

LENDING CLUB CASE STUDY



Submitted by

Kajal Singha Mahata

Krishnaveni Poluru

PROBLEM STATEMENT

- Minimizing financial losses from loan approval process
 - Losses occur when borrowers default on loans
- Objective: Reduce credit losses by identifying risky applicants
 - Approving loans for likely-to-repay applicants generates profit
 - Approving loans for likely-to-default applicants results in losses
- Exploratory Data Analysis to understand driving factors behind loan default
 - Knowledge used for portfolio and risk assessment

AGENDA

- Introduction
- Problem Statement
- Data Understanding
- Data Cleaning & Pre-processing
- Univariate Analysis
- Bivariate Analysis
- Multivariate Analysis
- Correlation Analysis
- Suggestions
- References & Useful Links
- Conclusion



DATA UNDERSTANDING, CLEANING, AND PRE-PROCESSING

LoanStatNew	Description
acc_now_delinq	The number of accounts on which the borrower is now delinquent.
acc_open_past_24mths	Number of trades opened in past 24 months.
addr_state	The state provided by the borrower in the loan application
all_util	Balance to credit limit on all trades
annual_inc	The self-reported annual income provided by the borrower during registration.
annual_inc_joint	The combined self-reported annual income provided by the co-borrowers during registration
application_type	Indicates whether the loan is an individual application or a joint application with two co-borrowers
avg_cur_bal	Average current balance of all accounts
bc_open_to_buy	Total open to buy on revolving bankcards.
bc_util	Ratio of total current balance to high credit/credit limit for all bankcard accounts.
chargeoff_within_12_mths	Number of charge-offs within 12 months
collection_recovery_fee	post charge off collection fee
collections_12_mths_ex_med	Number of collections in 12 months excluding medical collections
delinq_2yrs	The number of 30+ days past-due incidences of delinquency in the borrower's credit file for the past 2 years
delinq_amnt	The past-due amount owed for the accounts on which the borrower is now delinquent.
desc	Loan description provided by the borrower
dti	A ratio calculated using the borrower's total monthly debt payments on the total debt obligations, excluding mortgage and the requested LC loan, divided by the borrower's self-reported monthly income.
dti_joint	A ratio calculated using the co-borrowers' total monthly payments on the total debt obligations, excluding mortgages and the requested LC loan, divided by the co-borrowers' combined self-reported monthly income
earliest_cr_line	The month the borrower's earliest reported credit line was opened
emp_length	Employment length in years. Possible values are between 0 and 10 where 0 means less than one year and 10 means ten or more years.
emp_title	The job title supplied by the Borrower when applying for the loan.*
fico_range_high	The upper boundary range the borrower's FICO at loan origination belongs to.
fico_range_low	The lower boundary range the borrower's FICO at loan origination belongs to.
*del_amt	The total amount committed to that loan at that point in time.
*amt_inv	The total amount committed by investors for that loan at that point in time.
	LC assigned loan grade
*ship	The home ownership status provided by the borrower during registration. Our values are: RENT, OWN, MORTGAGE, OTHER.

DATA DESCRIPTION: IMAGE 1

DATA UNDERSTANDING

- Primary Attribute: Loan Status
 - Fully-Paid: Customers who have successfully repaid their loans
 - Charged-Off: Customers who have defaulted on their loans
 - Current: Customers whose loans are presently in progress
- Decision Matrix: Loan Acceptance Outcome
 - Fully Paid: Applicants who have successfully repaid both the principal and the interest rate of the loan
 - Current: Applicants in the process of making loan installments
 - Charged-off: Applicants who have failed to make timely installments
- Key Columns of Significance
 - Customer Demographics: Annual Income, Home Ownership, Employment Length, Debt to Income, State
- Excluded Columns: Customer Behavior Columns

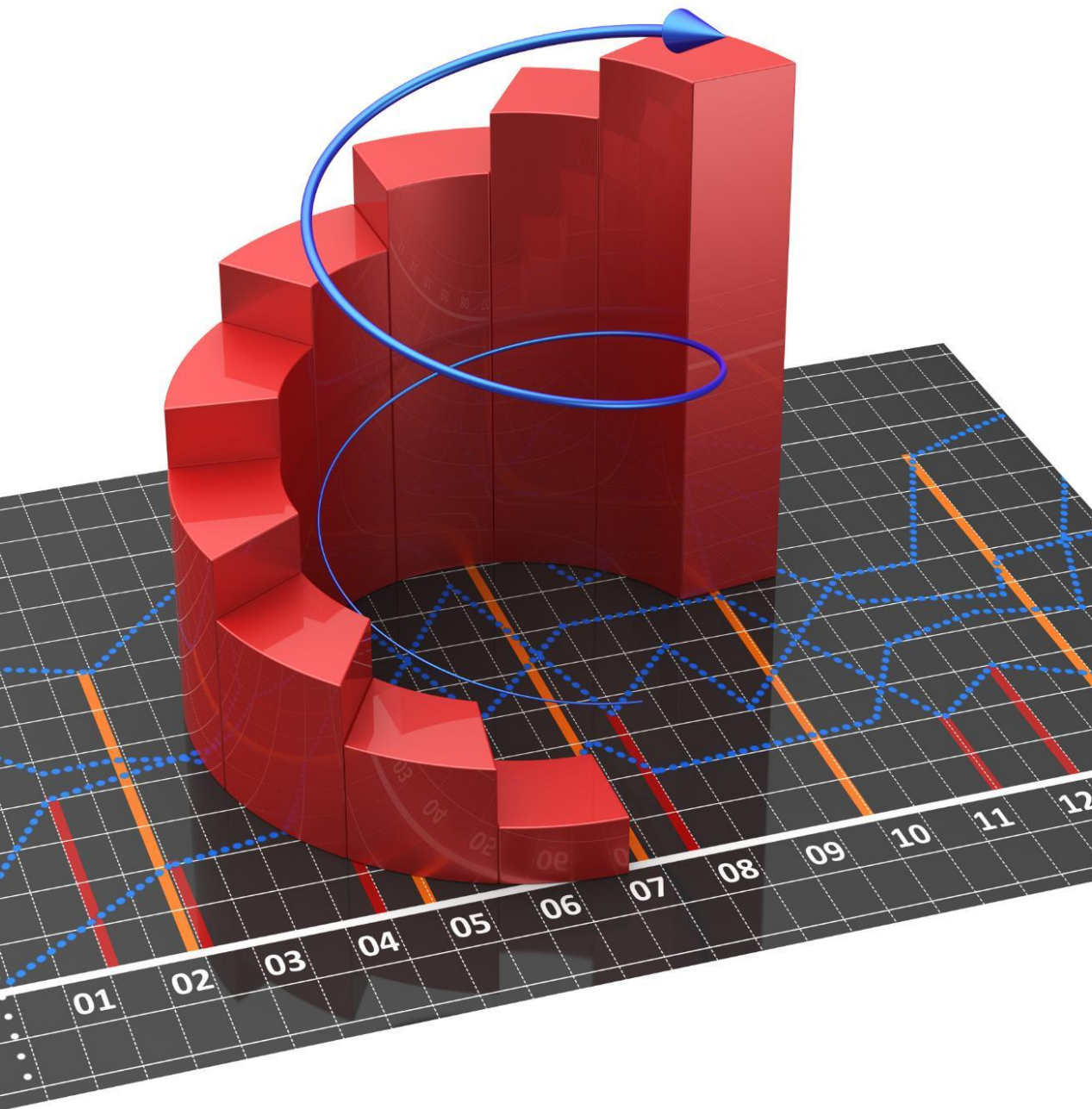
DATA UNDERSTANDING

- Granular Data
 - Columns with excessive detail will be omitted
 - Example: 'sub grade' column
- Columns with NA values
 - 54 columns contain only NA values
 - These columns will be removed
- Columns with only 0 values
 - These columns will also be dropped

DATA CLEANING & PRE-PROCESSING



- Loading data from loan CSV
 - Conversion of mixed data types
- Checking for null values in the dataset
 - 48% of columns with null values were dropped
- Checking for unique values
 - 9 columns with single unique values were removed
- Checking for duplicated rows in data
- Dropping Records & Columns
- Common Functions
- Data Conversion
- Outlier Treatment
- Imputing values in Columns



DATA CLEANING & PRE-PROCESSING

- Null Values for pub_rec_bankruptcies
 - 660 null values dropped
 - Cannot be imputed
- Post Data Cleaning and Preprocessing
 - 36094 rows × 18 columns left



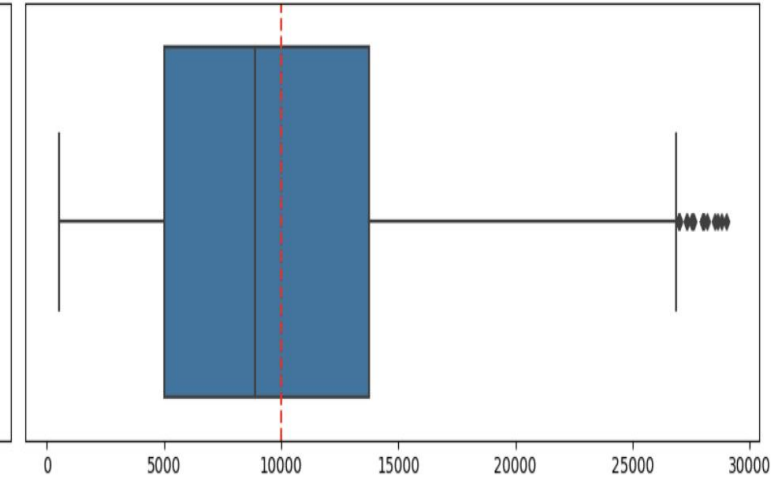
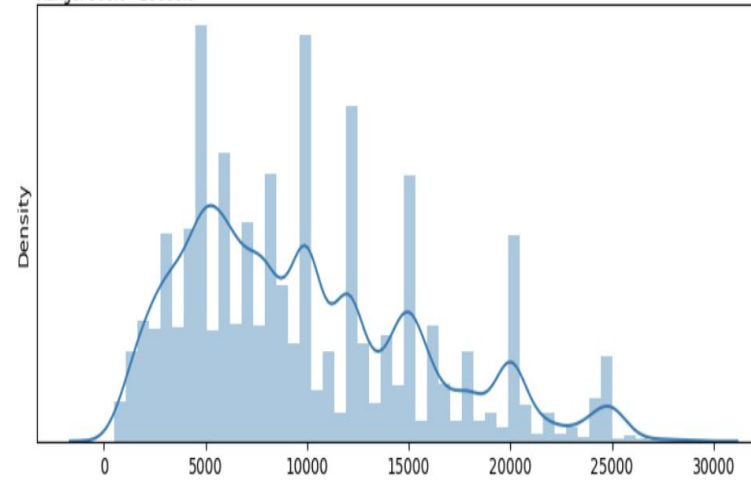
UNIVARIATE ANALYSIS

UNIVARIATE ANALYSIS:

1.Univariate analysis of Loan Amount

Most values between 5000.0 and 13750.0

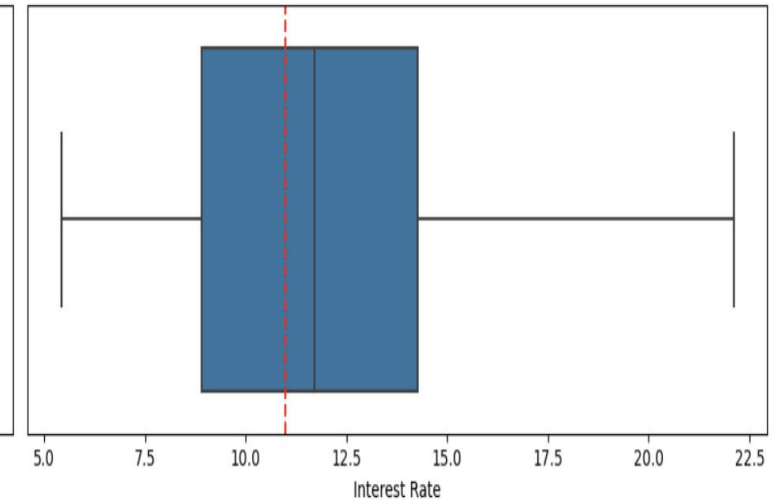
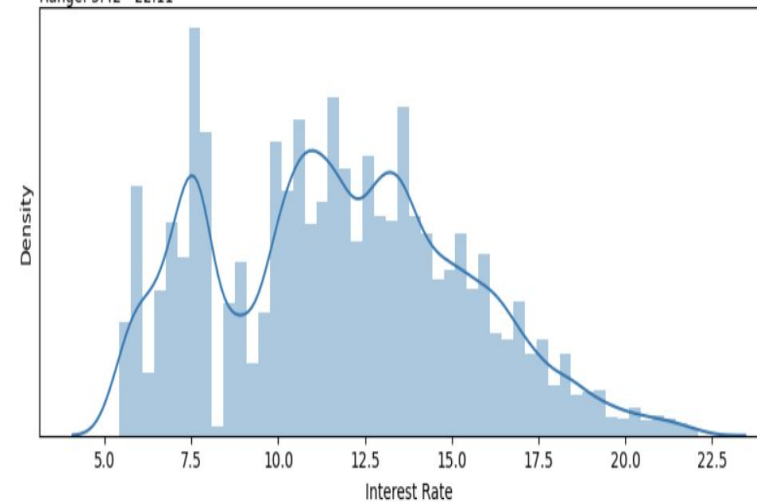
Range: 500.0 - 29000.0



2.Univariate analysis of Interest Rate

Most values between 8.9 and 14.26

Range: 5.42 - 22.11



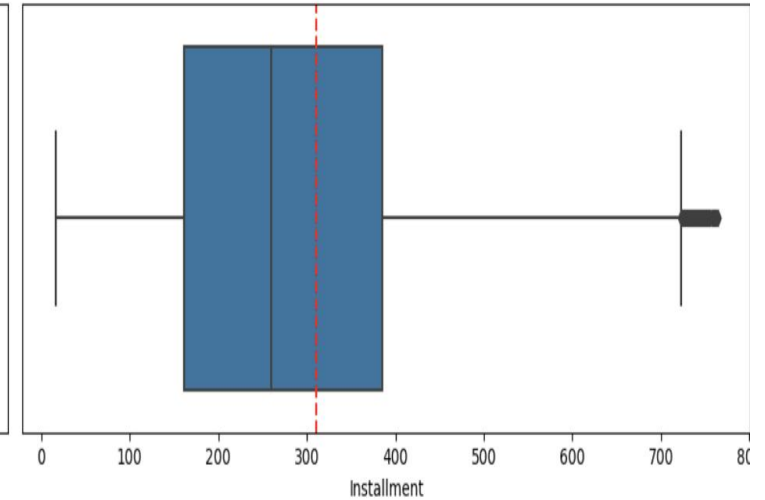
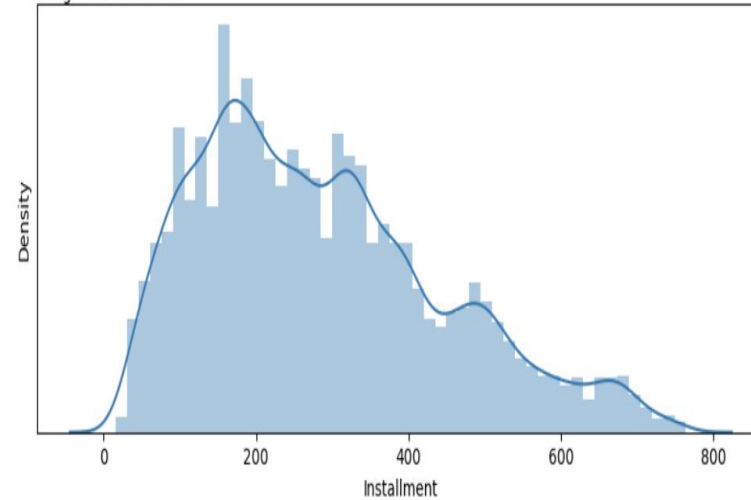
3.Univariate analysis of Installment

UNIVARIATE ANALYSIS:

3.Univariate analysis of Installment

Most values between 161.01500000000001 and 385.78

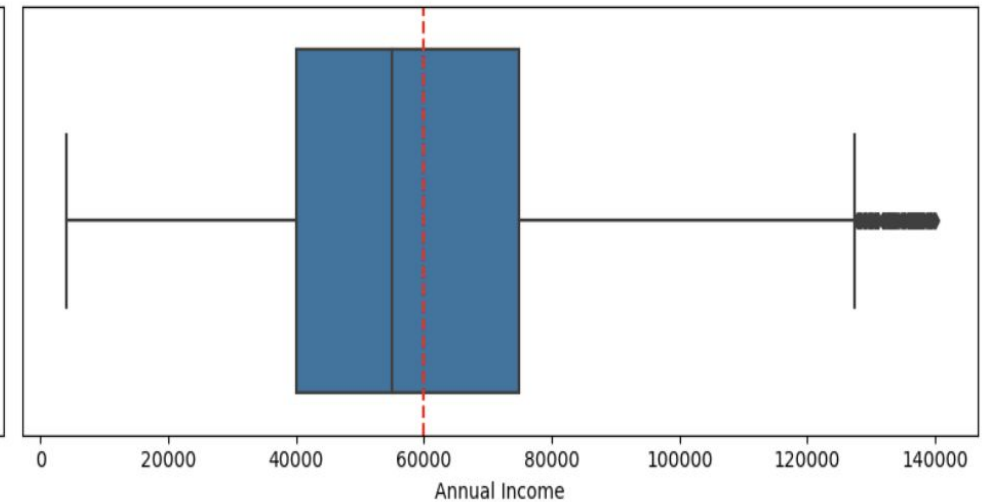
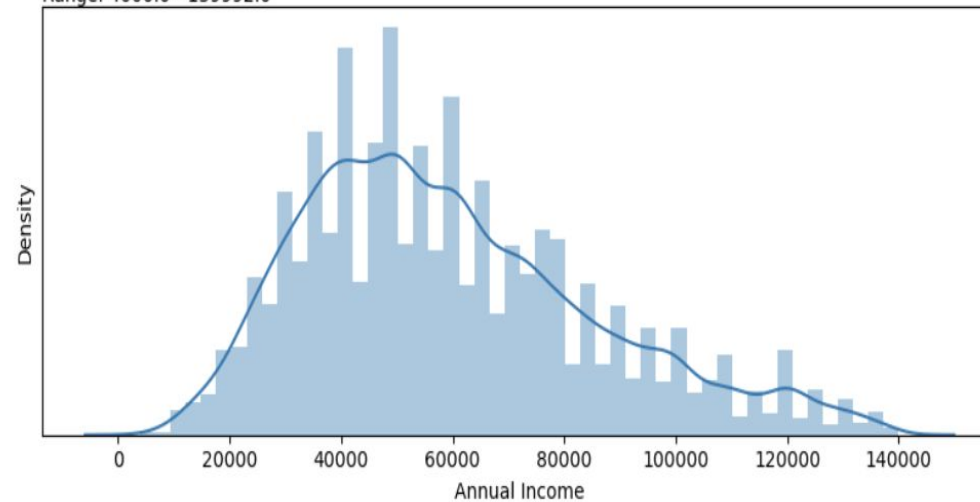
Range: 16.08 - 763.83



4.Univariate analysis of Annual Income

Most values between 40000.0 and 75000.0

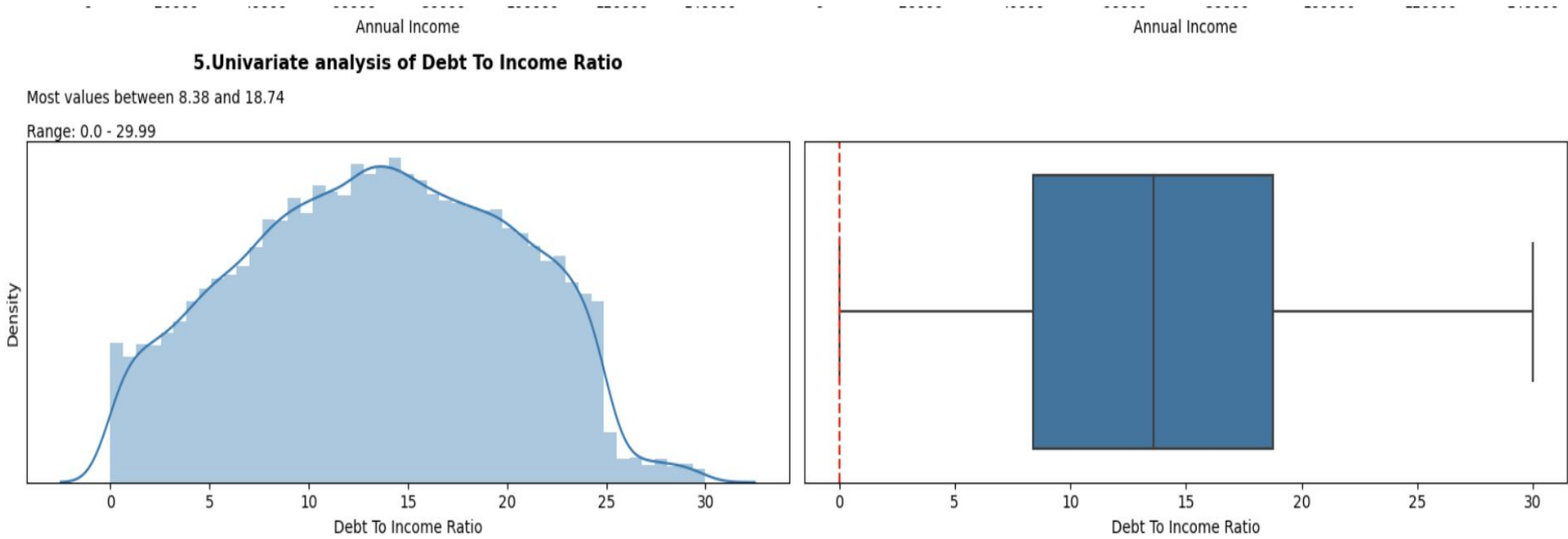
Range: 4000.0 - 139992.0



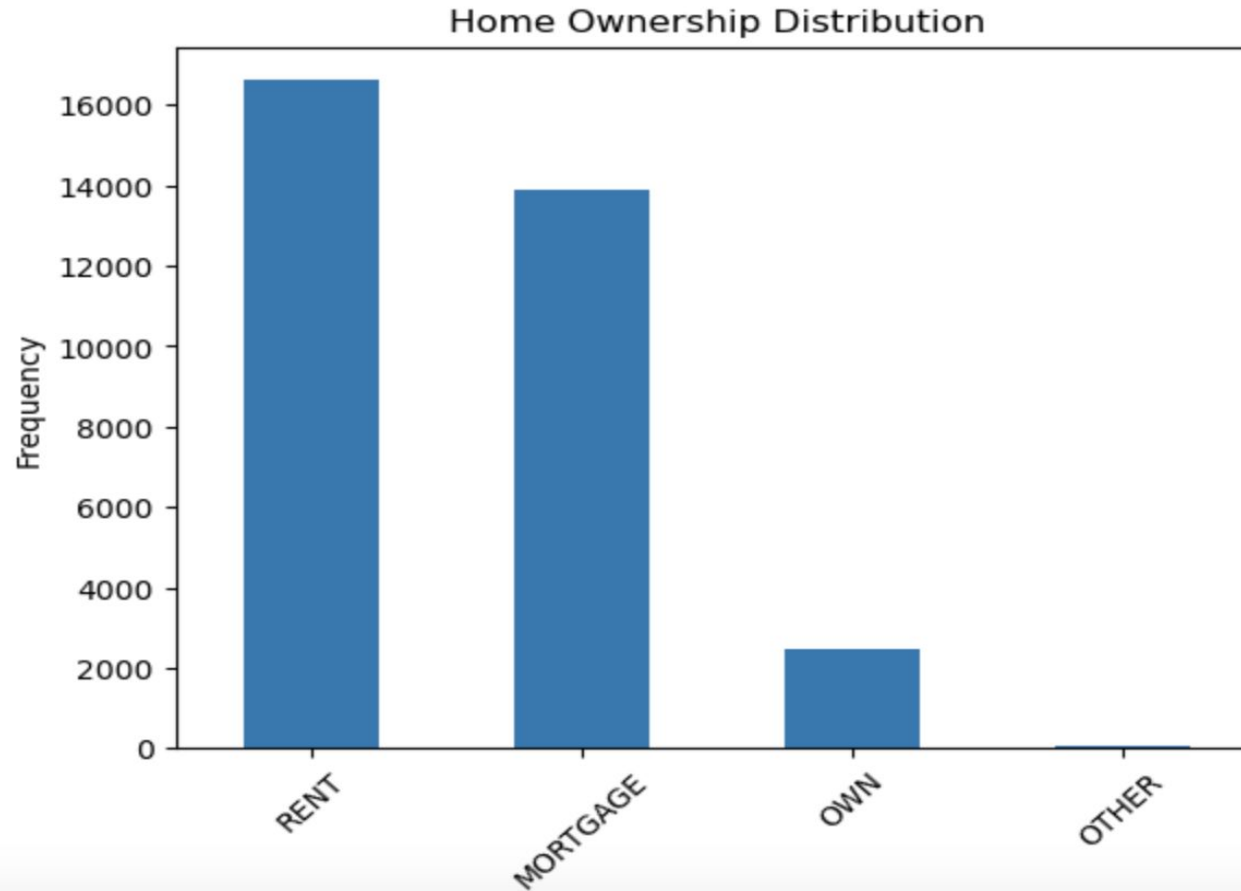
5.Univariate analysis of Debt To Income Ratio



UNIVARIATE ANALYSIS:



UNIVARIATE ANALYSIS:





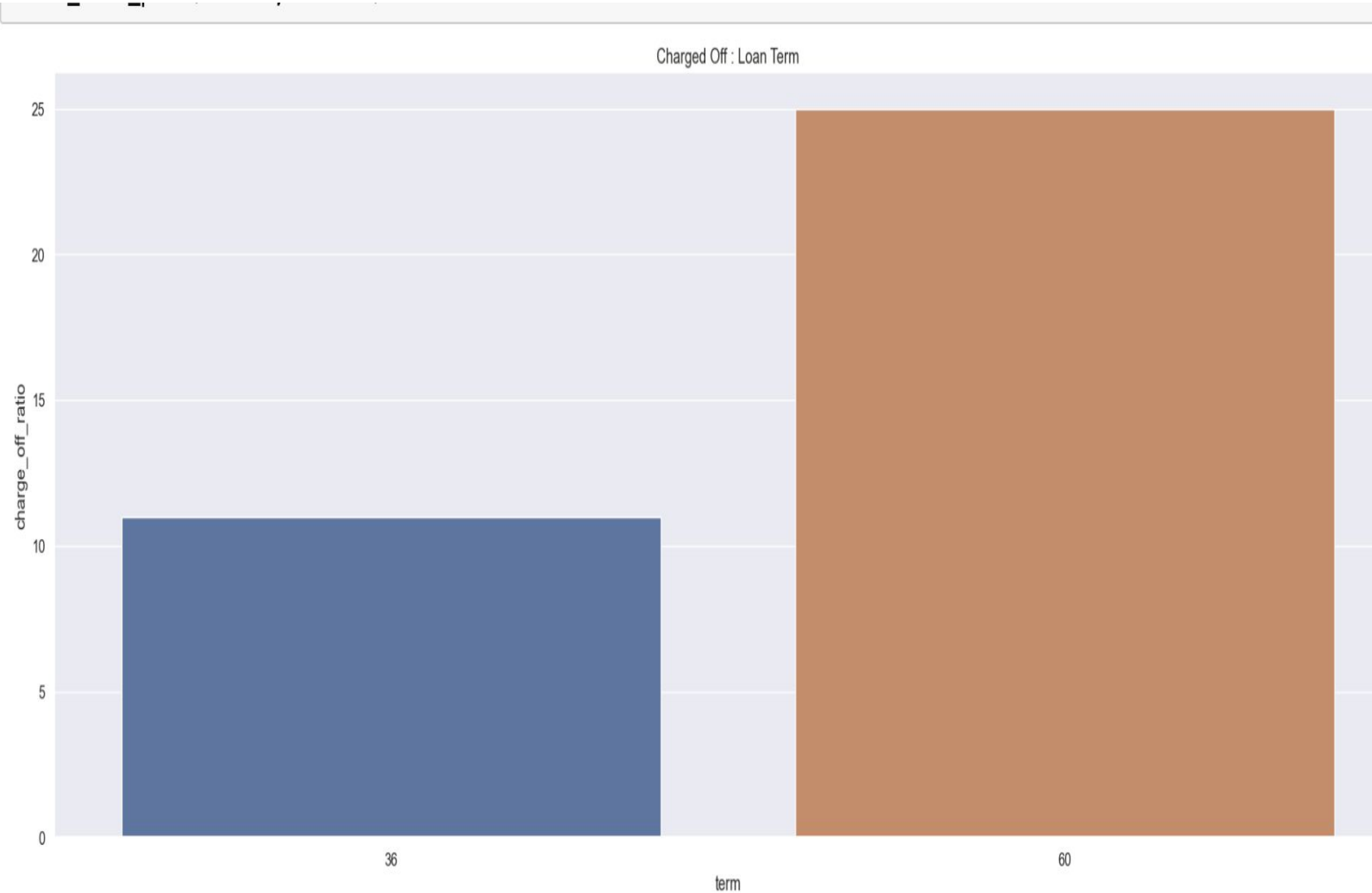
BIVARIATE, MULTIVARIATE, AND CORRELATION ANALYSIS

BIVARIATE ANALYSIS

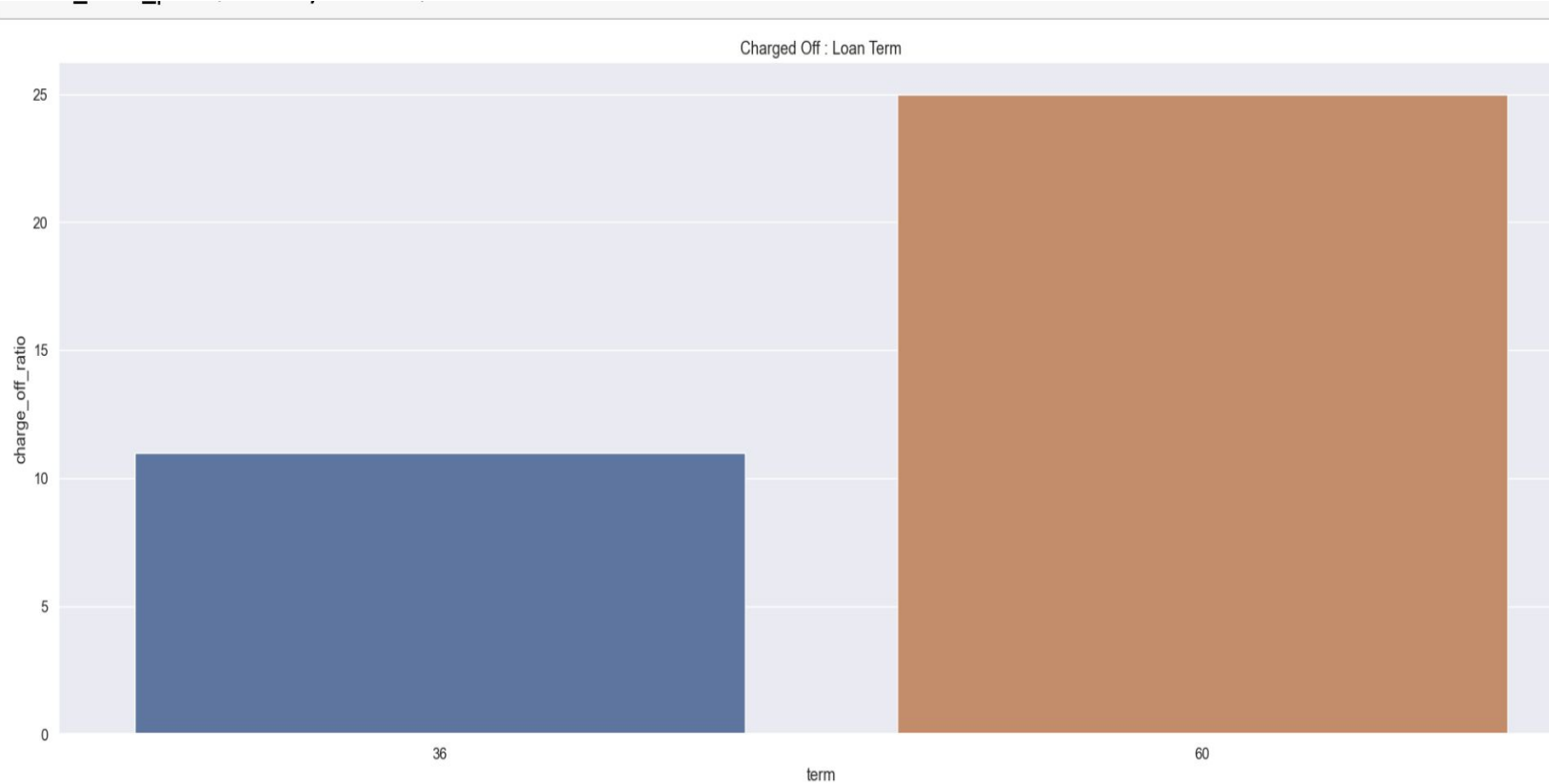
- Bivariate analysis is a statistical method that involves the simultaneous analysis of two variables (factors).
 - It aims to determine the empirical relationship between them.
 - The analysis can be used to test hypotheses, identify patterns, or explore relationships between the variables.
- It was carried out for both Categorical and Quantitative Variables
 - Categorical Variables: Ordered and Unordered
 - Quantitative Variables: Int Rate Bucket, Debt to Income Bucket, Annual Income Bucket, Funded Amount Bucket, Loan Amount Bucket
- Bivariate Analysis Observations
 - Ordered Categorical Variables: The loan applicants belonging to Grades B, C, and D contribute to most of the



BIIVARIATE ANALYSIS:



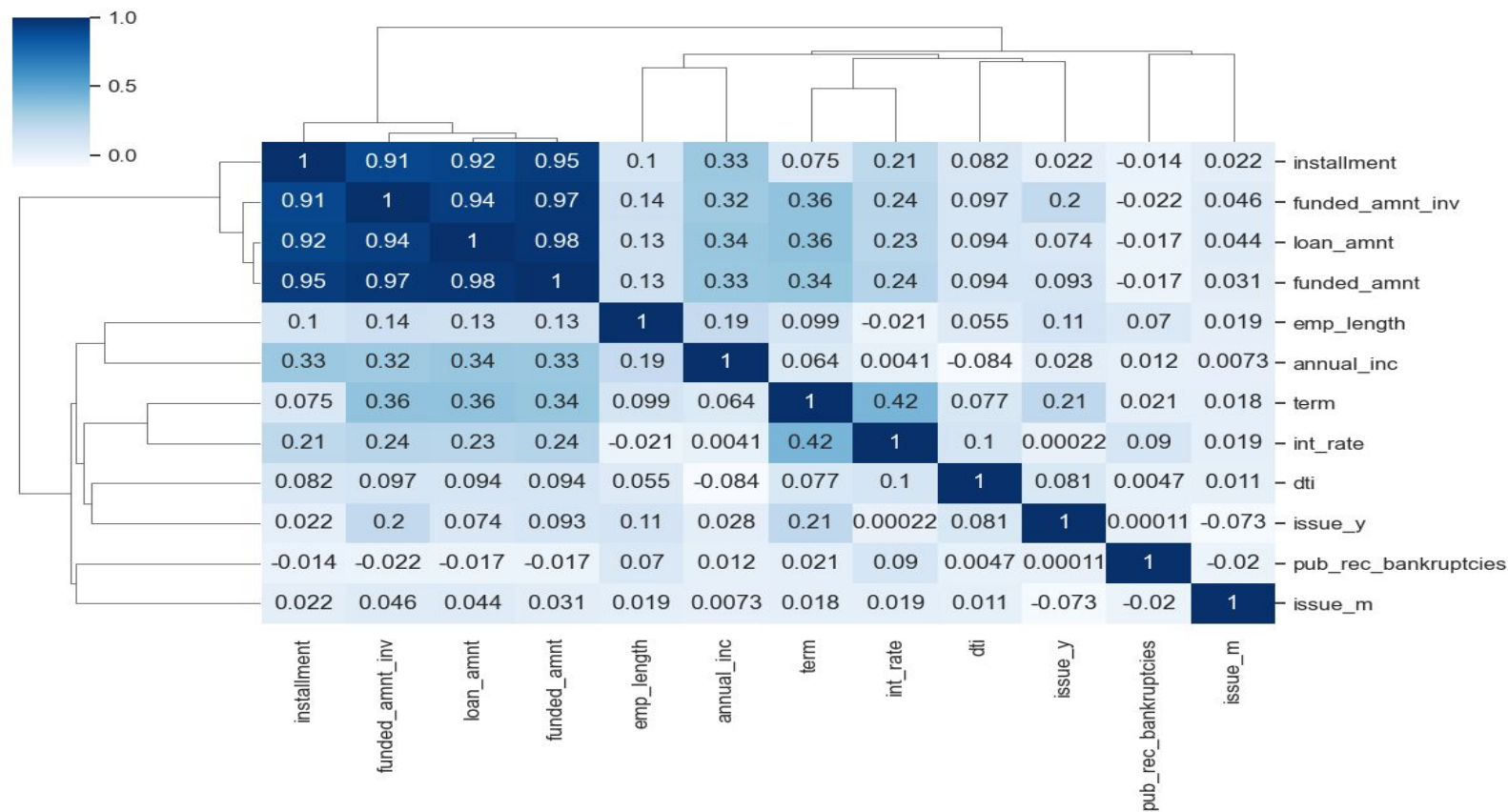
BIIVARIATE ANALYSIS:



MULTIVARIATE ANALYSIS

- Statistical technique used to analyze data involving more than two variables
 - Examines relationships between multiple variables simultaneously
- Widely used in various fields
 - Economics, social sciences, biology, marketing, and environmental science
- Can include different types of variables
 - Categorical, numerical, or a combination of both
- Observations and Inferences
 - Tendency to default the loan is likely with loan applicants belonging to B, C, D grades
 - Borrowers from sub grade B3, B4 and B5 have maximum tendency to default
 - Loan applicants with 10 years of experience has maximum tendency to default the loan
 - Borrowers from states CA, FL, NJ have maximum tendency to default the loan

MULTIVARIATE ANALYSIS:





CORRELATION ANALYSIS

- Statistical technique to measure strength and direction of relationship between variables
 - Quantifies degree of association between changes in variables
 - Widely used in various fields to understand patterns and relationships in data
- Correlation ranges from -1 to 1
 - $r=1$: perfect positive correlation
 - $r=-1$: perfect negative correlation
 - $r=0$: no correlation between variables

SUGGESTIONS, REFERENCES, AND CONCLUSION



SUGGESTIONS

- Implement Stricter Criteria for Grades B, C, and D
 - Minimize default risks with stricter risk assessment and underwriting criteria
- Focus on Subgrades B3, B4, and B5
 - Consider additional risk mitigation measures or lower loan amounts
- Evaluate and Limit 60-Month Loans
 - Decrease likelihood of defaults by limiting maximum term or adjusting interest rates
- Comprehensive Credit Scoring System
 - Incorporate various risk-related attributes for gauging creditworthiness
- Capitalizing on Market Growth
 - Maintain competitive edge while ensuring robust risk management practices
- Anticipate Peak Periods
 - Ensure efficient processing to meet customer demands during busy seasons

REFERENCES & USEFUL LINKS

- Technologies & Packages Used
- GitHub Repository Link:
<https://github.com/kajalmahata123/Lending-Club-Case-Study>
- Thank You!

Technology / Package	Version	Documentation
Python	3.11.4	https://www.python.org/
Matplotlib	3.7.1	https://matplotlib.org/
Numpy	1.24.3	https://numpy.org/
Pandas	1.5.3	https://pandas.pydata.org/
Seaborn	0.12.2	https://seaborn.pydata.org/