# EDA for loan defaults

Understand driving factors behind loan defaults for a corporate finance company using exponential data analysis

# Steps performed in EDA

- Data Loading
- Data Cleaning
- Univariate Analysis
  - Data Description
  - Segmented Analysis
- Bivariate Analysis
  - Categorical Attributes vs Loan Default Ratio
  - Continuous Attributes vs Loan Default Ratio
- Multivariate Analysis
  - Correlation Matrix

# Tools and Technologies Used

- Python 3
- Pandas
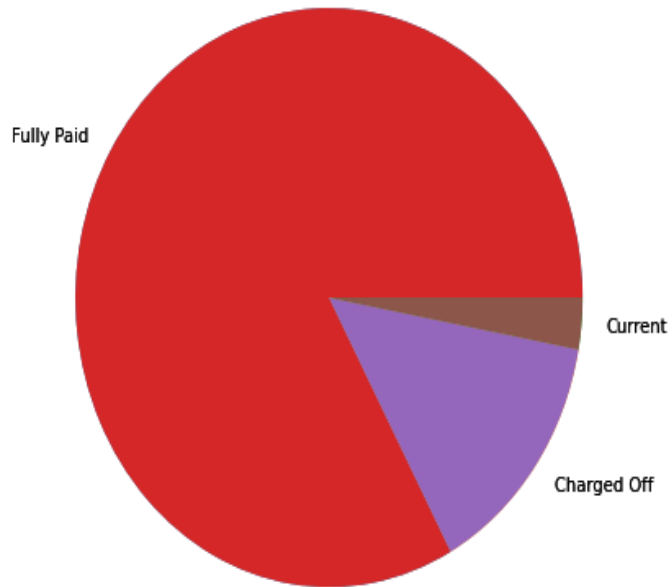- NumPy
- Matplotlib
- Seaborn
- Excel

# Dataset

- Data set is record of all loans issued by a corporate finance company from year 2007 to 2011
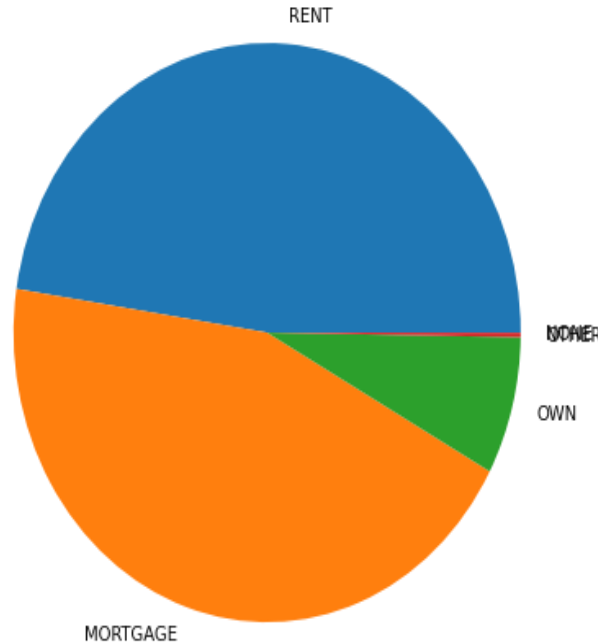
# Data Cleaning

- Cleaning empty rows and columns
- Standardising Data
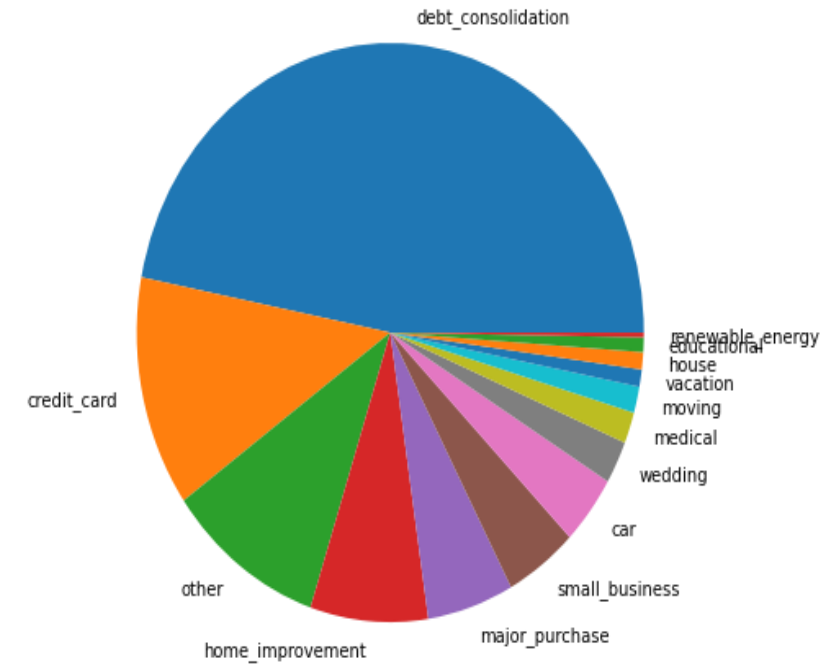- Removing attributes not relevant to analysis

# Univariate Analysis



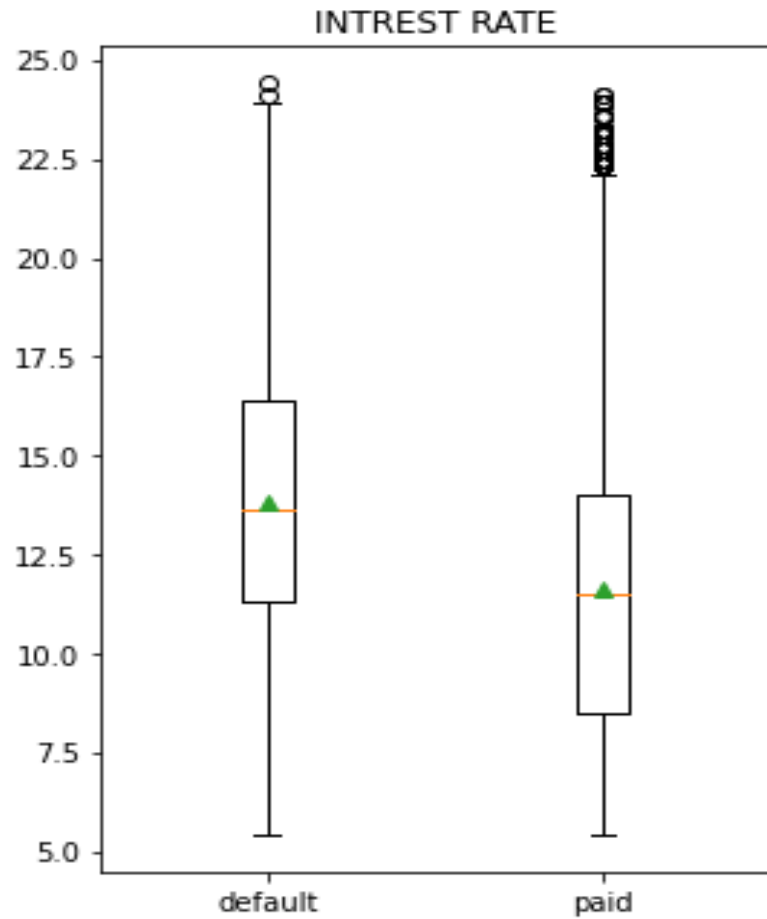LOAN STATUS COMPOSITION

HOME OWNERSHIP COMPOSITION

LOAN PURPOSE COMPOSITION
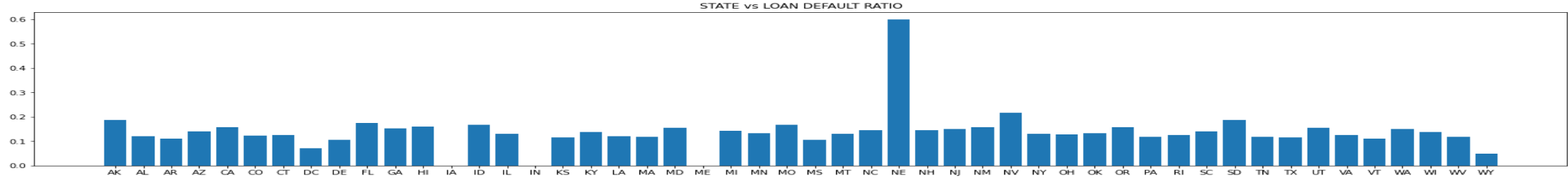
Univariate analysis of categorical variables reaveal
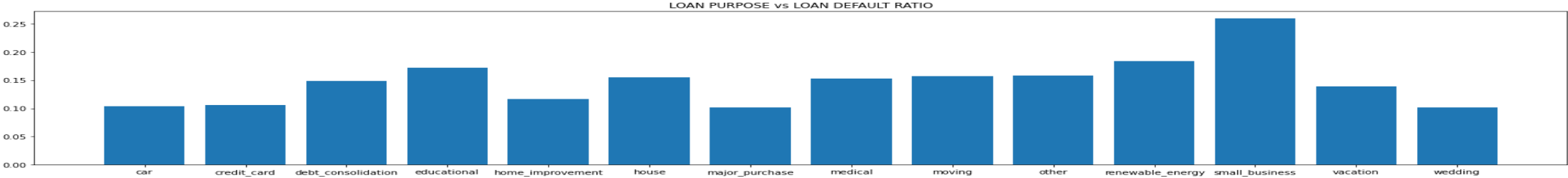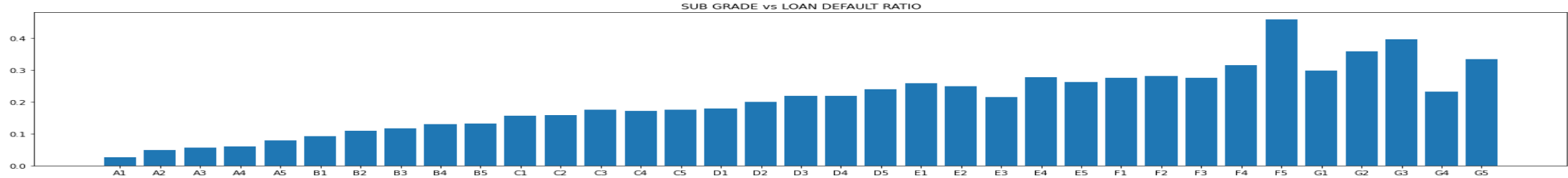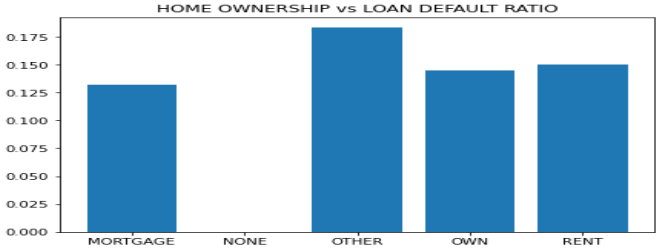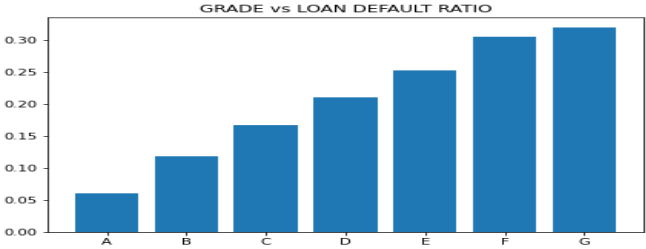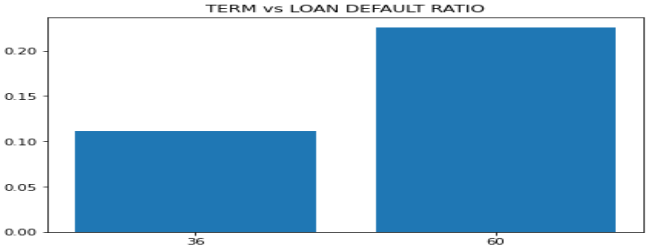- Majority of loan records are fully paid
- Most loan applicants are living on rent or have mortgaged house
- Debt Consolidation is most common purpose for taking loans

# Univariate Analysis - Segmented



Box plot of interest rate for paid vs defaulted loan reveals that defaulted loan on average have high interest rate
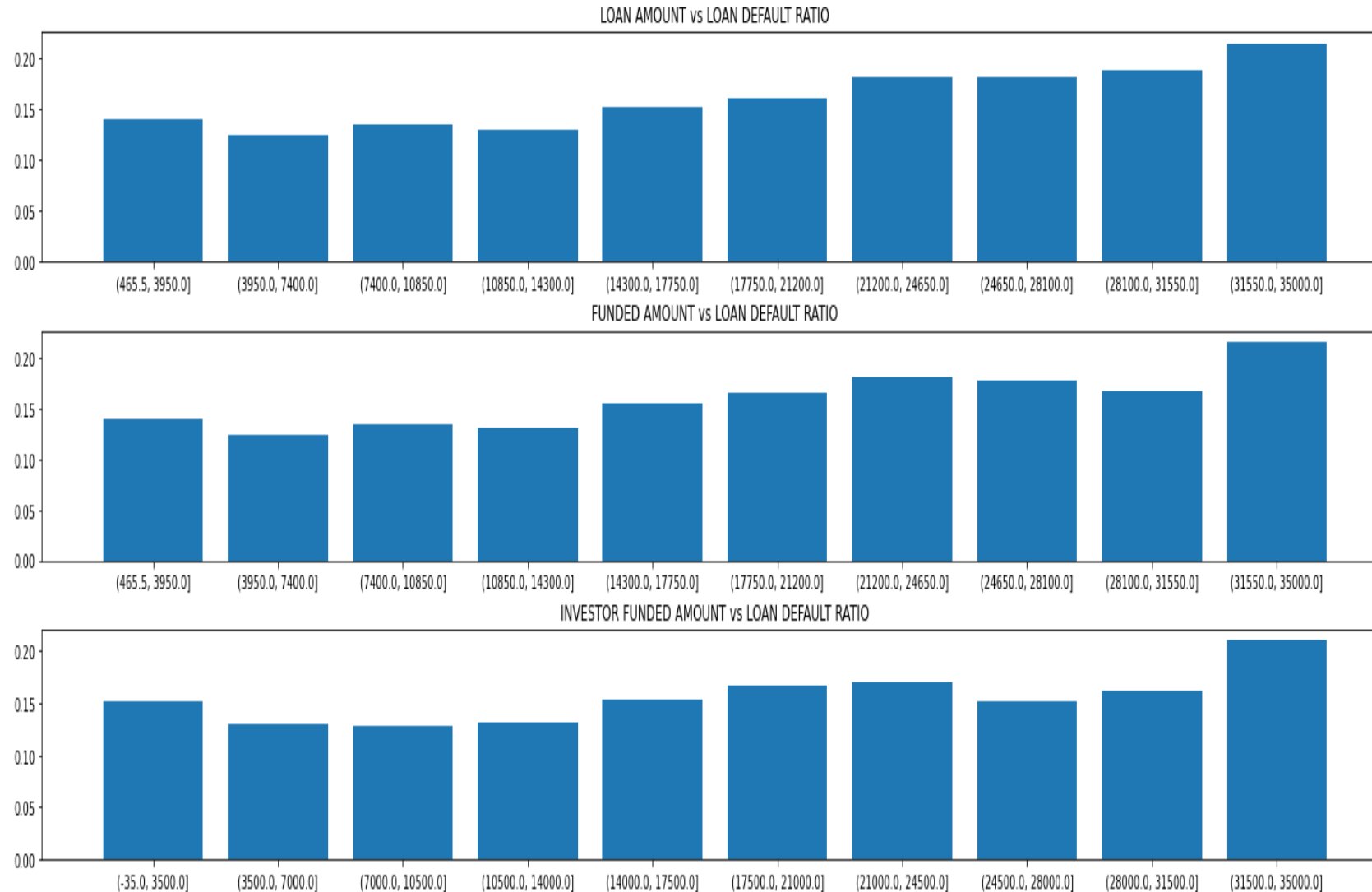
# Bivariate Analysis - Categorical
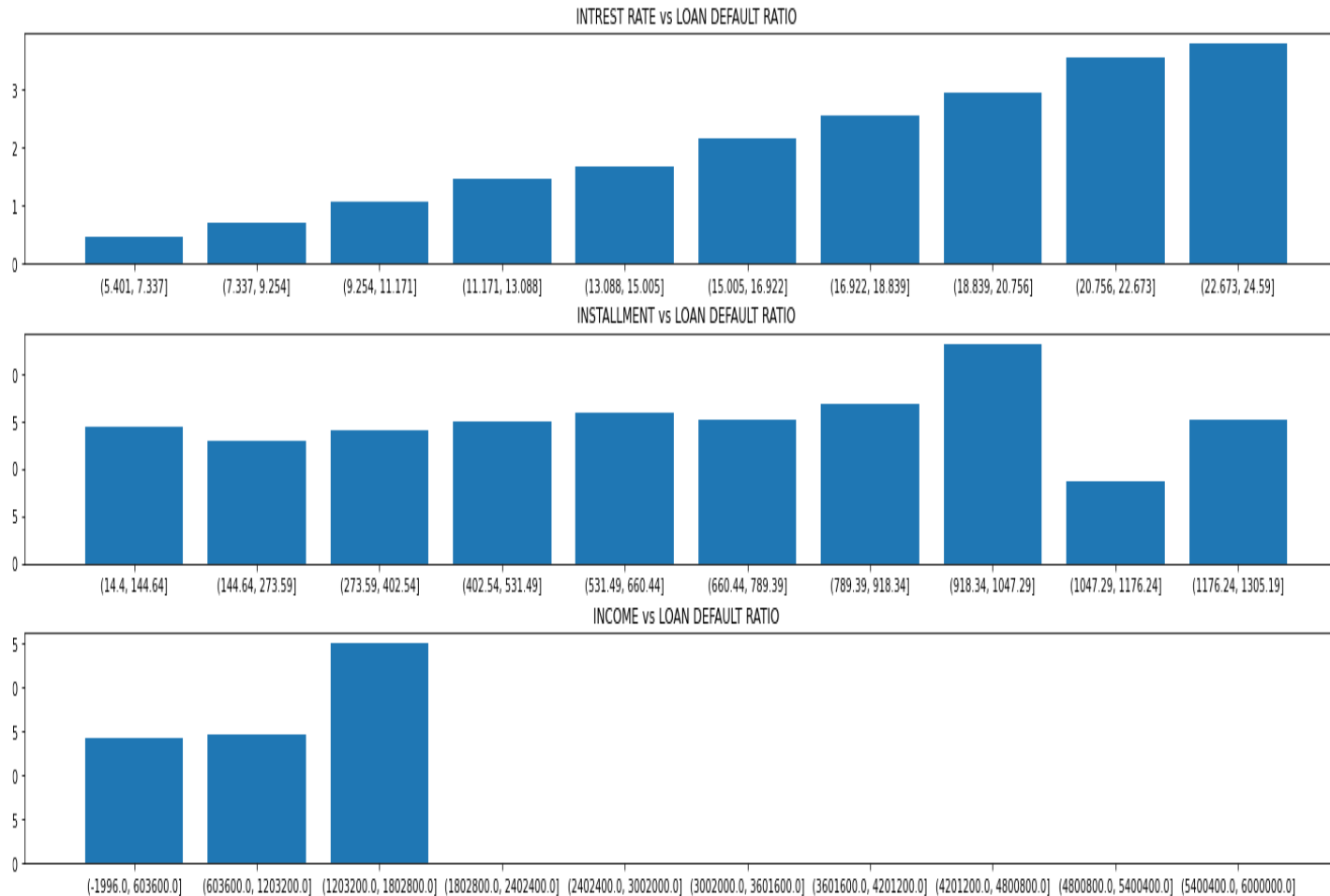
# Bivariate Analysis – Categorical(contd.)

- Long term loans in general have higher default rates. This can be attributed to chance of financial distress being higher in long duration.

- Grade and Subgrade both show positive correlation with loan default ratio

- Home Ownership doesn't have much influence on Loan Default.

- Loan taken for small businesses have significantly higher default rate. This can be attributed to unpredictability with successful running of small business. However rest of the loan purpose are close to average.
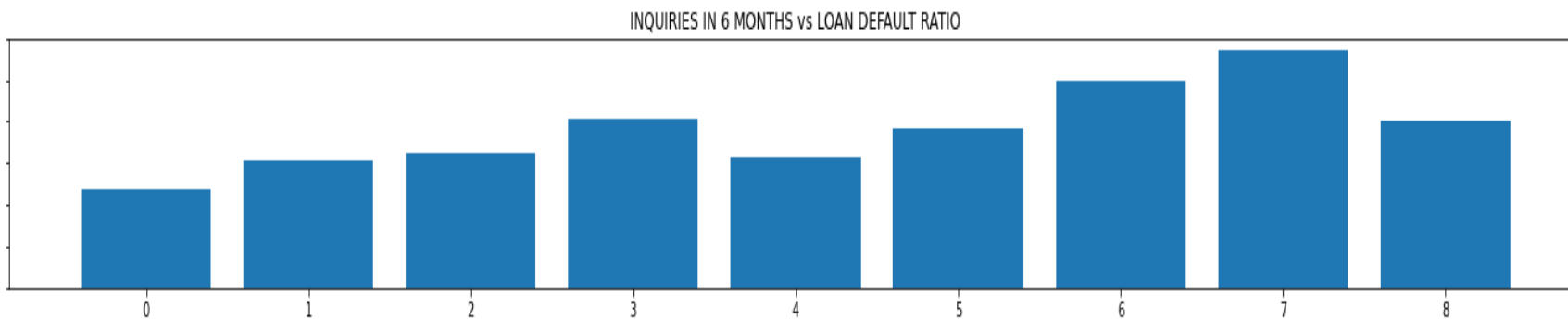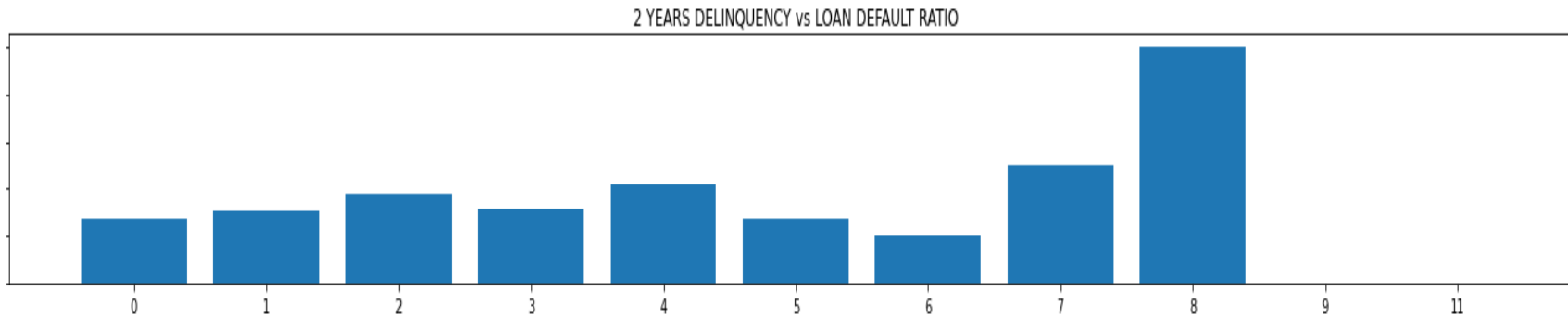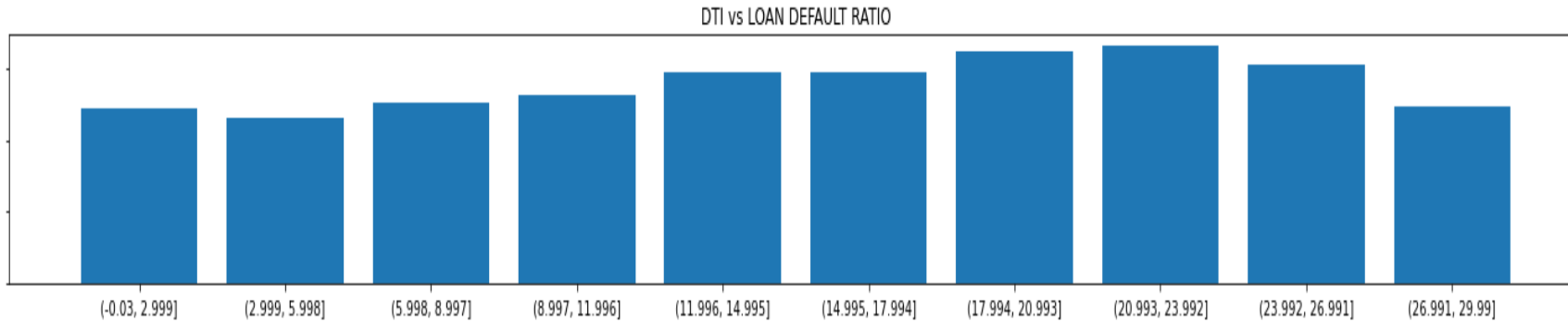
# Bivariate Analysis - Continuous



- Higher Loan Amount, Funded Amount, Investor Funded Amount result in higher default.

# Bivariate Analysis – Continuous(contd.)



INTREST RATE vs LOAN DEFAULT RATIO
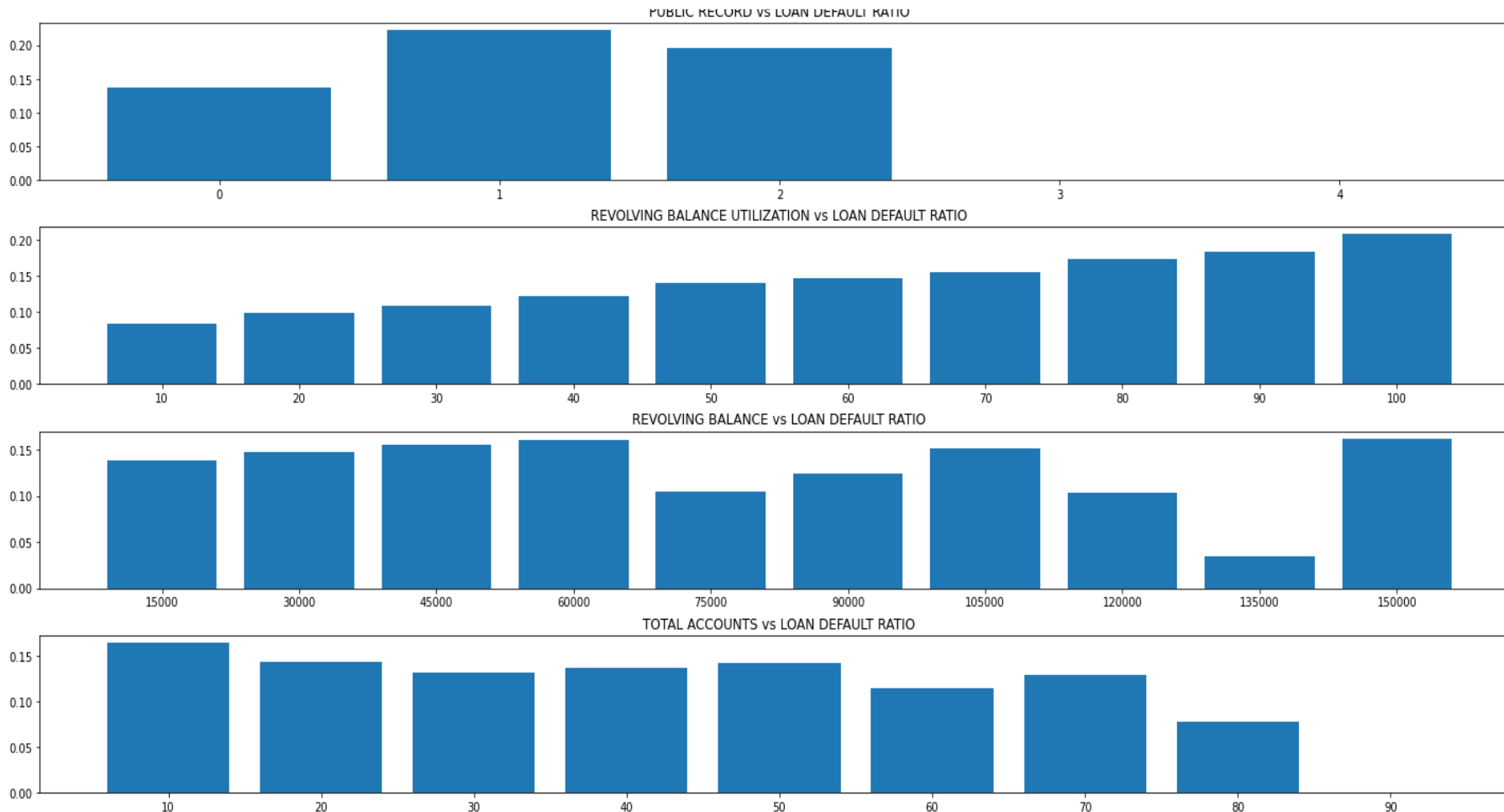INSTALLMENT vs LOAN DEFAULT RATIO
INCOME vs LOAN DEFAULT RATIO

- Loan Default Rate increases with interest rate
- Installment and loan default ratio show neutral relation
- Higher Income Groups have negligible loan defaults

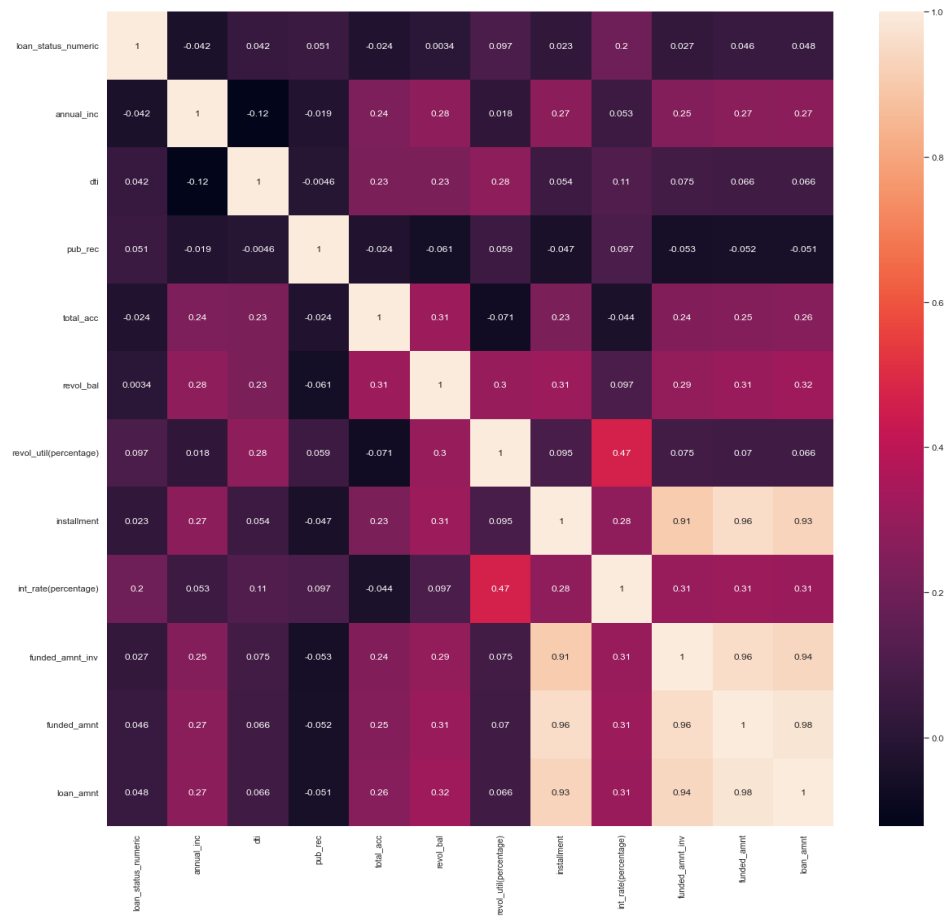# Bivariate Analysis – Continuous(contd.)



- DTI, 2 years delinquency and inquiries in 6 month doesn't show any relation with loan default ratio

# Bivariate Analysis – Continuous(contd.)



- Loan Default Rate increases with Revolving Balance Utilization.

- Public Record, Revolving Balance, Total accounts show no relation.

# Multivariate Analysis



Correlation matrix show higher positive correlation with interest rate, grade, subgrade and revolving balance utilization

# Conclusion

Most important deciding factors according to analysis are

- Loan Amount, Funded Amount, Investor Funded Amount
- Interest Rate
- Term
- Grade and Subgrade
- Revolving Balance Utilization