

A decorative graphic on the left side of the slide, consisting of a network of light blue lines and circles of varying sizes, resembling a circuit board or a neural network diagram.

# LENDING CLUB EDA CASE STUDY

BY

SUJEET KUMAR GOIT

VENUGOPAL PAMIDIPATI

# PROBLEM STATEMENT

- Financial companies faces the critical task of assessing loan applications of urban customers to make informed decisions. When evaluating loan applications, the company confronts two primary risks: denying a loan to a creditworthy applicant, resulting in a loss of business, or approving a loan to a borrower likely to default, thereby risking financial losses.
- The objective is to analyze the past loan applicant's data 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.

# ANALYSIS APPROACH

The approach is to perform the following activities in sequence to identify the patterns for defaults:-

- Data Cleaning
- Univariate Analysis
- Bivariate Analysis

# DATA CLEANING

Following Data Cleaning sequence applied on the data

1. Analyse dataset
  - Check columns for missing values
  - Process Data - Drop columns where more than 40% data is missing
2. Handling Incorrect data types
  - Remove special characters (%) from data
  - Convert data type to integer or float
3. Impute Missing Values
  - Remove rows where data is missing (less than 1% data) - As removing this will have no impact on the data set
4. Sanity checks
  - Check if Data is within the range as expected.
5. Outlier finding for numerical Data
  - Check for outlier data and process dataset by removing
6. Derived metrics
  - Create derived metrics from dataset, such as based on Annual Income, Interest rate, Loan amount

The background is a blue gradient. In the corners, there are decorative white lines resembling circuit traces or data paths, with small circles at various points.

# UNIVARIATE ANALYSIS

# APPROACH

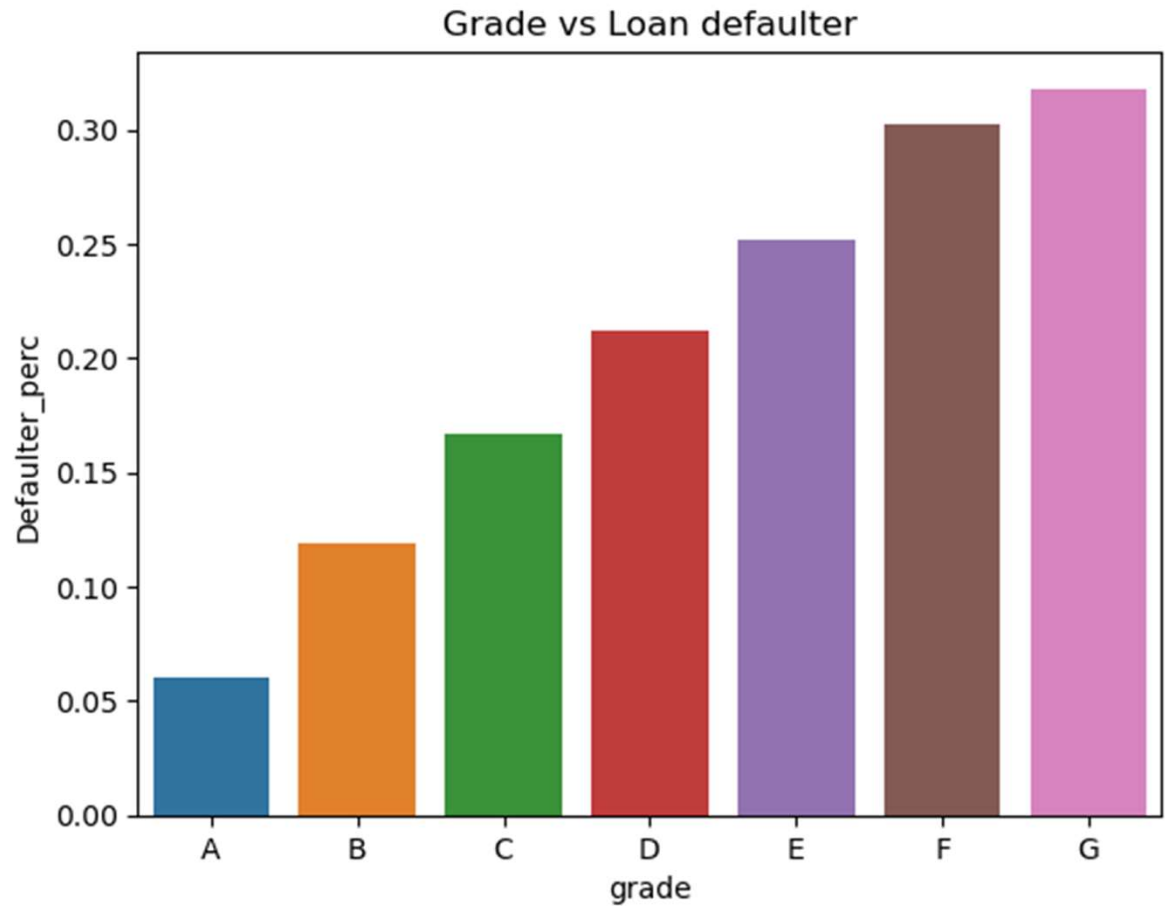
Grouping is done for Grade and Loan Status variables. Derived Defaulter\_perc as new variable in the group that shows percentages of defaulters for each Grade. By using bar plot, Segmented univariate analysis is done on defaulters across different segmented grades.

x-axis : grade

y-axis : Defaulter\_perc (Percentage of defaulters in each Grade)

# INFERENCE

In the bar plot, as the grade moves along the x-axis from A to G, the probability of default increases.



# APPROACH

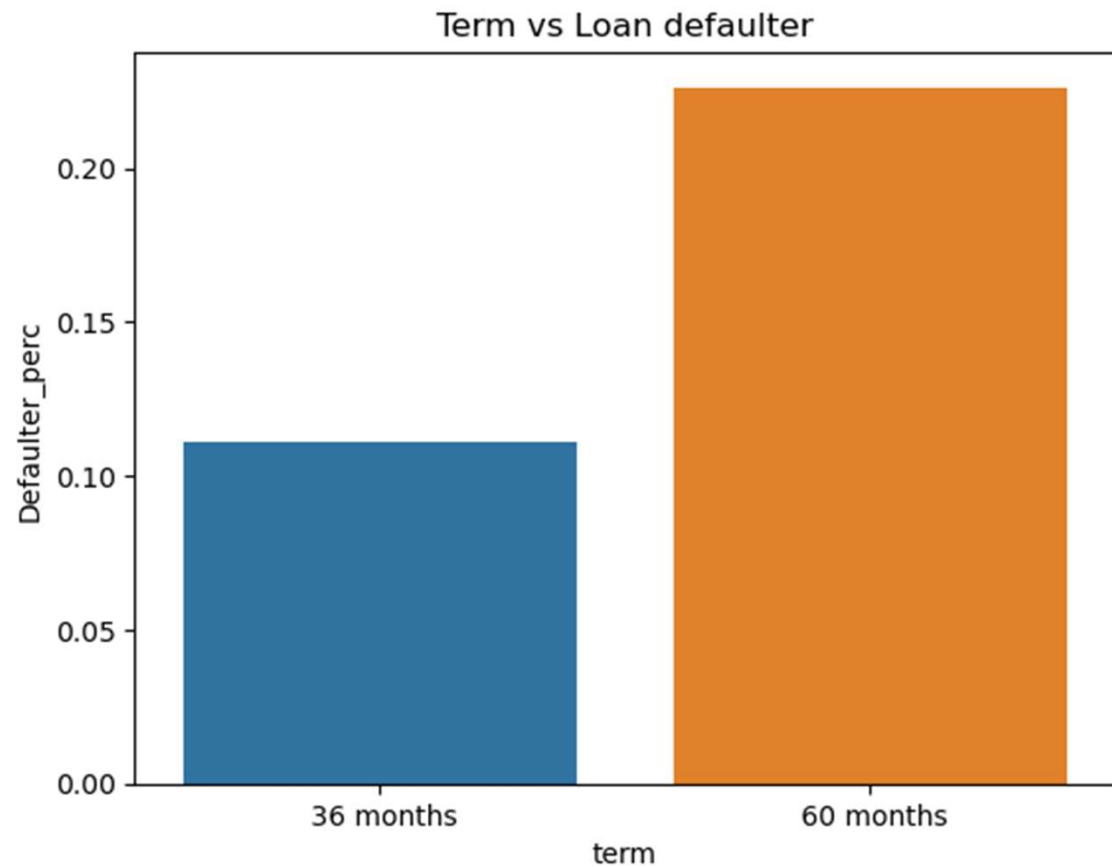
Grouping is done for term and Loan Status variables. Derived Defaulter\_perc as new variable in the group that shows percentages of defaulters for each Term. By using bar plot, Segmented univariate analysis is done on defaulters across different segmented terms.

x-axis : term

y-axis : Defaulter\_perc (Percentage of defaulters in each Term)

# INFERENCE

In the bar plot, the probability of default increases for loan having 60 months term.



# APPROACH

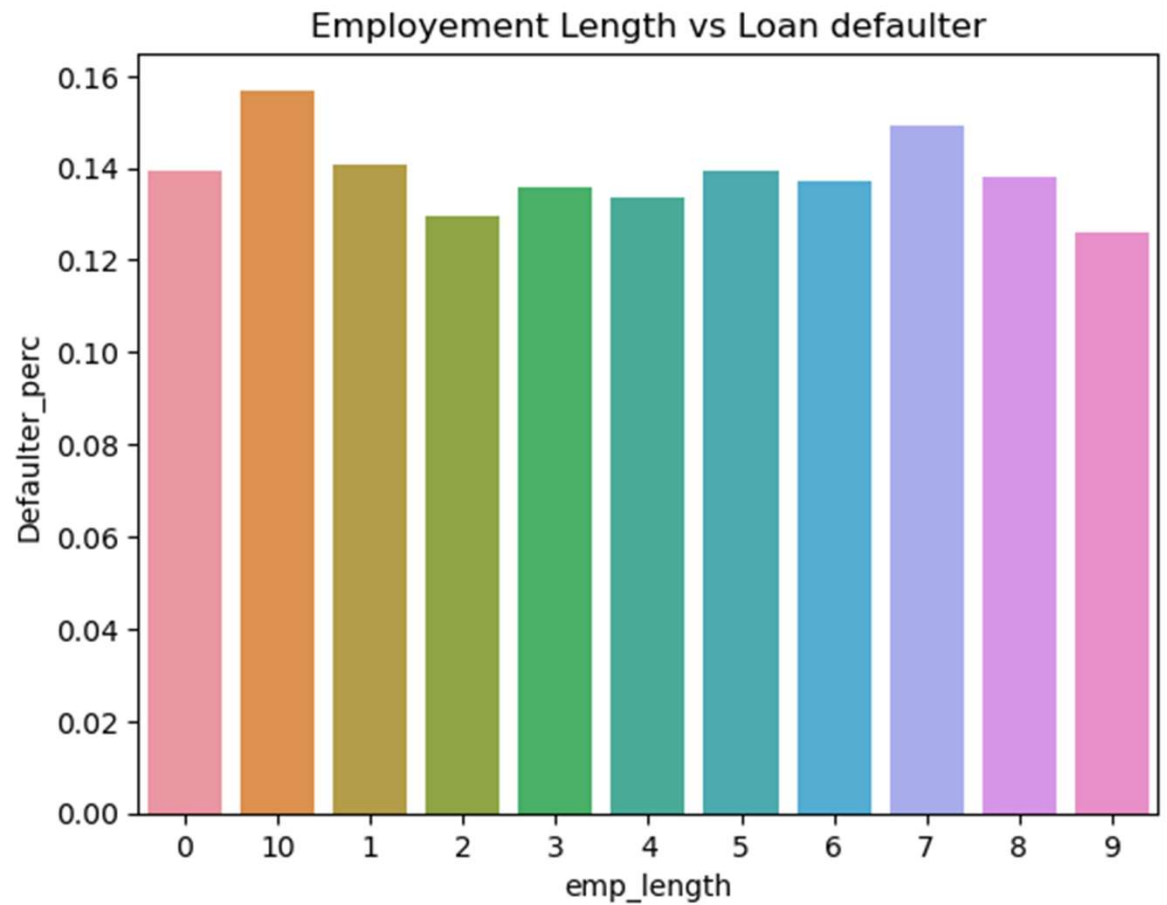
Grouping is done for Employee Length and Loan Status variables. Derived Defaulter\_perc as new variable in the group that shows percentages of defaulters for each Employee Length. By using bar plot, Segmented univariate analysis is done on defaulters across different segmented Employee Length.

x-axis : emp\_length (Employee Length)

y-axis : Defaulter\_perc (Percentage of defaulters in each Employee Length)

# INFERENCE

In the bar plot, No inference is concluded as all are almost in equal height.





# APPROACH

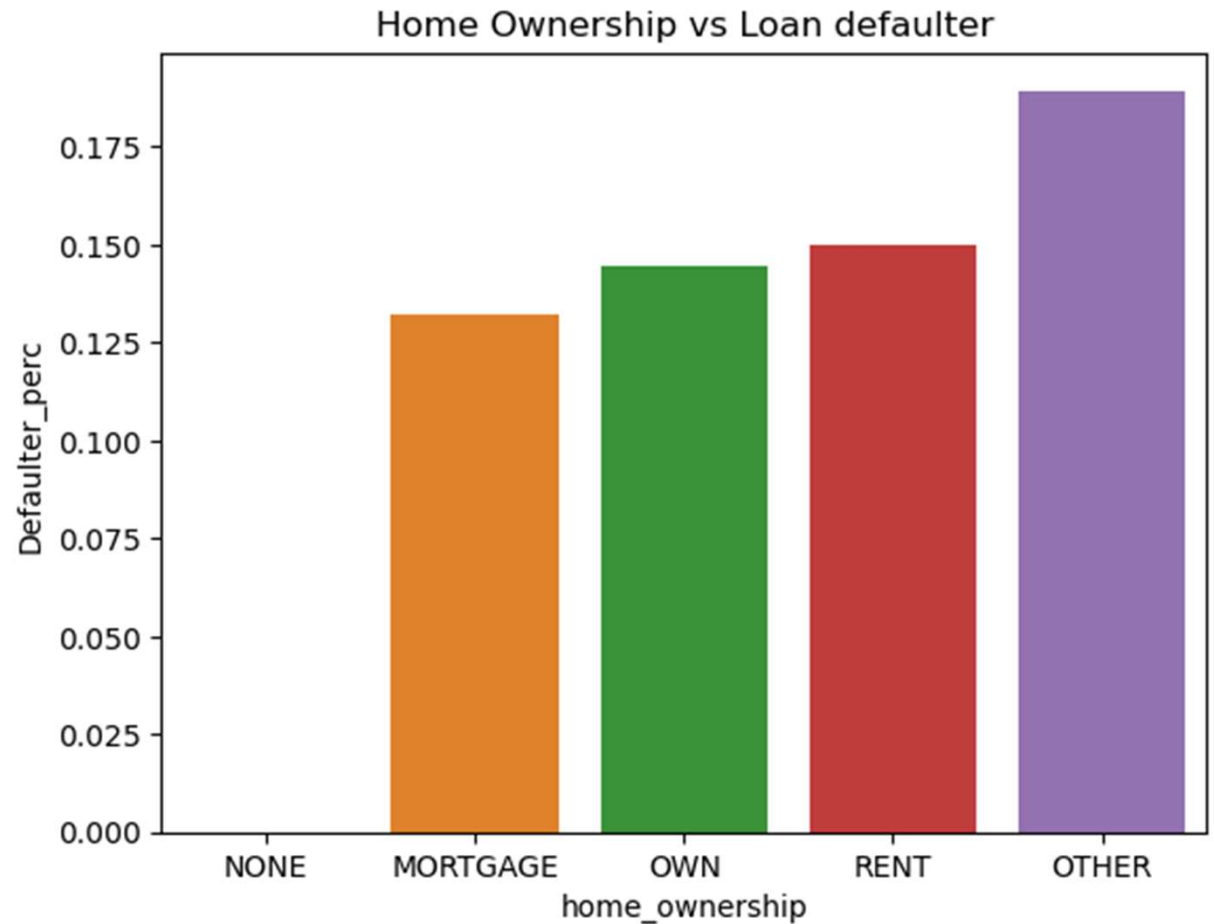
Grouping is done for Home Ownership and Loan Status variables . Derived Defaulter\_perc as new variable in the group that shows percentages of defaulters for each Home Ownership. By using bar plot, Segmented univariate analysis is done on defaulters across different segmented Home Ownership.

x-axis : home\_ownership (Home Ownership)

y-axis : Defaulter\_perc (Percentage of defaulters in each Home Ownership)

# INFERENCE

In the bar plot, The probability of default increases for Home Ownership as OTHER.



# APPROACH

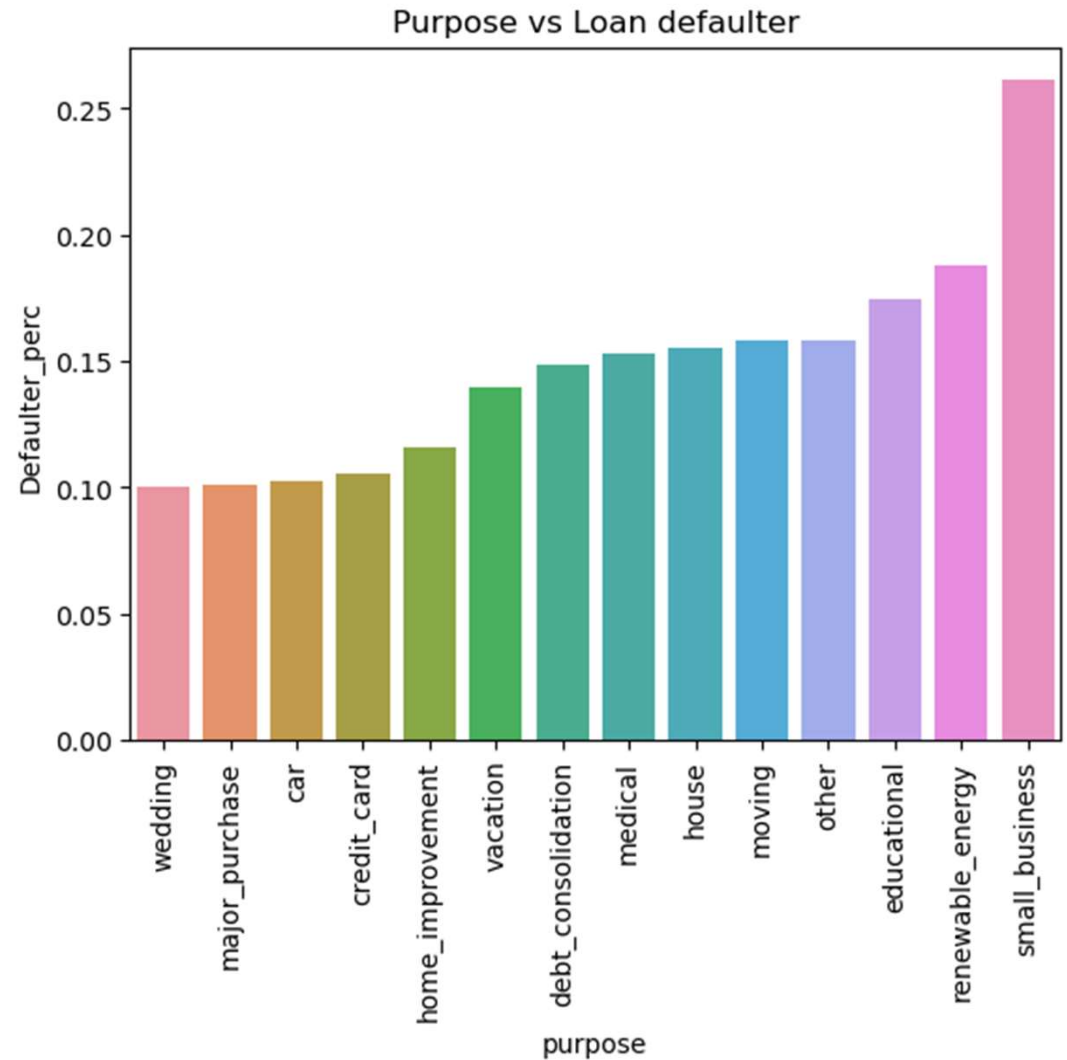
Grouping is done for Purpose and Loan Status variables. Derived Defaulter\_perc as new variable in the group that shows percentages of defaulters for each Purpose. By using bar plot, Segmented univariate analysis is done on defaulters across different segmented Purpose.

x-axis : purpose

y-axis : Defaulter\_perc (Percentage of defaulters in each Purpose)

# INFERENCE

In the bar plot, Small businesses are considered riskier, as their default percentage is notably high.



# APPROACH

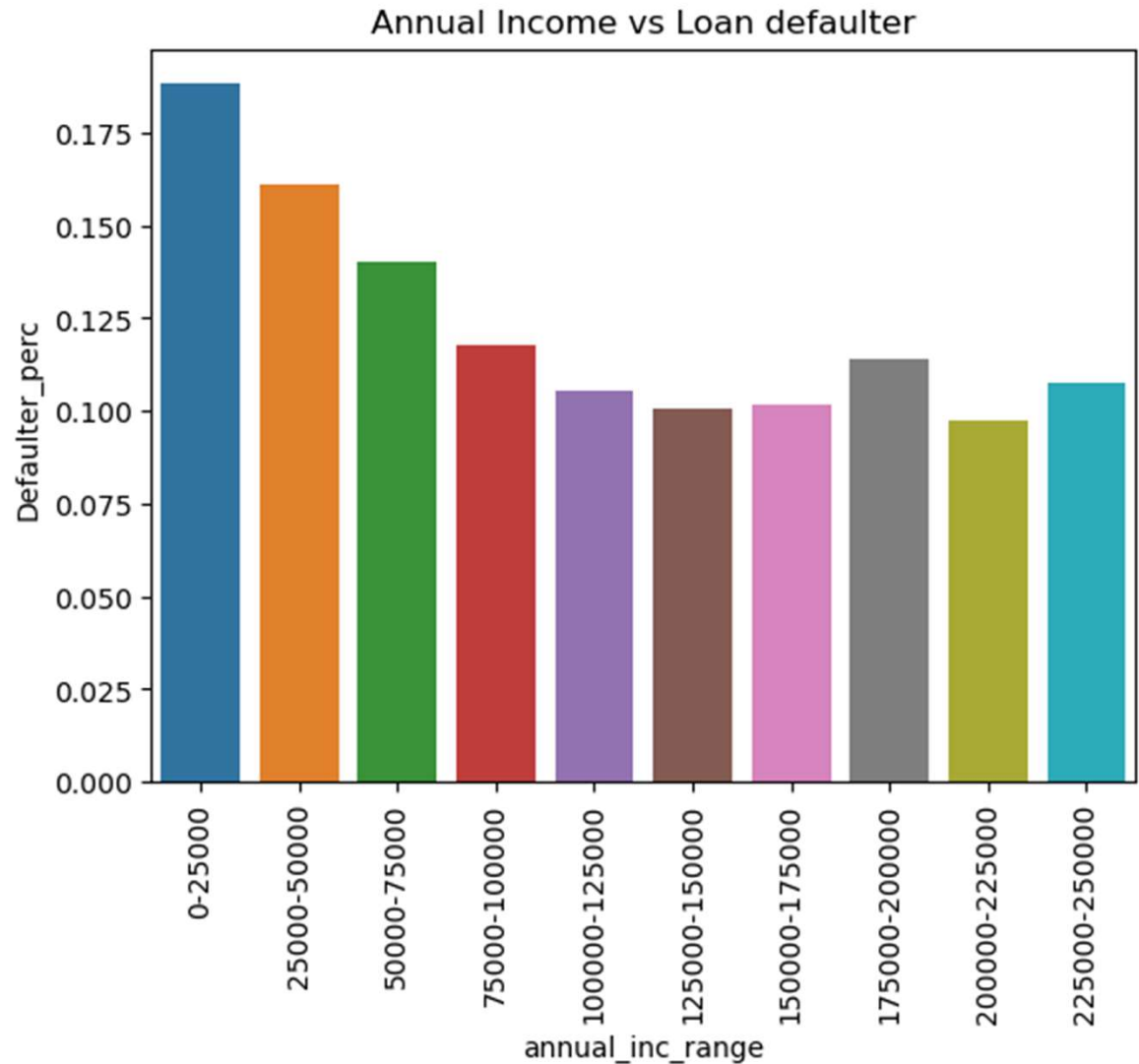
Grouping is done for Annual Income Range and Loan Status variables . Derived Defaulter\_perc as new variable in the group that shows percentages of defaulters for each Annual Income Range. By using bar plot, Segmented univariate analysis is done on defaulters across different segmented Annual Income Range.

x-axis : annual\_inc\_range (Annual Income Range)

y-axis : Defaulter\_perc (Percentage of defaulters in each Annual Income Range)

# INFERENCE

In the bar plot, There are almost 18% defaulters whose Annual Income Range is less than or equal to 25000 and 16% defaulters for Annual Income Range between 25000 – 50000. Bank should be more vigilant on giving loan to customers whose Annual Income Range is less than or equal to 50000.



# APPROACH

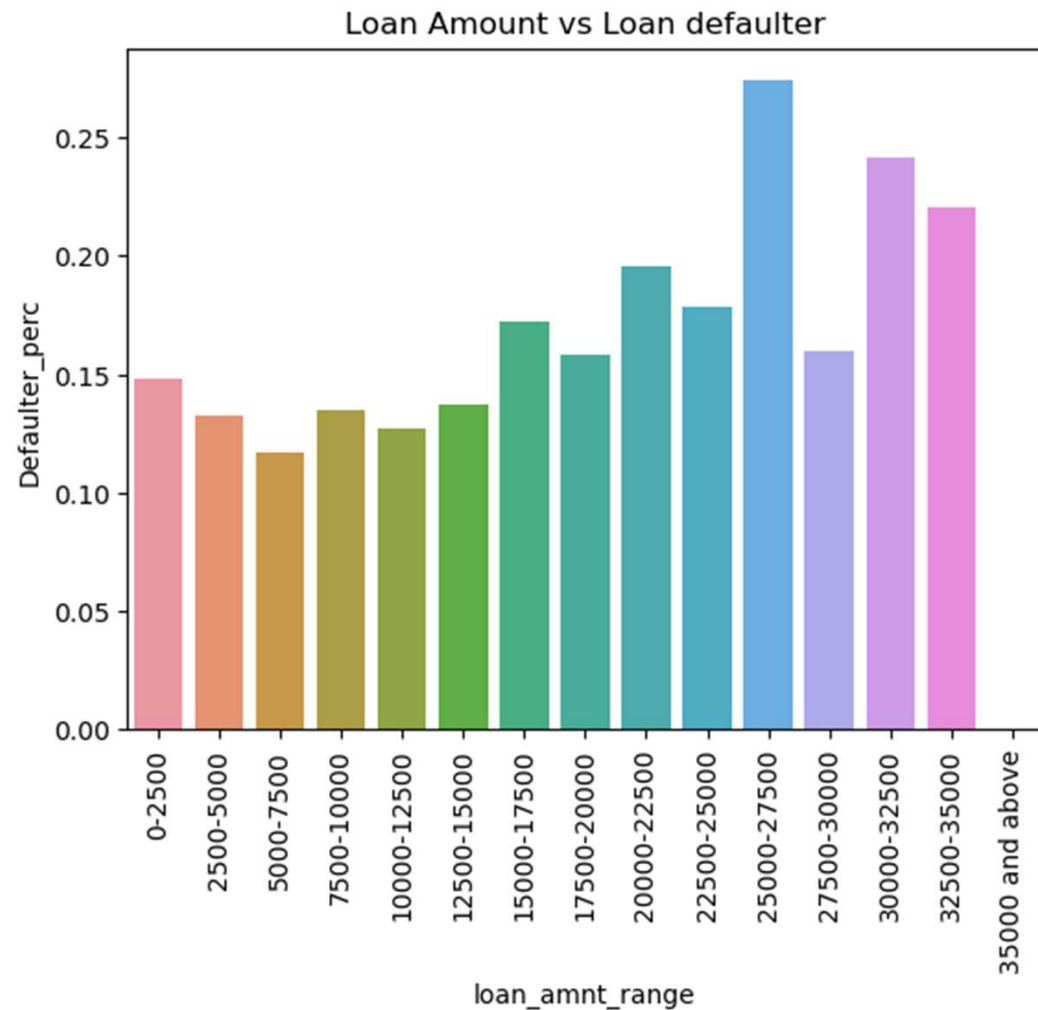
Grouping is done for Loan Amount and Loan Status variables . Derived Defaulter\_perc as new variable in the group that shows percentages of defaulters for each Loan Amount Range. By using bar plot, Segmented univariate analysis is done on defaulters across different segmented Loan Amount Range.

x-axis : loan\_amnt\_range (Loan Amount Range)

y-axis : Defaulter\_perc (Percentage of defaulters in each Loan Amount Range)

# INFERENCE

In the bar plot, More defaulters are observed in the loan amount range of 25000 – 27500 followed by 30000 and above.



# APPROACH

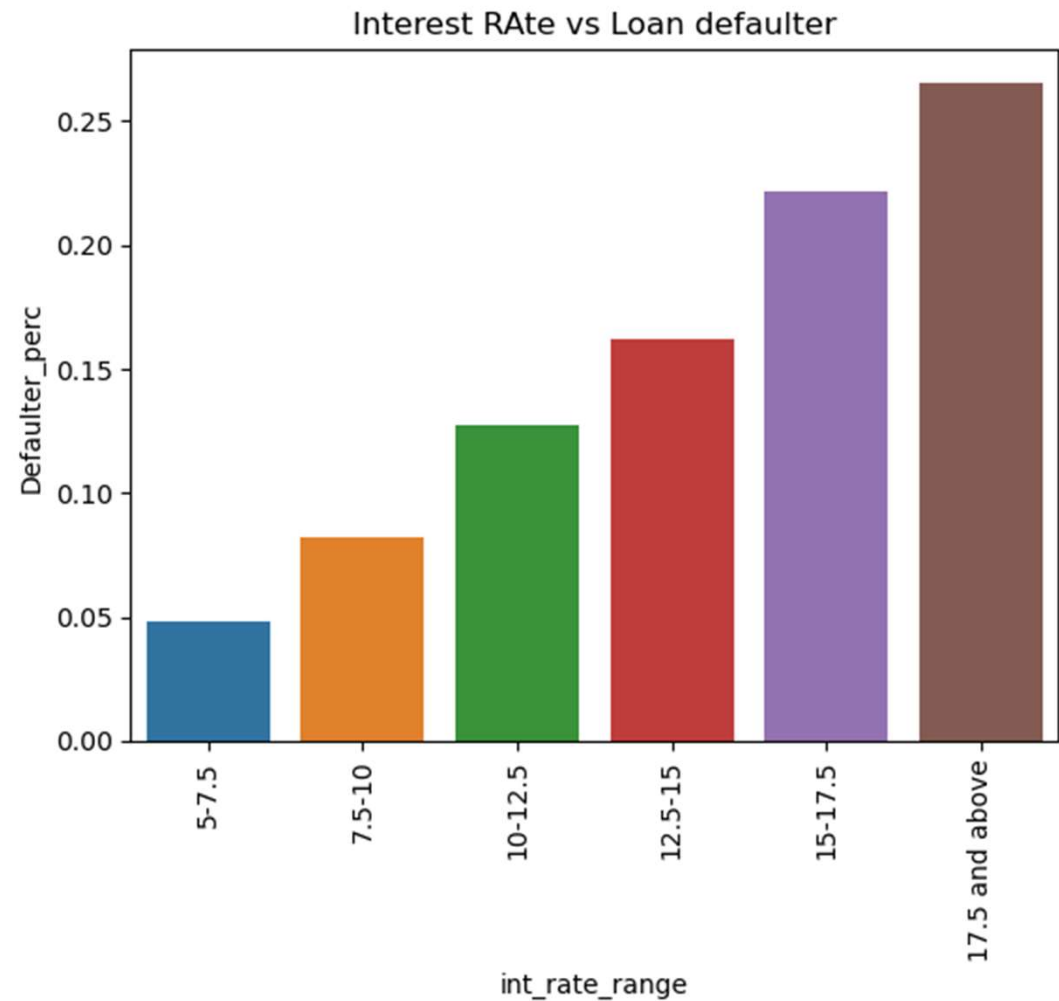
Grouping is done for Interest Rates and Loan Status variables. Derived Defaulter\_perc as new variable in the group that shows percentages of defaulters for each Interest Rates Range. By using bar plot, Segmented univariate analysis is done on defaulters across different segmented Interest Rates Range.

x-axis : int\_rate\_range (Interest Rates Range)

y-axis : Defaulter\_perc (Percentage of defaulters in each Interest Rates Range)

# INFERENCE

In the bar plot, Its likely that the number of defaulters will increase as interest rates rise.



# APPROACH

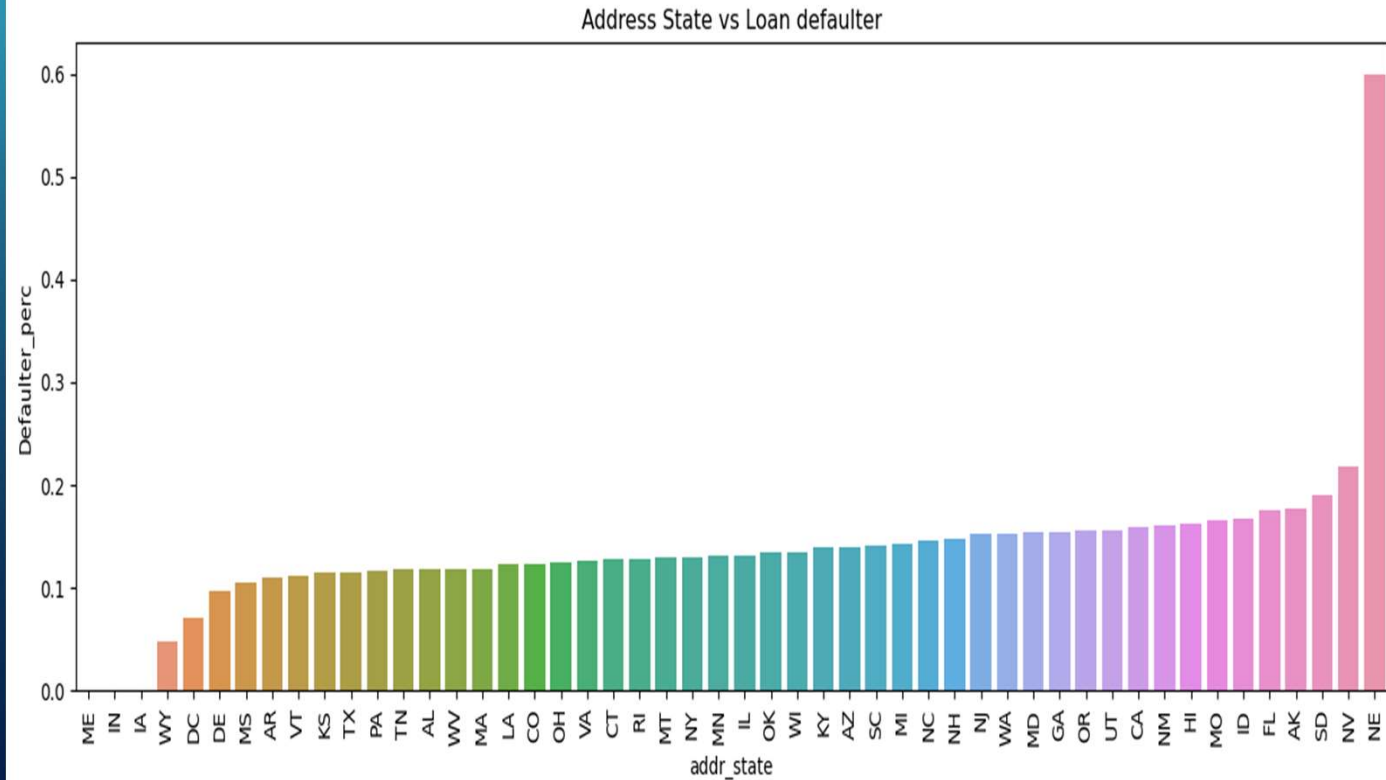
Grouping is done for Address State and Loan Status variables. Derived Defaulter\_perc as new variable in the group that shows percentages of defaulters for each Address State. By using bar plot, Segmented univariate analysis is done on defaulters across different segmented Address State.

x-axis : addr\_state (Address State)

y-axis : Defaulter\_perc (Percentage of defaulters in each Address State)

# INFERENCE

In the bar plot, NE state is showing high percentage of defaulters, but the total customers count is very minimal, so its not considerable. The NV state can be taken into consideration for having a high number of defaulters, as this state has reasonable amount of loan customers.



# APPROACH

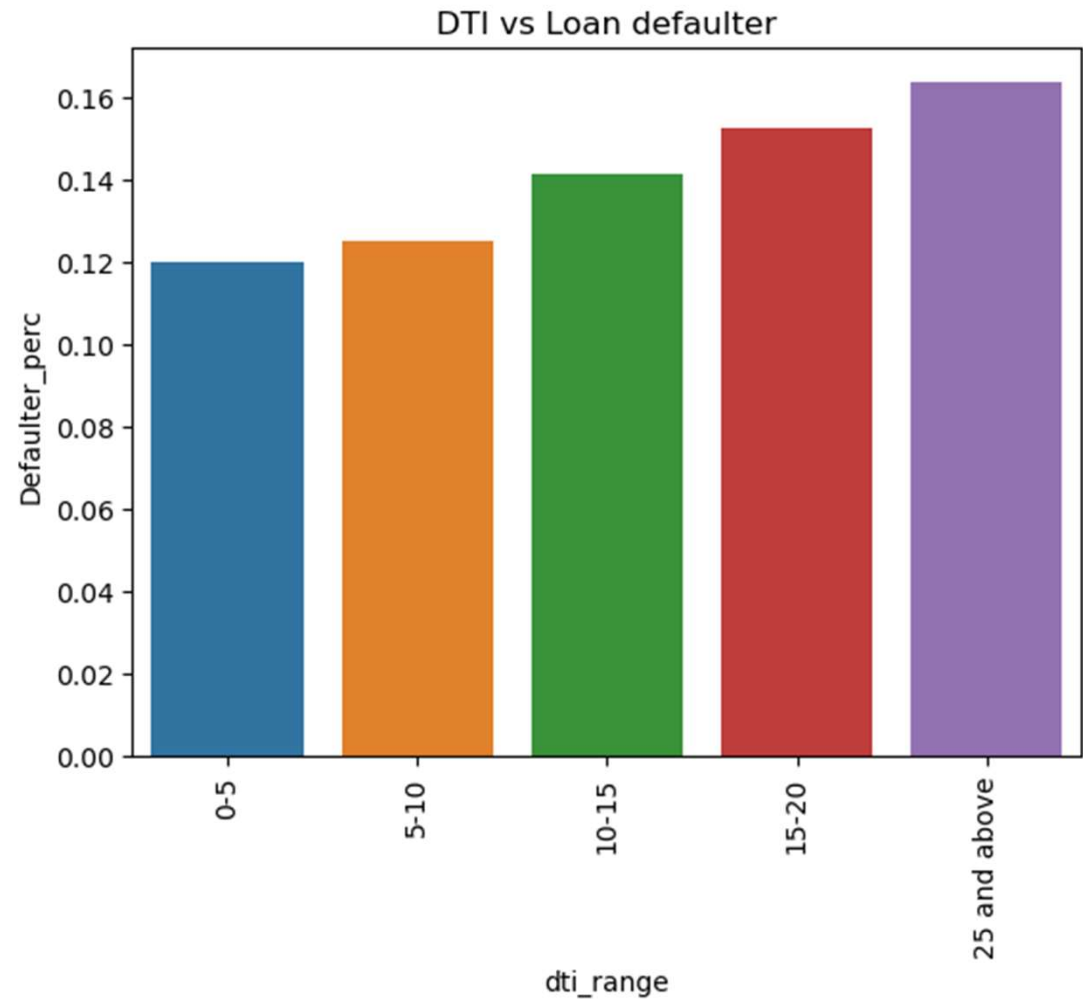
Grouping is done for DTI and Loan Status variables. Derived Defaulter\_perc as new variable in the group that shows percentages of defaulters for each DTI range. By using bar plot, Segmented univariate analysis is done on defaulters across different segmented DTI Range.

x-axis : dti\_range (dti range)

y-axis : Defaulter\_perc (Percentage of defaulters in each dti range)

# INFERENCE

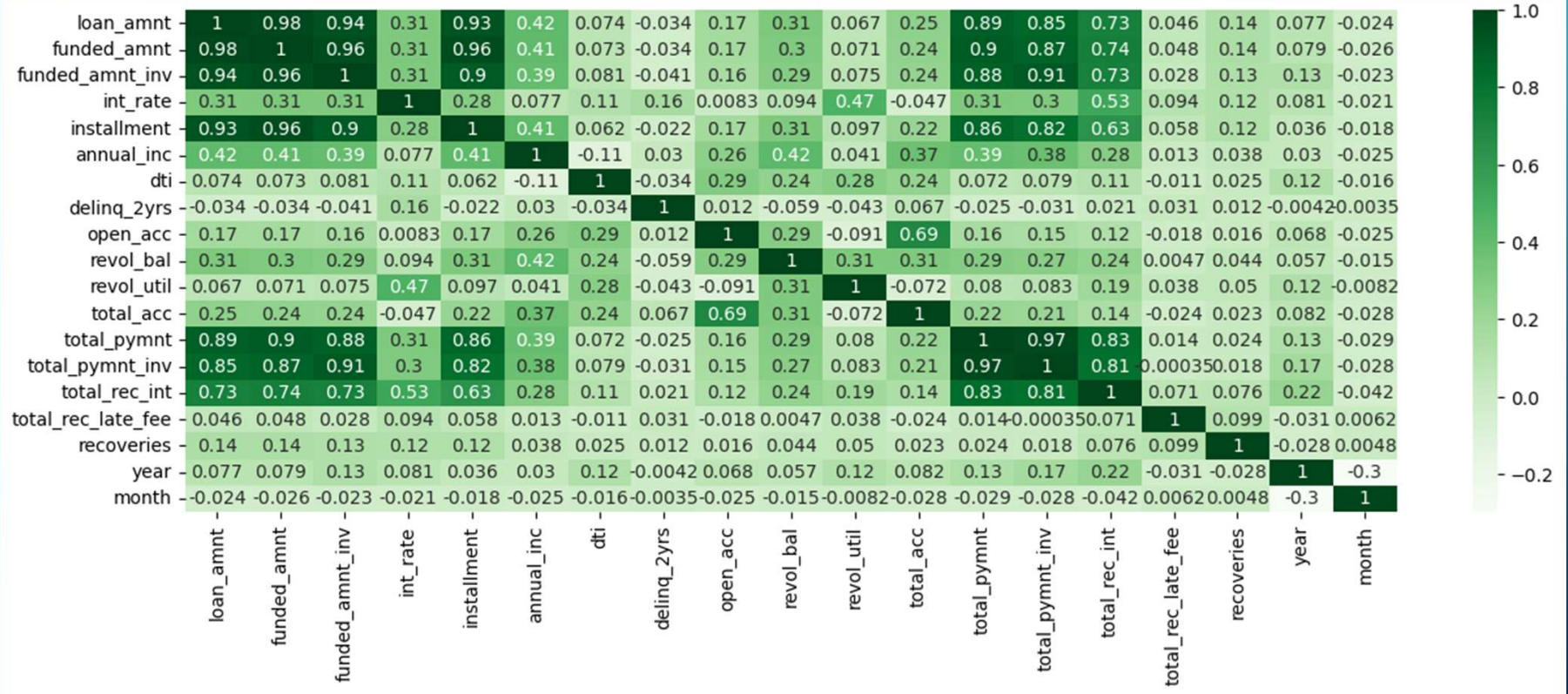
In the bar plot, Its likely that number of defaulters will increase as DTI value rise.



The background is a blue gradient with decorative circuit-like lines in the corners. These lines are composed of small circles connected by straight segments, resembling a stylized electronic circuit or data flow diagram.

# BIVARIATE ANALYSIS AMONG NUMERICAL VARIABLES





# CORRELATION AMONG NUMERICAL VARIABLES

Points to be concluded from the graph presented before.

- Loan Amount, Funded Amount, Funded Amount Inv, Total Payment, Total Payment Inv, Total Received Interest and Installments are formed a cluster of positive correlation. Means, all these variables have same kind of relation with Loan Status.

The background is a blue gradient with decorative circuit-like lines in the corners. These lines are composed of small circles connected by straight lines, resembling a stylized electronic circuit board. They are located in the top-left, top-right, bottom-left, and bottom-right corners.

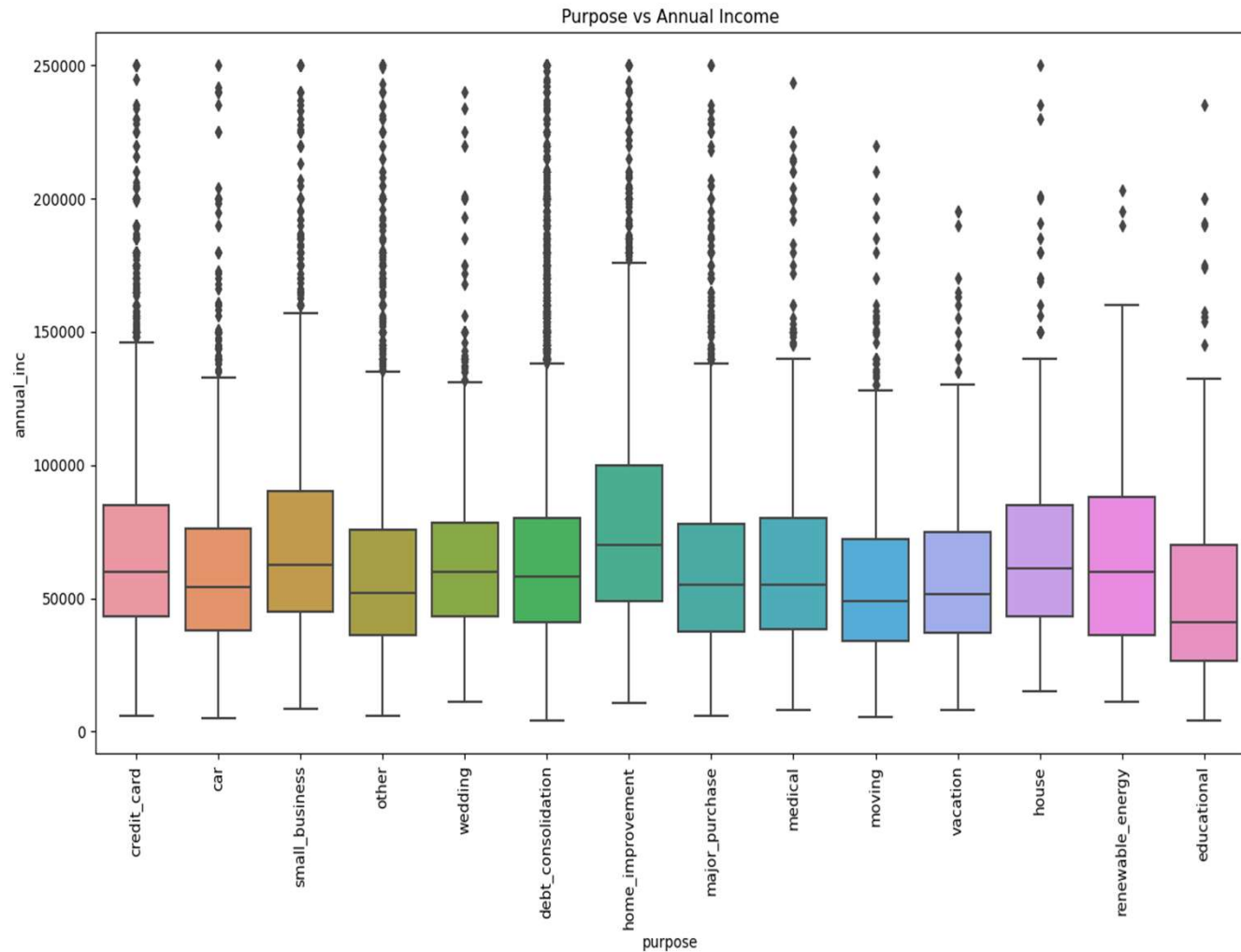
# BIVARIATE ANALYSIS BETWEEN CATEGORICAL VARIABLES

# APPROACH

IN the previous slides, Segmented Univariate Analysis showcased Small Business purpose has high default percentage.  
This Box Plot is used for Bivariate Analysis to showcase the relation between Purpose and other numerical variable like Annual Income.

# INFERENCE

IN the Box Plot, Its likely that Highly annual income people having small business purpose loan will fall into more defaults compare to lower annual income customers

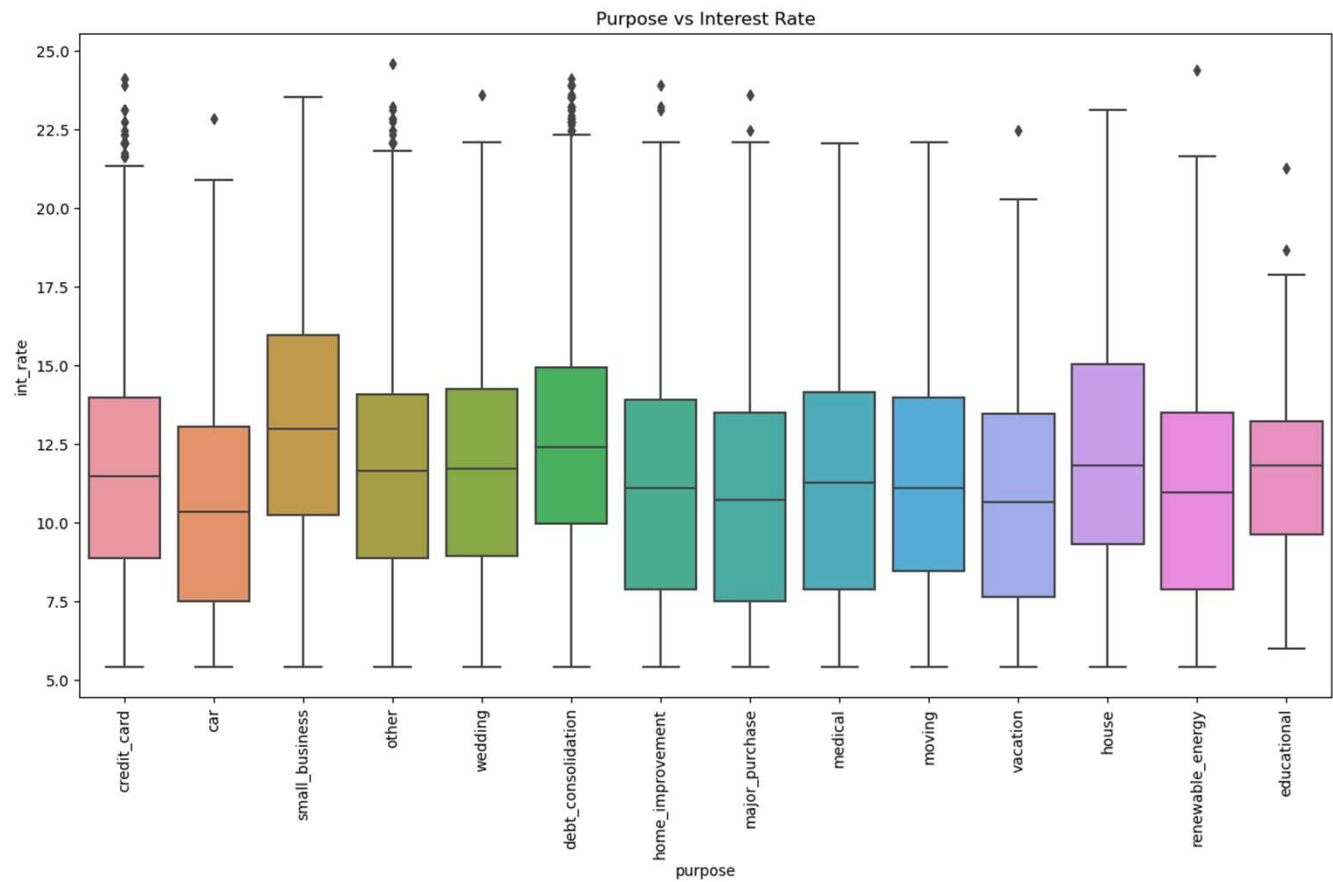


# APPROACH

IN the previous slides, Segmented Univariate Analysis showcased Small Business purpose has high default percentage. This Box Plot is used for Bivariate Analysis to showcase the relation between Purpose and other numerical variable like Interest rate.

# INFERENCE

IN the Box Plot, Its likely that Small Business purpose loans having High interest rates would fall into more defaults.

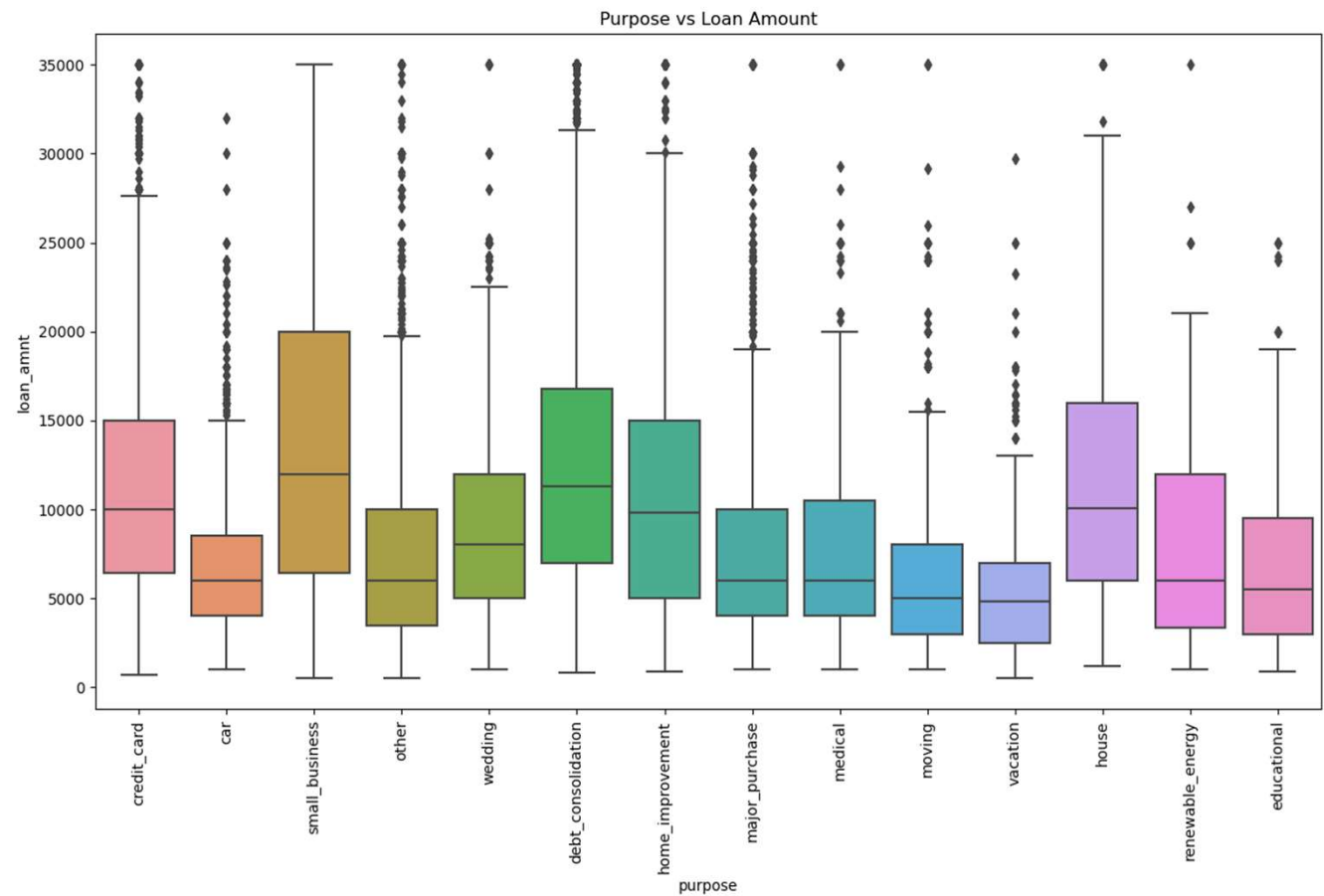


# APPROACH

IN the previous slides, Segmented Univariate Analysis showcased Small Business purpose has high default percentage. This Box Plot is used for Bivariate Analysis to showcase the relation between Purpose and other numerical variable like Loan Amounts.

# INFERENCE

IN the Box Plot, Its likely Small Business purpose loans having High loan amounts would fall into more defaults.



## CONCLUSION

1. Small business are the riskiest to give loan to.
2. Small business should be given loan at lower Interest Rate or less Loan amount, as these will reduce the risk of default.
3. As the Grades move from A to G, default is increasing. Bank can focus of this parameter.
4. Shorter duration loans(36 months) are paid more often



THANK YOU