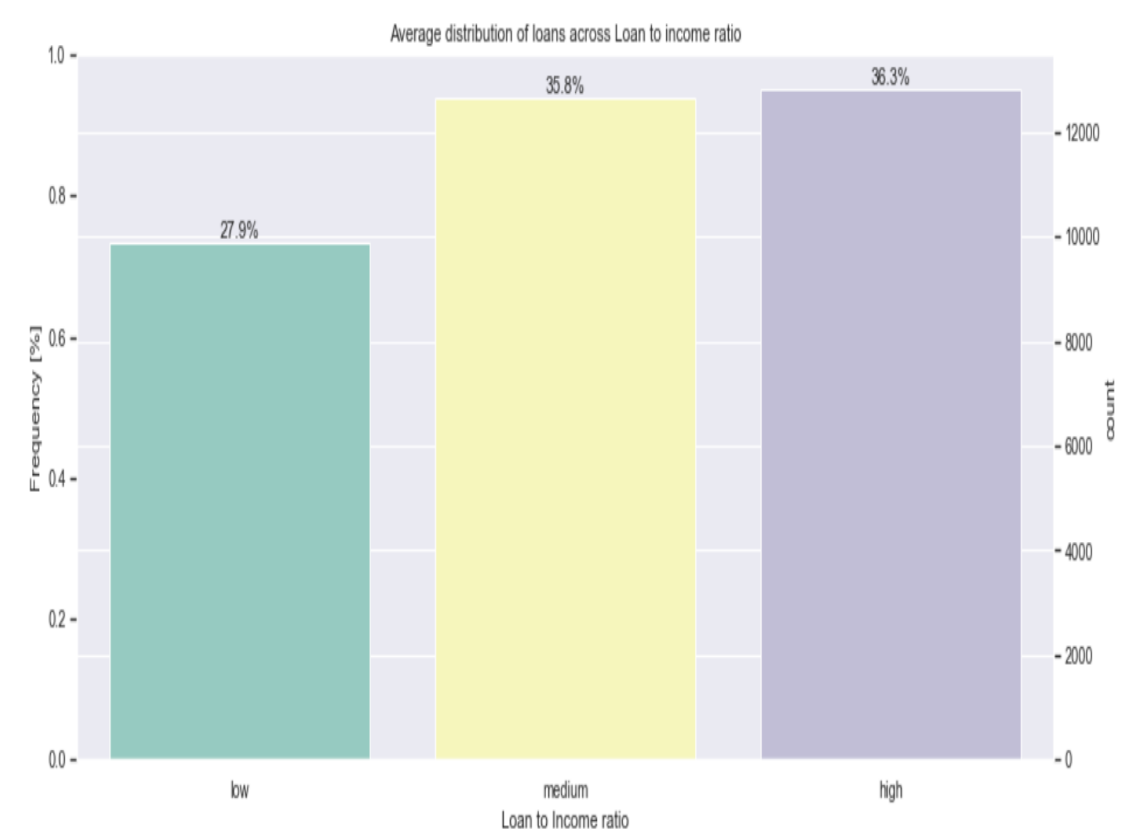
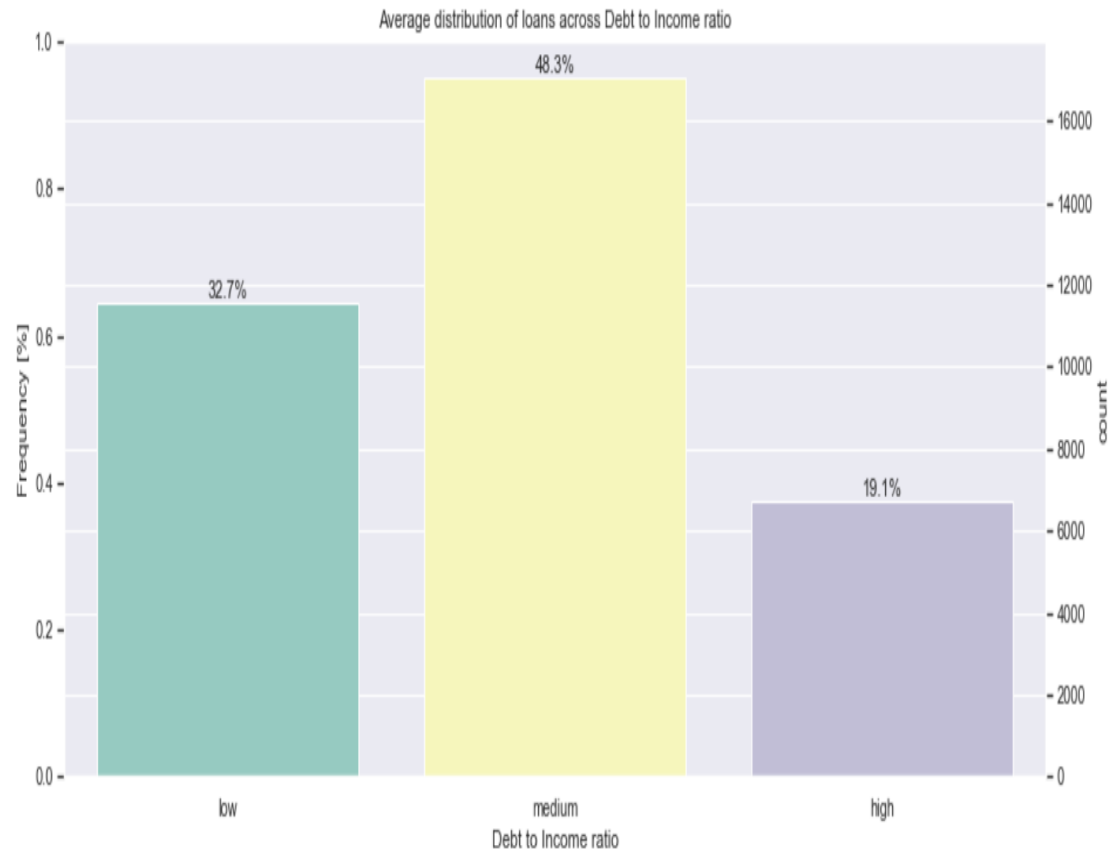


Average distribution of loans across various Grades



Average distribution of loans across various categories





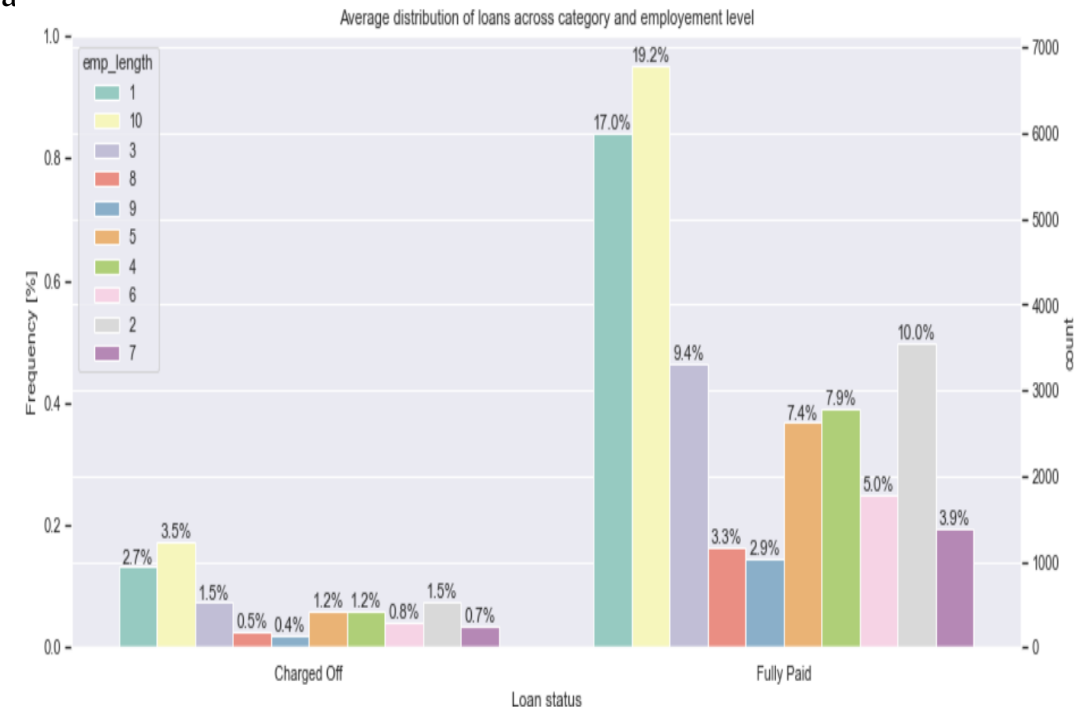
Segmented Univariate Analysis:

- On making an analysis for experience level of employee for different loan status we have found that, employees of junior level (≤ 3 Years) are having the more chances of default than other employees
- On making analysis on DTI(Debt to Income ratio) range, people with medium(>10.00 and ≤ 20.00) DTI range are having more chances of default than the other people. Similarly, people with High(>0.20) Loan to Income ratio range are having more chances of Default.

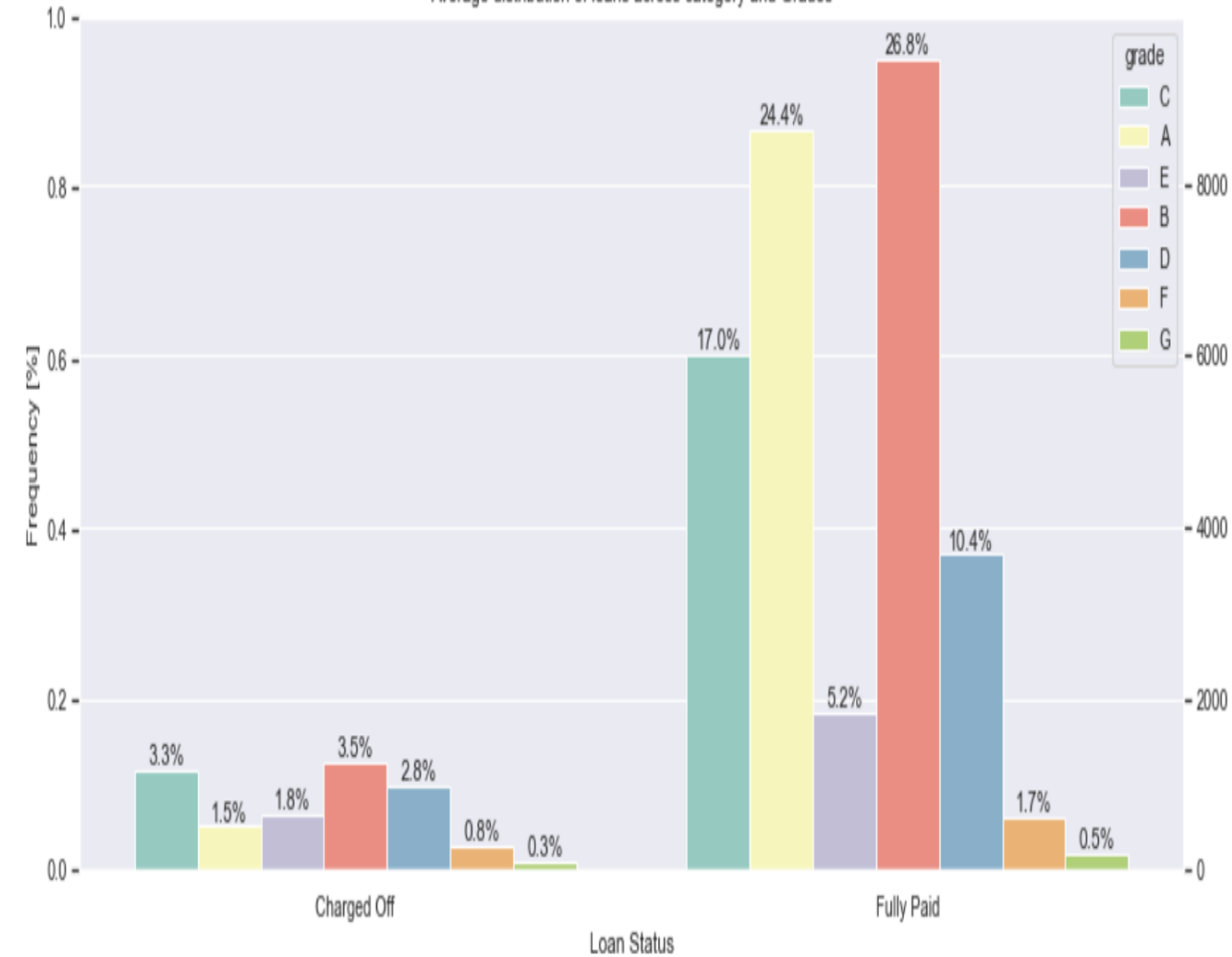
It means that people of having less employment experience will be having more DTI and LTI ratio since, there income will be less and expenditures including loan instalments will be high. Since, there will be more chance of defaulting the loan due to the burden of other expenditures.

- On making analysis for different grades of people, we found that people with grades B,C are having more chances of default
- When it comes to purpose of loan, people who took loan for the purpose of 'debt_consolidation' are having more chances of default than other purposes. It means, people who are taking loans to clear other debts are having more chances of default
- People who are having own houses are having less chances of default than people who are living in rented/mortgaged houses

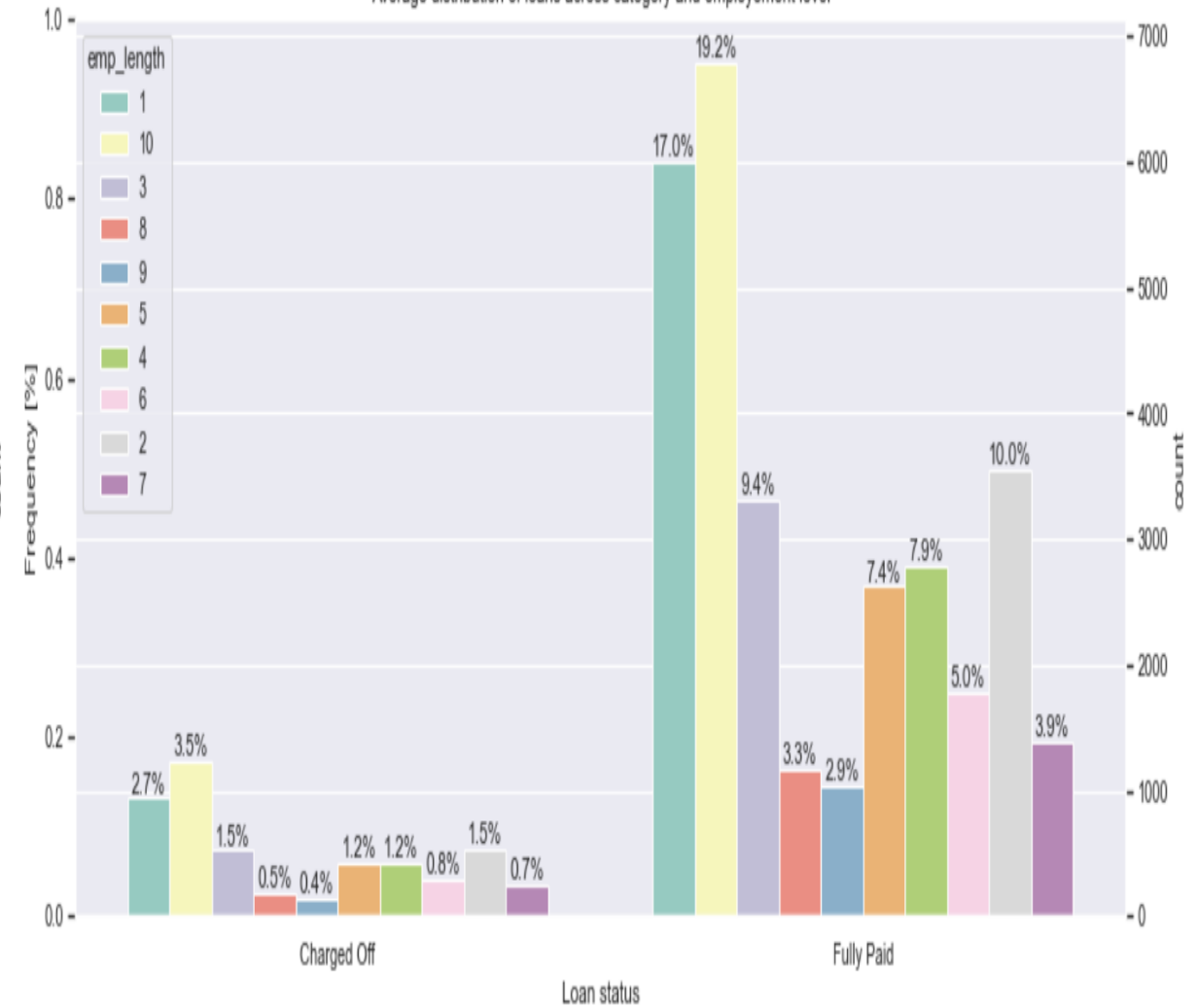
Following graphs helps in better understanding of above analysis easily

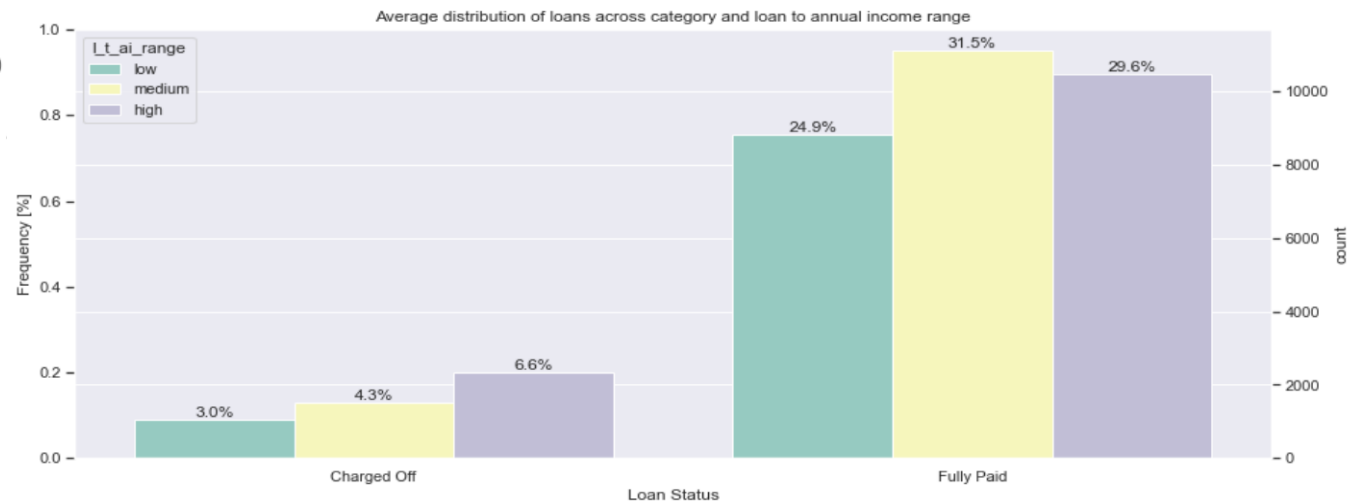
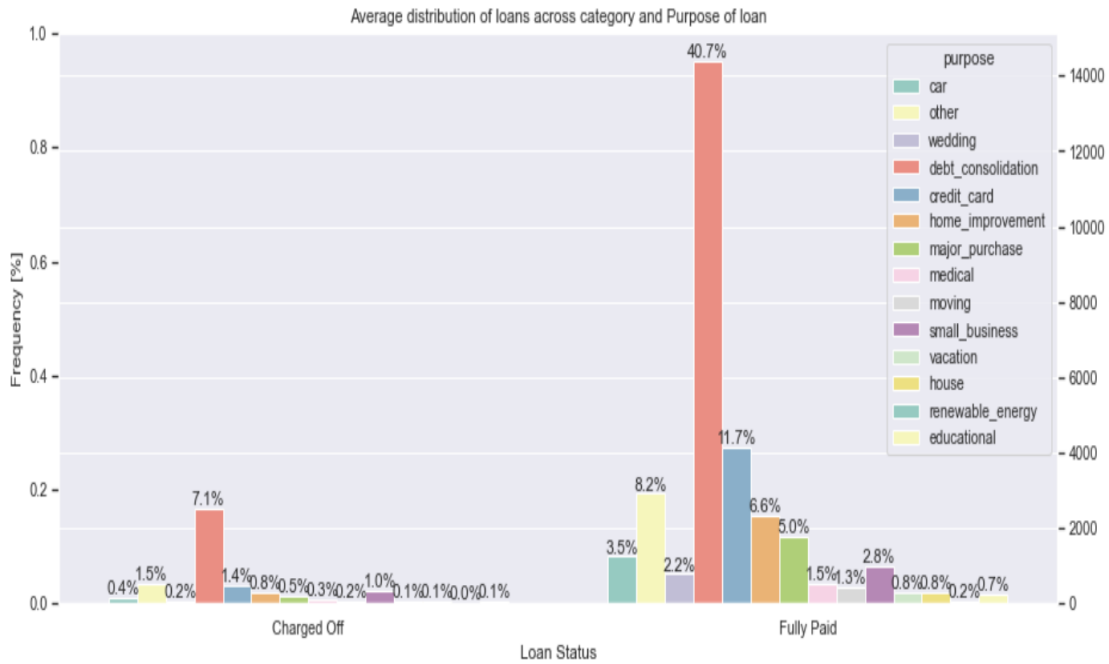
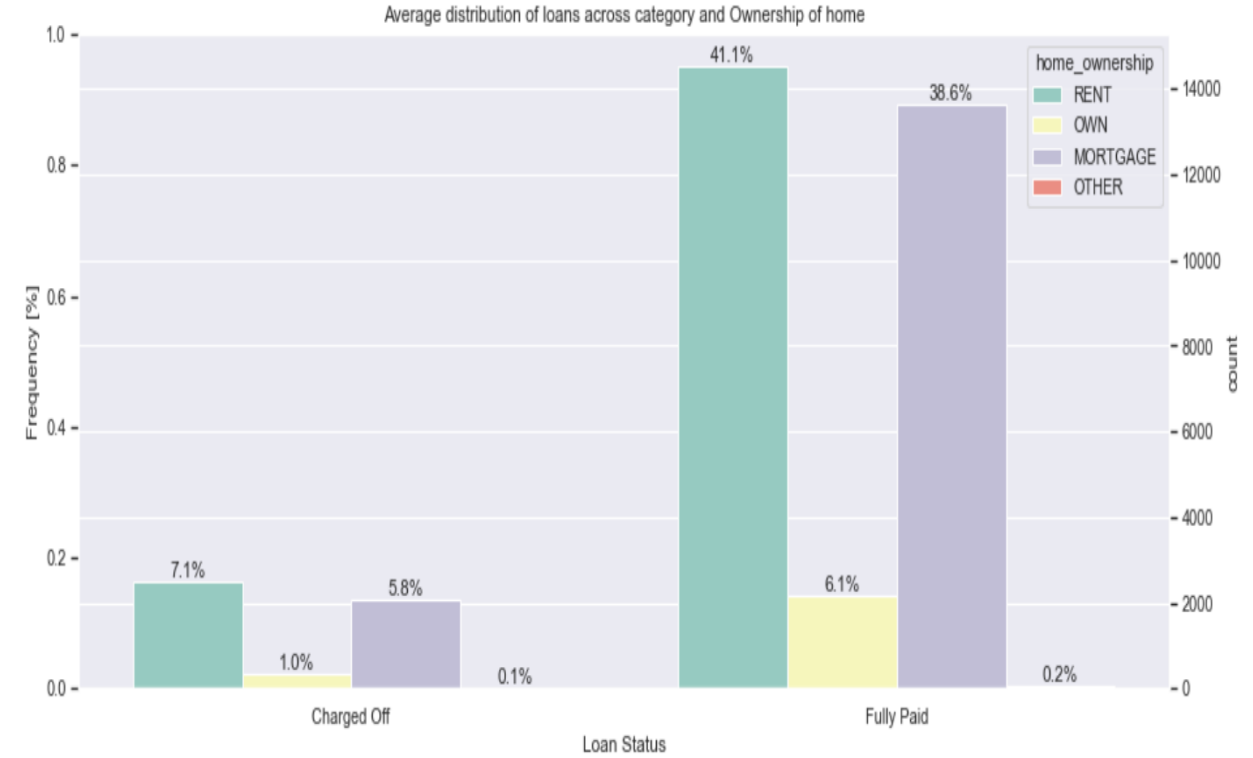
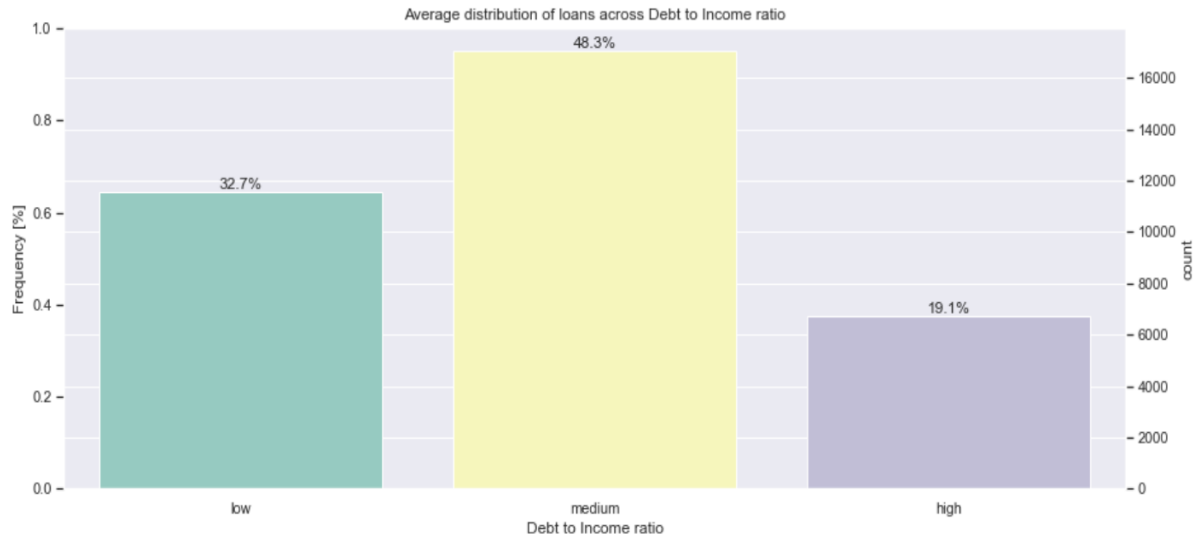


Average distribution of loans across category and Grades



Average distribution of loans across category and employment level





Bi Variet Analysis:

- In bivariant analysis we made analysis how ownership of home and verification status effects the loan status and made following outcomes from it
- **Ownership of house:**
 1. Rented House
 - A. Debt to Income ratio : Medium DTI is having more chances of default
 - B. Loan to Annual Income : High loan to Income ratio is having more chances of default
 - C. Experience level : Juniors are having more chances of default
 - D. Grades : B, C grade people are having more chances of default
 2. Mortgaged House
 - A. Debt to Income ratio : Medium DTI is having more chances of default
 - B. Loan to Annual Income : High loan to Income ratio is having more chances of default
 - C. Experience level : Specialist are having more chances of default
 - D. Grades : B, C grade people are having more chances of default
 3. Own House
 - A. Debt to Income ratio : Medium DTI is having more chances of default
 - B. Loan to Annual Income : High loan to Income ratio is having more chances of default
 - C. Experience level : Junior are having more chances of default
 - D. Grades : B, C grade people are having more chances of default

- **Verification Status**

1. Verified Status

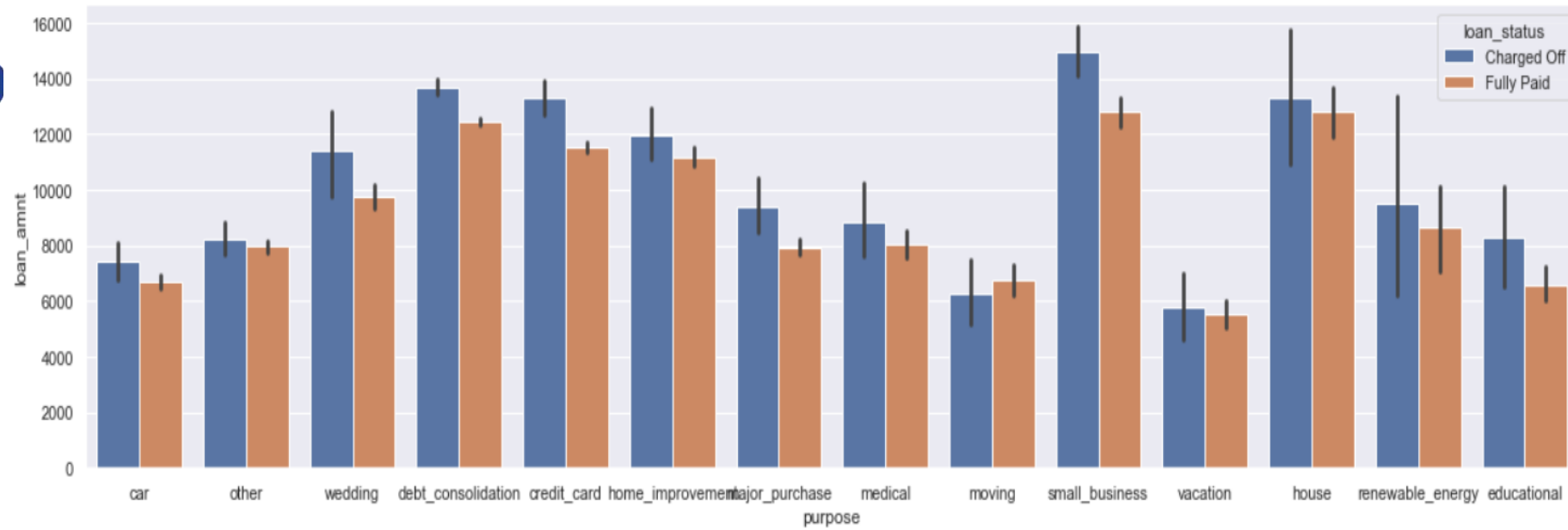
- A. Debt to Income ratio : Medium DTI is having more chances of default
- B. Loan to Annual Income : High loan to Income ratio is having more chances of default
- C. Experience level : Juniors are having more chances of default
- D. Grades : B, C grade people are having more chances of default
- E. No.Of Years Experience: 1, 10 Years of experience are having more chances of default

2. Not Verified Status

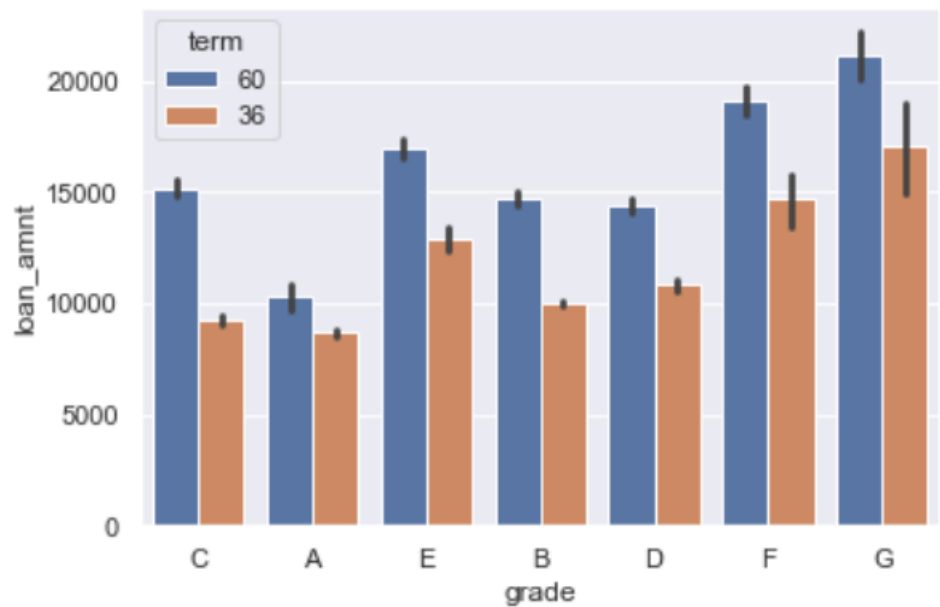
- A. Debt to Income ratio : Medium DTI is having more chances of default
- B. Loan to Annual Income : Medium loan to Income ratio is having more chances of default
- C. Experience level : Juniors are having more chances of default
- D. Grades : B, C grade people are having more chances of default
- E. No.Of Years Experience: 1, 10 Years of experience are having more chances of default

- Similarly, while making analysis between loan status, amount with different variables like as follows:
 - Verification Status: ***Loans with higher amount are verified but still they are more charged off***
 - Grades: **People with low grades takes more loan amount and are the people having more chances of defaults and people with low grades take higher amount of loans for longer duration**
 - Term: **People with high term are having more chances of default**
 - DTI Range: **People with medium & high debt to income ratio are having more chances of default. Since, they have more debts than income and fails ton pay loan**
 - LTI Range: **People with high income to loan ratio are having high chances of default by taking more amount of loan. Since there is less compared to their loan taken**
 - Employee Exp: **People with more experience level around 8-10 are having high chances of default by taking more amount of loan**
 - Purpose of loan: **People who took loan for the purpose of 'small_business', 'debt_consolidation', 'credit_card' are most likely to have chance of default by taking more amount of loan**
 - Ownership of home: **People with house of other and mortgage took more amount of loan and likely to default more**
- Lets see the above insights in graphs for better understanding

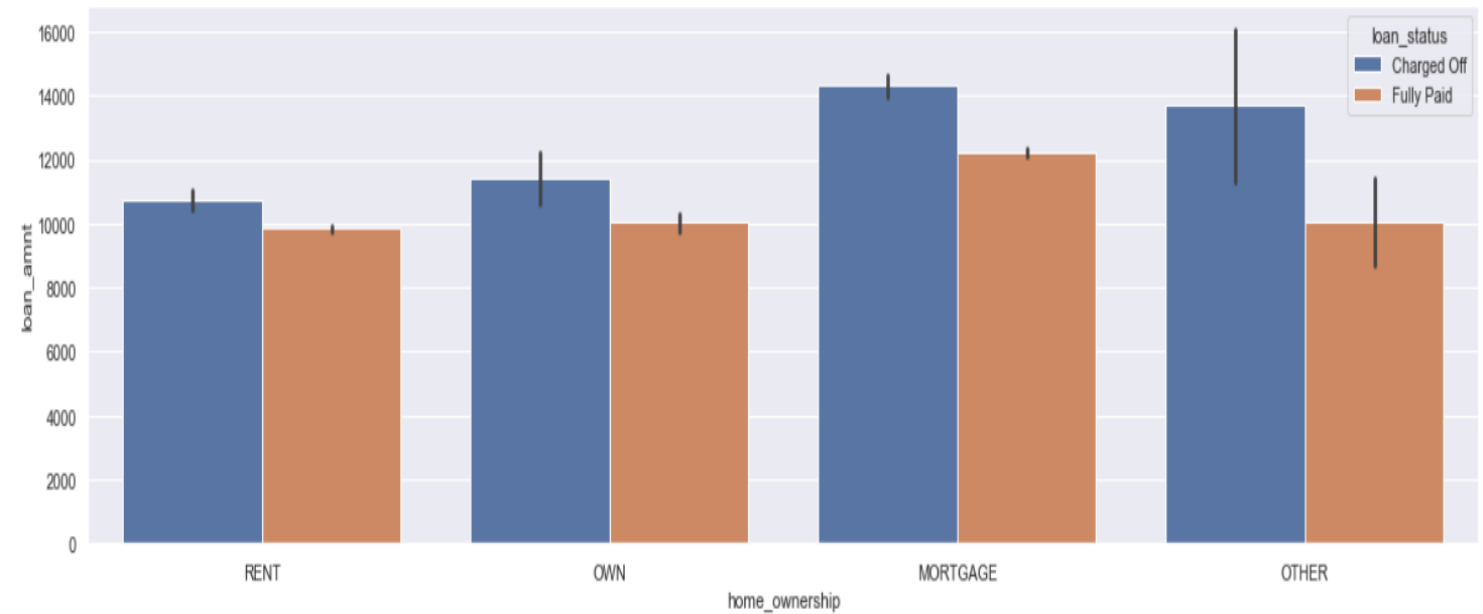
Different purposes of loan having different loan amounts and its loan statuses



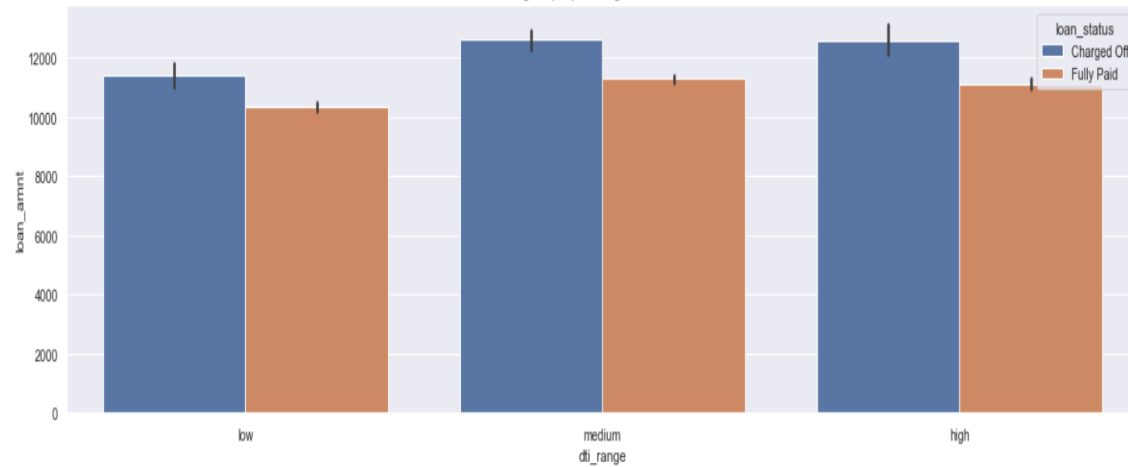
Different grades of people taking loan amount for different terms



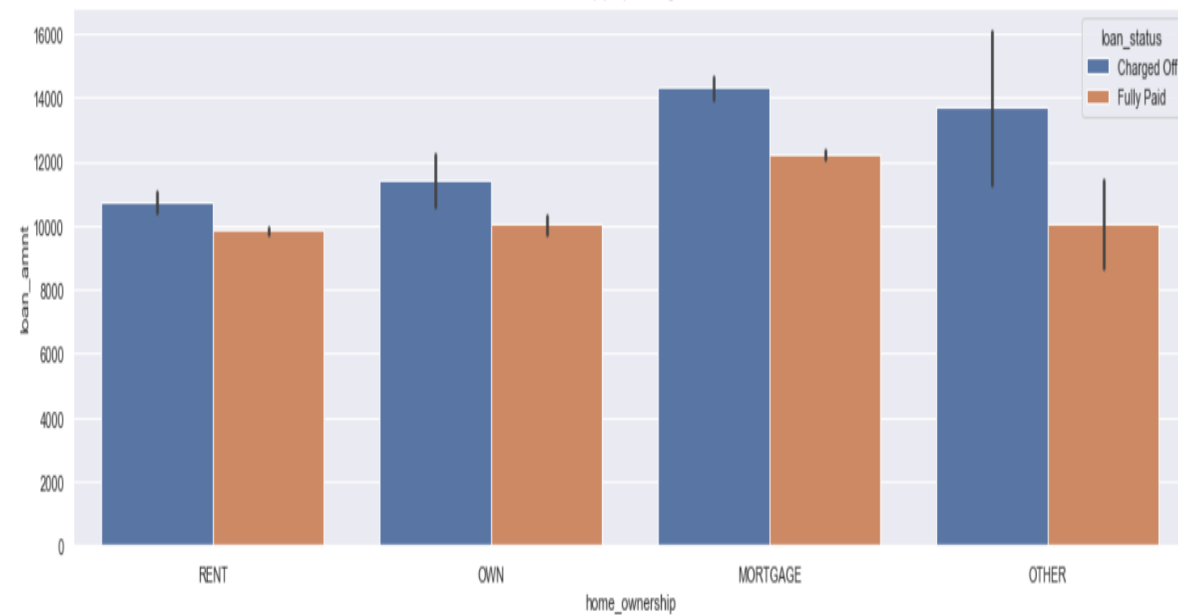
Different kinds of house ownership people taking loan amounts and its statuses



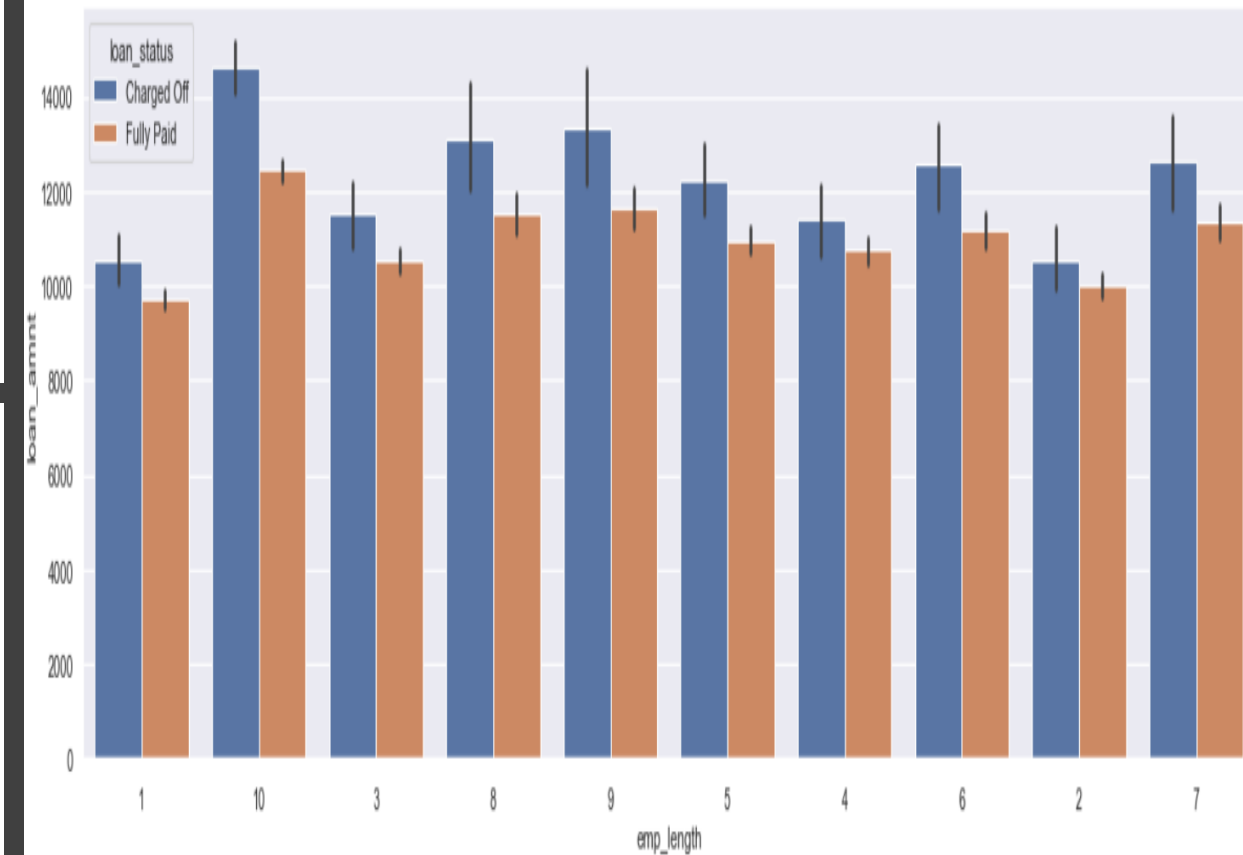
Different Dti range of people taking loan amount for loan status

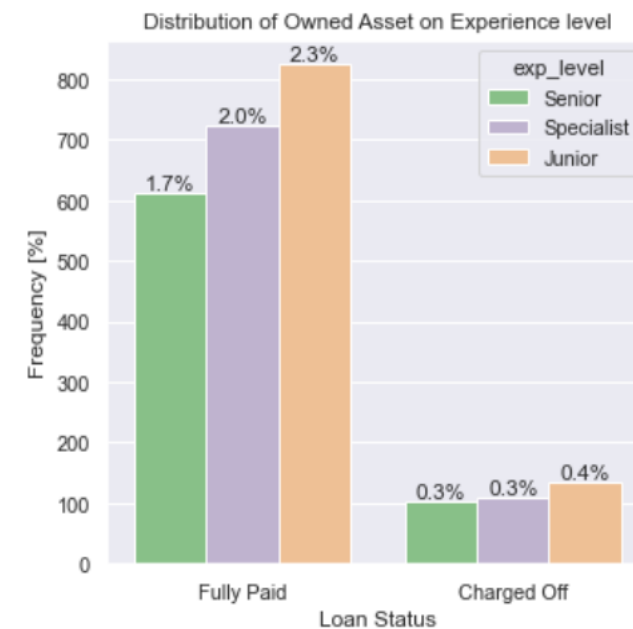
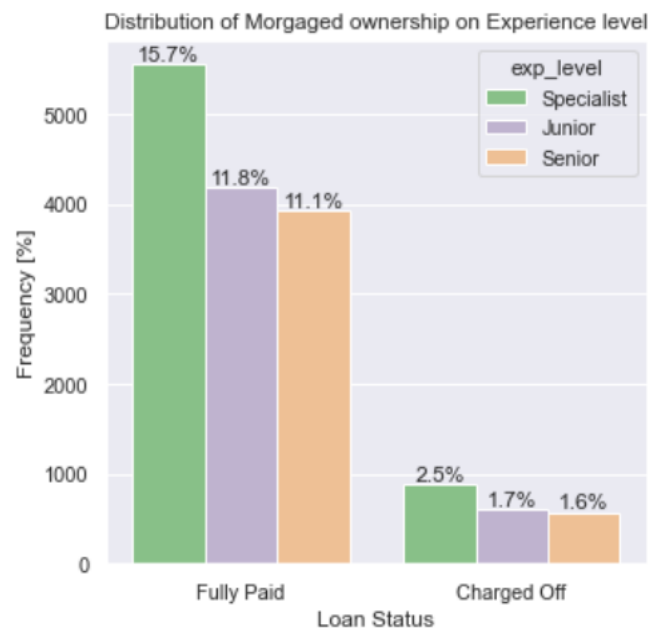
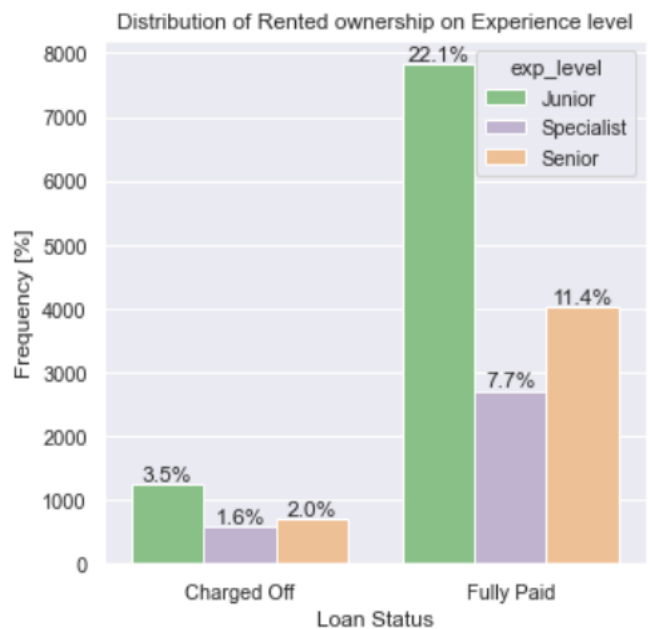
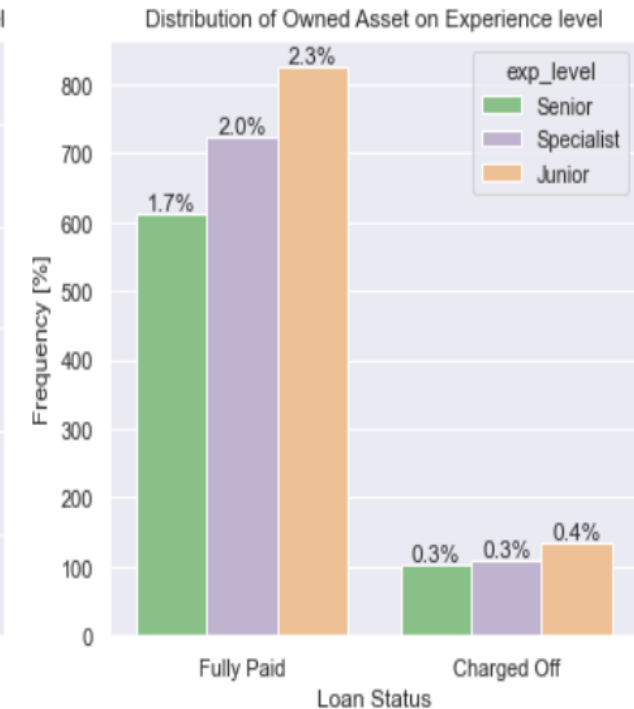
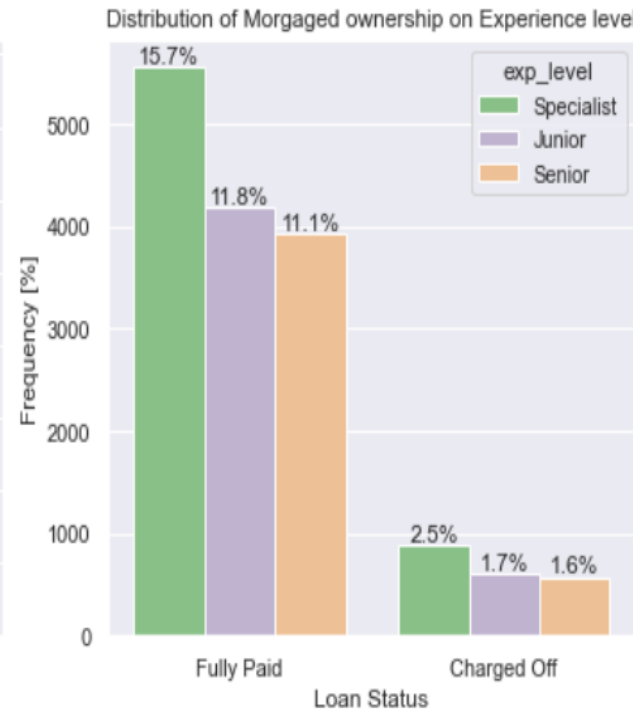
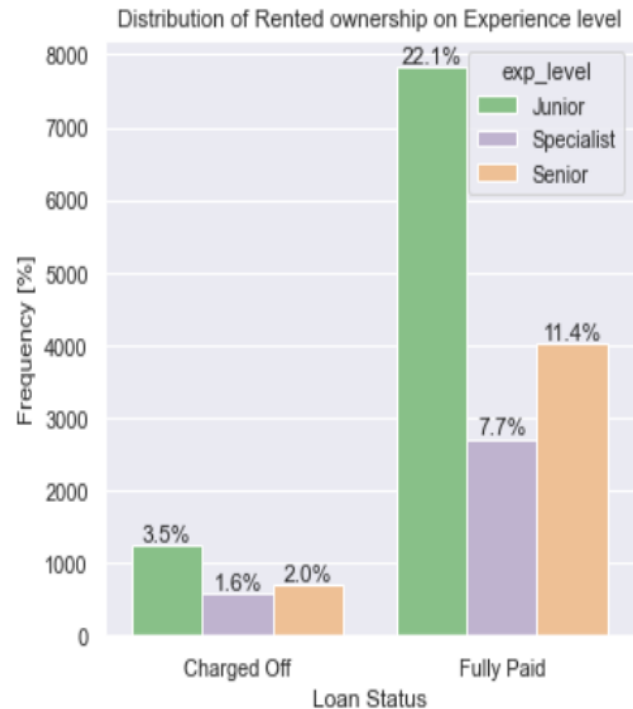


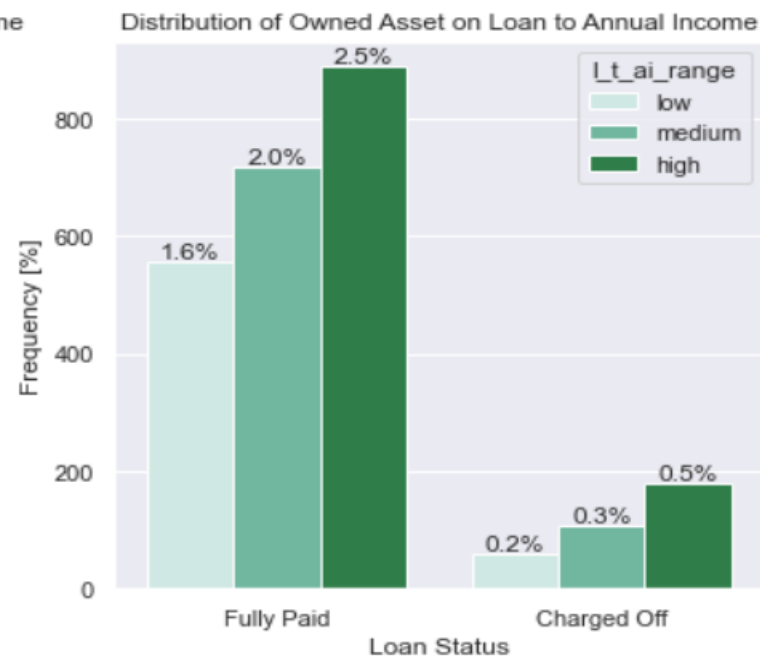
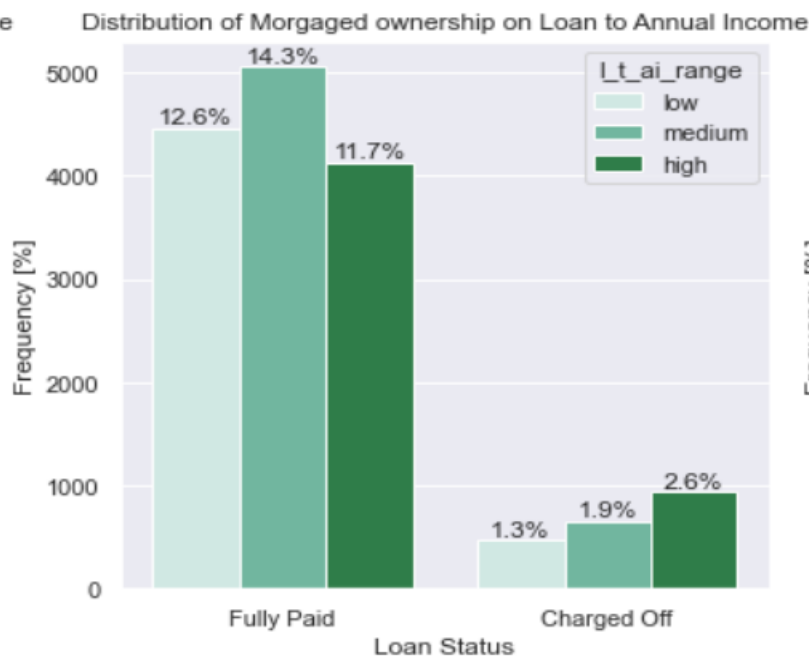
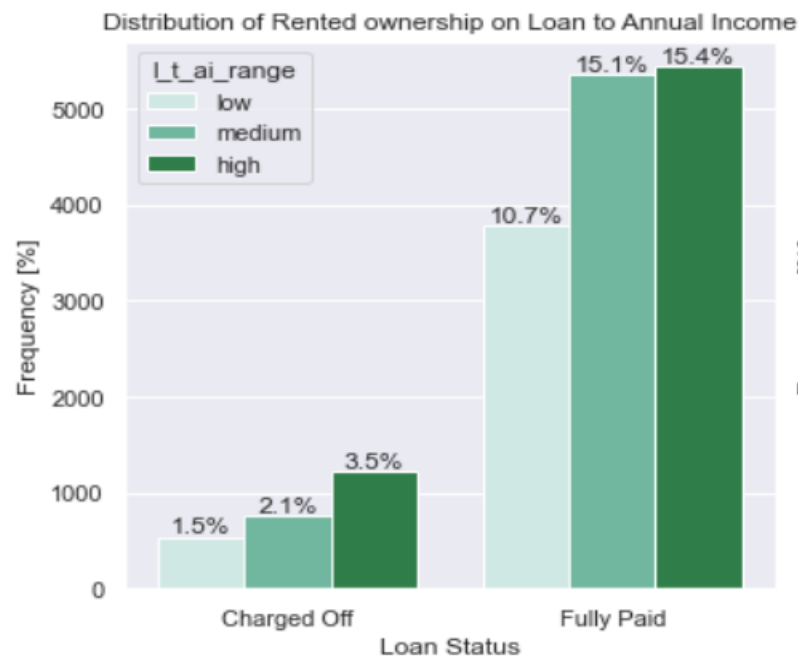
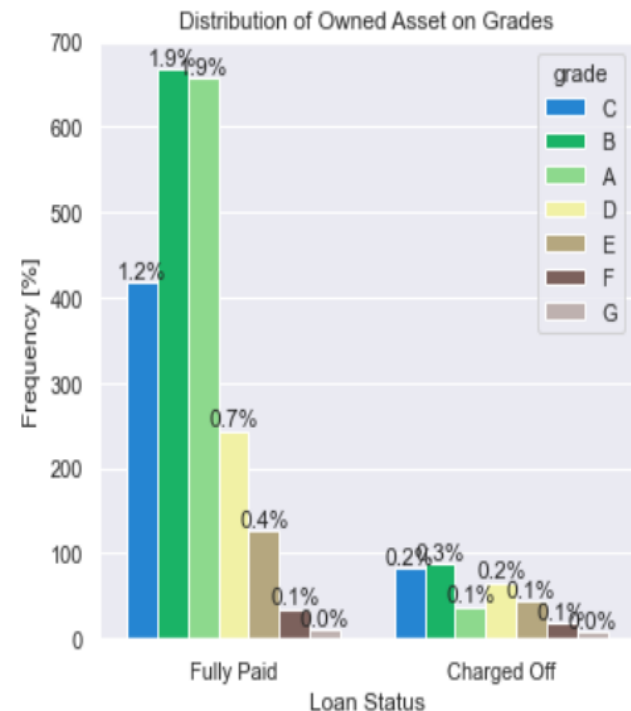
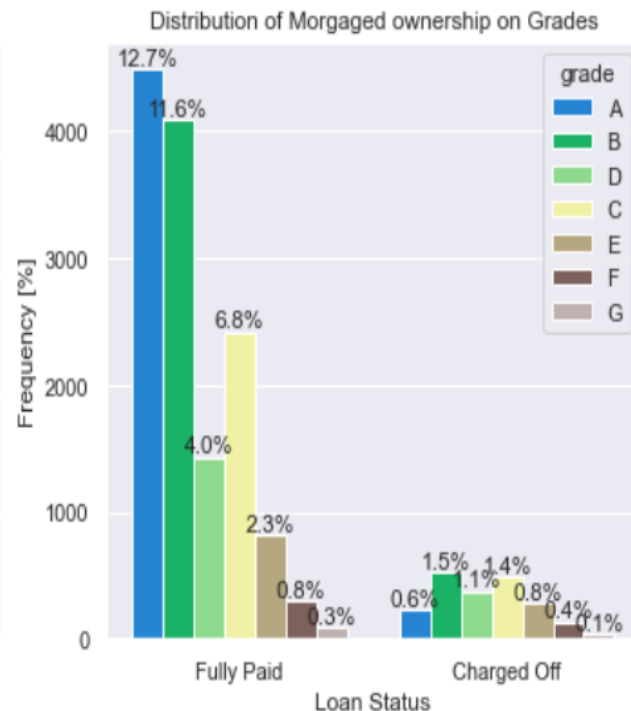
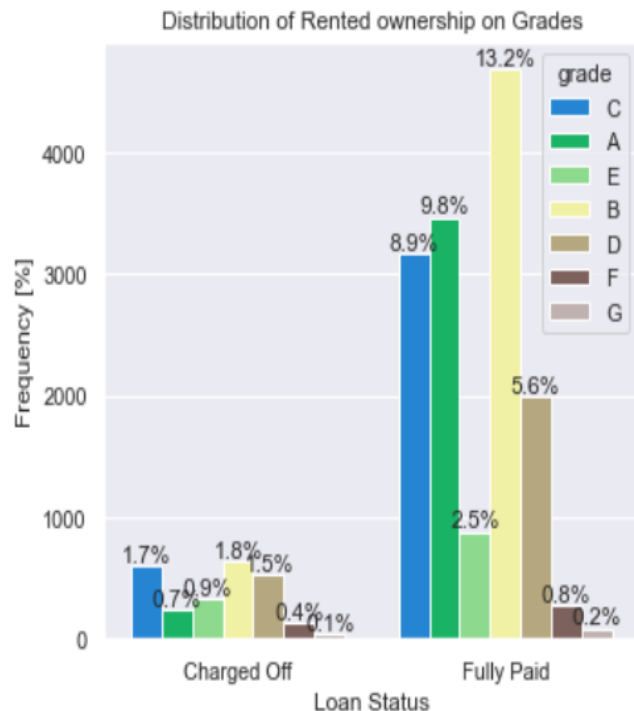
Different kinds of house ownership people taking loan amounts and its statuses



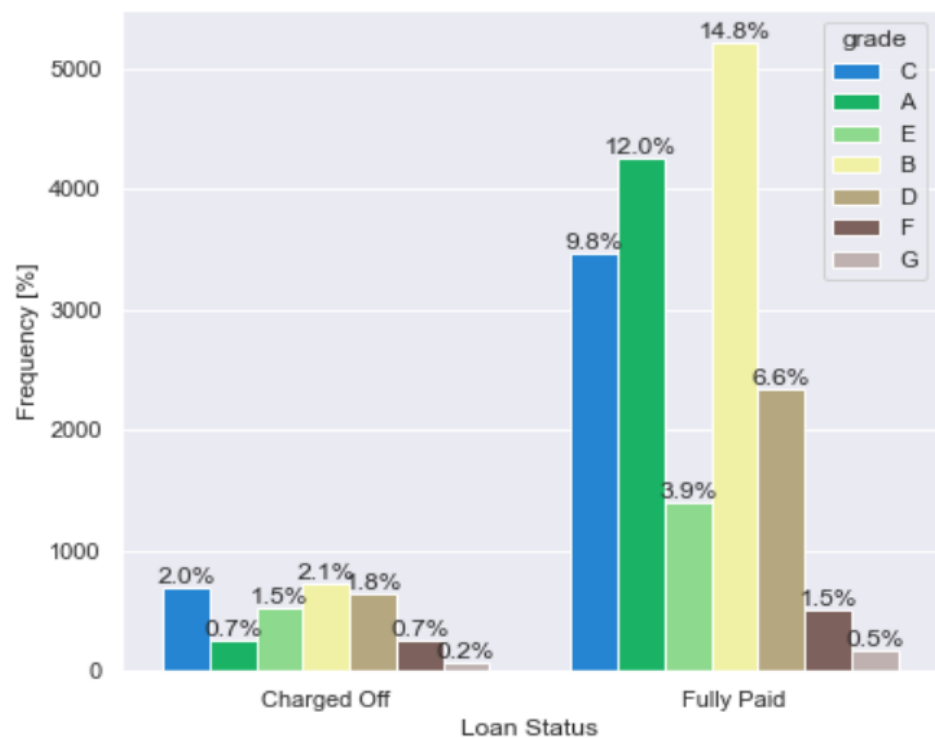
Different Employment length of people taking loan amount and their loan status



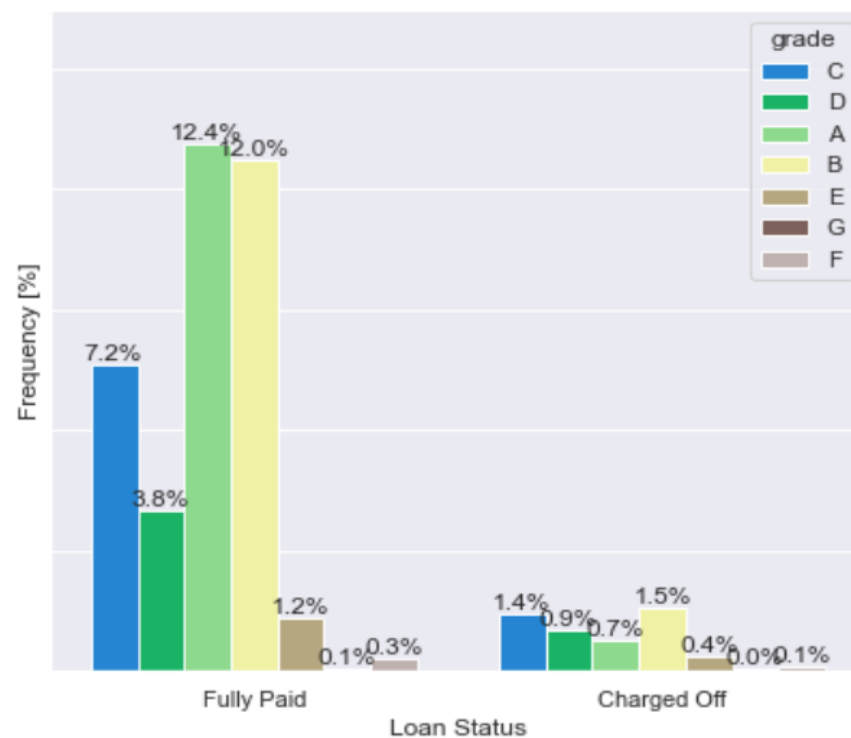




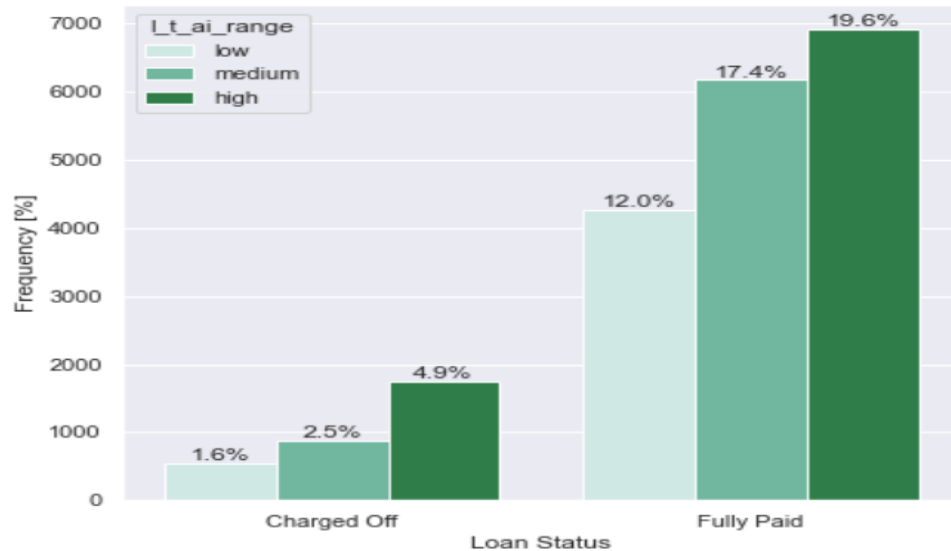
Distribution of Verified Source Income on Grades



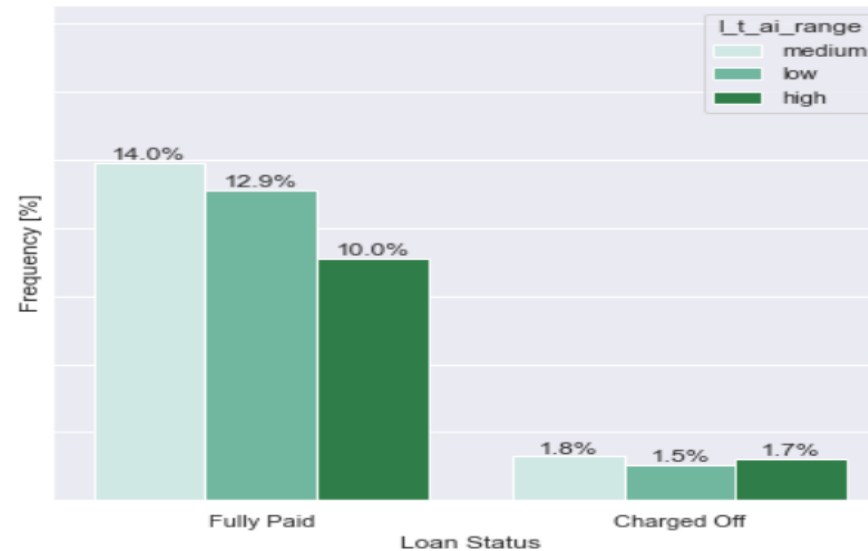
Distribution of Not Verified Source Income on Grades



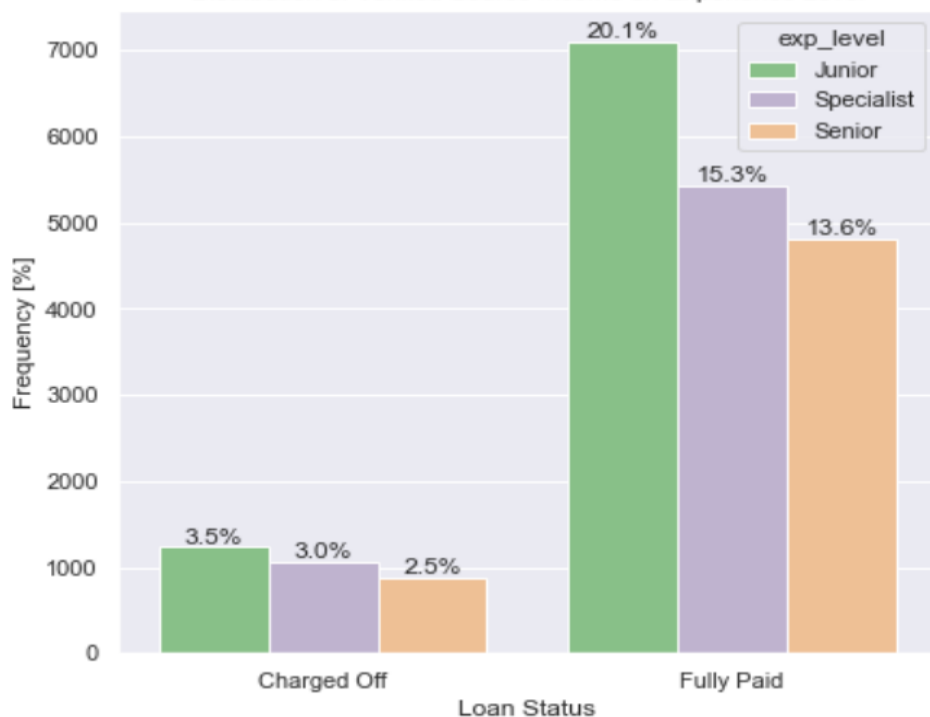
Distribution of Verified Source Income on Loan to Annual Income



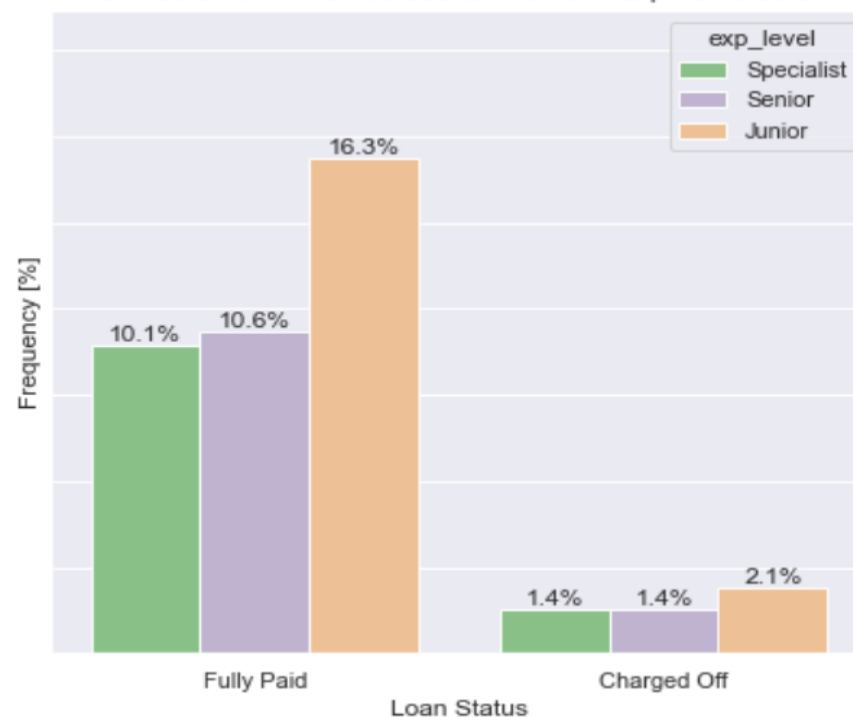
Distribution of Not Verified Source Income on Loan to Annual Income



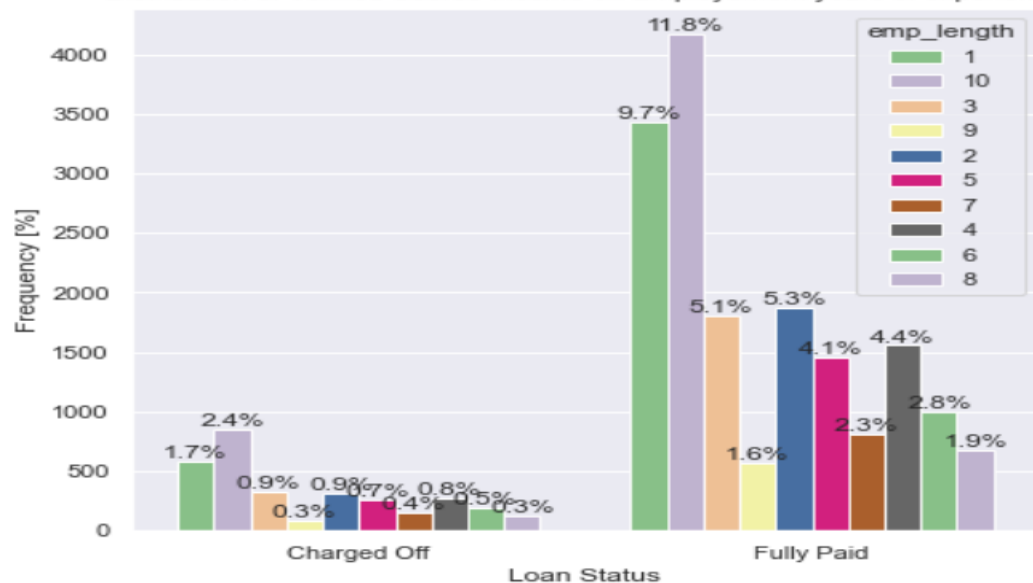
Distribution of Verified Source Income on Experience Level



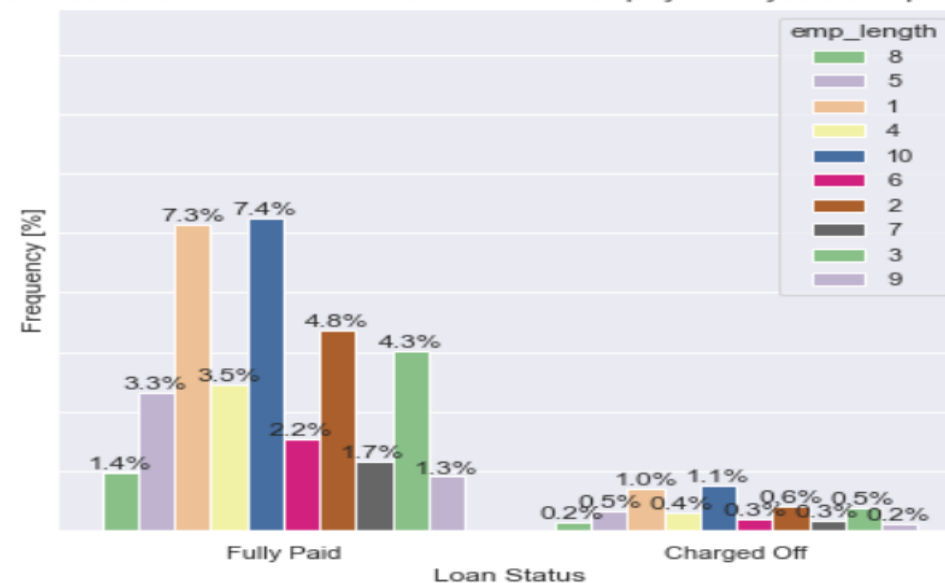
Distribution of Not Verified Source Income on Experience Level



Distribution of Verified Source Income on Employment years of experience



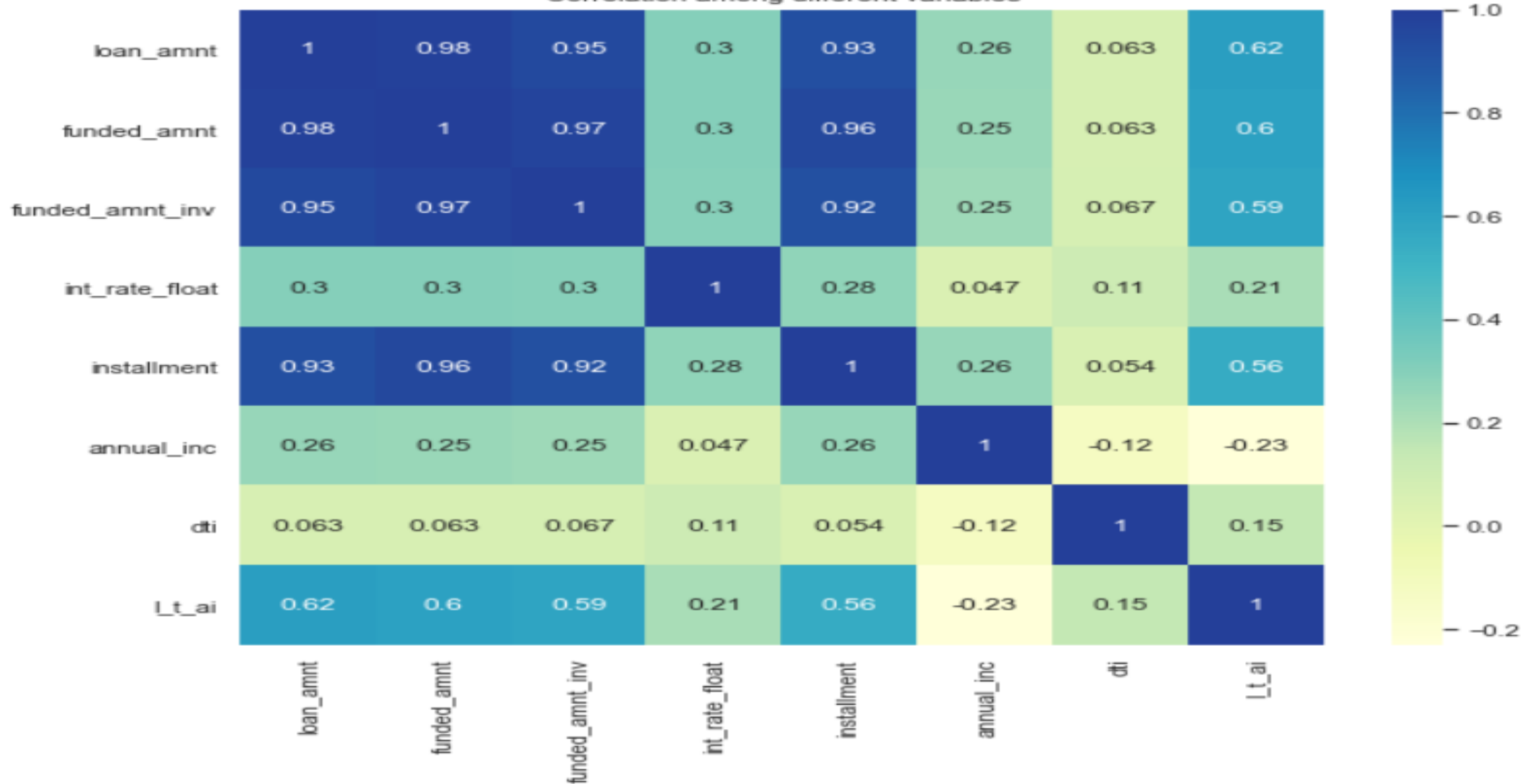
Distribution of Not Verified Source Income on Employment years of experience



Correlation:

- While making analysis using correlation, the following are the outcomes observed
 - High correlated:
 1. loan amount is highly correlated with funded_amt, funded_aqmt_inv, installment, income to amount respectively from high to low
 2. funded_amt_inv is highly correlated with funded_amt, loan_amt, installement, income to loan ratio respectively from high to low
 3. installment is correlated with funded_amt, loan_amt, funded_amt_inv, income to loan ratio respectively from high to low
 4. income to loan ratio is high correlated with loan_amt, funded_amt_inv, installement respectively from high to low
 - Low Correlated:
 1. Annual income is negative or less correlated with loan to income ratio and debt to income
- Lets see the graph for better understanding of above analysis

Correlation among different variables



Conclusion:

Lending needs to take care of following characters while considering loan:

- Debt, Loan to Income ratio
- Grade
- Purpose of loan
- Income

Following are conclusions that should be considered as per the above analysis while providing loan to the person:

- People with low grade are taking more amount of loans for long tenure are having more chances of default. So, better to give loan who are having the A grade.
- People who are taking loans with the purpose of debt_consideration, credit_card, others are having more chances of default
- People who are having own houses are having less chances of default
- People having less income and more debts are having high chances of default
- People with less delinquencies are having less chances of default