

---

---

# EDA CASE STUDY

By - Prabhat Kumar & Preetha Buddhan

---

---

# Problem Statement

- 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, etc.

# ABOUT THE DATASET

- This dataset has 2 files as explained below:
- 'loan.csv' contains It contains the complete loan data for all loans issued through
- the time period 2007 to 2011.
- 'Data\_Dictionary.csv' is data dictionary which describes the meaning of the
- variables.

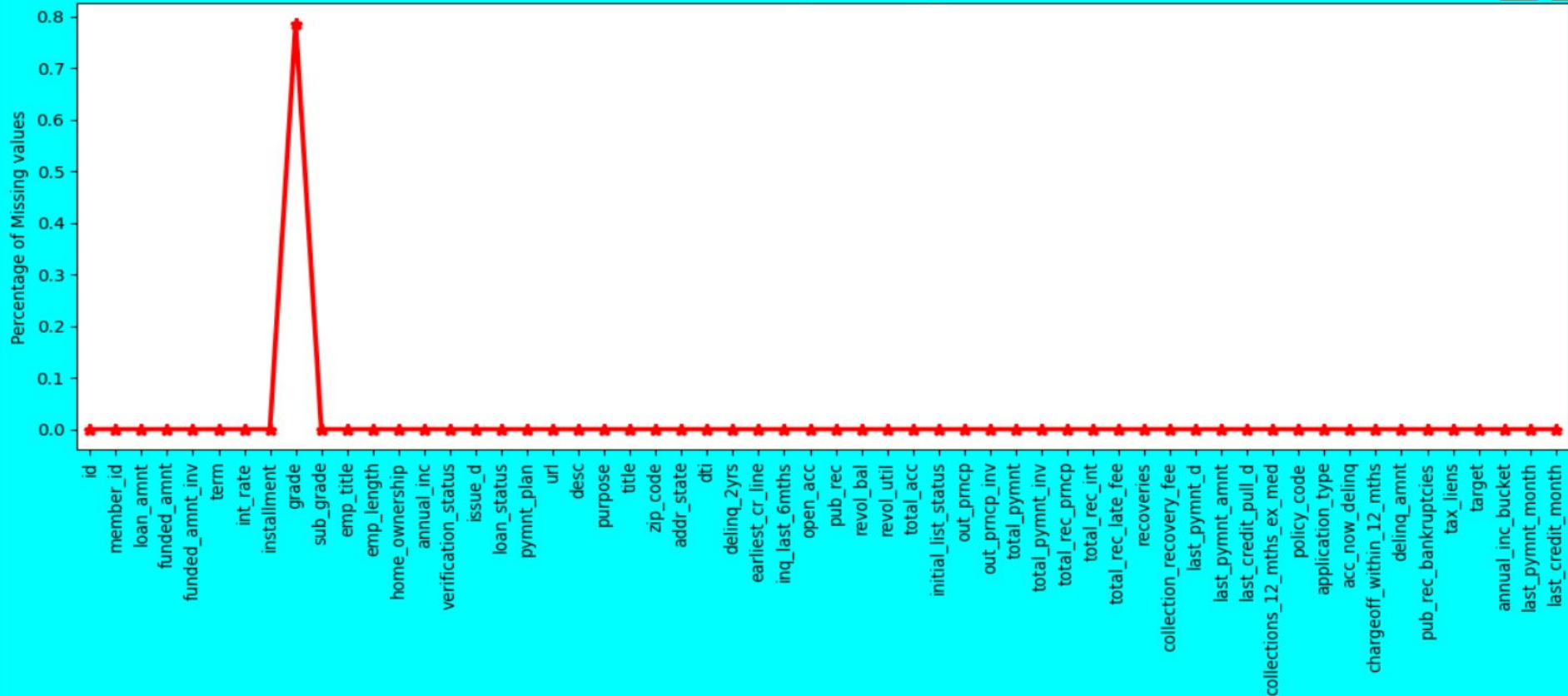
# MAJOR STEPS IN ANALYSIS

- ❖ Data Sourcing
- ❖ Data Understanding
- ❖ Checking and Handling Missing values in the data
- ❖ Handling Data Errors
- ❖ Outlier Identification and Analysis
- ❖ Univariate Analysis
- ❖ Bivariate and Multivariate Analysis
- ❖ Finding Top Correlated Features those Support Target Column.

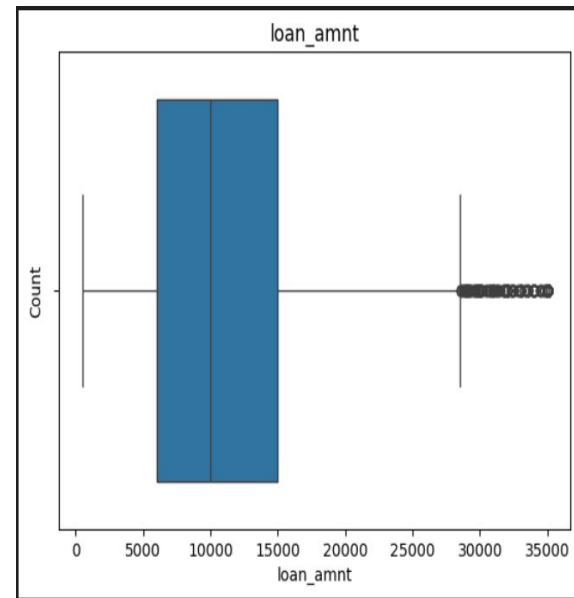
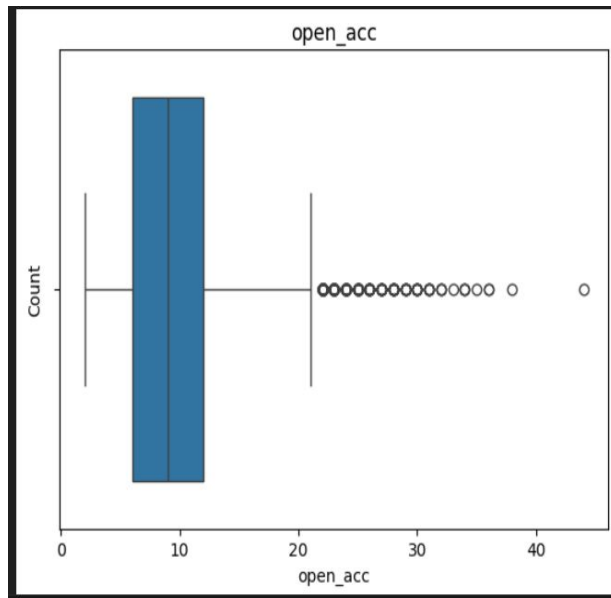
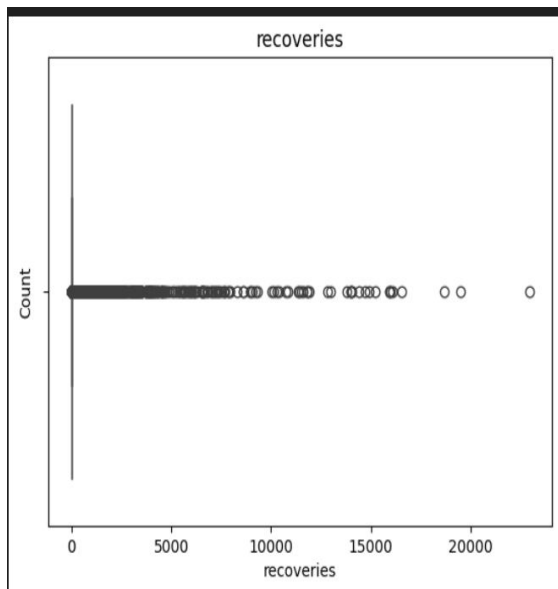
## Results on Loan.csv

# HANDLING MISSING DATA

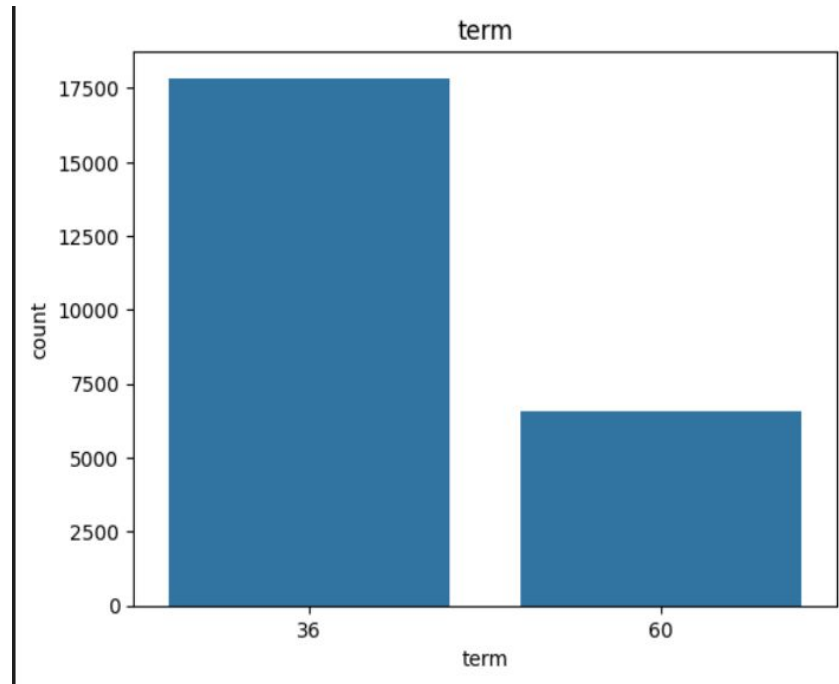
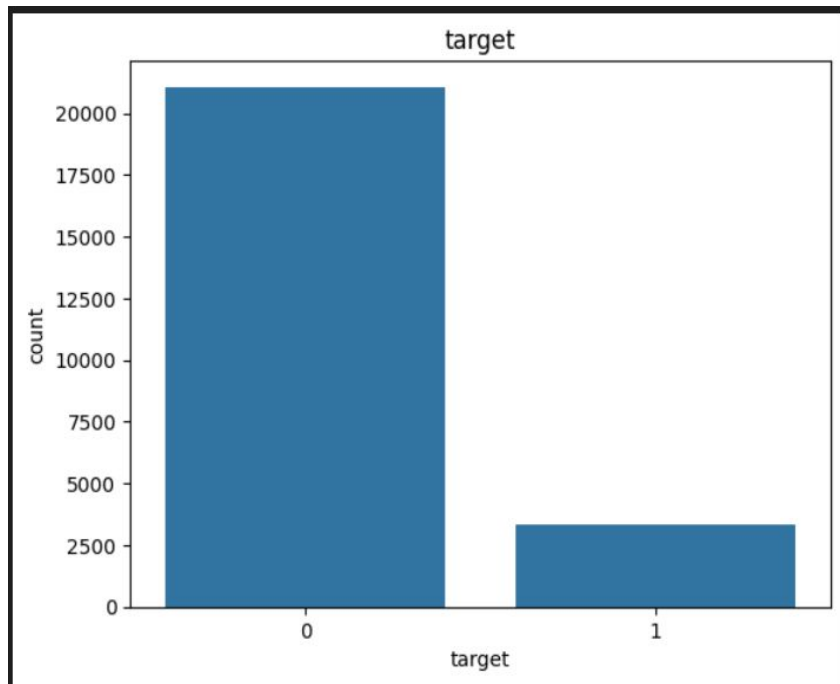
Plot for Percentage of Missing values



# Outlier Analysis

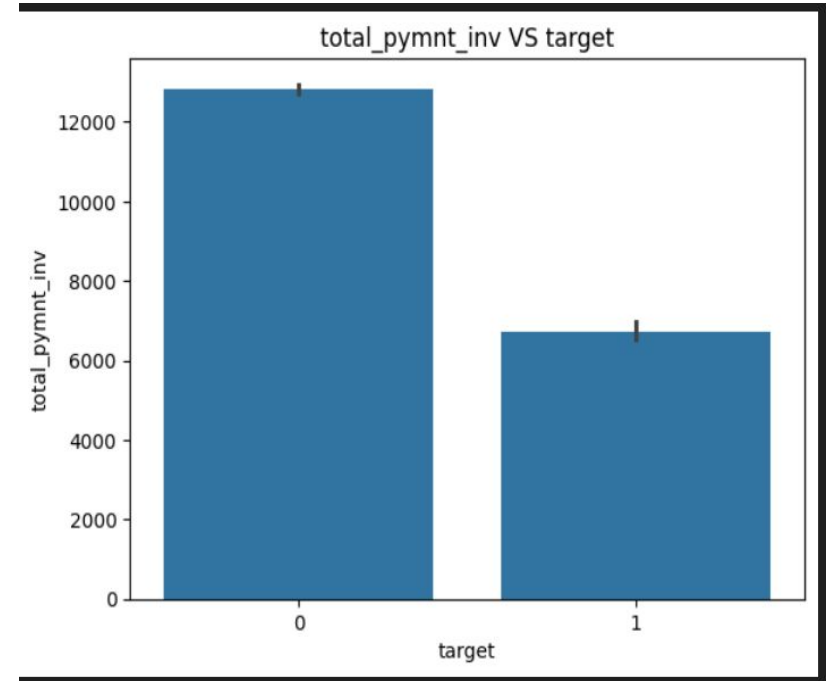
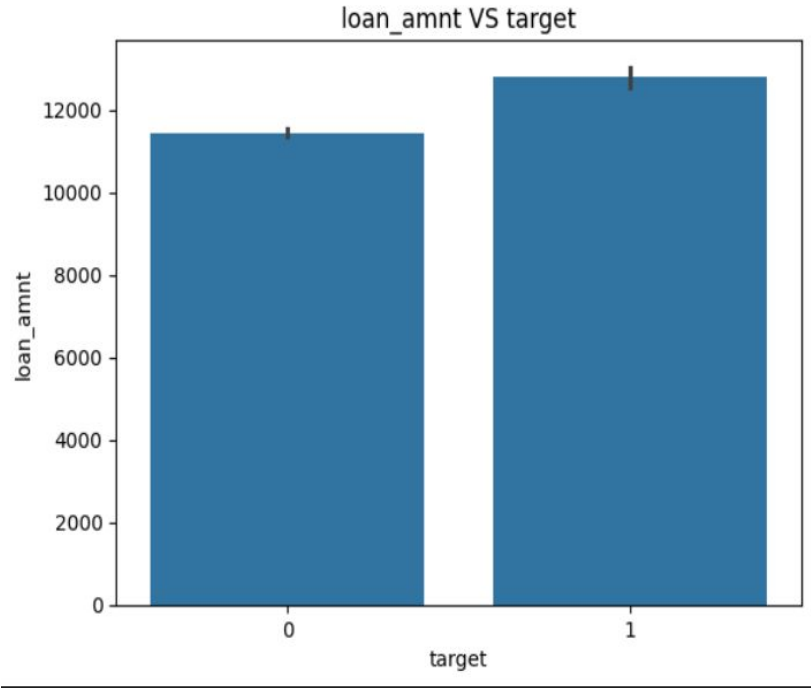


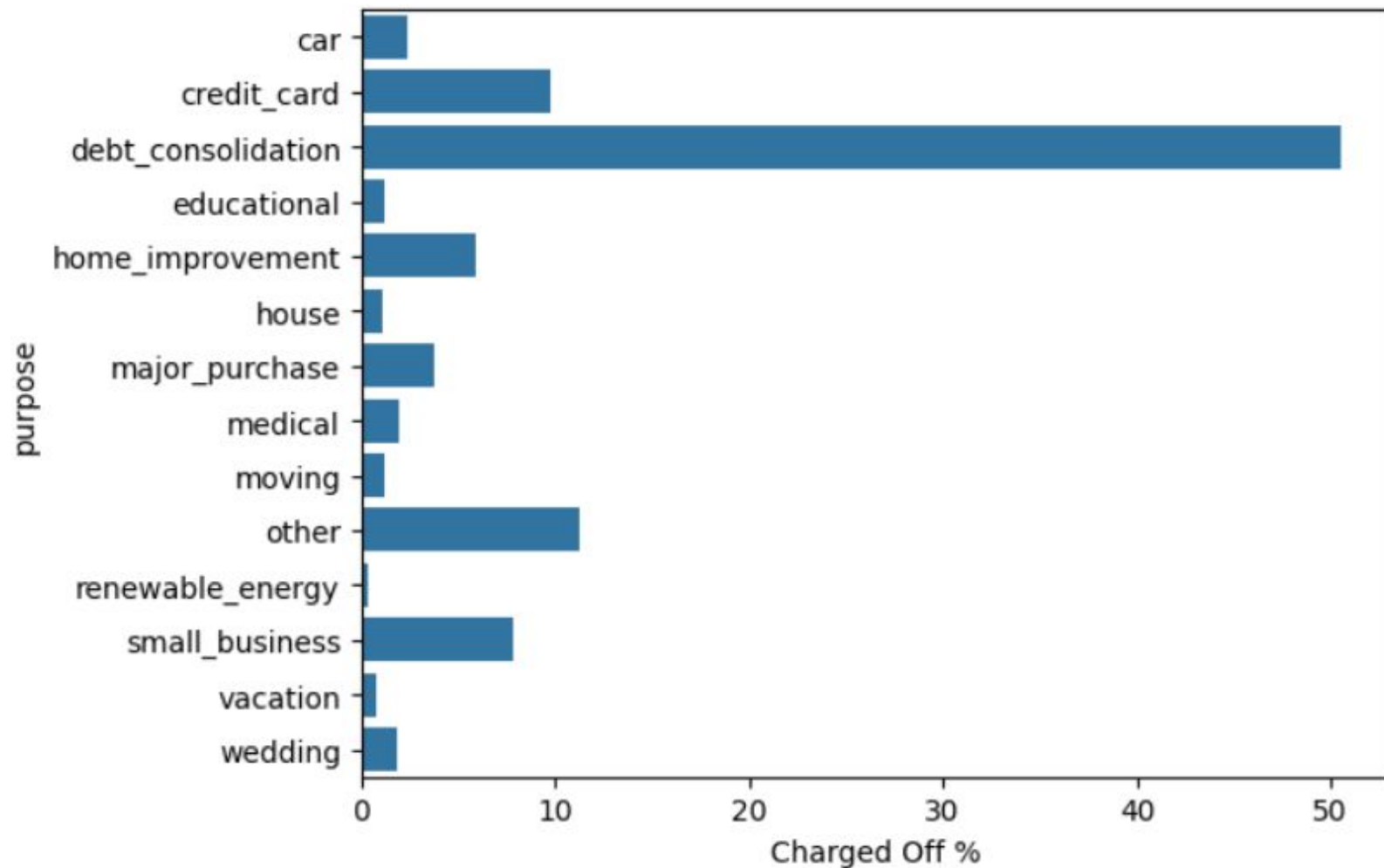
# Univariate Analysis



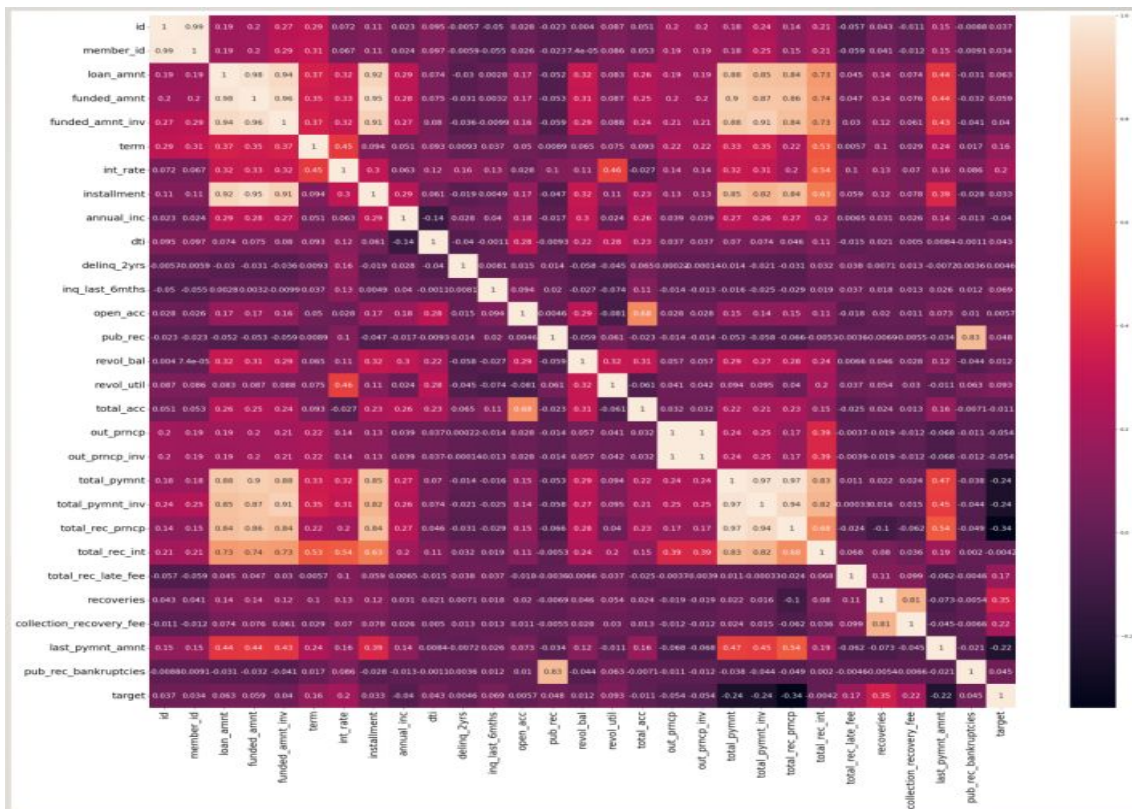


# Bivariate Analysis





# Cluster Map



# IMPORTANT OBSERVATIONS FROM THE EDA

- NO major differentiation seen in interest rates with income segment however, as income segment increases interest rates increase very slightly
- As grade changes from A to B and finally to F interest rates significantly increase. This means F are more risky customers as compared to A
- We can observe that the month May has more defaulted values compared to other in last\_credit\_pull\_d months.
- Median incomes of all three categories of customers are nearly similar. However, many Fully Paid customers have higher income levels than charged off and current customers.
- People take Higher Loan amount for long term loans and vice versa i.e. Higher Loan amount in 60 months loan tenure

# Major Observations

- Median incomes of all three categories of customers are nearly similar with increasing trend from fully paid to charged off and current customers.
- However, Many Fully Paid customers have higher installments than charged off and current customers.
- Almost 49% loan are charged off when taken for the purpose of debt consolidation which is very high
- As income segment increases installment also increases

# Take Keyaway:

From the above heatmap we can observe that the columns term, int\_rate, revol\_util has positive correlation with the target column and total\_payment, total\_payment\_inv, total\_rec\_prncp, total\_rec\_late\_fee, recoveries, collection\_recovery\_fee and last\_payment\_amount has negative correlation with the target column.

And we also observe that columns loan\_amnt, funded\_amnt, funded\_amnt\_inv, total\_payment, total\_payment\_inv, total\_rec\_prncp, total\_rec\_int has high correlation among them self.

[Detailed Analysis can be found in .ipynb file](#)

**THANK YOU**