



Module Code: CT050-3-M-DAP

Module Name: Data Analytical Programming

Individual Assignment

Loan Approval Prediction for LFI

Student Name: Muhammad Yousouf Ali Budullah

TP Number: TP086704

Intake Code: APDMF2501DSBA(BI)(PR)

Programme: MSc Data Science & Business Analytics

Module Lecturer: Mr. Dhason Padmakumar

Date of Submission: September 7, 2025

Contents

Part 1 — Foundations & Data Setup	4
Chapter 1: Introduction	4
Chapter 2: Background	5
Chapter 3: Assumptions and Justification	6
Chapter 4: Data Dictionary / Metadata	7
<i>Upload the Datasets to SAS Assignment Folder</i>	7
<i>Transfer the Datasets to ASGMLIB Permanent Library</i>	7
<i>Data Dictionary / Metadata Structure</i>	8
Part 2 — Literature Review & Variable Analysis	10
Chapter 5: Literature Review	10
<i>Loan Application Process</i>	10
<i>Financing Companies and Credit Assessment</i>	12
<i>Bad Debts and Risk Management</i>	13
Chapter 6: Analysis of the Variables	15
<i>Variable Overview in (ASGMLIB.TRAINING_DS)</i>	15
<i>Univariate: Categorical (ASGMLIB.TRAINING_DS)</i>	17
<i>Univariate: Numerical (ASGMLIB.TRAINING_DS)</i>	20
<i>Bivariate: Categorical (ASGMLIB.TRAINING_DS)</i>	24
<i>Bivariate: Categorical vs. Numerical (ASGMLIB.TRAINING_DS)</i>	27
<i>Variable Overview in (ASGMLIB.TESTING_DS)</i>	30
<i>Univariate: Categorical (ASGMLIB.TESTING_DS)</i>	32
<i>Univariate: Numerical (ASGMLIB.TESTING_DS)</i>	33
<i>Bivariate: Categorical (ASGMLIB.TESTING_DS)</i>	34
<i>Bivariate: Categorical vs. Numerical (ASGMLIB.TESTING_DS)</i>	38
Chapter 7: Data Cleaning and Imputation	41
<i>Backup Utilities (ASGMLIB)</i>	41
<i>Categorical Imputation (ASGMLIB)</i>	42
<i>Continuous Imputation (ASGMLIB)</i>	47
Chapter 8: Model Creation	52
<i>Model: Training (ASGMLIB)</i>	52
<i>Scoring: Testing (ASGMLIB)</i>	55
Part 3 — Report Generation, Data Visualization & Conclusion	57
Chapter 9: Report Generation	57
<i>Report: PDF Generation (ASGMLIB)</i>	57
<i>Report: Complex PDF Generation (ASGMLIB)</i>	60
Chapter 10: Data Visualisation	62

Chapter 11: Conclusion	66
<i>Loan Application Model Results.</i>	66
<i>Reflections as a Data Scientist</i>	66

Part 1

Foundations & Data Setup

Chapter 1

Introduction

In today's fast-paced financial environment, automation and data-driven decision-making are essential to maintain efficiency, accuracy, and competitiveness. As a data scientist at Lasiandra Finance Inc. (LFI), I have been tasked with developing a predictive model to improve the company's loan approval process. LFI is a leading private lender serving small and medium enterprises across the United States. To continue delivering customer-centric services, LFI must adopt innovative approaches to improve the current manual loan approval processes.

This project aims to build a predictive model using SAS analytics that can determine customer loan applications approvals and rejections with greater accuracy than the current process. By focusing on historical loan data, this project seeks to minimize human error in decision-making and improve the reliability of approvals. The model will serve as the foundation for automating LFI's loan eligibility screening, reduce processing time, support fair decision-making and help the business move toward a more scalable solution.

Chapter 2

Background

Lasiandra Finance Inc. is a private financing company based in New York that provides tailored loan services to SMEs. The business offers tailor-made financial solutions and flexibility to help small businesses access funding for growth and operations. LFI's strength lies in its personalised customer service and flexibility, but its current operations now faces challenges with its manual loan approval system.

Currently the manual process begins when an applicant submits their loan request, including personal and financial details. Loan officers then manually verify each application against internal criteria and past trends. They review each application one by one, checking details like gender, marital status, number of family members, income levels, employment type, and loan history. Each application is reviewed individually, requiring time-consuming cross-validation and decision-making.

This process causes delays and inconsistency. Without a data-driven model it also increases the risk of biases and errors in determining eligibility especially when dealing with a large number of loan applications. Therefore, there is an urgent need for an automated, predictive system that can promptly and accurately identify eligible applicants based on historical trends and attributes.

Chapter 3

Assumptions and Justification

Several assumptions have been made to guide the analysis and development of the predictive model. These assumptions have been made based on the provided dataset and the goal of the project to improve the loan approval process at LFI.

The following are assumed:

- The dataset provided is accurate
- The variables included are relevant to the decision-making process
- The data reflects real historical decisions made by LFI officers
- External factors like economic shifts or policy changes have not significantly affected the loan outcomes in the dataset
- Each record in the dataset represents a single loan application
- There is no major class imbalance that would distort model training
- All numeric and categorical variables are properly formatted and interpretable in SAS

Justification for the Use of SAS.

SAS was chosen because of its reliability, maturity, and strong support for statistical analysis. Its strengths include combining data management and analytics in one platform, producing reproducible reports through macros and ODS, and providing optimized procedures (e.g., PROC LOGISTIC, PROC FREQ, PROC MEANS) that simplify analysis. These features, along with its wide adoption in financial institutions, make SAS a natural choice for LFI's business needs.

At the same time, SAS has limitations. It is less flexible for modern machine learning integration, can be resource-heavy on large datasets, and is more costly than open-source alternatives. Compared with big-data pipelines such as NiFi or Spark, it also offers less scalability.

Despite these drawbacks, SAS is well suited for this project because the dataset is moderate in size, the focus is on interpretable statistical modeling, and the priority is transparent, auditable results. Using SAS ensures outputs that are reliable, consistent and aligned with the operational environment of the LFI.

Chapter 4

Data Dictionary / Metadata

Upload the Datasets to SAS Assignment Folder

Output

DAP_FT_JUN_2025_ASGMT_TP086704

TESTING_DS.csv

TRAINING DS.csv

Figure 1: Datasets visible in the SAS Assignment folder

The TRAINING_DS and TESTING_DS datasets were successfully uploaded to the SAS Assignment folder using the SAS Studio interface. The data is comprised of loan application records and are the basis for all future pre-processing and analysis steps.

Transfer the Datasets to ASGMLIB Permanent Library

Output

ASGMLIB

TESTING_DS

TRAINING_DS

Figure 2: TRAINING_DS and TESTING_DS saved under ASGMLIB

Using the Task Utilities in the SAS Studio interface, the data scientist selected the TRAINING_DS and TESTING_DS datasets from the uploaded files in the SAS assignment folder. By ensuring the library location was set to ASGMLIB and executing the import, the datasets were saved permanently. This ensures they remain accessible across SAS sessions, supporting consistency throughout the project.

Data Dictionary / Metadata Structure

Outputs					
Alphabetic List of Variables and Attributes					
#	Variable	Type	Len	Format	Informat
7	CANDIDATE_INCOME	Num	8	BEST12.	BEST32.
6	EMPLOYMENT	Char	3	\$3.	\$3.
4	FAMILY_MEMBERS	Char	2	\$2.	\$2.
2	GENDER	Char	6	\$6.	\$6.
8	GUARANTEE_INCOME	Num	8	BEST12.	BEST32.
9	LOAN_AMOUNT	Num	8	BEST12.	BEST32.
13	LOAN_APPROVAL_STATUS	Char	1	\$1.	\$1.
10	LOAN_DURATION	Num	8	BEST12.	BEST32.
11	LOAN_HISTORY	Num	8	BEST12.	BEST32.
12	LOAN_LOCATION	Char	7	\$7.	\$7.
3	MARITAL_STATUS	Char	11	\$11.	\$11.
5	QUALIFICATION	Char	14	\$14.	\$14.
1	SME_LOAN_ID_NO	Char	8	\$8.	\$8.

Listing 1: SAS code to describe and inspect the permanent table structure

Using PROC CONTENTS, the metadata of the TRAINING_DS dataset was obtained. The schema includes both categorical and numeric variables and confirms consistent formats/informats, supporting later preprocessing and analysis.

Part 2

Literature Review & Variable Analysis

Chapter 5

Literature Review

Loan Application Process

Traditionally the loan application process has been littered with inefficiencies due to manual workflows, departmental handoffs, and limited compliance validation. Traditionally, even processing a single loan could take 11–13 business days, with approval rates as low as 42% and error rates exceeding 25%. These delays are the result of fragmented data entry, inconsistent application of underwriting criteria, and incomplete documentation checks (Munnangi, 2024). Consequently, many institutions had high abandonment rates, reducing competitiveness and customer satisfaction.

Munnangi (2024) provides a breakdown of these inefficiencies. As illustrated in Figure 3, manual loan processing involved several departmental focal points, with error rates of up to 27%, and processing times took longer than 11 business days. These results illustrate why automation has become essential, as systematic problems, including compliance backlogs and inadequate paperwork checks, were more accountable for inefficiency than applicant quality.

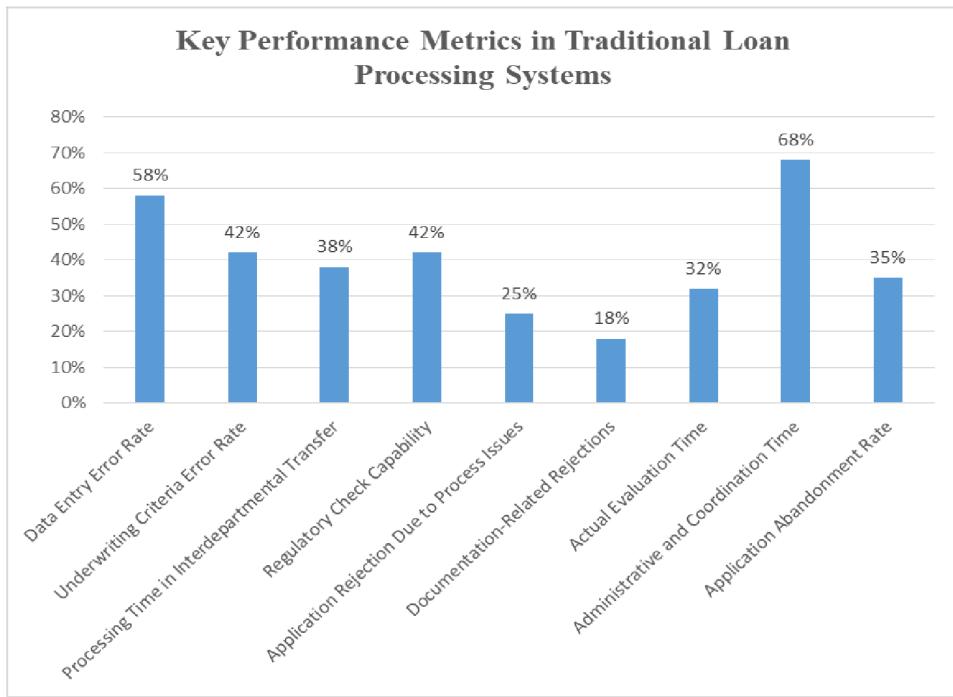


Figure 3: Analysis of manual loan processing inefficiencies in banking operations (adapted from (Munnangi, 2024)).

Recent developments illustrate how AI-driven document processing and intelligent automation are influencing this environment. Loan processing times decreased from almost two weeks to

less than 24 hours due to the introduction of a Loan Origination System (LOS) that combines intelligent routing, real-time compliance validation, and workflow automation (Munnangi, 2024). Customer satisfaction jumped to 89% and approval rates increased to 65%. When compared with manual processing, the LOS greatly reduced errors and delays, as seen in Figure 4, demonstrating visible efficiency gains.

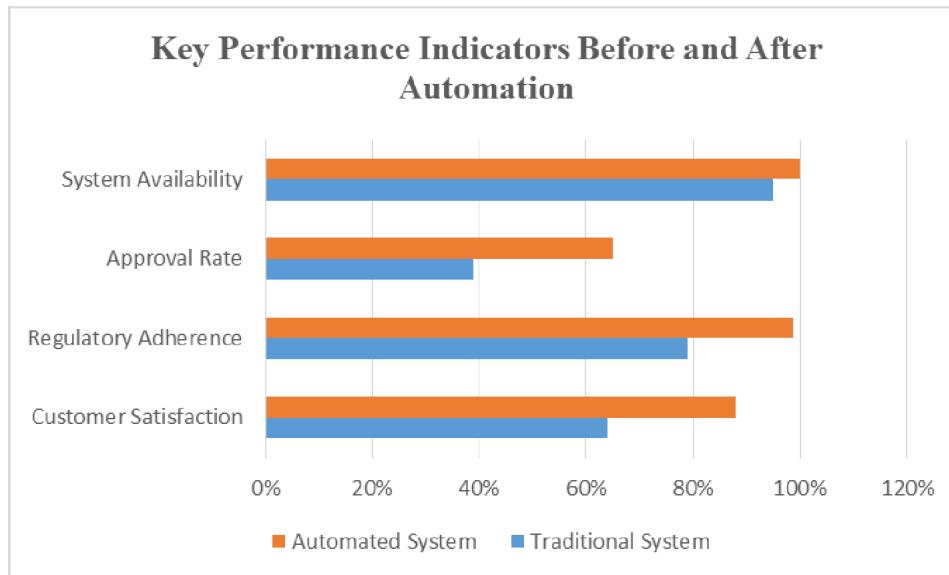


Figure 4: Before and after performance of loan processing systems, showing reduced delays and error rates through automation (adapted from (Munnangi, 2024)).

AI-driven Intelligent Document Processing (IDP) is a complementing invention. IDP uses robotic process automation, machine learning, and natural language processing to extract, validate, and filter applicant data in real time since loan applications involve numerous pieces of papers. According to Ramesh Pingili (2025), banks that adopted IDP have claimed 70% faster approvals, 50% higher fraud detection, and 40% lower compliance costs. Figure 5 reveals these cost reductions, comparing savings across categories including fraud detection and compliance.

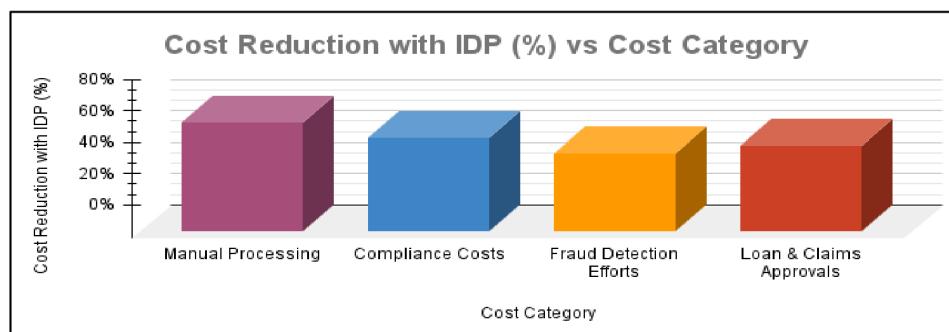


Figure 5: Cost savings achieved with AI-driven intelligent document processing in banking (adapted from (Ramesh Pingili, 2025)).

When used together, LOS and IDP offer an important change in loan applications from inef-

ficient, error prone systems to intelligent, flexible platforms that achieve a balance between accuracy, speed, and compliance. These tools enable financial institutions to deliver efficiency and confidence by reducing processing times.

Financing Companies and Credit Assessment

Predicting loan applications involves understanding how financing companies decide which applications are approved. Different criteria are used by different types of institutions, including banks, equity-based lenders, and non-bank financial institutions. Businesses often mix different forms of funding instead of depending only on one. For predictive modeling, this means that applicant qualities alone are not sufficient to explain loan approvals, institutional procedures and available funding options must also be taken into consideration (Santos et al., 2024).

Credit assessment has gone through a pretty substantial systematic change. Traditional approaches often introduced subjectivity and inconsistency since they rely on human judgment or basic rule-based scorecards. The use of machine learning is being employed by finance companies to produce faster and more accurate decisions. On a dataset of more than 1.3 million loans, Suhadolnik et al. (2023) examined ten methods and discovered that ensemble models in particular, XGBoost consistently performed better than decision trees and logistic regression. The most important variables were the loan period and interest rate, which together accounted for over half of the predictive power, according to their feature importance analysis (Figure 6).

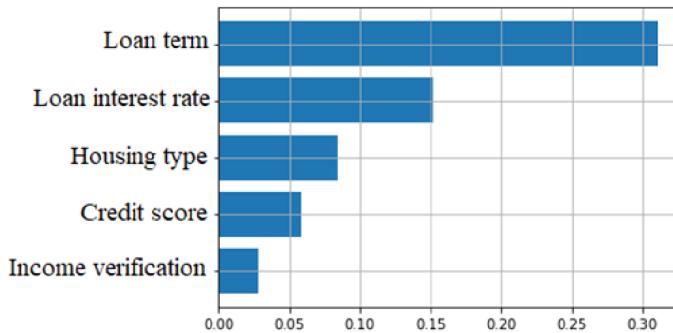


Figure 12. Relative importance of features.

Figure 6: Relative importance of features in loan application classification (adapted from (Suhadolnik et al., 2023)).

Model performance was also assessed using accuracy, precision, recall, F1, and AUC. As shown in Figure 7, ensemble models achieved higher predictive accuracy and better separation compared to traditional methods.

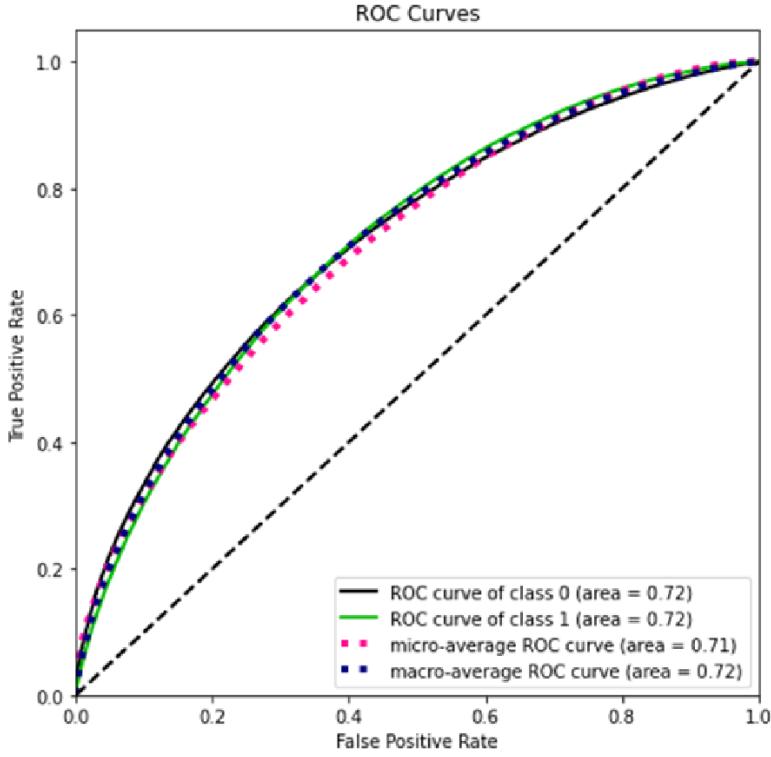


Figure 7: ROC curve comparison of classifiers in loan classification (adapted from (Suhadolnik et al., 2023)).

While financing firms offer a variety of tools, the factor that significantly impacts the prediction of loan applications is their increasing usage of machine learning. These methodological developments allow prediction models to more closely resemble the approval decisions taken by actual human lenders.

Bad Debts and Risk Management

Although loan approvals are a necessary step, financing companies also have to manage the default risk in order to ensure long-term sustainability of their decisions. Because they reduce profitability and threaten financial stability, bad debts, which are often referred to as non-performing loans (NPLs) are a significant problem to lenders. Consequently, it is vital to properly understand the factors that impact default risk to link loan application choices with future credit results.

Recent data from the European peer-to-peer (P2P) lending market, a sector of the financial system that has expanded substantially as part of FinTech lending, are presented by Nigmonov et al. (2024). Their study uses complex econometric techniques, such as endogeneity correction and the least absolute shrinkage and selection operator (LASSO), to analyse default risk using a loan-book dataset dating from 2014 to 2020. The results revealed that the chances of loan defaults are increased by both higher interest rates and larger stock market returns.

This is particularly relevant since it illustrates that credit risk is influenced by both borrowers' characteristics and the state of the financial market.

The study also reveals that the degree of FinTech adoption and the difficulties faced in the banking system have an impact on how severe these consequences are. Essentially, this means that default risk on alternative lending platforms can be increased by systemic stress and market volatility. Nigmonov et al. (2024) highlight the importance of incorporating macro-financial data into risk management frameworks by directly linking interest rates, stock market movements, and delinquency rates.

These findings illustrate how important the role of thorough screening during the application process is for predictive models of loan approvals. Although approval databases might not monitor final defaults, understanding the factors that contribute to bad debts adds substantial context. Effective approval models could help lower the long-term probability of non-performing loans in addition to predicting immediate results.

Chapter 6

Analysis of the Variables

Variable Overview in (ASGMLIB.TRAINING_DS)

This section provides an overview of the categorical, numerical, and bivariate variables in the ASGMLIB.TRAINING_DS dataset. Variables marked with a checkmark (✓) were selected for analysis, while those with a cross (✗) were excluded.

Categorical Variable	Analysed
EMPLOYMENT	✗
FAMILY_MEMBERS	✓
GENDER	✗
LOAN_APPROVAL_STATUS	✗
LOAN_LOCATION	✗
MARITAL_STATUS	✓
QUALIFICATION	✓
SME_LOAN_ID_NO	✗
LOAN_HISTORY	✗

Table 2: Univariate categorical variables in ASGMLIB.TRAINING_DS

Numerical Variable	Analysed
CANDIDATE_INCOME	✓
GUARANTEE_INCOME	✓
LOAN_AMOUNT	✓
LOAN_DURATION	✓
LOAN_HISTORY	✗

Table 3: Univariate numerical variables in ASGMLIB.TRAINING_DS

Variable 1	Variable 2	Analysed
MARITAL_STATUS	FAMILY_MEMBERS	✓
MARITAL_STATUS	QUALIFICATION	✓
GENDER	EMPLOYMENT	✓

Table 4: Bivariate categorical variable pairs in ASGMLIB.TRAINING_DS

Categorical	Numerical	Analysed
EMPLOYMENT	LOAN_AMOUNT	✓
QUALIFICATION	LOAN_DURATION	✓
GENDER	GUARANTEE_INCOME	✓

Table 5: Bivariate categorical vs numerical variable pairs in ASGMLIB.TRAINING_DS

Univariate Analysis of Categorical Variables in ASGMLIB.TRAINING_DS

Univariate Analysis of Categorical Variable: Family Members

SAS Code		Outputs				
Line	Code					
1	<i>/* Family Members: Frequency distribution and bar chart */</i>					
2	TITLE "Univariate Analysis - Frequency of Family Members";					
3	PROC FREQ DATA = ASGMLIB.					
	TRAINING_DS;					
4	TABLE FAMILY_MEMBERS;					
5	RUN;					
6						
7	TITLE "Univariate Analysis - Bar Chart of Family Members";					
8	ODS GRAPHICS / RESET WIDTH=3.0in HEIGHT=4.0in IMAGEMAP;					
9	PROC SGPLOT DATA = ASGMLIB.					
	TRAINING_DS;					
10	VBAR FAMILY_MEMBERS;					
11	RUN;					

Listing 2: SAS code for frequency table and bar chart of FAMILY_MEMBERS

Univariate Analysis - Bar Chart of FAMILY_MEMBERS

FAMILY_MEMBERS	Frequency
0	345
1	102
2	101
3+	51

Figure 8: Bar chart of FAMILY_MEMBERS

Most applicants have no family members (57.6%), followed by one (17.0%) and two (16.9%). Only 8.5% lies in the category three-or-more, indicating that the distribution is skewed toward smaller households. This makes sense since many SME loan applicants apply on their own rather than as part of a larger household. There are also 15 missing records which should be flagged in the table so they don't get overlooked later.

Univariate Analysis of Categorical Variable: Marital Status

SAS Code	Outputs																				
<pre> 1 /* Marital Status: Frequency distribution and bar chart */ 2 TITLE "Univariate Analysis - Frequency of Marital Status"; 3 PROC FREQ DATA = ASGMLIB. TRAINING_DS; TABLE MARITAL_STATUS; 5 RUN; 6 7 TITLE "Univariate Analysis - Bar Chart of Marital Status"; 8 ODS GRAPHICS / RESET WIDTH=3.0in HEIGHT=4.0in IMAGEMAP; 9 PROC SGPlot DATA = ASGMLIB. TRAINING_DS; VBAR MARITAL_STATUS; 11 RUN; </pre> <p>Listing 3: SAS code for frequency table and bar chart of MARITAL_STATUS</p>	<table border="1"> <thead> <tr> <th>MARITAL_STATUS</th><th>Frequency</th><th>Percent</th><th>Cumulative Frequency</th><th>Cumulative Percent</th></tr> </thead> <tbody> <tr> <td>Married</td><td>398</td><td>65.14</td><td>398</td><td>65.14</td></tr> <tr> <td>Not Married</td><td>213</td><td>34.86</td><td>611</td><td>100.00</td></tr> <tr> <td colspan="5">Frequency Missing = 3</td></tr> </tbody> </table> <p>Table 7: Frequency distribution of MARITAL_STATUS</p> <p>Figure 9: Bar chart of MARITAL_STATUS</p>	MARITAL_STATUS	Frequency	Percent	Cumulative Frequency	Cumulative Percent	Married	398	65.14	398	65.14	Not Married	213	34.86	611	100.00	Frequency Missing = 3				
MARITAL_STATUS	Frequency	Percent	Cumulative Frequency	Cumulative Percent																	
Married	398	65.14	398	65.14																	
Not Married	213	34.86	611	100.00																	
Frequency Missing = 3																					

The majority of applicants are married (65.1%) while 34.9% are not married. The distribution is skewed toward married applicants, as they make up almost two-thirds of the dataset. A possible reason is that married individuals may be seen as more financially stable or more likely to apply for loans with family responsibilities in mind. There are also 3 missing cases that need to be highlighted for future modeling.

Univariate Analysis of Categorical Variable: Qualification

SAS Code	Outputs															
<pre> 1 /* Qualification: Frequency distribution and bar chart */ 2 TITLE "Univariate Analysis - Frequency of Qualification"; 3 PROC FREQ DATA = ASGMLIB. TRAINING_DS; TABLE QUALIFICATION; 5 RUN; 6 7 TITLE "Univariate Analysis - Bar Chart of Qualification"; 8 ODS GRAPHICS / RESET WIDTH=3.0in HEIGHT=4.0in IMAGEMAP; 9 PROC SGPlot DATA = ASGMLIB. TRAINING_DS; VBAR QUALIFICATION; 11 RUN; </pre> <p>Listing 4: SAS code for frequency table and bar chart of QUALIFICATION</p>	<table border="1"> <thead> <tr> <th>QUALIFICATION</th><th>Frequency</th><th>Percent</th><th>Cumulative Frequency</th><th>Cumulative Percent</th></tr> </thead> <tbody> <tr> <td>Graduate</td><td>480</td><td>78.18</td><td>480</td><td>78.18</td></tr> <tr> <td>Under Graduate</td><td>134</td><td>21.82</td><td>614</td><td>100.00</td></tr> </tbody> </table> <p>Table 8: Frequency distribution of QUALIFICATION</p> <p>Figure 10: Bar chart of QUALIFICATION</p>	QUALIFICATION	Frequency	Percent	Cumulative Frequency	Cumulative Percent	Graduate	480	78.18	480	78.18	Under Graduate	134	21.82	614	100.00
QUALIFICATION	Frequency	Percent	Cumulative Frequency	Cumulative Percent												
Graduate	480	78.18	480	78.18												
Under Graduate	134	21.82	614	100.00												

Most of the applicants are graduates (78.2%) while undergraduates make up only 21.8%. The distribution shows a strong bias toward higher education. It's likely because graduates have more access to opportunities, resources, and financial literacy which makes them more inclined to apply for SME loans. No missing values were recorded here, so the variable is complete and easy to work with.

Univariate Analysis of Numerical Variables in ASGMLIB.TRAINING.DS

Univariate Analysis of Numerical Variable: Candidate Income

SAS Code

```
1 /* Candidate Income: Summary  
   statistics and histogram */  
2 TITLE "Univariate Analysis - Summary  
      of Candidate Income";  
3 PROC MEANS DATA = ASGMLIB.  
      TRAINING_DS N NMISS MIN MAX MEAN  
      MEDIAN STD;  
4      VAR CANDIDATE_INCOME;  
5      RUN;  
6  
7 TITLE "Univariate Analysis -  
      Histogram of Candidate Income";  
8 PROC SGLOT DATA = ASGMLIB.  
      TRAINING_DS;  
9      HISTOGRAM CANDIDATE_INCOME;  
10     RUN;
```

Listing 5: SAS code for descriptive statistics and histogram of CANDIDATE_INCOME

Outputs

Analysis Variable : CANDIDATE_INCOME						
N	N Miss	Minimum	Maximum	Mean	Median	Std Dev
614	0	150.0000000	81000.00	5403.46	3812.50	6109.04

Table 9: Descriptive statistics of CANDIDATE_INCOME

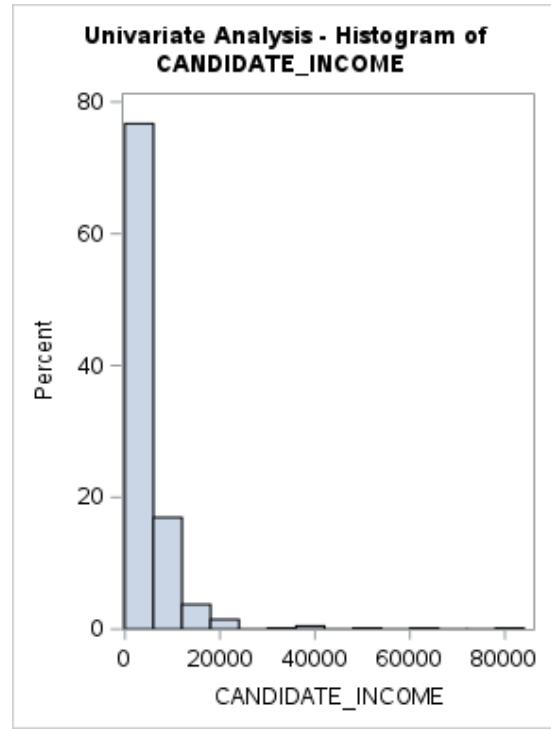
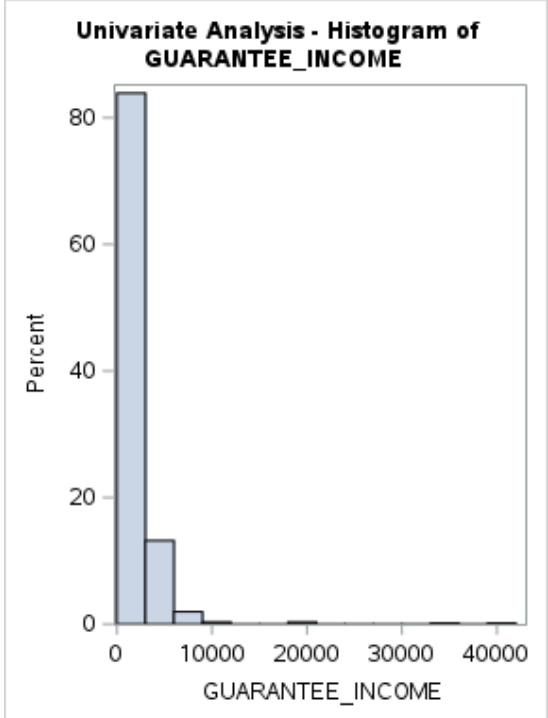


Figure 11: Histogram of CANDIDATE_INCOME

The histogram (Figure 11) shows that the majority of the applicants have a low income, with a very small number of them having a significantly higher income. This trend is also supported by the summary statistics table (Table 9), which displays a mean income of 5,403.46, which is higher than the median income of 3,812.50. This difference indicates that the distribution is skewed to the right, with the few high-income applicants acting as outliers.

Univariate Analysis of Numerical Variable: Guarantee Income

SAS Code	Outputs														
<pre> 1 /* Guarantee Income: Summary statistics and histogram */ 2 TITLE "Univariate Analysis - Summary of Guarantee Income"; 3 PROC MEANS DATA = ASGMLIB. TRAINING_DS N NMISS MIN MAX MEAN MEDIAN STD; 4 VAR GUARANTEE_INCOME; 5 RUN; 6 7 TITLE "Univariate Analysis - Histogram of Guarantee Income"; 8 PROC SGPLOT DATA = ASGMLIB. TRAINING_DS; 9 HISTOGRAM GUARANTEE_INCOME; 10 RUN; </pre> <p>Listing 6: SAS code for descriptive statistics and histogram of GUARANTEE_INCOME</p>	<p>Analysis Variable : GUARANTEE_INCOME</p> <table border="1"> <thead> <tr> <th>N</th><th>N Miss</th><th>Minimum</th><th>Maximum</th><th>Mean</th><th>Median</th><th>Std Dev</th></tr> </thead> <tbody> <tr> <td>614</td><td>0</td><td>0</td><td>41667.00</td><td>1621.25</td><td>1188.50</td><td>2926.25</td></tr> </tbody> </table> <p>Table 10: Descriptive statistics of GUARANTEE_INCOME</p>  <p>The histogram displays the distribution of GUARANTEE_INCOME. The x-axis is labeled 'GUARANTEE_INCOME' and ranges from 0 to 40,000 with major ticks every 10,000 units. The y-axis is labeled 'Percent' and ranges from 0 to 80 with major ticks every 20 units. The distribution is highly right-skewed, with the highest frequency occurring at the lowest income levels (around 0-10,000). A single bar at approximately 10,000 represents about 10% of the data. The distribution tapers off significantly as income increases, with very few individuals having incomes above 20,000.</p> <p>Figure 12: Histogram of GUARANTEE_INCOME</p>	N	N Miss	Minimum	Maximum	Mean	Median	Std Dev	614	0	0	41667.00	1621.25	1188.50	2926.25
N	N Miss	Minimum	Maximum	Mean	Median	Std Dev									
614	0	0	41667.00	1621.25	1188.50	2926.25									

GUARANTEE_INCOME indicates a highly right-skewed distribution, meaning that most applicants have a low guarantee income, with a few having very high incomes. The histogram (Figure 12) shows a large number of values at the lowest income levels, with a long tail extending to the right. The summary table (Table 10) confirms this skewness, as the mean (\$1,621.25) is higher than the median (\$1,188.50), pulled upward by high-income outliers.

Univariate Analysis of Numerical Variable: Loan Amount

SAS Code

```

1 /* Loan Amount: Summary statistics
   and histogram */
2 TITLE "Univariate Analysis - Summary
   of Loan Amount";
3 PROC MEANS DATA = ASGMLIB.
   TRAINING_DS N NMISS MIN MAX MEAN
   MEDIAN STD;
4   VAR LOAN_AMOUNT;
5   RUN;
6
7 TITLE "Univariate Analysis -
   Histogram of Loan Amount";
8 PROC SGPLOT DATA = ASGMLIB.
   TRAINING_DS;
9   HISTOGRAM LOAN_AMOUNT;
10  RUN;

```

Listing 7: SAS code for descriptive statistics and histogram of LOAN_AMOUNT

Outputs

Analysis Variable : LOAN_AMOUNT						
N	N Miss	Minimum	Maximum	Mean	Median	Std Dev
592	22	9.000000	700.000000	146.4121622	128.0000000	85.5873252

Table 11: Descriptive statistics of LOAN_AMOUNT

Univariate Analysis - Histogram of LOAN_AMOUNT

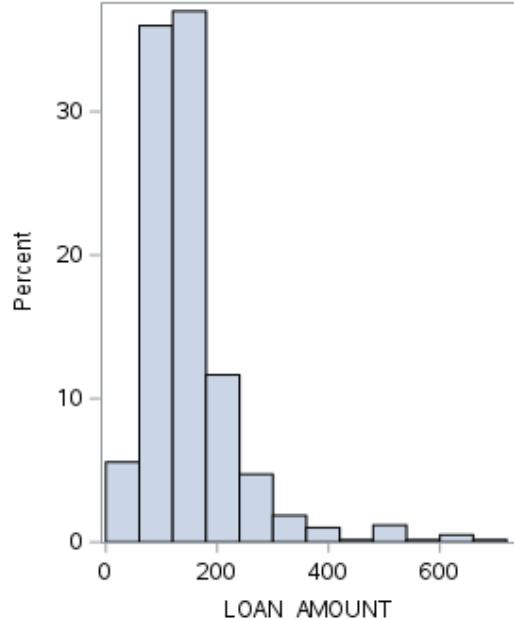
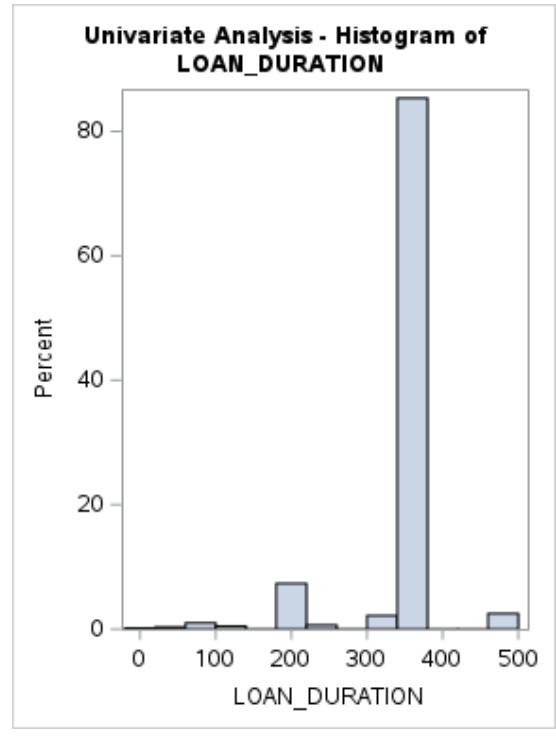


Figure 13: Histogram of LOAN_AMOUNT

LOAN_AMOUNT shows that the distribution is skewed to the right, with the majority of loan amounts falling between 100 and 200. This is clearly seen in the histogram (Figure 13), which shows a clear peak in this range. The summary table (Table 11) further supports this, as the mean (\$146.41) is slightly higher than the median (\$128.00), a sign of a right-skewed distribution caused by a few larger loan amounts. It also worth noting that there are 22 missing values LOAN_AMOUNT.

Univariate Analysis of Numerical Variable: Loan Duration

SAS Code	Outputs																					
<pre> 1 /* Loan Duration: Summary statistics and histogram */ 2 TITLE "Univariate Analysis - Summary of Loan Duration"; 3 PROC MEANS DATA = ASGMLIB. TRAINING_DS N NMISS MIN MAX MEAN MEDIAN STD; 4 VAR LOAN_DURATION; 5 RUN; 6 7 TITLE "Univariate Analysis - Histogram of Loan Duration"; 8 PROC SGPLOT DATA = ASGMLIB. TRAINING_DS; 9 HISTOGRAM LOAN_DURATION; 10 RUN; </pre> <p>Listing 8: SAS code for descriptive statistics and histogram of LOAN_DURATION</p>	<table border="1"> <thead> <tr> <th colspan="7">Analysis Variable : LOAN_DURATION</th> </tr> <tr> <th>N</th><th>N Miss</th><th>Minimum</th><th>Maximum</th><th>Mean</th><th>Median</th><th>Std Dev</th></tr> </thead> <tbody> <tr> <td>600</td><td>14</td><td>12.000000</td><td>480.000000</td><td>342.000000</td><td>360.000000</td><td>65.1204099</td></tr> </tbody> </table> <p>Table 12: Descriptive statistics of LOAN_DURATION</p>  <p>The histogram displays the distribution of loan durations. The x-axis is labeled 'LOAN_DURATION' and ranges from 0 to 500. The y-axis is labeled 'Percent' and ranges from 0 to 80. There are several bars, with the highest bar reaching approximately 85% at a duration of about 360 days. Other bars are much smaller, indicating a long tail of longer loans.</p> <p>Figure 14: Histogram of LOAN_DURATION</p>	Analysis Variable : LOAN_DURATION							N	N Miss	Minimum	Maximum	Mean	Median	Std Dev	600	14	12.000000	480.000000	342.000000	360.000000	65.1204099
Analysis Variable : LOAN_DURATION																						
N	N Miss	Minimum	Maximum	Mean	Median	Std Dev																
600	14	12.000000	480.000000	342.000000	360.000000	65.1204099																

LOAN_DURATION shows a strong focus of loans at around 360 days , with over 80% of the loans serving this duration. This is clearly seen by the tall bar in the histogram (Figure 14). The summary table (Table 12) shows that the median is 360 days, reflecting this concentration. The mean is 342 days, which is slightly lower than the median, suggesting a small number of short duration loans are pulling the mean down. It also worth noting that there are 14 missing values in LOAN_DURATION.

Bivariate Analysis of Categorical Variables in ASGMLIB.TRAINING_DS

Bivariate Analysis of Categorical Variable: Gender vs Employment

SAS Code

```

1 /* Relationship between Gender and
   Employment */
2 TITLE1 "Bivariate Analysis of
   Categorical Variables";
3 TITLE2 "Gender vs Employment";
4 PROC FREQ DATA = ASGMLIB.
   TRAINING_DS;
5 TABLE GENDER * EMPLOYMENT /
   PLOTS = FREQPLOT(TWOWAY=
   STACKED SCALE=GROUPPCT);
7 RUN;

```

Listing 9: SAS code for cross-tabulation and stacked bar chart of GENDER vs EMPLOYMENT

Outputs

Table of GENDER by EMPLOYMENT			
GENDER	EMPLOYMENT		
	No	Yes	Total
Female	89 15.64	15 2.64	104 18.28
	85.58	14.42	
	18.13	19.23	
Male	402 70.65	63 11.07	465 81.72
	86.45 81.87	13.55 80.77	
Total	491 86.29	78 13.71	569 100.00

Frequency Missing = 45

Table 13: Crosstab of GENDER by EMPLOYMENT

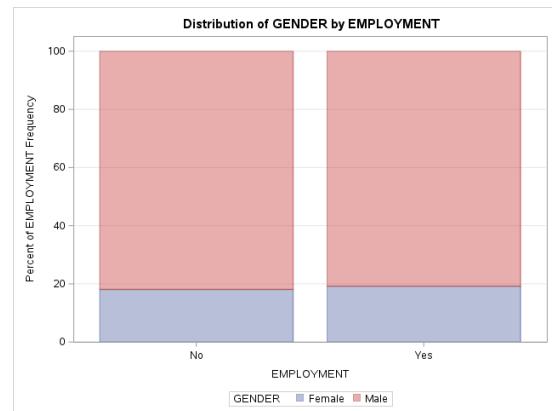


Figure 15: Stacked bar chart of GENDER by EMPLOYMENT

The bivariate analysis of GENDER and EMPLOYMENT show that employment rates are consistent between genders. Given on the table's row percentages (Table 13), 14.42% of females are employed, which is only marginally higher than 13.55% for males. The stacked bar chart (Figure 15) validates this, as the proportion of employed applicants (the "Yes" bar) has an equal gender composition to the non-employed applicants (the "No" bar). However, due to the much larger number of males present in ASGMLIB.TRAINING_DS, they represent a larger percentage of the total employed population at 80.77% compared to females at 19.23%. It is important to note that the analysis was conducted on 569 applicants due to 45 missing values, a significant amount.

Bivariate Analysis of Categorical Variable: Marital Status vs Family Members

SAS Code		Outputs																																			
<pre> 1 /* Relationship between Marital Status and Family Members */ 2 TITLE1 "Bivariate Analysis of Categorical Variables"; 3 TITLE2 "Marital Status vs Family Members"; 4 PROC FREQ DATA = ASGMLIB. TRAINING_DS; TABLE MARITAL_STATUS * FAMILY_MEMBERS / PLOTS = FREQPLOT(TWOWAY= STACKED SCALE=GROUPPCT); 7 RUN;</pre>		<p>Table 14: Crosstab of MARITAL_STATUS by FAMILY_MEMBERS</p> <table border="1"> <thead> <tr> <th rowspan="2">MARITAL_STATUS</th> <th colspan="5">FAMILY_MEMBERS</th> </tr> <tr> <th>0</th> <th>1</th> <th>2</th> <th>3+</th> <th>Total</th> </tr> </thead> <tbody> <tr> <td>Married</td> <td>174 29.05 44.62 50.43</td> <td>79 13.19 20.26 77.45</td> <td>93 15.53 23.85 92.08</td> <td>44 7.35 11.28 86.27</td> <td>390 65.11</td> </tr> <tr> <td>Not Married</td> <td>171 28.55 81.82 49.57</td> <td>23 3.84 11.00 22.55</td> <td>8 1.34 3.83 7.92</td> <td>7 1.17 3.35 13.73</td> <td>209 34.89</td> </tr> <tr> <td>Total</td> <td>345 57.60</td> <td>102 17.03</td> <td>101 16.86</td> <td>51 8.51</td> <td>599 100.00</td> </tr> <tr> <td colspan="5">Frequency Missing = 15</td></tr> </tbody> </table> <p>Figure 16: Stacked bar chart of MARITAL_STATUS by FAMILY_MEMBERS</p>		MARITAL_STATUS	FAMILY_MEMBERS					0	1	2	3+	Total	Married	174 29.05 44.62 50.43	79 13.19 20.26 77.45	93 15.53 23.85 92.08	44 7.35 11.28 86.27	390 65.11	Not Married	171 28.55 81.82 49.57	23 3.84 11.00 22.55	8 1.34 3.83 7.92	7 1.17 3.35 13.73	209 34.89	Total	345 57.60	102 17.03	101 16.86	51 8.51	599 100.00	Frequency Missing = 15				
MARITAL_STATUS	FAMILY_MEMBERS																																				
	0	1	2	3+	Total																																
Married	174 29.05 44.62 50.43	79 13.19 20.26 77.45	93 15.53 23.85 92.08	44 7.35 11.28 86.27	390 65.11																																
Not Married	171 28.55 81.82 49.57	23 3.84 11.00 22.55	8 1.34 3.83 7.92	7 1.17 3.35 13.73	209 34.89																																
Total	345 57.60	102 17.03	101 16.86	51 8.51	599 100.00																																
Frequency Missing = 15																																					

Through the bivariate analysis we can see MARITAL_STATUS and FAMILY_MEMBERS has a clear relationship. Where, as the number of family members increases, so does the chance of a person being married. The stacked bar chart (Figure 16) demonstrates this trend, with the “Married” (blue) section growing as the FAMILY_MEMBERS category increases from 0 to 3+. The table’s row percentages (Table 14) validates this, as applicants who are married rises from 50.43% for those with 0 family members to 86.27% for those with 3 or more family members. It is important to note that the analysis was conducted on 599 individuals due to 15 missing values, a moderate amount.

Bivariate Analysis of Categorical Variable: Marital Status vs Qualification

SAS Code	Outputs																			
<pre> 1 /* Relationship between Marital Status and Qualification */ 2 TITLE1 "Bivariate Analysis of Categorical Variables"; 3 TITLE2 "Marital Status vs Qualification"; 4 PROC FREQ DATA = ASGMLIB. TRAINING_DS; TABLE MARITAL_STATUS * QUALIFICATION / PLOTS = FREQPLOT(TWOWAY= STACKED SCALE=GROUPPCT); 7 RUN; </pre> <p>Listing 11: SAS code for cross-tabulation and stacked bar chart of MARITAL_STATUS vs QUALIFICATION</p>	<p>Table 15: Crosstab of MARITAL_STATUS by QUALIFICATION</p> <table border="1"> <thead> <tr> <th rowspan="2">MARITAL_STATUS</th> <th colspan="3">QUALIFICATION</th> </tr> <tr> <th>Graduate</th> <th>Under Graduate</th> <th>Total</th> </tr> </thead> <tbody> <tr> <td>Married</td> <td>309 50.57 77.64 64.78</td> <td>89 14.57 22.36 66.42</td> <td>398 65.14</td> </tr> <tr> <td>Not Married</td> <td>168 27.50 78.87 35.22</td> <td>45 7.36 21.13 33.58</td> <td>213 34.86</td> </tr> <tr> <td>Total</td> <td>477 78.07</td> <td>134 21.93</td> <td>611 100.00</td> </tr> </tbody> </table> <p style="text-align: right;">Frequency Missing = 3</p> <p>Figure 17: Stacked bar chart of MARITAL_STATUS by QUALIFICATION</p>	MARITAL_STATUS	QUALIFICATION			Graduate	Under Graduate	Total	Married	309 50.57 77.64 64.78	89 14.57 22.36 66.42	398 65.14	Not Married	168 27.50 78.87 35.22	45 7.36 21.13 33.58	213 34.86	Total	477 78.07	134 21.93	611 100.00
MARITAL_STATUS	QUALIFICATION																			
	Graduate	Under Graduate	Total																	
Married	309 50.57 77.64 64.78	89 14.57 22.36 66.42	398 65.14																	
Not Married	168 27.50 78.87 35.22	45 7.36 21.13 33.58	213 34.86																	
Total	477 78.07	134 21.93	611 100.00																	

The bivariate analysis of MARITAL_STATUS and QUALIFICATION shows a weak relationship. The stacked bar chart (Figure 17) illustrates that the proportion of married and not married individuals is nearly the same for both “Graduate” and “Under Graduate” qualifications. The table’s row percentages (Table 15) validates this, as 77.64% of married individuals are graduates and 78.87% of not married individuals are graduates. This could imply that a person’s qualification is not a determining factor of their marital status. It is important to note that the analysis was conducted on 611 individuals due to 3 missing values, a minimal amount.

Bivariate Analysis of Categorical vs. Numerical Variables in ASGMLIB.TRAINING_DS

Bivariate Analysis of Categorical Variable & Continuous: Guarantee Income vs Gender

SAS Code		Outputs																																	
<pre>1 /* Compare Guarantee Income across Gender categories */ 2 TITLE1 "Bivariate Analysis of Variables"; 3 TITLE2 "Gender vs Guarantee Income"; 4 PROC MEANS DATA = ASGMLIB. TRAINING_DS; 5 CLASS GENDER; /* Grouping by Gender (categorical) */ 6 VAR GUARANTEE_INCOME; /* Summarizing Guarantee Income (numeric) */ 7 RUN;</pre>		<table border="1"><thead><tr><th colspan="7">Analysis Variable : GUARANTEE_INCOME</th></tr><tr><th>GENDER</th><th>N Obs</th><th>N</th><th>Mean</th><th>Std Dev</th><th>Minimum</th><th>Maximum</th></tr></thead><tbody><tr><td>Female</td><td>112</td><td>112</td><td>1108.01</td><td>4094.60</td><td>0</td><td>41667.00</td></tr><tr><td>Male</td><td>489</td><td>489</td><td>1742.93</td><td>2606.51</td><td>0</td><td>33837.00</td></tr></tbody></table>						Analysis Variable : GUARANTEE_INCOME							GENDER	N Obs	N	Mean	Std Dev	Minimum	Maximum	Female	112	112	1108.01	4094.60	0	41667.00	Male	489	489	1742.93	2606.51	0	33837.00
Analysis Variable : GUARANTEE_INCOME																																			
GENDER	N Obs	N	Mean	Std Dev	Minimum	Maximum																													
Female	112	112	1108.01	4094.60	0	41667.00																													
Male	489	489	1742.93	2606.51	0	33837.00																													

Listing 12: SAS code for summarizing GUARANTEE_INCOME by GENDER

Table 16: Descriptive statistics of GUARANTEE_INCOME by GENDER

The bivariate analysis of GUARANTEE_INCOME by GENDER shows a significant difference in income between the two genders. The mean guarantee income for males is 1,742.93, which is greater than the mean of 1,108.01 for females. The standard deviation for females (4,094.60) is also significantly larger than for males (2,606.51), which possibly means a greater spread and more extreme outliers in female guarantee incomes (Table 16).

Bivariate Analysis of Categorical Variable & Continuous: Loan Amount vs Employment

SAS Code	Outputs																					
<pre> 1 /* Compare Loan Amount across Employment categories */ 2 TITLE1 "Bivariate Analysis of Variables"; 3 TITLE2 "Employment vs Loan Amount"; 4 PROC MEANS DATA = ASGMLIB. TRAINING_DS; 5 CLASS EMPLOYMENT; /* Grouping by Employment (categorical) */ 6 VAR LOAN_AMOUNT; /* Summarizing Loan Amount (numeric) */ 7 RUN; </pre> <p>Listing 13: SAS code for summarizing LOAN_AMOUNT by EMPLOYMENT</p>	<p style="text-align: center;">Analysis Variable : LOAN_AMOUNT</p> <table border="1"> <thead> <tr> <th>EMPLOYMENT</th> <th>N Obs</th> <th>N</th> <th>Mean</th> <th>Std Dev</th> <th>Minimum</th> <th>Maximum</th> </tr> </thead> <tbody> <tr> <td>No</td> <td>500</td> <td>482</td> <td>141.7489627</td> <td>79.7809192</td> <td>9.0000000</td> <td>700.0000000</td> </tr> <tr> <td>Yes</td> <td>82</td> <td>79</td> <td>172.0000000</td> <td>108.6346500</td> <td>25.0000000</td> <td>650.0000000</td> </tr> </tbody> </table>	EMPLOYMENT	N Obs	N	Mean	Std Dev	Minimum	Maximum	No	500	482	141.7489627	79.7809192	9.0000000	700.0000000	Yes	82	79	172.0000000	108.6346500	25.0000000	650.0000000
EMPLOYMENT	N Obs	N	Mean	Std Dev	Minimum	Maximum																
No	500	482	141.7489627	79.7809192	9.0000000	700.0000000																
Yes	82	79	172.0000000	108.6346500	25.0000000	650.0000000																

Table 17: Descriptive statistics of LOAN_AMOUNT by EMPLOYMENT

The bivariate analysis of LOAN_AMOUNT by EMPLOYMENT illustrates that employed applicants are more likely to take out larger loan amounts. The mean loan amount for employed applicants is 172, which is greater than the mean loan amount of 141.75 for unemployed applicants. Furthermore, the standard deviation is higher for employed individuals, at 108.63, compared to 79.78 for unemployed applicants. This indicates a wider range of loan amounts within the employed group compared to the unemployed group,(Table 17).

Bivariate Analysis of Categorical Variable & Continuous: Loan Duration vs Qualification

SAS Code	Outputs																					
<pre> 1 /* Compare Loan Duration across Qualification categories */ 2 TITLE1 "Bivariate Analysis of Variables"; 3 TITLE2 "Qualification vs Loan Duration"; 4 PROC MEANS DATA = ASGMLLIB. TRAINING_DS; 5 CLASS QUALIFICATION; /* Grouping by Qualification (categorical) */ 6 VAR LOAN_DURATION; /* Summarizing Loan Duration (numeric) */ 7 RUN; </pre> <p>Listing 14: SAS code for summarizing LOAN_DURATION by QUALIFICATION</p>	<p style="text-align: center;">Analysis Variable : LOAN_DURATION</p> <table border="1"> <thead> <tr> <th>QUALIFICATION</th> <th>N Obs</th> <th>N</th> <th>Mean</th> <th>Std Dev</th> <th>Minimum</th> <th>Maximum</th> </tr> </thead> <tbody> <tr> <td>Graduate</td> <td>480</td> <td>472</td> <td>344.6694915</td> <td>61.2996615</td> <td>12.0000000</td> <td>480.0000000</td> </tr> <tr> <td>Under Graduate</td> <td>134</td> <td>128</td> <td>332.1562500</td> <td>77.0796129</td> <td>36.0000000</td> <td>480.0000000</td> </tr> </tbody> </table>	QUALIFICATION	N Obs	N	Mean	Std Dev	Minimum	Maximum	Graduate	480	472	344.6694915	61.2996615	12.0000000	480.0000000	Under Graduate	134	128	332.1562500	77.0796129	36.0000000	480.0000000
QUALIFICATION	N Obs	N	Mean	Std Dev	Minimum	Maximum																
Graduate	480	472	344.6694915	61.2996615	12.0000000	480.0000000																
Under Graduate	134	128	332.1562500	77.0796129	36.0000000	480.0000000																

The bivariate analysis of LOAN_DURATION by QUALIFICATION reveals that there is a minor difference in the loan duration between graduates and under graduates. The mean loan duration for graduate candidates is 344.67 days, which is only marginally longer than the mean of 332.16 days for under graduate candidates. The standard deviations are also nearly identical, indicating that both groups have a similar spread in their loan durations (Table 18).

Variable Overview in (ASGMLIB.TESTING_DS)

This section provides an overview of the categorical, numerical, and bivariate variables in the ASGMLIB.TESTING_DS dataset. Variables marked with a checkmark (✓) were analysed, while those with a cross (✗) were excluded due to time constraints.

Categorical Variable	Analysed
EMPLOYMENT	✗
FAMILY_MEMBERS	✓
GENDER	✓
LOAN_APPROVAL_STATUS	✗
LOAN_LOCATION	✗
MARITAL_STATUS	✓
QUALIFICATION	✓
SME_LOAN_ID_NO	✗
LOAN_HISTORY	✗

Table 19: Univariate categorical variables in ASGMLIB.TESTING_DS

Numerical Variable	Analysed
CANDIDATE_INCOME	✓
GUARANTEE_INCOME	✓
LOAN_AMOUNT	✓
LOAN_DURATION	✓
LOAN_HISTORY	✗

Table 20: Univariate numerical variables in ASGMLIB.TESTING_DS

Variable 1	Variable 2	Analysed
MARITAL_STATUS	FAMILY_MEMBERS	✓
MARITAL_STATUS	QUALIFICATION	✓
GENDER	EMPLOYMENT	✓

Table 21: Bivariate categorical variable pairs in ASGMLIB.TESTING_DS

Categorical	Numerical	Analysed
EMPLOYMENT	LOAN_AMOUNT	✓
QUALIFICATION	LOAN_DURATION	✓
GENDER	GUARANTEE_INCOME	✓

Table 22: Bivariate categorical vs numerical variable pairs in ASGMLIB.TESTING_DS

Univariate Analysis of Categorical Variables in ASGMLIB.TESTING_DS

Univariate Analysis of Categorical Variables: Family Members, Marital Status, Qualification & Gender

SAS Code

```

1 /* Enable compile-time notes for macros
   */
2 OPTIONS MCOMPILENOTE=ALL;
3
4 /* Macro to generate frequency table for
   a categorical variable */
5 %MACRO UVACATE_VAR(ptitle, pds, pvar);
6   TITLE "&PTITLE";
7   PROC FREQ DATA=&PDS;
8     TABLE &PVAR;
9   RUN;
10 %MEND UVACATE_VAR;
11
12 /* Macro calls: categorical variables */
13 %UVACATE_VAR(UNIVARIATE ANALYSIS OF THE
14   CATEGORICAL VARIABLE: Family_Members,
15   ASGMLIB.TESTING_DS, FAMILY_MEMBERS)
16 ;
17 %UVACATE_VAR(UNIVARIATE ANALYSIS OF THE
18   CATEGORICAL VARIABLE: Marital_Status,
19   ASGMLIB.TESTING_DS, MARITAL_STATUS)
20 ;
21 %UVACATE_VAR(UNIVARIATE ANALYSIS OF THE
22   CATEGORICAL VARIABLE: Qualification,
23   ASGMLIB.TESTING_DS, QUALIFICATION);
24 %UVACATE_VAR(UNIVARIATE ANALYSIS OF THE
25   CATEGORICAL VARIABLE: Gender,
26   ASGMLIB.TESTING_DS, GENDER);

```

Listing 15: SAS macro and calls for categorical frequencies in ASGMLIB.TESTING_DS

Outputs

FAMILY_MEMBERS	Frequency	Percent	Cumulative Frequency	Cumulative Percent
0	200	56.02	200	56.02
1	58	16.25	258	72.27
2	59	16.53	317	88.80
3+	40	11.20	357	100.00
Frequency Missing = 10				

Table 23: Frequency distribution of FAMILY_MEMBERS

GENDER	Frequency	Percent	Cumulative Frequency	Cumulative Percent
Female	70	19.66	70	19.66
Male	286	80.34	356	100.00
Frequency Missing = 11				

Table 24: Frequency distribution of GENDER

MARITAL_STATUS	Frequency	Percent	Cumulative Frequency	Cumulative Percent
Married	233	63.49	233	63.49
Not Married	134	36.51	367	100.00

Table 25: Frequency distribution of MARITAL_STATUS

QUALIFICATION	Frequency	Percent	Cumulative Frequency	Cumulative Percent
Graduate	283	77.11	283	77.11
Under Graduate	84	22.89	367	100.00

Table 26: Frequency distribution of QUALIFICATION

The analysis of the applicants shows a clear trend of concentration within a few categories. A significant amount of the applicants are male, making up over 80% of the testing dataset (Table 24). Applicants are predominantly married, accounting for nearly two-thirds of the candidates (Table 25), and a large majority hold a Graduate qualification, representing 77.11% (Table 26). The frequency table for family members also shows a concentration, with more than half the candidates (56.02%) having no family members (Table 23). It is important to note that a small number of missing values exist for both the Gender and Family Members variables.

Univariate Analysis of Numerical Variables in ASGMLIB.TESTING_DS

Univariate Analysis of Numerical Variables: Candidate Income, Guarantee Income, Loan Amount & Loan Duration

SAS Code

```

1 /* Enable compile-time notes for macros
   */
2 OPTIONS MCOMPILENOTE=ALL;
3 /* Macro to generate summary statistic
   for a continuous variable */
4 %MACRO UVA_CONTI_VAR(ptitle, pds, pvar);
5   TITLE "&PTITLE";
6   PROC MEANS DATA=&pds N NMISS MIN MAX
7     MEAN MEDIAN STD;
8     VAR &pvar;
9   RUN;
10 %MEND UVA_CONTI_VAR;
11 /* Macro calls: continuous variables */
12 %UVA_CONTI_VAR(UNIVARIATE ANALYSIS OF THE
13   CONTINUOUS VARIABLE:
14   Candidate_Income, ASGMLIB.TESTING_DS
15   , CANDIDATE_INCOME);
16 %UVA_CONTI_VAR(UNIVARIATE ANALYSIS OF THE
17   CONTINUOUS VARIABLE:
18   Guarantee_Income, ASGMLIB.TESTING_DS
19   , GUARANTEE_INCOME);
20 %UVA_CONTI_VAR(UNIVARIATE ANALYSIS OF THE
21   CONTINUOUS VARIABLE: Loan_Amount,
22   ASGMLIB.TESTING_DS, LOAN_AMOUNT);
23 %UVA_CONTI_VAR(UNIVARIATE ANALYSIS OF THE
24   CONTINUOUS VARIABLE: Loan_Duration,
25   ASGMLIB.TESTING_DS, LOAN_DURATION
26 );

```

Listing 16: SAS macro and calls for continuous variables in ASGMLIB.TESTING_DS

Outputs

Analysis Variable : CANDIDATE_INCOME						
N	N Miss	Minimum	Maximum	Mean	Median	Std Dev
367	0	0	72529.00	4805.60	3786.00	4910.69

Table 27: Descriptive statistics of CANDIDATE_INCOME

Analysis Variable : GUARANTEE_INCOME						
N	N Miss	Minimum	Maximum	Mean	Median	Std Dev
367	0	0	24000.00	1569.58	1025.00	2334.23

Table 28: Descriptive statistics of GUARANTEE_INCOME

Analysis Variable : LOAN_AMOUNT						
N	N Miss	Minimum	Maximum	Mean	Median	Std Dev
362	5	28.000000	550.000000	136.1325967	125.0000000	61.3666524

Table 29: Descriptive statistics of LOAN_AMOUNT

Analysis Variable : LOAN_DURATION						
N	N Miss	Minimum	Maximum	Mean	Median	Std Dev
361	6	6.000000	480.000000	342.5373961	360.0000000	65.1566434

Table 30: Descriptive statistics of LOAN_DURATION

Based on the descriptive statistics, the mean candidate income is 4,805.60, which is significantly higher than the average guarantor income of 1,569.58, which implies that applicants tend to have a greater financial capacity than their guarantors (Table 27, Table 28). The average loan amount is 136.13, with a typical loan duration of about 342.54 days (Table 29, Table 30). It is also important to note that data for LOAN_AMOUNT and LOAN_DURATION is not complete, with 5 and 6 missing values respectively, while the income data for both candidates and guarantors are complete with 367 observations each.

Bivariate Analysis of Categorical Variables in ASGMLIB.TESTING_DS

SAS Code

```
1 /* Enable compile-time notes for macros */
2 OPTIONS MCOMPILENOTE=ALL;
3
4 /* Macro: BVA_CATE_CATE
5 Purpose : Cross-tabulation and stacked frequency plot for two categorical
       variables
6 Parameters:
7     PTITLE1 - First line of title
8     PTITLE2 - Second line of title (variables being compared)
9     PDS      - Dataset
10    PVAR1   - First categorical variable
11    PVAR2   - Second categorical variable */
12 %MACRO BVA_CATE_CATE(ptitle1, ptitle2, pds, pvar1, pvar2);
13   TITLE1 "&ptitle1";
14   TITLE2 "&ptitle2";
15   PROC FREQ DATA=&pds;
16     TABLE &pvar1 * &pvar2 /
17       PLOTS=FREQPLOT(TWOWAY=STACKED SCALE=GROUPTPCT);
18   RUN;
19 %MEND BVA_CATE_CATE;
20
21 /* Macro calls */
22 %BVA_CATE_CATE(BIVARIATE ANALYSIS OF THE VARIABLES:, Categorical Variable -
23   Marital_Status vs Categorical Variable - Family_Members, ASGMLIB.TESTING_DS
24   , MARITAL_STATUS, FAMILY_MEMBERS);
23 %BVA_CATE_CATE(BIVARIATE ANALYSIS OF THE VARIABLES:, Categorical Variable -
24   Marital_Status vs Categorical Variable - Qualification, ASGMLIB.TESTING_DS,
25   MARITAL_STATUS, QUALIFICATION);
24 %BVA_CATE_CATE(BIVARIATE ANALYSIS OF THE VARIABLES:, Categorical Variable -
25   Gender vs Categorical Variable - Employment, ASGMLIB.TESTING_DS, GENDER,
26   EMPLOYMENT);
```

Listing 17: SAS macro for bivariate analysis of categorical variables in ASGMLIB.TESTING_DS

Bivariate Analysis of Categorical Variable: Marital Status vs Family Members

Output

Table of MARITAL_STATUS by FAMILY_MEMBERS					
MARITAL_STATUS	FAMILY_MEMBERS				
	0	1	2	3+	Total
Married	95 26.61 41.67 47.50	45 12.61 19.74 77.59	53 14.85 23.25 89.83	35 9.80 15.35 87.50	228 63.87
Not Married	105 29.41 81.40 52.50	13 3.64 10.08 22.41	6 1.68 4.65 10.17	5 1.40 3.88 12.50	129 36.13
Total	200 56.02	58 16.25	59 16.53	40 11.20	357 100.00
Frequency Missing = 10					

Table 31: Cross-tabulation of MARITAL_STATUS by FAMILY_MEMBERS (TESTING_DS)

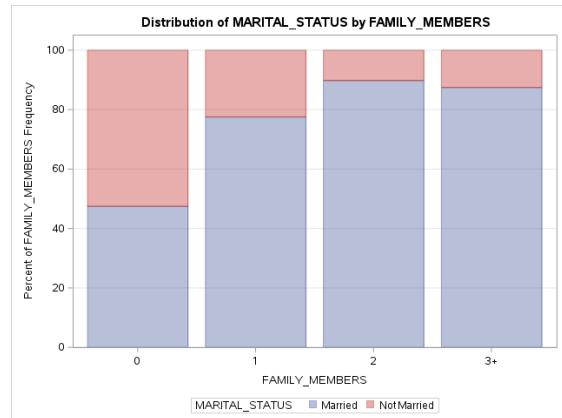


Figure 18: Stacked bar chart of MARITAL_STATUS vs FAMILY_MEMBERS (TESTING_DS)

Through the bivariate analysis we can see MARITAL_STATUS and FAMILY_MEMBERS has a clear relationship similar to the training dataset. The stacked bar chart (Figure 18) validates this, demonstrating that as the number of family members increases, the numbers of applicants who are married also increases . The table (Table 31) validates this as well, showing that the percentage of married individuals rises from 47.50% for those with 0 family members to 87.50% for those with 3 or more. It is important to note that the analysis was conducted on 357 applicants due to 10 missing values, a minimal amount.

Bivariate Analysis of Categorical Variable: Marital Status vs Qualification

Output

Table of MARITAL_STATUS by QUALIFICATION			
MARITAL_STATUS	QUALIFICATION		
	Graduate	Under Graduate	Total
Married	176 47.96 75.54 62.19	57 15.53 24.46 67.86	233 63.49
	107 29.16 79.85 37.81	27 7.36 20.15 32.14	134 36.51
	283 77.11	84 22.89	367 100.00

Table 32: Cross-tabulation of MARITAL_STATUS by QUALIFICATION (TESTING_DS)

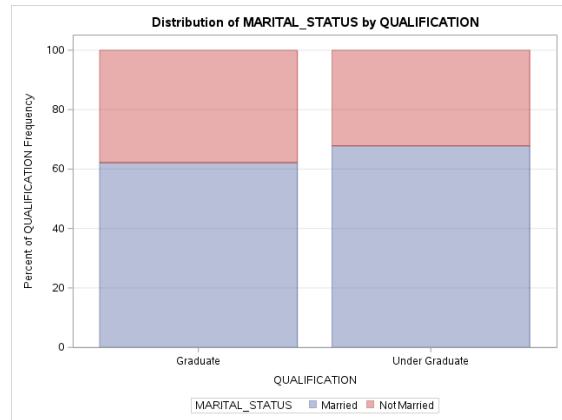


Figure 19: Stacked bar chart of MARITAL_STATUS vs QUALIFICATION (TESTING_DS)

The bivariate analysis of MARITAL_STATUS and QUALIFICATION displays a weak relationship. The stacked bar chart (Figure 19) shows that the percentage of married and not married individuals is nearly the same for both "Graduate" and "Under Graduate" qualifications. The table (Table 32) also reveals the same, with 75.54% of married applicants being graduates and 79.85% of not married applicants also being graduates. This implies that a person's qualification does not have a strong influence on their marital status. This analysis was conducted on the all 367 individuals, with no missing values.

Bivariate Analysis of Categorical Variable: Gender vs Employment

Output

Table of GENDER by EMPLOYMENT			
GENDER	EMPLOYMENT		
	No	Yes	Total
Female	63 18.92 94.03 21.14	4 1.20 5.97 11.43	67 20.12
Male	235 70.57 88.35 78.86	31 9.31 11.65 88.57	266 79.88
Total	298 89.49	35 10.51	333 100.00
Frequency Missing = 34			

Table 33: Cross-tabulation of GENDER by EMPLOYMENT (TESTING_DS)

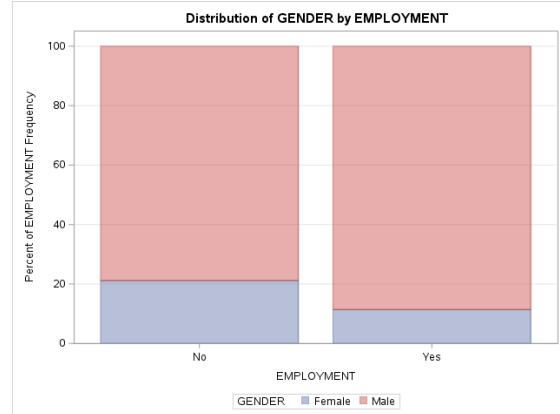


Figure 20: Stacked bar chart of GENDER vs EMPLOYMENT (TESTING_DS)

The bivariate analysis of GENDER and EMPLOYMENT displays a difference in employment rates between male and female applicants. The stacked bar chart (Figure 20) reveals that employed applicants have a much higher share of males than females compared to the unemployed applicants. Based on the cross-tabulation (Table 33), the employment rate for males (11.65%) is significantly higher than for females (5.97%). It is also important to note that males represent the majority of the total employed population with 88.57%, compared to 11.43% for females. It is important to note that the analysis was conducted on 333 applicants due to 34 missing values, a substantial amount.

Bivariate Analysis of Categorical vs. Numerical Variables in ASGMLIB.TESTING_DS

SAS Code

```
1 /* Enable compile-time notes for macros */
2 OPTIONS MCOMPILENOTE=ALL;
3
4 /* Macro: BVA_CATE_CONTI
5 Purpose : Summary statistics of a continuous variable grouped by a
categorical variable
6 Parameters:
7     PTITLE1 - First line of title
8     PTITLE2 - Second line of title (variables being compared)
9     PDS      - Dataset
10    PCATE    - Categorical variable
11    PCONTI   - Continuous variable */
12 %MACRO BVA_CATE_CONTI(ptitle1, ptitle2, pds, pcate, pconti);
13     TITLE1 "&ptitle1";
14     TITLE2 "&ptitle2";
15     PROC MEANS DATA=&pds;
16         CLASS &pcate;
17         VAR &pconti;
18     RUN;
19 %MEND BVA_CATE_CONTI;
20
21 /* Macro calls */
22 %BVA_CATE_CONTI(BIVARIATE ANALYSIS OF THE VARIABLES:, Categorical Variable -
Employment vs Continuous Variable - Loan_Amount, ASGMLIB.TESTING_DS,
EMPLOYMENT, LOAN_AMOUNT);
23 %BVA_CATE_CONTI(BIVARIATE ANALYSIS OF THE VARIABLES:, Categorical Variable -
Qualification vs Continuous Variable - Loan_Duration, ASGMLIB.TESTING_DS,
QUALIFICATION, LOAN_DURATION);
24 %BVA_CATE_CONTI(BIVARIATE ANALYSIS OF THE VARIABLES:, Categorical Variable -
Gender vs Continuous Variable - Guarantee_Income, ASGMLIB.TESTING_DS,
GENDER, GUARANTEE_INCOME);
```

Listing 18: SAS macro and calls for categorical vs. numerical analysis in ASGMLIB.TESTING_DS

Bivariate Analysis of Categorical Variable & Continuous: Guarantee Income vs Gender

Output						
Analysis Variable : GUARANTEE_INCOME						
GENDER	N Obs	N	Mean	Std Dev	Minimum	Maximum
Female	70	70	1171.96	1979.82	0	11666.00
Male	286	286	1670.87	2433.94	0	24000.00

Table 34: Descriptive statistics of GUARANTEE_INCOME by GENDER in ASGMLIB.TESTING_DS

The bivariate analysis of GUARANTEE_INCOME by GENDER (Table 34) displays a difference in income between the male and female. Males have a larger mean guarantee income of 1,670.87 compared to 1,171.96 for females. However, the standard deviation for males is also larger, at 2,433.94, implying a wider range of income among males compared to females, who only have a standard deviation of 979.82.

Bivariate Analysis of Categorical Variable & Continuous: Loan Amount vs Employment

Output						
Analysis Variable : LOAN_AMOUNT						
EMPLOYMENT	N Obs	N	Mean	Std Dev	Minimum	Maximum
No	307	302	133.7218543	57.9927432	28.0000000	460.0000000
Yes	37	37	150.1891892	84.6797358	50.0000000	550.0000000

Table 35: Descriptive statistics of LOAN_AMOUNT by EMPLOYMENT in ASGMLIB.TESTING_DS

The bivariate analysis of LOAN_AMOUNT by EMPLOYMENT (Table 35) reveals that employed applicants tend to take out much larger loans. The mean loan amount for employed applicants is 150.19, compared to a mean loan amount of 133.72 for unemployed applicants. In addition, the standard deviation is higher for employed applicants (84.68) compared to the unemployed applicants (57.99), which could imply a greater variation in loan amounts for employed applicants.

Bivariate Analysis of Categorical Variable & Continuous: Loan Duration vs Qualification

Output						
Analysis Variable : LOAN_DURATION						
QUALIFICATION	N Obs	N	Mean	Std Dev	Minimum	Maximum
Graduate	283	279	340.2867384	66.9216513	6.0000000	480.0000000
Under Graduate	84	82	350.1951220	58.4884581	36.0000000	480.0000000

Table 36: Descriptive statistics of LOAN_DURATION by QUALIFICATION in ASGMLIB.TESTING_DS

The bivariate analysis of LOAN_DURATION by QUALIFICATION (Table 36) demonstrate a marginal difference in the loan duration between graduates and under graduates. The mean loan duration for graduate applicants is 340.29 days, which is only a little bite smaller than the mean of 350.20 days for under graduate applicant. The standard deviations are also nearly identical, implying that both have a similar spread in their loan durations.

Chapter 7

Data Cleaning and Imputation

Backup & Restore Utilities for ASGMLIB.TRAINING_DS and ASGMLIB.TESTING_DS

SAS Code

```
1 /* Enable compile-time notes for macros */
2 OPTIONS MCOMPILENOTE=ALL;
3 /* Macro: BKP_MAKE
4 Purpose : Create a backup copy of a dataset ( _BK)
5 Param    : PDS - Dataset */
6 %MACRO BKP_MAKE(PDS);
7 PROC SQL;
8   CREATE TABLE &PDS._BK AS
9   SELECT * FROM &PDS;
10 QUIT;
11 %MEND BKP_MAKE;
12 /* Macro: BKP_RESTORE
13 Purpose : Restore the dataset from its _BK copy
14 Param    : PDS - Dataset */
15 %MACRO BKP_RESTORE(PDS);
16 PROC SQL;
17   CREATE TABLE &PDS AS
18   SELECT * FROM &PDS._BK;
19 QUIT;
20 %MEND BKP_RESTORE;
21 /* Make backups before any imputation runs --- */
22 %BKP_MAKE(ASGMLIB.TRAINING_DS);
23 %BKP_MAKE(ASGMLIB.TESTING_DS);
```

Listing 19: SAS macros for creating and restoring dataset backups

Output

Table: ASGMLIB.TRAINING_DS_BK						
Columns	Total rows: 614	Total columns: 13	View:	Column names	Filter: (none)	Rows 1-100
<input checked="" type="checkbox"/> Select all				SME_LOAN_ID...	GEND...	MARITAL_STA...
<input checked="" type="checkbox"/> SME_LOAN_ID_NO	1	LP001002	Male	Not Married	0	Graduate
<input checked="" type="checkbox"/> GENDER	2	LP001003	Male	Married	1	Graduate
<input checked="" type="checkbox"/> MARITAL_STATUS	3	LP001005	Male	Married	0	Graduate

Figure 21: Created backup table
ASGMLIB.TRAINING_DS_BK

Table: ASGMLIB.TESTING_DS_BK						
Columns	Total rows: 367	Total columns: 13	View:	Column names	Filter: (none)	Rows 1-100
<input checked="" type="checkbox"/> Select all				SME_LOAN_ID...	GEND...	MARITAL_STA...
<input checked="" type="checkbox"/> SME_LOAN_ID_NO	1	LP001015	Male	Married	0	Graduate
<input checked="" type="checkbox"/> GENDER	2	LP001022	Male	Married	1	Graduate
<input checked="" type="checkbox"/> MARITAL_STATUS	3	LP001031	Male	Married	2	Graduate

Figure 22: Created backup table
ASGMLIB.TESTING_DS_BK

Before any imputation, both the training and testing datasets were backed up into separate copies (_BK). This ensures that the original datasets remain intact and can be restored if needed. Figures 21 and 22 confirm the successful creation of backup tables.

Categorical Imputation of Variables in ASGMLIB.TRAINING_DS and ASGMLIB.TESTING_DS

SAS Code

```
1 /* Enable compile-time notes for macros */
2 OPTIONS MCOMPILENOTE=ALL;
3 /* Macro: IMPUTE_CATE_MODE
4 Purpose : Impute a categorical variable's missing/blank values with its MODE.
5 Params   : PTITLEVAR - Pretty variable name for TITLES
6          PDS      - Dataset
7          PVAR     - Categorical variable name
8          PSTATS   - Output dataset to store counts/mode */
9 %MACRO IMPUTE_CATE_MODE(PTITLEVAR, PDS, PVAR, PSTATS);
10 /*Imputing values - STEP 1*/
11 TITLE1 "STEP 1: LIST THE DETAILS OF THE LOAN APPLICANTS";
12 TITLE2 "WHO SUBMITTED THEIR LOAN APPLICATIONS WITHOUT &PTITLEVAR DATA";
13 FOOTNOTE "----End----";
14 PROC SQL;
15   SELECT *
16   FROM &PDS t
17   WHERE (t.&PVAR IS MISSING OR t.&PVAR = '');
18 QUIT;
19 /*Imputing values - STEP 2*/
20 TITLE1 "STEP 2: COUNT THE TOTAL NUMBER OF THE LOAN APPLICANTS";
21 TITLE2 "WHO SUBMITTED THEIR LOAN APPLICATIONS WITHOUT &PTITLEVAR DATA";
22 FOOTNOTE "----End----";
23 PROC SQL;
24   SELECT COUNT(*) LABEL = "Total number of loan applicants"
25   FROM &PDS t
26   WHERE (t.&PVAR IS MISSING OR t.&PVAR = '');
27 QUIT;
28 /*Imputing values - STEP 3 */
29 TITLE1 "STEP 3: FIND THE MOD VALUE AND SAVE THE STATISTICS IN A TEMPORARY DATASET";
30 FOOTNOTE "----End----";
31 PROC SQL;
32   CREATE TABLE &PSTATS AS
33   SELECT t.&PVAR AS &PVAR, COUNT(*) AS counts
34   FROM &PDS t
35   WHERE ( (t.&PVAR IS NOT MISSING) OR (t.&PVAR NE '') )
36   GROUP BY t.&PVAR;
37 QUIT;
38 /*Imputing values - STEP 4: impute the missing values found...*/
39 TITLE1 "STEP 4: IMPUTE THE MISSING VALUES FOUND...";
40 FOOTNOTE "----End----";
41 PROC SQL;
42   UPDATE &PDS
43   SET &PVAR = (
44     SELECT to.&PVAR
45     FROM &PSTATS to
46     WHERE to.counts = (SELECT MAX(ti.counts) FROM &PSTATS ti)
47   )
48   WHERE (&PVAR IS MISSING OR &PVAR = '');
49 QUIT;
50 /*Imputing values - STEP 5(AI)*/
51 TITLE1 "STEP 5(AI): LIST THE DETAILS OF THE LOAN APPLICANTS";
52 TITLE2 "WHO SUBMITTED THEIR LOAN APPLICATIONS WITHOUT &PTITLEVAR DATA";
53 FOOTNOTE "----End----";
54 PROC SQL;
55   SELECT *
56   FROM &PDS t
57   WHERE (t.&PVAR IS MISSING OR t.&PVAR = '');
58 QUIT;
59 /*Imputing values - STEP 6(AI)*/
60 TITLE1 "STEP 6(AI): COUNT THE TOTAL NUMBER OF THE LOAN APPLICANTS";
61 TITLE2 "WHO SUBMITTED THEIR LOAN APPLICATIONS WITHOUT &PTITLEVAR DATA";
62 FOOTNOTE "----End----";
63 PROC SQL;
64   SELECT COUNT(*) LABEL = "Total number of loan applicants"
65   FROM &PDS t
66   WHERE (t.&PVAR IS MISSING OR t.&PVAR = '');
67 QUIT;
68 %MEND IMPUTE_CATE_MODE;
```

Listing 20: Mode-based imputation macro for (ASGMLIB.TRAINING_DS and ASGMLIB.TESTING_DS)

SAS Code

```
1 /* TRAINING: MARITAL_STATUS */
2 %IMPUTE_CATE_MODE(Marital_Status, ASGMLIB.TRAINING_DS, MARITAL_STATUS,
ASGMLIB.TRAINING_STATS_DS);

3
4 /* TRAINING: FAMILY_MEMBERS */
5 %IMPUTE_CATE_MODE(Family_Members, ASGMLIB.TRAINING_DS, FAMILY_MEMBERS,
ASGMLIB.TRAINING_STATS_DS);

6
7 /* TESTING: MARITAL_STATUS */
8 %IMPUTE_CATE_MODE(Family_Members, ASGMLIB.TESTING_DS, FAMILY_MEMBERS,
ASGMLIB.TESTING_STATS_DS);
```

Listing 21: Macro calls for imputing categorical variables for (ASGMLIB.TRAINING_DS and ASGMLIB.TESTING_DS

Categorical Imputation

Categorical variables were imputed using the mode strategy. The SAS macro followed these steps: (i) list applicants with missing or blank values, (ii) count total missing, (iii) build a temporary statistics table to identify the mode, (iv) replace missing values with the mode using an UPDATE query, and (v) re-list and re-count post-imputation to confirm that no missing values remained.

Imputation of MARITAL_STATUS in ASGMLIB.TRAINING_DS

Output

STEP 1: LIST THE DETAILS OF THE LOAN APPLICANTS WHO SUBMITTED THEIR LOAN APPLICATIONS WITHOUT Marital_Status DATA												
SME_LOAN_ID_NO	GENDER	MARITAL_STATUS	FAMILY_MEMBERS	QUALIFICATION	EMPLOYMENT	CANDIDATE_INCOME	GUARANTEE_INCOME	LOAN_AMOUNT	LOAN_DURATION	LOAN_HISTORY	LOAN_LOCATION	LOAN_APPROVAL_STATUS
LP001357	Male			Graduate	No	3816	754	160	360	1	City	Y
LP001760	Male			Graduate	No	4758	0	158	480	1	Town	Y
LP002393	Female			Graduate	No	10047	0	.	240	1	Town	Y

—End—

STEP 2: COUNT THE TOTAL NUMBER OF THE LOAN APPLICANTS WHO SUBMITTED THEIR LOAN APPLICATIONS WITHOUT Marital_Status DATA												
Total number of loan applicants												
3												

—End—

STEP 5(A): LIST THE DETAILS OF THE LOAN APPLICANTS WHO SUBMITTED THEIR LOAN APPLICATIONS WITHOUT Marital_Status DATA												
Total number of loan applicants												
0												

—End—

Figure 23: Training set: Pre- and post-imputation checks for MARITAL_STATUS (list+count before; list+count after)

Table: ASGMLIB.TRAINING_STATS_DS		View: Column names	CODE	LOG	RESULTS	OUTPUT DATA
Columns		Total rows: 2 Total columns: 2				
<input checked="" type="checkbox"/> Select all		MARIT	Total rows: 2 Total columns: 2			
<input checked="" type="checkbox"/> ▲ MARITAL_STATUS		1	Married			
<input checked="" type="checkbox"/> i counts		2	Not Married			
				counts		
				398		
				213		
						NOTE: 3 rows were updated in ASGMLIB.TRAINING_DS.

Figure 24: Training set: Temporary stats table for MARITAL_STATUS (mode counts)

Figure 25: Training set: SAS log confirming successful imputation (“3 rows were updated”).

Based on the outputs, MARITAL_STATUS was missing for 3 applicants (as shown in the pre-imputation checks, Fig. 23). The mode of this variable was identified as Married in the temporary statistics table (Table 24). A total of 3 rows were updated, confirmed in the SAS log (Fig. 25), and no applicants remained with missing values after imputation (Fig. 23).

Imputation of FAMILY_MEMBERS in ASGMLIB.TRAINING_DS

Output

STEP 1: LIST THE DETAILS OF THE LOAN APPLICANTS WHO SUBMITTED THEIR LOAN APPLICATIONS WITHOUT Family_Members DATA													
SME_LOAN_ID_NO	GENDER	MARITAL_STATUS	FAMILY_MEMBERS	QUALIFICATION	EMPLOYMENT	CANDIDATE_INCOME	GUARANTEE_INCOME	LOAN_AMOUNT	LOAN_DURATION	LOAN_HISTORY	LOAN_LOCATION	LOAN_APPROVAL_STATUS	
LP001350	Male	Married		Graduate	No	13650	0	.	360	1	City	Y	
LP001357	Male	Married		Graduate	No	3816	754	160	360	1	City	Y	
LP001426	Male	Married		Graduate	No	5667	2687	180	360	1	Village	Y	
LP001754	Male	Married		Under Graduate	Yes	4735	0	138	360	1	City	N	
LP001760	Male	Married		Graduate	No	4758	0	158	480	1	Town	Y	
LP001945	Female	Not Married		Graduate	No	5417	0	143	480	0	City	N	
LP001972	Male	Married		Under Graduate	No	2875	1750	105	360	1	Town	Y	
LP002100	Male	Not Married		Graduate	No	2833	0	71	360	1	City	Y	
LP002106	Male	Married		Graduate	Yes	5503	4490	70	.	1	Town	Y	
LP002130	Male	Married		Under Graduate	No	3223	3230	152	360	0	Village	N	
LP002144	Female	Not Married		Graduate	No	3813	0	116	180	1	City	Y	
LP002393	Female	Married		Graduate	No	10047	0	.	240	1	Town	Y	
LP002682	Male	Married		Under Graduate	No	3074	1800	123	360	0	Town	N	
LP002847	Male	Married		Graduate	No	5116	1451	165	360	0	City	N	
LP002943	Male	Not Married		Graduate	No	2987	0	88	360	0	Town	N	

—End—

STEP 2: COUNT THE TOTAL NUMBER OF THE LOAN APPLICANTS WHO SUBMITTED THEIR LOAN APPLICATIONS WITHOUT Family_Members DATA													
Total number of loan applicants													
15													

—End—

STEP 5(A): LIST THE DETAILS OF THE LOAN APPLICANTS WHO SUBMITTED THEIR LOAN APPLICATIONS WITHOUT Family_Members DATA													
Total number of loan applicants													
0													

—End—

STEP 6(A): COUNT THE TOTAL NUMBER OF THE LOAN APPLICANTS WHO SUBMITTED THEIR LOAN APPLICATIONS WITHOUT Family_Members DATA													
Total number of loan applicants													
0													

—End—

Figure 26: Training set: Pre- and post-imputation checks for FAMILY_MEMBERS (list+count before; list+count after)

Table: _TEMP1.TRAINING_STATS_DS		View: Column names	Filter: (none)	CODE		LOG	RESULTS	OUTPUT DATA
Columns	⑤	Total rows: 4 Total columns: 2		Rows 1-4	Counts			
⑤ Select all		Total rows: 4 Total columns: 2						
<input checked="" type="checkbox"/> FAMILY_MEMBERS	1	0			345			
<input checked="" type="checkbox"/> counts	2	1			102			
	3	2			101			
	4	3+			51			

Figure 27: Training set: Temporary stats table for MARITAL_STATUS (mode counts)

Figure 28: Training set: SAS log confirming successful imputation (“rows were updated”).

Based on the outputs, FAMILY_MEMBERS was missing for 15 applicants (as indicated in the pre-imputation checks, Fig. 26). The mode was found to be 0 family members in the statistics table (Table 27). The SAS log confirmed that 15 rows were updated (Fig. 28), and the after imputation (See Fig. 26) verified that no records remained with missing values.

Imputation of FAMILY_MEMBERS in ASGMT.TESTING_DS

Output

STEP 1: LIST THE DETAILS OF THE LOAN APPLICANTS WHO SUBMITTED THEIR LOAN APPLICATIONS WITHOUT Family_Members DATA

SME_LOAN_ID_NO	GENDER	MARITAL_STATUS	FAMILY_MEMBERS	QUALIFICATION	EMPLOYMENT	CANDIDATE_INCOME	GUARANTEE_INCOME	LOAN_AMOUNT	LOAN_DURATION	LOAN_HISTORY	LOAN_LOCATION	LOAN_APPROVAL_STATUS
LP001237	Male	Married		Under Graduate	No	4163	1475	162	360	1	City	
LP001366	Female	Not Married		Graduate	No	3250	0	95	360	1	Town	
LP001587	Male	Married		Graduate	No	4082	0	93	360	1	Town	
LP001789		Not Married		Graduate	No	3333	1250	110	360	1	Town	
LP002111	Male	Married		Graduate	No	3016	1300	100	360	-	City	
LP002360	Male	Married		Graduate	No	10000	0	-	360	1	City	
LP002385	Male	Married		Graduate	No	3863	0	70	300	1	Town	
LP002441	Male	Not Married		Graduate	No	3579	3308	138	360	-	Town	
LP002654	Female	Not Married		Graduate	Yes	14987	0	177	360	1	Village	
LP002754	Male	Not Married		Graduate	No	2066	2108	104	84	1	City	

—End—

STEP 2: COUNT THE TOTAL NUMBER OF THE LOAN APPLICANTS WHO SUBMITTED THEIR LOAN APPLICATIONS WITHOUT Family_Members DATA

Total number of loan applicants	10
---------------------------------	----

—End—

STEP 5(A): LIST THE DETAILS OF THE LOAN APPLICANTS WHO SUBMITTED THEIR LOAN APPLICATIONS WITHOUT Family_Members DATA

—End—

STEP 6(A): COUNT THE TOTAL NUMBER OF THE LOAN APPLICANTS WHO SUBMITTED THEIR LOAN APPLICATIONS WITHOUT Family_Members DATA

Total number of loan applicants	0
---------------------------------	---

—End—

Figure 29: Testing set: Pre- and post-imputation checks for FAMILY_MEMBERS (list+count before; list+count after)

Table: TEMP1.TESTING_STATS_DS | View: Column names | Filter: (none)

Columns

- Select all
- FAMILY_MEMBERS
- counts

Total rows: 4 Total columns: 2

FAMILY_MEMBERS counts

FAMILY_MEMBERS	counts
1 0	200
2 1	58
3 2	59
4 3+	40

ROWS 1-4

CODE LOG RESULTS OUTPUT DATA

▼ Errors, Warnings, Notes

► Errors (0)

► Warnings (0)

► Notes (9)

NOTE: 10 rows were updated in ASGMLIB.TESTING_DS.

Figure 30: Testing set: Temporary stats table for MARITAL_STATUS (mode counts)

Figure 31: Testing set: SAS log confirming successful imputation (“rows were updated”).

Based on the outputs, FAMILY_MEMBERS was missing for 5 applicants (as shown in the pre-imputation steps, Fig. 29). The mode of this variable was again 0, identified in the statistics table (Table 30). The SAS log reported that 5 rows were successfully updated (Fig. 31), and this was further confirmed in the after imputation check (Fig. 29).

Continuous Imputation of Variables in ASGMLIB.TRAINING_DS and ASGMLIB.TESTING_DS

SAS Code

```
1  /* Enable compile-time notes for macros */
2  OPTIONS MCOMPILENOTE=ALL;
3
4  /* Macro: IMPUTE_CONT_MEAN
5   Purpose : Impute missing values of a continuous variable using MEAN
6   Params  : PTITLEVAR - Pretty variable name for titles (e.g., Loan_Amount)
7   PDS      - Dataset (libref.member)
8   PVAR     - Continuous variable name
9 */
10 %MACRO IMPUTE_CONT_MEAN(PTITLEVAR, PDS, PVAR);
11
12 *****
13 STEP 1: List applicants with missing values (continuous)
14 *****
15 TITLE1 "STEP 1: LIST THE DETAILS OF THE LOAN APPLICANTS";
16 TITLE2 "WHO SUBMITTED THEIR LOAN APPLICATIONS WITHOUT &PTITLEVAR DATA";
17 FOOTNOTE "----End----";
18 PROC SQL;
19   SELECT *
20   FROM &PDS t
21   WHERE (t.&PVAR IS MISSING OR t.&PVAR = .);
22 QUIT;
23
24 *****
25 STEP 2: Count applicants with missing values (continuous)
26 *****
27 TITLE1 "STEP 2: COUNT THE TOTAL NUMBER OF THE LOAN APPLICANTS";
28 TITLE2 "WHO SUBMITTED THEIR LOAN APPLICATIONS WITHOUT &PTITLEVAR DATA";
29 FOOTNOTE "----End----";
30 PROC SQL;
31   SELECT COUNT(*) LABEL='Number of Loan Applicants'
32   FROM &PDS t
33   WHERE (t.&PVAR IS MISSING OR t.&PVAR = .);
34 QUIT;
35
36 *****
37 STEP 3: Impute missing values with MEAN
38 *****
39 PROC STDIZE DATA=&PDS REONLY METHOD=MEAN OUT=&PDS;
40   VAR &PVAR;
41 QUIT;
42
43 *****
44 STEP 4(AI): Post-imputation listing (should be empty)
45 *****
46 TITLE1 "STEP 4(AI): LIST THE DETAILS OF THE LOAN APPLICANTS";
47 TITLE2 "WHO SUBMITTED THEIR LOAN APPLICATIONS WITHOUT &PTITLEVAR DATA";
48 FOOTNOTE "----End----";
49 PROC SQL;
50   SELECT *
51   FROM &PDS t
52   WHERE (t.&PVAR IS MISSING OR t.&PVAR = .);
53 QUIT;
54
55 *****
56 STEP 5(AI): Post-imputation count (should be zero)
57 *****
58 TITLE1 "STEP 5(AI): COUNT THE TOTAL NUMBER OF THE LOAN APPLICANTS";
59 TITLE2 "WHO SUBMITTED THEIR LOAN APPLICATIONS WITHOUT &PTITLEVAR DATA";
60 FOOTNOTE "----End----";
61 PROC SQL;
62   SELECT COUNT(*) LABEL='Number of Loan Applicants'
63   FROM &PDS t
64   WHERE (t.&PVAR IS MISSING OR t.&PVAR = .);
65 QUIT;
66
67 %MEND IMPUTE_CONT_MEAN;
```

Listing 22: Mean-based imputation macro for (ASGMLIB.TRAINING_DS and ASGMLIB.TESTING_DS)

SAS Code

```
1 /* TRAINING: LOAN_AMOUNT */
2 %IMPUTE_CONT_MEAN(Loan_Amount, ASGMLIB.TRAINING_DS, LOAN_AMOUNT);
3
4 /* TRAINING: LOAN_DURATION */
5 %IMPUTE_CONT_MEAN(Loan_Duration, ASGMLIB.TRAINING_DS, LOAN_DURATION);
6
7 /* TESTING: LOAN_AMOUNT */
8 %IMPUTE_CONT_MEAN(Loan_Amount, ASGMLIB.TESTING_DS, LOAN_AMOUNT);
```

Listing 23: Macro calls for imputing continuos variables for (ASGMLIB.TRAINING_DS and ASGMLIB.TESTING_DS

Continuous Imputation

Continuous variables were imputed using the mean strategy. The SAS macro followed these steps: (i) list applicants with missing values, (ii) count total missing, (iii) replace missing values with the mean using PROC STDIZE, and (iv) re-list and re-count post-imputation to confirm that all missing values had been handled.

Imputation of LOAN_AMOUNT in ASGMLIB.TRAINING_DS

Output

STEP 1: LIST THE DETAILS OF THE LOAN APPLICANTS WHO SUBMITTED THEIR LOAN APPLICATIONS WITHOUT Loan_Amount DATA													
SME_LOAN_ID_NO	GENDER	MARITAL_STATUS	FAMILY_MEMBERS	QUALIFICATION	EMPLOYMENT	CANDIDATE_INCOME	GUARANTEE_INCOME	LOAN_AMOUNT	LOAN_DURATION	LOAN_HISTORY	LOAN_LOCATION	LOAN_APPROVAL_STATUS	
LP001002	Male	Not Married	0	Graduate	No	5849	0	.	360	1	City	Y	
LP001106	Male	Married	0	Graduate	No	2275	2067	.	360	1	City	Y	
LP001213	Male	Married	1	Graduate	No	4945	0	.	360	0	Village	N	
LP001266	Male	Married	1	Graduate	Yes	2395	0	.	360	1	Town	Y	
LP001326	Male	Not Married	0	Graduate	.	6782	0	.	360	.	City	N	
LP001350	Male	Married	0	Graduate	No	13650	0	.	360	1	City	Y	
LP001356	Male	Married	0	Graduate	No	4652	3583	.	360	1	Town	Y	
LP001392	Female	Not Married	1	Graduate	Yes	7451	0	.	360	1	Town	Y	
LP001449	Male	Not Married	0	Graduate	No	3865	1640	.	360	1	Village	Y	
LP001682	Male	Married	3+	Under Graduate	No	3992	0	.	180	1	City	N	
LP001922	Male	Married	0	Graduate	No	20667	0	.	360	1	Village	N	
LP001990	Male	Not Married	0	Under Graduate	No	2000	0	.	360	1	City	N	
LP002054	Male	Married	2	Under Graduate	No	3601	1590	.	360	1	Village	Y	
LP002113	Female	Not Married	3+	Under Graduate	No	1830	0	.	360	0	City	N	
LP002243	Male	Married	0	Under Graduate	No	3010	3136	.	360	0	City	N	
LP002393	Female	Married	0	Graduate	No	10047	0	.	240	1	Town	Y	
LP002401	Male	Married	0	Graduate	No	2213	1125	.	360	1	City	Y	
LP002533	Male	Married	2	Graduate	No	2947	1603	.	360	1	City	N	
LP002697	Male	Not Married	0	Graduate	No	4690	2087	.	360	1	Town	N	
LP002778	Male	Married	2	Graduate	Yes	6033	0	.	360	0	Village	N	
LP002784	Male	Married	1	Under Graduate	No	2492	2375	.	360	1	Village	Y	
LP002960	Male	Married	0	Under Graduate	No	2400	3800	.	180	1	City	N	

—End—

STEP 2: COUNT THE TOTAL NUMBER OF THE LOAN APPLICANTS WHO SUBMITTED THEIR LOAN APPLICATIONS WITHOUT Loan_Amount DATA													
Number of Loan Applicants	22												

—End—

STEP 4(AI): LIST THE DETAILS OF THE LOAN APPLICANTS WHO SUBMITTED THEIR LOAN APPLICATIONS WITHOUT Loan_Amount DATA													
--	--	--	--	--	--	--	--	--	--	--	--	--	--

—End—

STEP 5(AI): COUNT THE TOTAL NUMBER OF THE LOAN APPLICANTS WHO SUBMITTED THEIR LOAN APPLICATIONS WITHOUT Loan_Amount DATA													
Number of Loan Applicants	0												

—End—

Figure 32: Training: pre-/post-imputation checks for LOAN_AMOUNT (Steps 1–5(AI))

Based on the outputs, LOAN_AMOUNT was missing for 22 applicants (as shown in the pre-imputation checks, Fig. 32). The variable was imputed using the mean, replacing all missing values. After imputation, no applicants remained with missing loan amount values (Fig. 32).

Imputation of LOAN_DURATION in ASGMLIB.TRAINING_DS

Output

STEP 1: LIST THE DETAILS OF THE LOAN APPLICANTS WHO SUBMITTED THEIR LOAN APPLICATIONS WITHOUT Loan_Duration DATA													
SME_LOAN_ID_NO	GENDER	MARITAL_STATUS	FAMILY_MEMBERS	QUALIFICATION	EMPLOYMENT	CANDIDATE_INCOME	GUARANTEE_INCOME	LOAN_AMOUNT	LOAN_DURATION	LOAN_HISTORY	LOAN_LOCATION	LOAN_APPROVAL_STATUS	
LP001041	Male	Married	0	Graduate		2600	3500	115	.	1	City	Y	
LP001109	Male	Married	0	Graduate	No	1828	1330	100	.	0	City	N	
LP001136	Male	Married	0	Under Graduate	Yes	4695	0	96	.	1	City	Y	
LP001137	Female	Not Married	0	Graduate	No	3410	0	88	.	1	City	Y	
LP001250	Male	Married	3+	Under Graduate	No	4755	0	95	.	0	Town	N	
LP001391	Male	Married	0	Under Graduate	No	3572	4114	152	.	0	Village	N	
LP001574	Male	Married	0	Graduate	No	3707	3166	182	.	1	Village	Y	
LP001669	Female	Not Married	0	Under Graduate	No	1907	2365	120	.	1	City	Y	
LP001749	Male	Married	0	Graduate	No	7578	1010	175	.	1	Town	Y	
LP001770	Male	Not Married	0	Under Graduate	No	3189	2598	120	.	1	Village	Y	
LP002106	Male	Married	0	Graduate	Yes	503	4490	70	.	1	Town	Y	
LP002188	Male	Not Married	0	Graduate	No	5124	0	124	.	0	Village	N	
LP002357	Female	Not Married	0	Under Graduate	No	2720	0	80	.	0	City	N	
LP002362	Male	Married	1	Graduate	No	7250	1667	110	.	0	City	N	

—End—

STEP 2: COUNT THE TOTAL NUMBER OF THE LOAN APPLICANTS WHO SUBMITTED THEIR LOAN APPLICATIONS WITHOUT Loan_Duration DATA													
Number of Loan Applicants													
14													

—End—

STEP 4(A): LIST THE DETAILS OF THE LOAN APPLICANTS WHO SUBMITTED THEIR LOAN APPLICATIONS WITHOUT Loan_Duration DATA													
Number of Loan Applicants													
0													

—End—

STEP 5(A): COUNT THE TOTAL NUMBER OF THE LOAN APPLICANTS WHO SUBMITTED THEIR LOAN APPLICATIONS WITHOUT Loan_Duration DATA													
Number of Loan Applicants													
0													

—End—

Figure 33: Training: pre-/post-imputation checks for LOAN_DURATION (Steps 1–5(AI))

Based on the outputs, LOAN_DURATION was missing for 14 applicants (as shown in the pre-imputation checks, Fig. 33). The variable was imputed using the mean, ensuring all missing values were replaced. After imputation, no applicants remained with missing loan duration values (Fig. 33).

Imputation of LOAN_AMOUNT in ASGMLIB.TESTING_DS

Output

STEP 1: LIST THE DETAILS OF THE LOAN APPLICANTS WHO SUBMITTED THEIR LOAN APPLICATIONS WITHOUT Loan_Amount DATA													
SME_LOAN_ID_NO	GENDER	MARITAL_STATUS	FAMILY_MEMBERS	QUALIFICATION	EMPLOYMENT	CANDIDATE_INCOME	GUARANTEE_INCOME	LOAN_AMOUNT	LOAN_DURATION	LOAN_HISTORY	LOAN_LOCATION	LOAN_APPROVAL_STATUS	
LP001415	Male	Married	1	Graduate	No	3413	4053	.	360	1	Town		
LP001542	Female	Married	0	Graduate	No	2262	0	.	480	0	Town		
LP002057	Male	Married	0	Under Graduate	No	13083	0	.	360	1	Village		
LP002360	Male	Married	0	Graduate	No	10000	0	.	360	1	City		
LP002593	Male	Married	1	Graduate	No	8333	4000	.	360	1	City		

—End—

STEP 2: COUNT THE TOTAL NUMBER OF THE LOAN APPLICANTS WHO SUBMITTED THEIR LOAN APPLICATIONS WITHOUT Loan_Amount DATA													
Number of Loan Applicants													
5													
—End—													

STEP 4(AI): LIST THE DETAILS OF THE LOAN APPLICANTS WHO SUBMITTED THEIR LOAN APPLICATIONS WITHOUT Loan_Amount DATA													
Number of Loan Applicants													
0													
—End—													

STEP 5(AI): COUNT THE TOTAL NUMBER OF THE LOAN APPLICANTS WHO SUBMITTED THEIR LOAN APPLICATIONS WITHOUT Loan_Amount DATA													
Number of Loan Applicants													
0													
—End—													

Figure 34: Testing: pre-/post-imputation checks for LOAN_AMOUNT (Steps 1–5(AI))

Based on the outputs, LOAN_AMOUNT was missing for 5 applicants (as shown in the pre-imputation checks, Fig. 34). The mean was used to impute these values, and after imputation, no applicants remained with missing values in loan amount (Fig. 34).

Chapter 8

Model Creation

Logistic Regression Model Creation in ASGMLIB.TRAINING_DS

SAS Code

```
1  /* Train logistic regression model */
2  TITLE "Logistic Regression - Model Creation (Training)";
3  PROC LOGISTIC DATA=ASGMLIB.TRAINING_DS
4          OUTMODEL=ASGMLIB.TRAINING_DS_LR_MODEL;      /* save model */
5  CLASS
6      GENDER
7      MARITAL_STATUS
8      FAMILY_MEMBERS
9      QUALIFICATION
10     EMPLOYMENT
11     LOAN_HISTORY
12     LOAN_LOCATION
13     ;
14 MODEL LOAN_APPROVAL_STATUS =
15     GENDER
16     MARITAL_STATUS
17     FAMILY_MEMBERS
18     QUALIFICATION
19     EMPLOYMENT
20     LOAN_HISTORY
21     LOAN_LOCATION
22     CANDIDATE_INCOME
23     GUARANTEE_INCOME
24     LOAN_AMOUNT
25     LOAN_DURATION
26     ;
27     /* Save predicted probabilities to an output dataset */
28     OUTPUT OUT=ASGMLIB.TRAINING_OUT_DS
29         P=PRED_PROB;    /* predicted probability of modeled event */
30 RUN;
```

Listing 24: SAS code for training the logistic regression model on ASGMLIB.TRAINING_DS

Output

Model Information	
Data Set	ASGMLIB.TRAINING_DS
Response Variable	LOAN_APPROVAL_STATUS
Number of Response Levels	2
Model	binary logit
Optimization Technique	Fisher's scoring

Figure 35: Model Information for logistic regression in ASGMLIB.TRAINING_DS

Figure 35 summarizes the modeling setup: binary logistic regression on LOAN_APPROVAL_STATUS using Fisher's scoring. Categorical predictors were specified in the CLASS statement, ensuring correct handling of nominal variables.

Output

Model Convergence Status	
Convergence criterion (GCONV=1E-8) satisfied.	

Figure 36: Convergence status of logistic regression (criterion satisfied)

The optimization converged successfully as shown in Figure 36, where the convergence criterion (GCONV=1E-8) is satisfied. This ensures stable coefficient estimates and valid inference.

Output

Number of Observations Read	614
Number of Observations Used	526

Figure 37: Number of observations read and used in model estimation

Figure 37 shows that 614 observations were read, of which 526 were used. This implies that 88 records were excluded, most likely due to missing values. Minimizing exclusions allows us to preserve statistical power.

Output

Model Fit Statistics		
Criterion	Intercept Only	Intercept and Covariates
AIC	663.975	501.945
SC	668.240	565.925
-2 Log L	661.975	471.945

Figure 38: Model fit statistics (AIC, SC, and $-2 \log L$)

Model fit improved markedly as shown in Figure 38. AIC decreased from 663.98 to 501.95, SC from 668.24 to 565.93, and $-2 \log L$ from 661.98 to 471.95. These reductions indicate that the full model fits substantially better than the intercept-only model.

Output

Type 3 Analysis of Effects			
Effect	DF	Wald Chi-Square	Pr > ChiSq
GENDER	1	0.4878	0.4849
MARITAL_STATUS	1	3.7881	0.0516
FAMILY_MEMBERS	3	1.7620	0.6232
QUALIFICATION	1	3.3770	0.0661
EMPLOYMENT	1	0.1640	0.6855
LOAN_HISTORY	1	81.0825	<.0001
LOAN_LOCATION	2	12.5217	0.0019
CANDIDATE_INCOME	1	0.0010	0.9743
GUARANTEE_INCOME	1	1.2406	0.2654
LOAN_AMOUNT	1	2.3906	0.1221
LOAN_DURATION	1	0.0705	0.7906

Figure 39: Type 3 analysis of effects (Wald chi-square tests)

The Wald chi-square tests in Figure 39 show that LOAN_HISTORY ($p < 0.0001$) and LOAN_LOCATION ($p = 0.0019$) are significant predictors. MARITAL_STATUS ($p = 0.0516$) and QUALIFICATION ($p = 0.0661$) are borderline significant, while all other predictors were not statistically significant at $\alpha = 0.05$.

Results Summary

The logistic regression converged and showed improved fit compared to the intercept-only baseline (AIC 663.98 → 501.95; SC 668.24 → 565.93; $-2 \log L$ 661.98 → 471.95). LOAN_HISTORY and LOAN_LOCATION emerged as strong predictors of approval, while marital status and qualification were borderline. Other demographic or financial variables did not reach significance, indicating potential improvement from further cleaning.

Logistic Regression Scoring on ASGMLIB.TESTING_DS

SAS Code — Scoring the Testing Set

```

1 /*Predict the Loan Approval Status using the model created*/
2 TITLE "Logistic Regression - Scoring (Testing)";
3 PROC LOGISTIC INMODEL=ASGMLIB.TRAINING_DS_LR_MODEL; /* load saved model */
4     SCORE DATA=ASGMLIB.TESTING_DS
5             OUT=ASGMLIB.TESTING_LAS_PRED_TP086704_DS;
6 QUIT;
7
8 TITLE "Report - Loan Approval Status Predicted";
9 FOOTNOTE "----End----";
10 /* Display the details of the loan approval status predicted */
11 PROC SQL;
12     SELECT *
13     FROM ASGMLIB.TESTING_LAS_PRED_TP086704_DS;
14 QUIT;

```

Listing 25: SAS code for scoring and predicting loan approval status on ASGMLIB.TESTING_DS

Outputs

Columns		Total rows: 367 Total columns: 17	View: Column names	Filter: (none)	Rows 1-100	P_N	P_Y
<input checked="" type="checkbox"/>	Select all						
<input checked="" type="checkbox"/>	SME_LOAN_ID_NO	1	City			0.1681370992	0.8318629008
<input checked="" type="checkbox"/>	GENDER	1	City			0.231497332	0.768502668
<input checked="" type="checkbox"/>	MARITAL_STATUS	1	City			0.1814165065	0.8185834935
<input checked="" type="checkbox"/>	FAMILY_MEMBERS	. City					
<input checked="" type="checkbox"/>	QUALIFICATION	1	City			0.3504984817	0.6495015183
<input checked="" type="checkbox"/>	EMPLOYMENT	1	City			0.3406633502	0.6593366498
<input checked="" type="checkbox"/>	CANDIDATE_INCOME	1	Town			0.2585814233	0.7414185767
<input checked="" type="checkbox"/>	GUARANTEE_INCOME	0	Village			0.9322651935	0.0677348065
<input checked="" type="checkbox"/>	LOAN_AMOUNT	1	City				
<input checked="" type="checkbox"/>	LOAN_DURATION	1	Town				
						0.2189515295	0.7810484705

Figure 40: Predicted loan approval status results for ASGMLIB.TESTING_DS

Report - Loan Approval Status Predicted																
SME_LOAN_ID_NO	GENDER	MARITAL_STATUS	FAMILY_MEMBERS	QUALIFICATION	EMPLOYMENT	CANDIDATE_INCOME	GUARANTEE_INCOME	LOAN_AMOUNT	LOAN_DURATION	LOAN_HISTORY	LOAN_LOCATION	From: LOAN_APPROVAL_STATUS	Ind: LOAN_APPROVAL_STATUS	Predicted Probability: LOAN_APPROVAL_STATUS=N	Predicted Probability: LOAN_APPROVAL_STATUS=Y	
LP001015	Male	Married	0	Graduate	No	5720	0	110	360	1	City		Y	0.1681370992	0.8318629008	
LP001022	Male	Married	1	Graduate	No	3076	1500	126	360	1	City		Y	0.231497332	0.768502668	
LP001029	Male	Married	2	Graduate	No	2000	1600	208	360	1	City		Y	0.1814165065	0.8185834935	
LP001035	Male	Married	2	Under Graduate	No	2340	2548	100	360	1	City		Y	0.3504984817	0.6495015183	
LP001051	Male	Not Married	0	Under Graduate	No	3276	0	78	360	1	City		Y	0.3406633502	0.6593366498	
LP001054	Male	Married	0	Under Graduate	Yes	2165	3422	152	360	1	City		Y	0.2585814233	0.7414185767	
LP001058	Female	Not Married	1	Under Graduate	No	2228	0	59	360	1	Town		Y	0.9322651935	0.0677348065	
LP001059	Male	Married	2	Under Graduate	No	3881	0	147	360	0	Village		N	0.0677348065	0.9322651935	
LP001069	Male	Not Married	2	Under Graduate	No	13933	0	200	360	1	City		Y	0.2189515295	0.7810484705	
LP001077	Male	Not Married	0	Under Graduate	No	2400	2400	123	360	1	Town		Y	0.2189515295	0.7810484705	
LP001078	Male	Not Married	0	Under Graduate	No	3091	0	90	360	1	City		Y	0.3504984817	0.6495015183	
LP001082	Male	Married	1	Graduate	No	2188	1516	162	360	1	Town		Y	0.3406633502	0.6593366498	
LP001083	Male	Not Married	3+	Graduate	No	4196	0	40	180	1	City		Y	0.2585814233	0.7414185767	
LP001094	Male	Married	2	Graduate	No	12173	0	166	360	1	Town		Y	0.3504984817	0.6495015183	
LP001098	Female	Not Married	0	Graduate	No	4666	0	124	360	1	Town		Y	0.157248	0.842154	

Table 37: Detailed prediction output for ASGMLIB.TESTING_DS

The logistic regression model was applied to the testing dataset (ASGMLIB.TESTING_DS), and predicted loan approval statuses were generated. Figure 40 summarizes the predicted outcomes for each applicant, while Table 37 provides a more detailed view of the predictions. These results confirm that the trained logistic regression model can be successfully deployed to score new applications and support decision making for loan applications.

Part 3

Report Generation, Data Visualization & Conclusion

Chapter 9

Report Generation

Report Generation of Predicted Loan Approval Status in ASGMLIB.TESTING_DS

SAS Code

```
1  /* Close any open ODS destinations (clean start) */
2  ods html close;
3  ods pdf close;
4  ods listing close;
5
6  /* PDF output location & basic options */
7  ods pdf file="/home/u64179868/sasuser.v94/DAP_FT_JUN_2025_ASGMT_TP086704/LAS_REPORT_TP086704.pdf";
8  options nodate;
9
10 /* Report titles/footers */
11 title1 "Predicted Bank Loan Approval Status Report";
12 title2 "APU, TPM";
13 footnote " End of Report ";
14
15 /* Main report */
16 ods proclabel "Loan Approval Predictions";
17 proc report data= ASGMLIB.TESTING_LAS_PRED_TP086704_DS nowindows;
18   by SME_LOAN_ID_NO;
19
20 /* Column definitions (labels chosen for clean PDF headings) */
21 define SME_LOAN_ID_NO / group "Loan ID";
22 define GENDER / group "Gender";
23 define MARITAL_STATUS / group "Marital Status";
24 define FAMILY_MEMBERS / group "Family Members";
25 define CANDIDATE_INCOME / group "Monthly Income";
26 define GUARANTEE_INCOME / group "Co-Applicant's Income";
27 define LOAN_AMOUNT / group "Loan Amount";
28 define LOAN_DURATION / group "Loan Duration";
29 define LOAN_HISTORY / group "Loan History";
30 define LOAN_LOCATION / group "Loan Location";
31 run;
```

Listing 26: SAS code used to generate the PDF report with PROC REPORT

Output

DAP_FT_JUN_2025_ASGMT_TP086704.pdf
LAS_REPORT_TP086704.pdf
my_dap_project_tp086704.sas

Figure 41: Folder view showing the generated PDF LAS_REPORT_TP086704.pdf

CODE LOG RESULTS
+ Errors, Warnings, Notes
+ Errors
+ Warnings(45)
+ Notes(49)
76 ODS LISTING CLOSED;
77 /* PDF output location & basic options */
78 ods pdf file='/home/u417986/reuser.v94/DAP_FT_JUN_2025_ASGMT_TP086704/LAS_REPORT_TP086704.pdf';
NOTE: Writing ODS PDF output to DISK destination
79 '/home/u417986/reuser.v94/DAP_FT_JUN_2025_ASGMT_TP086704/LAS_REPORT_TP086704.pdf', printer 'PDF'.
80
options nodate;

Figure 42: ODS log confirming PDF was written to disk (highlighted path)

Output

Predicted Bank Loan Approval Status Report APU, TPM SME_LOAN_ID_NO=LP001022																
Loan ID	Gender	Marital Status	Family Members	QUALIFICATION	EMPLOYMENT	Monthly Income	Co-Applicant's Income	Loan Amount	Loan Duration	Loan History	Loan Location	LOAN_APPROVAL_STATUS	From:	Info:	Predicted Probability: LOAN_APPROVAL_STATUS=N	Predicted Probability: LOAN_APPROVAL_STATUS=Y
LP001022	Male	Married	1	Graduate	No	3076	1500	126	360	1	City	N	Y	0.2214973	0.7685027	

Figure 43: Results table (sample row) with predicted probabilities

Predicted Bank Loan Approval Status Report APU, TPM SME_LOAN_ID_NO=LP001031																
Loan ID	Gender	Marital Status	Family Members	QUALIFICATION	EMPLOYMENT	Monthly Income	Co-Applicant's Income	Loan Amount	Loan Duration	Loan History	Loan Location	LOAN_APPROVAL_STATUS	From:	Info:	Predicted Probability: LOAN_APPROVAL_STATUS=N	Predicted Probability: LOAN_APPROVAL_STATUS=Y
LP001031	Male	Married	2	Graduate	No	5000	1800	208	360	1	City	N	Y	0.1814165	0.8185835	

Figure 44: Results table (second sample row) with predicted probabilities

Output

Predicted Bank Loan Approval Status Report APU, TPM SME_LOAN_ID_NO=LP001022																
Loan ID	Gender	Marital Status	Family Members	QUALIFICATION	EMPLOYMENT	Monthly Income	Co-Applicant's Income	Loan Amount	Loan Duration	Loan History	Loan Location	LOAN_APPROVAL_STATUS	From:	Info:	Predicted Probability: LOAN_APPROVAL_STATUS=N	Predicted Probability: LOAN_APPROVAL_STATUS=Y
LP001022	Male	Married	1	Graduate	No	3076	1500	126				N	Y	0.234973	0.765027	

Predicted Bank Loan Approval Status Report APU, TPM SME_LOAN_ID_NO=LP001031																
Loan ID	Gender	Marital Status	Family Members	QUALIFICATION	EMPLOYMENT	Monthly Income	Co-Applicant's Income	Loan Amount	Loan Duration	Loan History	Loan Location	LOAN_APPROVAL_STATUS	From:	Info:	Predicted Probability: LOAN_APPROVAL_STATUS=N	Predicted Probability: LOAN_APPROVAL_STATUS=Y
LP001031	Male	Married	2	Graduate	No	5000	1800	208				N	Y	0.1814165	0.8185835	

Figure 45: PDF report page showing SME_LOAN_ID_NO = LP001022

Figure 46: PDF report page showing SME_LOAN_ID_NO = LP001031

The code in Listing 26 closes open ODS destinations, opens an ODS PDF target, and renders a grouped PROC REPORT with clear column labels. The folder view and ODS note (Figures 41–42) verify that LAS_REPORT_TP086704.pdf was created at the specified path. Sample rows from the Results window (Figures 43–44) illustrate the predicted probabilities stored in the scored dataset. Figures 45–46 show the final PDF layout per loan ID, including both probabilities P_N and P_Y.

Complex Report Generation of Predicted Loan Approval Status in ASGMLIB.TESTING_DS

SAS Code

```
1  /* Enable compile-time notes for macros */
2  OPTIONS MCOMPILENOTE=ALL;
3
4  /* Macro: MACRO_RPT_3
5   Purpose : Generate a complex report filtered by loan location (CITY/VILLAGE/TOWN)
6   and create a location-specific output table with derived fields.
7   Params  : PLOAN_LOCATION - Location filter (CITY / VILLAGE / TOWN)
8   */
9  %MACRO MACRO_RPT_3(PLOAN_LOCATION);
10
11 /* Map location parameter to dataset suffix */
12 %LOCAL SFX OUTDS;
13 %IF %UPCASE(&PLOAN_LOCATION) = CITY    %THEN %LET SFX = CTY;
14 %ELSE %IF %UPCASE(&PLOAN_LOCATION) = VILLAGE %THEN %LET SFX = VLG;
15 %ELSE %IF %UPCASE(&PLOAN_LOCATION) = TOWN    %THEN %LET SFX = TWN;
16 %ELSE %DO;
17     %PUT ERROR: (MACRO_RPT_3) INVALID LOAN_LOCATION=&PLOAN_LOCATION. USE CITY, VILLAGE, OR TOWN. ;
18     %RETURN;
19 %END;
20
21 %LET OUTDS = ASGMLIB.TESTING_LAS_PRED_TP086704_&SFX;
22
23 /* Drop output table if it already exists */
24 %IF %SYSFUNC(EXIST(&OUTDS)) %THEN %DO;
25     PROC SQL;
26         DROP TABLE &OUTDS;
27     QUIT;
28 %END;
29
30 /* Create specific location output table */
31 PROC SQL;
32     CREATE TABLE &OUTDS AS
33     SELECT
34         L.SME_LOAN_ID_NO AS Loan_ID,
35
36         /* Gender short name */
37         CASE
38             WHEN UPCASE(L.GENDER) = 'MALE' THEN 'M'
39             WHEN UPCASE(L.GENDER) = 'FEMALE' THEN 'F'
40             ELSE 'N/A'
41         END AS Gender_short_name,
42
43         /* Remarks based on gender + loan history */
44         CASE
45             WHEN UPCASE(L.GENDER) = 'MALE' AND L.LOAN_HISTORY = 1 THEN
46                 'He is a good applicant. For he has settled his past loan(s) on time.'
47             WHEN UPCASE(L.GENDER) = 'MALE' AND L.LOAN_HISTORY = 0 THEN
48                 'He is not a good applicant. For he has not settled his past loan(s) on time.'
49             WHEN UPCASE(L.GENDER) = 'FEMALE' AND L.LOAN_HISTORY = 1 THEN
50                 'She is a good applicant. For she has settled her past loan(s) on time.'
51             WHEN UPCASE(L.GENDER) = 'FEMALE' AND L.LOAN_HISTORY = 0 THEN
52                 'She is not a good applicant. For she has not settled her past loan(s) on time.'
53             ELSE 'N/A'
54         END AS Remarks
55     FROM ASGMLIB.TESTING_LAS_PRED_TP086704_DS L
56     WHERE UPCASE(L.LOAN_LOCATION) = "%UPCASE(&PLOAN_LOCATION)" ;
57     QUIT;
58
59 /* Print Final Report */
60 TITLE1 "COMPLEX REPORT: %SYSFUNC(DATE(),WORDDATE.)";
61 TITLE2 "DETAILS OF THE LOAN APPLICANTS  &PLOAN_LOCATION";
62 PROC PRINT DATA=&OUTDS;
63     VAR Loan_ID Gender_short_name Remarks;
64     RUN;
65
66 %MEND MACRO_RPT_3;
67
68 /*Location: City*/
69 %MACRO_RPT_3(CITY);
70 /*Location: Village*/
71 %MACRO_RPT_3(VILLAGE);
72 /*Location: Town*/
73 %MACRO_RPT_3(TOWN);
```

Output

COMPLEX REPORT: September 3, 2025 DETAILS OF THE LOAN APPLICANTS — CITY

Obs	Loan_ID	Gender_short_name	Remarks
1	LP001015	M	He is a good applicant. For he has settled his past loan(s) on time.
2	LP001022	M	He is a good applicant. For he has settled his past loan(s) on time.
3	LP001031	M	He is a good applicant. For he has settled his past loan(s) on time.
4	LP001035	M	N/A
5	LP001051	M	He is a good applicant. For he has settled his past loan(s) on time.

Figure 47: Complex report filtered to CITY: derived fields *Gender_short_name* and *Remarks* populated from GENDER and LOAN_HISTORY.

COMPLEX REPORT: September 3, 2025 DETAILS OF THE LOAN APPLICANTS — TOWN

Obs	Loan_ID	Gender_short_name	Remarks
1	LP001055	F	She is a good applicant. For she has settled her past loan(s) on time.
2	LP001067	M	He is a good applicant. For he has settled his past loan(s) on time.
3	LP001082	M	He is a good applicant. For he has settled his past loan(s) on time.
4	LP001094	M	He is not a good applicant. For he has not settled his past loan(s) on time.
5	LP001096	F	She is a good applicant. For she has settled her past loan(s) on time.

Figure 48: Complex report filtered to TOWN: shows both positive and negative *Remarks* narratives reflecting LOAN_HISTORY = 1 vs. 0.

COMPLEX REPORT: September 3, 2025 DETAILS OF THE LOAN APPLICANTS — VILLAGE

Obs	Loan_ID	Gender_short_name	Remarks
1	LP001056	M	He is not a good applicant. For he has not settled his past loan(s) on time.
2	LP001153	M	He is not a good applicant. For he has not settled his past loan(s) on time.
3	LP001317	F	She is a good applicant. For she has settled her past loan(s) on time.
4	LP001347	F	She is not a good applicant. For she has not settled her past loan(s) on time.
5	LP001361	M	He is not a good applicant. For he has not settled his past loan(s) on time.

Figure 49: Complex report filtered to VILLAGE: concise listing with derived fields consistent with CITY/TOWN logic.

The complex reports for CITY, TOWN, and VILLAGE (Figs. 47–49) share the same structure, listing loan applicants with derived fields *Gender_short_name* (M/F) and narrative-style *Remarks* based on loan history. Applicants with LOAN_HISTORY=1 are described as “good applicants” who repaid on time, while those with LOAN_HISTORY=0 are flagged as “not good applicants.” The TOWN report (Fig. 48) clearly illustrates both positive and negative cases across genders, demonstrating how the reporting logic captures different applicant profiles.

Chapter 10

Data Visualisation

Stacked Bar Chart: Family Members by Loan Location

SAS Code

```
1  ****
2  STACKED BAR CHART
3  - Shows number of FAMILY_MEMBERS by LOAN_LOCATION
4  - Groups are stacked to compare location composition
5  ****
6  TITLE "Number of Family Members by Loan Location";
7  PROC SGPILOT DATA=ASGMTLIB.TESTING_LAS_PRED_TP086704_DS;
8    VBAR FAMILY_MEMBERS /
9      GROUP=LOAN_LOCATION
10     GROUPDISPLAY=STACK /* stack groups above each other */;
11   LABEL FAMILY_MEMBERS = "Number of Family Members"
12     LOAN_LOCATION = "Loan Location";
13  RUN;
```

Output

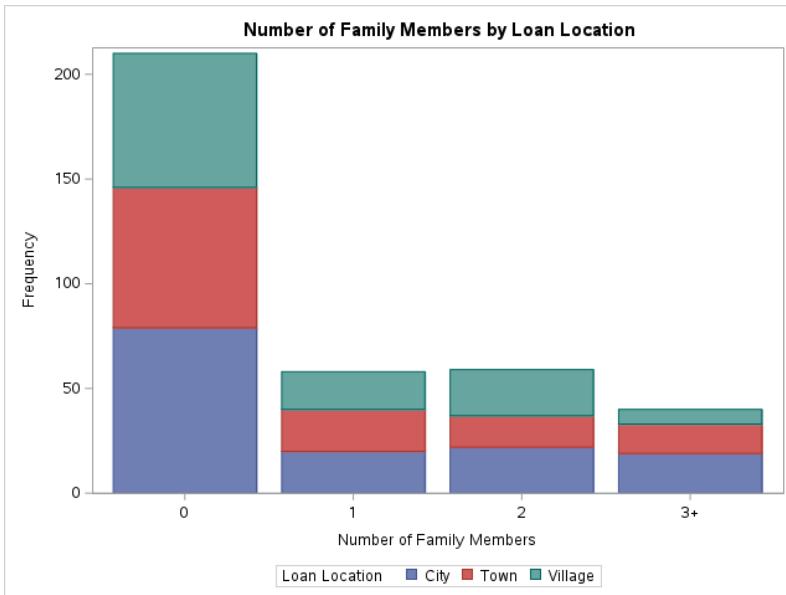


Figure 50: Stacked counts of FAMILY_MEMBERS by LOAN_LOCATION.

The stacked bar chart displays the number of family members across loan applicants by location (City, Town, Village). Applicants with no additional family members make up the largest group, with smaller proportions having one, two, or three-plus members. The distribution pattern is consistent across all three locations. This reveals that single applicants dominate the dataset regardless of their location.

Pie Chart: Loan Approval Status Predicted

SAS Code

```
1 ****
2 PIE CHART
3 - Overall distribution of predicted LOAN APPROVAL STATUS
4 - Quick view of class proportions
5 ****
6 TITLE "Loan Approval Status Predicted";
7 PROC GCHART DATA=ASGMLIB.TESTING_LAS_PRED_TP086704_DS;
8   PIE3D I_LOAN_APPROVAL_STATUS;
9   LABEL I_LOAN_APPROVAL_STATUS = "Loan Approval Status";
10 RUN;
11 QUIT;
```

Output

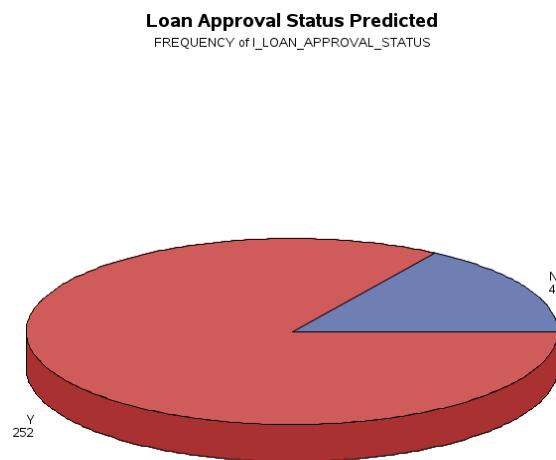


Figure 51: Overall split of predicted I_LOAN_APPROVAL_STATUS.

Based on the results, the pie chart shows the distribution of loan approval status. Out of 300 applications, 252 were approved (Y) while 48 were rejected (N). The large approval slice implies that the large majority of applicants satisfied the loan requirements. This shows that most of the applicants in the dataset are eligible to secure loan approval.

Box Plot: Loan Amount by Approval Status

SAS Code

```
1 ****
2 BOX PLOT
3 - Compares LOAN_AMOUNT distributions across
4 LOAN APPROVAL STATUS categories
5 - Highlights median, spread, and outliers
6 ****
7 TITLE "Loan Amount by Loan Approval Status";
8 PROC SGPILOT DATA=ASGMLIB.TESTING_LAS_PRED_TP086704_DS;
9 VBOX LOAN_AMOUNT / CATEGORY=I_LOAN_APPROVAL_STATUS;
10 LABEL LOAN_AMOUNT      = "Loan Amount"
11          I_LOAN_APPROVAL_STATUS = "Loan Approval Status";
12 RUN;
```

Output

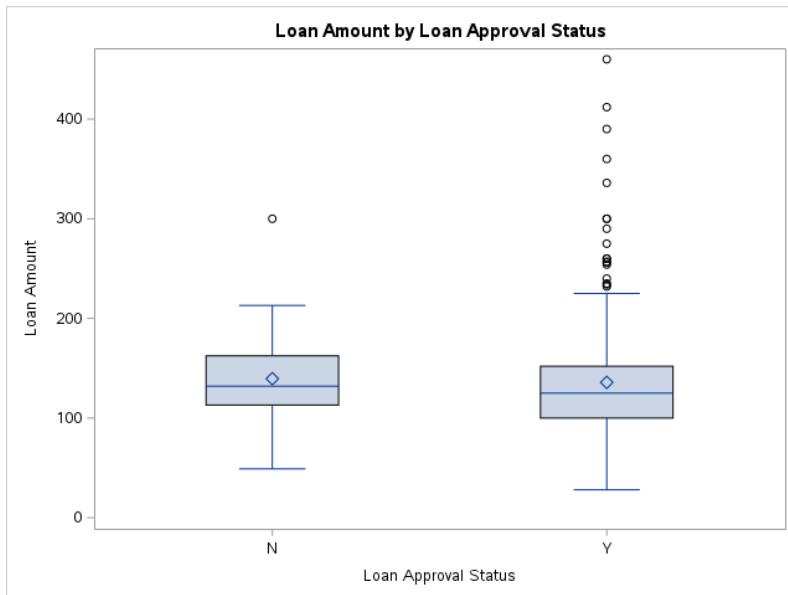


Figure 52: Distribution of LOAN_AMOUNT across predicted approval classes, showing median, spread, and outliers.

The box plot displays the distribution of loan amounts by approval status. Approved applications (Y) show a wider spread, with several outliers exceeding 400, while rejected ones (N) remain mostly below 200. Both groups have similar medians, but the variability is greater for approved loans. This indicates that larger loan amounts are more likely to be approved.

Clustered Bar Chart: Loan Approval Status by Qualification

SAS Code

```
1 ****CLUSTERED BAR CHART*****
2 - Shows LOAN APPROVAL STATUS across QUALIFICATION
3 - Side-by-side bars for clean comparison
4 ****/
5
6 TITLE "Loan Approval Status by Qualification";
7 PROC SGPLOT DATA=ASGMLIB.TESTING_LAS_PRED_TP086704_DS;
8 VBAR QUALIFICATION /
9     GROUP=I_LOAN_APPROVAL_STATUS
10    GROUPDISPLAY=CLUSTER; /* side-by-side grouping */
11    LABEL QUALIFICATION = "Qualification"
12    I_LOAN_APPROVAL_STATUS = "Loan Approval Status";
13 RUN;
```

Output



Figure 53: Side-by-side counts of predicted approval status by QUALIFICATION.

The clustered bar chart compares loan approval status by qualification (Graduate vs. Under Graduate). Graduates had nearly 200 approvals with relatively few rejections, while undergraduates showed about 50 approvals and 15 rejections. The difference suggests that higher education is associated with greater loan approval success. This highlights an applicant qualification as a key factor in loan decision.

Chapter 11

Conclusion

Loan Application Model Results.

The logistic regression model achieved stable convergence and demonstrated a significant improvement in fit compared to the intercept-only baseline, with AIC decreasing from 663.98 to 501.95 and SC from 668.24 to 565.93. Among the predictors, LOAN_HISTORY and LOAN_LOCATION emerged as the strongest and most statistically significant factors influencing loan approval. MARITAL_STATUS and QUALIFICATION were borderline significant, while other demographic and financial variables showed little impact at the 5% level. When applied to the testing dataset, the model generated consistent predictions of approval status, confirming its usefulness in supporting loan decision-making for new applications.

Reflections as a Data Scientist

Working on this project provided significant opportunities for technical and professional growth. I gained hands-on experience with SAS macro functions, logistic regression modeling, and ODS-based reporting, which were new to me. At the same time, I was able to deepen my knowledge of univariate and bivariate analysis, building on prior experience. My background in other programming languages made it easier to adapt to SAS, and I now see it as an additional pathway to perform statistical modeling and reporting efficiently.

The project reinforced the central role of data quality in predictive modeling. Although I already understood the importance of handling missing values, I learned new imputation techniques, such as applying the mode for categorical variables and the mean for continuous variables. This broadened my understanding of how imputation strategies are tailored to different data types.

Beyond technical skills, I became more familiar with the SAS environment and its ecosystem learning to navigate the interface, interpret logs, debug issues, and verify outputs. This process also indirectly improved my proficiency with LaTeX, as I had to design a structured, professional report that balanced code, outputs, and narrative in a clear format. Communicating findings in a way that is accessible to both technical and non-technical audiences was a key part of this learning process.

There were also challenges present. Repetition in code and formatting was a persistent issue, especially when presenting SAS code, outputs, and descriptions in a neat, non-overbearing way. I discovered how macros could reduce redundancy and simplify workflows, and I would use this strategy more extensively in future projects. Similarly, I would prepare more effectively by templating the document, standardising outputs, and minimizing unnecessary code.

Overall, this project strengthened my ability to apply data science tools critically, document findings effectively, and continuously refine workflows for greater clarity and efficiency.

References

- Munnangi, S. (2024). Revolutionizing Loan Systems Through Intelligent Automation. *International Journal of Research In Computer Applications and Information Technology (IJRCAIT)*, 7(2), 1508. <https://doi.org/10.5281/zenodo.14220828>
- Nigmonov, A., Shams, S., & Urbonas, P. (2024). Estimating probability of default via delinquencies? Evidence from European P2P lending market. *Global Finance Journal*, 63. <https://doi.org/10.1016/j.gfj.2024.101050>
- Ramesh Pingili. (2025). AI-driven intelligent document processing for banking and finance. *International Journal of Management & Entrepreneurship Research*, 7(2), 98–109. <https://doi.org/10.51594/ijmer.v7i2.1802>
- Santos, A. M., Cincera, M., & Cerulli, G. (2024). Sources of financing: Which ones are more effective in innovation–growth linkage? *Economic Systems*, 48(2). <https://doi.org/10.1016/j.ecosys.2023.101177>
- Suhadolnik, N., Ueyama, J., & Da Silva, S. (2023). Machine Learning for Enhanced Credit Risk Assessment: An Empirical Approach. *Journal of Risk and Financial Management*, 16(12). <https://doi.org/10.3390/jrfm16120496>