



**A • P • U**

---

**ASIA PACIFIC UNIVERSITY**  
**OF TECHNOLOGY & INNOVATION**

|                         |   |  |
|-------------------------|---|--|
| <b>Module Code</b>      | : | CT127-3-2-PFDA (Programming for Data Analysis)   |
| <b>Coursework Title</b> | : | Credit Risk Classification   |
| <b>Lecturer Name</b>    | : | Farhana Illiani Binti Hassan   |
| <b>Hand Out Date</b>    | : | 6 October 2024   |
| <b>Hand In Date</b>     | : | 9 December 2024  |
| <b>Intake</b>           | : | APD2F2409IT(IOT), APU2F2409IT(IOT), APU2F2409IT(BIS)   |
| <b>Group No.</b>        | : | 13   |
| <b>Name, TP Number</b>  | : | <ol style="list-style-type: none"><li>1. Chong Zhi Yue (TP067869)</li><li>2. Eleanor Permata Fry (TP072606)</li><li>3. Lee Jun Keat (TP067856)</li><li>4. Tham Wing Hein (TP67080)</li></ol> |

## Table of Contents

|  |    |
|--|----|
| 1.0 Introduction.....  | 4  |
| 1.1 Data Description .....   | 4  |
| 1.2 Assumptions.....   | 5  |
| 1.3 Hypothesis & Objectives .....  | 5  |
| 2.0 Data Preparation.....  | 6  |
| 2.1 Install Required Packages and Load The Libraries, Then Import Data.....  | 6  |
| 2.2 Data Pre-processing and Cleaning .....   | 6  |
| 2.3 Data Validation .....  | 8  |
| 3.0 Data Analysis .....  | 10 |
| Objective 3.1: To Investigate The Relationship Between Customer Demographics And Credit Risk Classification. Lee Jun Keat (TP067856) .....                   | 10 |
| 3.1.1: Is There A Significant Relationship Between Personal Status And Credit Risk Classification?.....  | 10 |
| 3.1.2: Do Job Skills Significantly Influence Credit Risk Classification? .....   | 11 |
| 3.1.3: Foreign Workers Got A Better Credit Risk Class?.....  | 13 |
| 3.1.4: Is There A Significant Relationship Between Customer Gender And Credit Risk Classification?.....  | 14 |
| 3.1.5: What Are The Combined Effects Of Gender, Job Skills, On Credit Risk Classification?.....  | 15 |
| 3.1.6 Additional Features .....  | 16 |
| 3.2 Objective 3-2: To Analyse The Effect Of Loan Related Factors On Credit Risk Classification – Tham Wing Hein (TP067080) .....                             | 17 |
| 3.2.1 Is there any relationship between loan purpose and credit risk classification? .....   | 17 |
| 3.2.2 Does the duration of the loan (installment duration) influence the credit risk classification?.....  | 19 |
| 3.2.3 Are loan amounts a strong predictor for credit risk classification? .....  | 20 |
| 3.2.4 Does employment status interacts with loan duration to influence credit risk classification?.....  | 22 |
| 3.2.5 Does the proportion of income committed to loan payments (installment commitment) influence the likelihood of higher credit risk classification? ..... | 24 |

|   |    |
|---|----|
| 3.2.6 Conclusion of Objective 3-2 .....   | 25 |
| 3.2.7 Additional Features.....  | 26 |
| 3.3 Objective 3-3: To Assess the Role of Financial Commitments on Credit Risk Classification. Eleanor Permata Fry (TP072606) .....                          | 27 |
| 3.3.1 Is There A Relationship Between The Instalment Commitments And Credit Risk Classification?.....   | 27 |
| 3.3.2 Does Income Level Predict Credit Risk Classification When Considered Alongside Instalment Commitment?.....  | 30 |
| 3.3.3 How Does The Presence Of the Financial Commitments (e.g., Other Payment Plans) Influence Credit Risk Classification? .....                            | 33 |
| 3.3.4 What Demographic Factors (e.g., Age, Marital Status) Moderate The Relationship Between Instalment Commitments and Credit Risk Classification? .....   | 35 |
| 3.3.5 How Does The Duration Of Residence Impact The Relationship Between Financial Commitments And Credit Risk Classification?.....                         | 37 |
| 3.3.6 Additional Features .....   | 39 |
| 3.4 Objective 4: To explore the impact of customer credit history on credit risk classification. - CHONG ZHI YUE (TP067869) .....                           | 41 |
| 3.4.1 Is there a significant relationship between the history of past loans and current credit risk classification?.....                                    | 41 |
| 3.4.2 Does the number of previous loan defaults predict the likelihood of being classified as a high-risk customer?.....                                    | 43 |
| 3.4.3 How does the length of credit history impact the credit risk classification? .....  | 45 |
| 3.4.4 What are the external factors that interact with past credit history (e.g., income, employment status) to influence credit risk classification? ..... | 49 |
| 3.4.5 Is there any relationship between the purpose of having the loan and the credit risk classification?.....   | 55 |
| 3.4.6 Is there any relationship between the categorical variables (e.g., credit history, purpose, and employment status)? .....                             | 57 |
| 4.0 Conclusion .....  | 60 |
| 4.1 Results.....  | 60 |
| 4.2 Recommendations.....  | 61 |
| 4.3 Limitations and Future Direction.....   | 61 |

## **1.0 Introduction**

The bank sector act as an important role in assessing credit risks to ensure financial stability and minimize losses. The case study of this assignment mainly focuses on implementing the application of data analytics techniques to identify the credit risks based on customer demographics and their credit status. Inside this dataset, it contains various attributes such as gender, marital status, employment status, loan purpose, and instalment's duration, providing a comprehensive view of customer profiles and their credit-related activities.

The objective of this assignment is to conduct an analysis of the dataset to identify the factors between high-risk and low-risk individuals. This will involve applying data exploration, manipulation, transformation, and visualization techniques to pre-process the data, and reflect our important insights to develop models for credit risk classification. Advanced analytical concepts will also be incorporated to enhance the effectiveness of the analysis.

## **1.1 Data Description**

In this assignment, the team conducted data analysis on the “3. credit\_risk\_classification.csv,” which contains approximately 6000 records, where each row corresponds to a unique customer’s record. The dataset includes various key variables:

A. Customer Demographics:

- Gender
- Marital Status
- Employment Status
- Age

B. Credit Behaviour:

- Credit Amount
- Duration of Instalment Commitment
- Purpose of Loan
- Credit History
- Account Type

C. Credit Risk:

- Credit Class

To determine the main variables impacting credit risk, a dataset with numerous numerical and categorical elements will be analysed. Before doing the analysis of each objective, the data will be pre-processed, cleaned, and explored to find patterns and relationships that can help with predictive modelling.

## **1.2 Assumptions**

For this assignment, each group member chose one goal to watch and evaluate to discover how it affected the classification of credit risk. We have done analyse the missing information during data cleaning and handle categorical and continuous variable properly. We make our predictions that the features featured are relevant to the categorization of credit risk and that it can predict the connections between customer demographics, credit behaviour, and credit risk classification.

## **1.3 Hypothesis & Objectives**

### **Hypothesis:**

1. Classification of credit risk is influenced by customer demographics and credit behaviours (e.g., credit amount, loan term, and job status).
2. Higher credit amounts and longer loan durations are associated with a higher probability of being categorized as a “Bad” credit risk.

### **Objectives:**

1. To investigate the relationship between customer demographics and credit risk classification.
2. To analyse the effect of loan-related factors on credit risk classification.
3. To assess the role of financial commitments on credit risk classification.
4. To explore the impact of customer credit history on credit risk classification.

## **2.0 Data Preparation**

### **2.1 Install Required Packages and Load The Libraries, Then Import Data**

```
#Install packages and load libraries needed
install.packages("dplyr")
library(dplyr)
install.packages("ggplot2")
library(ggplot2)
install.packages("ggridges")
library(ggridges)
install.packages("igraph")
library(igraph)
install.packages("plotrix")
library(plotrix)
install.packages("vcd")
library(vcd)
```

*Figure: Install Packages and Load Libraries*

The code loads several packages that are commonly used for data manipulation, visualization and analysis. ‘dplyr’ was being used to perform data manipulations such as filtering, selecting and summarizing data frames. ‘ggplot2’ was being used to create visualizations such as plots. ‘ggridges’ was being used to create ridge plots. ‘igraph’ was being used for network analysis and graph manipulation which can be used to create, manipulate, and visualize graphs. ‘plotrix’ was being used to do various plotting functions like 3d plots and others. ‘vcd’ is being used for visual display of categorical data.

```
AData <- read.csv("C:\\\\Users\\\\Itsuki\\\\Documents\\\\Degree Year 2 Semester 1\\\\Programming For Data Analysis\\\\Assignment\\\\3. credit_risk_classification.csv", sep=",", header=TRUE)
AData
#Table View
View(AData)
```

*Figure 1: Import Dataset*

The code above reads the .csv file named “3. Credit\_risk\_classification.csv” from the mentioned directory into an R data frame called ‘AData’ using ‘read.csv’ function. Then the data frame is being displayed into table form using ‘View’ function.

### **2.2 Data Pre-processing and Cleaning**

```
#view Missing Data
colsums(is.na(AData))
sum(complete.cases(AData))
which(is.na(AData))
sum(is.na(AData))
```

*Figure 2: View Missing Data in Dataset*

- `colsums(is.na(AData))` was used to showing how many missing values exist in each columns.
- `sum(complete.cases(AData))` will tell how many rows in AData have no missing data.

- `which(is.na(AData))` gives the exact positions of the missing values in rows and columns.
- `sum(is.na(AData))` show how many total missing values that are across the entire dataset.

```
#Remove Unused Column
AData<- AData %>% select(-X)

#Remove Duplicated Data
AData <- AData %>% distinct()
sum(duplicated(AData))
```

Figure 3: Remove Unused Column & Duplicated Data

- “-X” here means unnecessary or irrelevant columns. “-X” here will be removed.
- `distinct()` is a function from “dplyr” package that removes duplicated rows based on all columns.
- `sum()` will then check for any remaining duplicates and returns the total count of duplicated rows.

```
#Remove Error In age
AData <- AData[AData$age >= 18 & AData$age <= 100, ]
AData$age <- as.integer(AData$age)

#Remove Error In duration
AData$duration <- round(AData$duration)
AData <- AData[AData$duration >= 4 & AData$duration <= 72, ]

#Remove Error In installment commitment
AData$installment_commitment <- round(AData$installment_commitment)
AData <- AData[AData$installment_commitment >= 1 & AData$installment_commitment <= 4, ]

#Remove Error In residence since
AData$residence_since <- round(AData$residence_since)
AData <- AData[AData$residence_since >= 1 & AData$residence_since <= 4, ]

#Remove Error In existing credits
AData$existing_credits <- round(AData$existing_credits)
AData <- AData[AData$existing_credits >= 1 & AData$existing_credits <= 4, ]

#Remove Error In number dependents
AData$num_dependents <- round(AData$num_dependents)
AData <- AData[AData$num_dependents >= 1 & AData$num_dependents <= 2, ]
```

Figure 4: Remove Errors in Specified Column

The codes above filters out invalid or out-of-range values for “age”, “duration”, “installment\_commitment”, “residence\_since”, “existing\_credits”, and “num\_dependents”. It then rounds continuous numeric variables to the nearest appropriate integer ensuring the dataset only contains valid and logical entries for the columns.

```
#Filling other_payment_plans
AData$other_payment_plans[AData$other_payment_plans == ""] <- NA
mode_value <- names(sort(table(AData$other_payment_plans), decreasing = TRUE))[1]
print(mode_value) # shows stores
AData$other_payment_plans[is.na(AData$other_payment_plans)] <- mode_value
```

Figure 5: Filling Blank Cell and Replace with Mode Value

This code was being used to find and fill in the missing values (NA) in “other\_payment\_plans” with the most frequent value (mode).

```
#Change Label Name
AData <- AData %>%
  mutate(
    personal_status = case_when(
      grepl("male", personal_status) & grepl("single", personal_status) ~ "male single",
      grepl("female", personal_status) & grepl("divorced/dependent/married", personal_status) ~ "female divorced/dependent/married",
      grepl("male", personal_status) & grepl("divorced", personal_status) ~ "male divorced/separated",
      grepl("male", personal_status) & grepl("married/widowed", personal_status) ~ "male married/widowed",
      TRUE ~ personal_status
    )
  )
AData <- AData %>%
  mutate(
    job = case_when(
      job == "skilled" ~ "skilled",
      job == "unskilled resident" ~ "unskilled, resident",
      job == "high qualif/self emp/mgmt" ~ "high qualification, self-employed, management",
      job == "unemp/unskilled non res" ~ "unemployed, unskilled, non-resident",
      TRUE ~ job
    )
  )
AData <- AData %>%
  mutate(
    employment = case_when(
      employment == "<1" ~ "T or more years",
      employment == "1<x<4" ~ "1 to 3 years",
      employment == "4<x<7" ~ "4 to 6 years",
      employment == "c1" ~ "less than 1 year",
      employment == "unemployed" ~ "unemployed",
      TRUE ~ "other"
    )
  )
```

Figure 6: Change the Data Label Names in The Columns

```
#add column "gender"
AData$gender <- ifelse(grepl("female", AData$personal_status, ignore.case = TRUE), "female",
                      ifelse(grepl("male", AData$personal_status, ignore.case = TRUE), "male", NA))
```

Figure 7: Adding extra column called "gender"

## 2.3 Data Validation

```
#conversion
AData$age <- as.numeric(AData$age)
AData$credit_amount <- as.numeric(AData$credit_amount)
AData$installment_commitment <- as.numeric(AData$installment_commitment)
AData$residence_since <- as.numeric(AData$residence_since)
AData$existing_credits <- as.numeric(AData$existing_credits)
AData$num_dependents <- as.numeric(AData$num_dependents)

#credit history conversion
AData$credit_history <- gsub("delayed previously", "delays", AData$credit_history)
AData$credit_history <- gsub("critical/order existing credit", "critical accounts", AData$credit_history)
AData$credit_history <- gsub("no credits/all paid", "all paid", AData$credit_history)
```

Figure 8: Converting Data Types

```
data.Frame': 1400 obs. of  21 variables:
$ checking_status : chr  "+0" "0xx<x<200" "no checking" "+0" ...
$ duration       : num  6 41 12 42 24 36 24 36 18 ...
$ credit_history : chr  "delayed previously" "delays" "critical/accounts" "existing paid" ...
$ purpose        : chr  "radio/tv" "radio/tv" "education" "furniture/equipment" ...
$ credit_amount  : num  1169 5951 2098 7884 4870 ...
$ savings_status : chr  "1<x<6" "1<x<6" "1<x<6" "1<x<6" "1<x<6" ...
$ employment      : chr  ">w?" "1<x<6" "4xx<x?" "4xx<x?" ...
$ installmnt_commitment: num  4 2 2 3 2 3 2 2 4 ...
$ gender          : chr  "male" "male" "male single" "male single" "male single" ...
$ other_parties   : chr  "none" "none" "none" "guarantor" ...
$ residence_since : num  4 2 3 4 4 4 2 4 2 ...
$ property_magnitude: chr  "real estate" "real estate" "real estate" "life insurance" ...
$ age             : num  67 22 49 45 53 35 53 35 61 28 ...
$ other_payment_plan: chr  "stores" "stores" "stores" "stores" ...
$ existing_credits: num  2 1 1 1 2 1 1 1 1 2 ...
$ job             : chr  "skilled" "skilled" "unskilled, resident" "skilled" ...
$ num_dependents : num  1 1 2 1 2 1 1 1 1 ...
$ phone           : chr  "yes" "no" "no" "no" ...
$ foreign_worker  : chr  "yes" "yes" "yes" "yes" ...
$ class            : chr  "good" "bad" "good" "good" ...
```

Figures 9 & 10: Output of Corrected Data Types

```
#Save Cleaned Dataset
write.csv(AData, "C:\\\\Users\\\\Dev\\\\Downloads\\\\Assignment\\\\Assignment\\\\3. cleaned_credit_risk_classification.csv", row.names = FALSE)

#Read Cleaned Dataset
AData <- read.csv("C:\\\\Users\\\\Dev\\\\Downloads\\\\Assignment\\\\Assignment\\\\3. cleaned_credit_risk_classification.csv")
View(AData)
```

Figure 11: Save the processed data and read again (refresh)

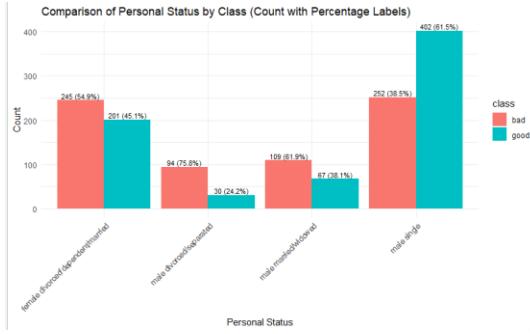
## 3.0 Data Analysis

### Objective 3.1: To Investigate The Relationship Between Customer Demographics And Credit Risk Classification. Lee Jun Keat (TP067856)

#### 3.1.1: Is There A Significant Relationship Between Personal Status And Credit Risk Classification?

The first step is to combine all the count within the personal status and classify it with the class to determine the credit risk score.

```
#analysis 1
personal_status_frequency <- Adata %>%
  count(personal_status) %>%
  as_tibble()
personal_status_frequency
# Display the tibble
personal_status_frequency_tibble
status_class_counts <- as.data.frame(table(Adata$personal_status, Adata$class))
colnames(status_class_counts) <- c("personal_status", "class", "count")
status_class_counts <- status_class_counts %>
  group_by(personal_status) %>%
  mutate(percentage = count / sum(count) * 100)
ggplot(status_class_counts, aes(x = personal_status, y = count, fill = class)) +
  geom_bar(stat = "identity", position = "dodge") +
  geom_text(aes(label = paste0(count, " (", round(percentage, 1), "%)"),
                position = position_dodge(width = 0.9),
                vjust = -0.3,
                size = 3) +
  labs(title = "Comparison of Personal Status by class (Count with Percentage Labels)",
       x = "Personal Status", y = "Count") +
  theme_minimal() +
  theme(axis.text.x = element_text(angle = 45, hjust = 1))
```



Figures 12 & 13: Bar Chart of Analysis 1

```
> personal_status_frequency
# A tibble: 4 × 2
  personal_status      n
  <chr>              <int>
1 female divorced/divorced/dependent    446
2 male divorced/separated            124
3 male married/widowed             176
4 male single                      654
```

Figure 14: Tibble table for Analysis 1

The analysis took place with 1400 individuals in the credit risk classification with the current personal status. Each bar displays the count of “bad” and “good” credit class for every category

#### Category:

Female divorced/divorced/dependent: There are 446 individuals inside this category, there are 201 (45.1%) individuals got a “good” class for the credit risk and 245 (54.9%) individuals are classified as “bad” class.

Male divorced/separated: Among 124 individuals, it shows that only 30 (24.2%) individuals got a “good” class and most of them which are 94 (75.8%) individuals got a “bad” class. The significance shows that it skews into the bad credit class in this category.

Male married/widowed: The result returns that only 67 (38.1%) individuals got “good” credit class during the analysis. While 109 of out 176 individuals is got a “bad” credit score. Same as the previous category, it shows that majority of them are towards higher financial risk.

Male singles: This category got the most count of individual among 4 of the categories. Which are 654 out of 1400 individuals. It shows that majority of them shows that they had a “good” credit class which are 402 (61.5%) individuals, and the rest 252 (38.5%) individual are “bad” credit class.

```
> chi_table <- table(AData$personal_status, AData$class)
>
> ctest <- chisq.test(chi_table)
>
> ctest

Pearson's Chi-squared test

data: chi_table
X-squared = 81.799, df = 3, p-value < 2.2e-16
```

Figure 15: Chi Squared Statistic for Analysis 1

In the chi square table, it shows that the chi squared statistic is on the value of 81.799 with the degree of freedom = 3 and a p value <2.2e-16.

### 3.1.2: Do Job Skills Significantly Influence Credit Risk Classification?

```
#Analysis 2
library(tidyverse)
job_frequency <- AData %>%
  count(job.class) %>%
  mutate(job_frequency = n / sum(n))
job_frequency

#library(tibble)
skilled.data <- job_frequency %>% filter(job == "skilled")
Labels <- paste0(skilled.data$class, "- (" , skilled.data$n, ") ", "(" , round(skilled.data$n / sum(skilled.data$n)) * 100, 1), "%")
print(skilled.data,
  labels = Labels,
  main = "Credit class distribution for job skills: skilled",
  explode = TRUE,
  col = c("blue", "red"),
  Tablelex = 1)
Tablelex = 1

unskilled.data <- job_frequency %>% filter(job == "unskilled", resident)
Labels <- paste0(unskilled.data$class, "- (" , unskilled.data$n, ") ", "(" , round(unskilled.data$n / sum(unskilled.data$n)) * 100, 1), "%")
print(unskilled.data,
  labels = Labels,
  main = "Credit Class Distribution for job skills: unskilled/resident",
  explode = TRUE,
  col = c("blue", "red"),
  Tablelex = 1)
Tablelex = 1

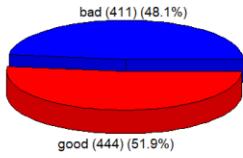
highskilled.data <- job_frequency %>% filter(job == "high qualification, self-employed, management")
Labels <- paste0(highskilled.data$class, "- (" , highskilled.data$n, ") ", "(" , round(highskilled.data$n / sum(highskilled.data$n)) * 100, 1), "%")
print(highskilled.data,
  labels = Labels,
  main = "Credit Class Distribution for job skills: high qualification/ self-employed/ management",
  explode = TRUE,
  col = c("blue", "red"),
  Tablelex = 1)
Tablelex = 1
```

| job   | class | n   |
|---|-------|-----|
| 1 high qualification, self-employed, management | bad   | 148 |
| 2 high qualification, self-employed, management | good  | 97  |
| 3 skilled                                       | bad   | 411 |
| 4 skilled                                       | good  | 444 |
| 5 unemployed, unskilled, non-resident           | bad   | 60  |
| 6 unemployed, unskilled, non-resident           | good  | 15  |
| 7 unskilled, resident                           | bad   | 81  |
| 8 unskilled, resident                           | good  | 144 |

Figure 16 & 17: Code Sampling and Tibble Table for Analysis 2

An analysis about a credit risk score of an individual within a dataset about the job skills with its credit risk score. 3 categories which are “skilled”, “unskilled/resident” and “high qualification/self-employed/management”.

Credit Class Distribution for Job Skills: Skilled



Credit Class Distribution for Job Skills: High qualification/ self-employed/ management

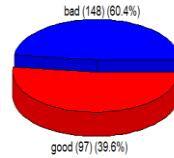


Figure 18 & 19: 3D Pie Chart- Job Skills:Skilled / Job Skills-High Qualifications/Self-Employed/Management

Credit Class Distribution for Job Skills: Unskilled/resident

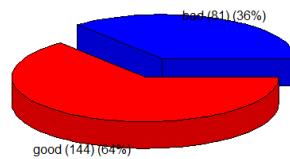


Figure 20 & 21: 3D Pie Chart- Unskilled/Resident

## **1. Skilled Workers:**

**Good Credit Class: 444 individuals (51.9%)**

**Bad Credit Class: 411 individuals (48.1%)**

In this category, both sides are nearly split and there is less difference between two sides, with a slight majority in the good credit class. It means that most of them is able to handle their financial status, but it is close that they might fall into bad credit class.

The analysis provided explores the impact of job skills on credit risk classification, examining three job categories: "skilled," "unskilled, resident," and "high qualification, self-employed, management." For each job category, the credit class distribution is divided into "good" and "bad" classifications and visualized using pie charts. Here's a breakdown:

## **2. Unskilled, Resident Workers:**

**Good Credit Class: 144 individuals (16.8%)**

**Bad Credit Class: 81 individuals (9.5%)**

Unskilled/resident workers show a larger proportion in the good credit class (64%) compared to the bad (36%). It shows that unskilled workers have a higher chance of being classified as good credit than bad as they don't have financial issues, but their overall distribution still tell us that unskilled status alone does not prevent a good credit classification.

## **3. High Qualification, Self-Employed, Management**

- **Good Credit Class:** 97 individuals (39.6%)
- **Bad Credit Class:** 148 individuals (60.4%)

In this category, individuals with high qualification or management roles show a larger proportion in the bad credit category 148 individuals (60.4%), with only 97 individuals (39.6%) in the good credit class. This reflects that even with high qualifications or self-employment status, there is a risk of getting bad credit classification due to income instability or other financial risks associated with self-employment or high-level management roles.

In conclusion, job skills do influence credit classification, but other underlying factors likely contribute to credit risk. High-skill or high-management roles do not guarantee a good credit classification, and some unskilled individuals can still achieve favourable credit ratings.

### **3.1.3: Foreign Workers Got A Better Credit Risk Class?**

```
#Analysis 3
library(dplyr)
foreign_worker_frequency <- Adata %>%
  count(foreign_worker.class) %>%
  as_tibble()
foreign_worker_frequency

foreign_worker_frequency2 <- Adata %>%
  count(foreign_worker) %>%
  as_tibble()
foreign_worker_frequency2
library(ggplot2)

ggplot(Adata, aes(x = foreign_worker, fill = class)) +
  geom_bar(position = "fill") +
  labs(title = "Proportion of credit class by Foreign worker status",
       x = "Foreign worker Status", y = "Proportion") +
  scale_fill_manual(values = c("good" = "blue", "bad" = "red")) +
  scale_y_continuous(labels = scales::percent)
```

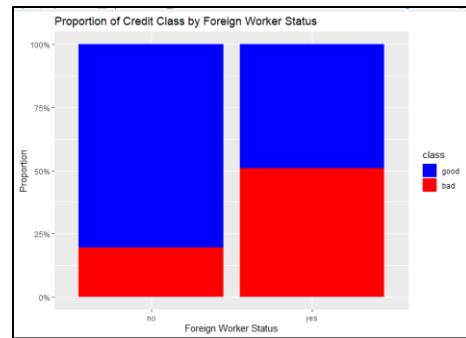


Figure 22 & 23: Sample code and Stacked Bar Graph (Proportion) for Analysis 3

This graph shows the proportion of credit class based on the foreign worker status. For the category who are not foreign workers, it shows that a large proportion falls into the good credit class and the bad credit class only takes a small proportion, which indicates that most of the individuals who are not foreign workers have a good credit rating.

Meanwhile, for the individuals who are foreign workers, the proportion of bad is higher than good class. We can see that there are barely any small differences between bad credit class and good credit class. We can understand that individuals who are from other countries are more likely to have a bad credit class.

In conclusion, individuals who are not foreign workers got a better credit class, while individuals who are foreign workers mostly got a bad credit class. This may relate that underlying factors associated with foreign worker status that impact creditworthiness, potentially due to economic, employment, or stability differences.

### 3.1.4: Is There A Significant Relationship Between Customer Gender And Credit Risk Classification?

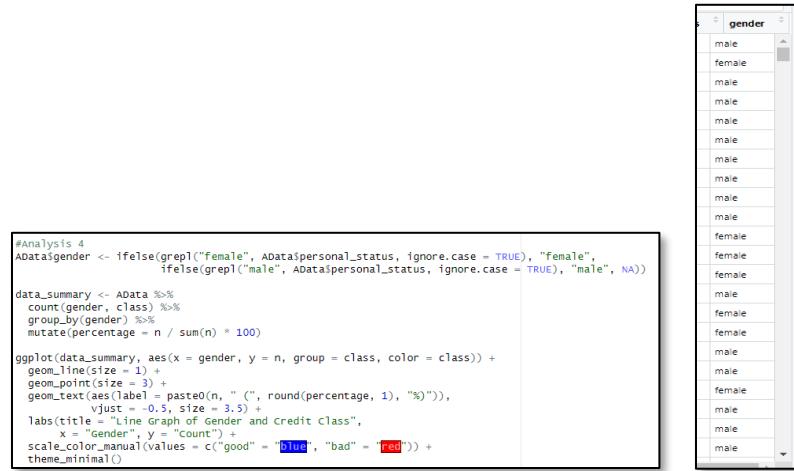


Figure 24 &amp; 25: Sampling Code and New Column for Analysis 4

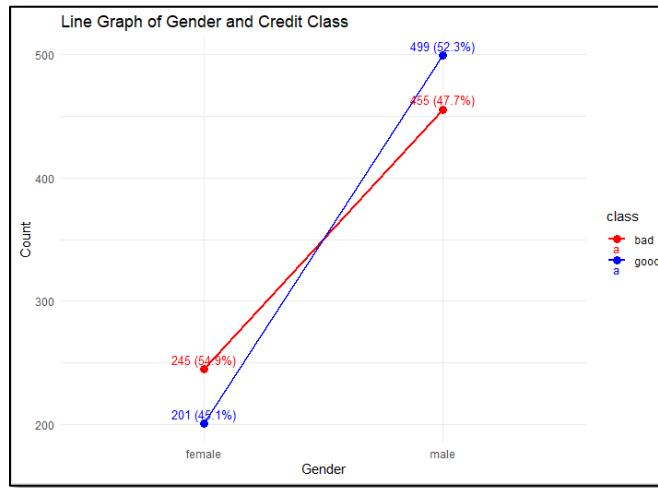


Figure 26: Line Graph for Analysis 4

The line graph above shows the relationship between the credit risk classification and their gender. Among the 1400 individuals, we found that there are 446 females and 954 males in our research. In 446 individuals who are female, it indicates that the one who get a bad credit score in credit risk classification is higher than the good one. Where 245 (54.9%) individuals got a "bad" credit score and the rest 201 (45.1%) got a "good" credit score. In contrast, males show a slightly better credit classification, with 499 (52.3%) individuals were classified as having "good" credit class, and the rest of 455(47.7%) individuals having "bad" credit class. These figures show a higher prevalence of good credit among males compared to females in this dataset. And there is potential influence of gender on credit risk classification, with males generally having a better credit profile than females.

### 3.1.5: What Are The Combined Effects Of Gender, Job Skills, On Credit Risk Classification?

```

Analysis 5
library(dplyr)
library(ggplot2)

ggplot(Adata, aes(x = "", fill = class)) +
  geom_bar(width = 1, stat = "count") +
  facet_wrap(~gender, scales = "free_y") +
  scale_fill_manual(values = c("bad" = "#E63333", "good" = "#3CB371")) +
  labs(title = "Credit risk distribution by Gender",
       x = "",
       y = "",
       fill = "Credit Risk") +
  geom_text(aes(label = ..count..), stat = "count", position = position_stack(vjust = 0.5)) +
  theme_void() +
  theme(legend.position = "bottom")

ggplot(Adata, aes(x = factor(1), fill = class)) +
  geom_bar(width = 1, stat = "count") +
  facet_wrap(~job, scales = "free_y") +
  scale_fill_manual(values = c("bad" = "#E63333", "good" = "#3CB371")) +
  geom_text(aes(label = ..count..), stat = "count", position = position_stack(vjust = 0.5)) +
  labs(title = "Credit risk distribution by job skills",
       x = "",
       y = "",
       fill = "Credit Risk") +
  theme_void() +
  theme(legend.position = "bottom")

ggplot(Adata, aes(x = factor(1), fill = class)) +
  geom_bar(width = 1, stat = "count") +
  facet_wrap(~foreign_worker, labeller = label_both, scales = "free_y") +
  scale_fill_manual(values = c("bad" = "#E63333", "good" = "#3CB371")) +
  geom_text(aes(label = ..count..), stat = "count", position = position_stack(vjust = 0.5)) +
  labs(title = "Credit risk distribution by Foreign status",
       x = "",
       y = "",
       fill = "Credit Risk") +
  theme_void() +
  theme(legend.position = "bottom")

```

Figure 27: Code Sampling for Analysis 5

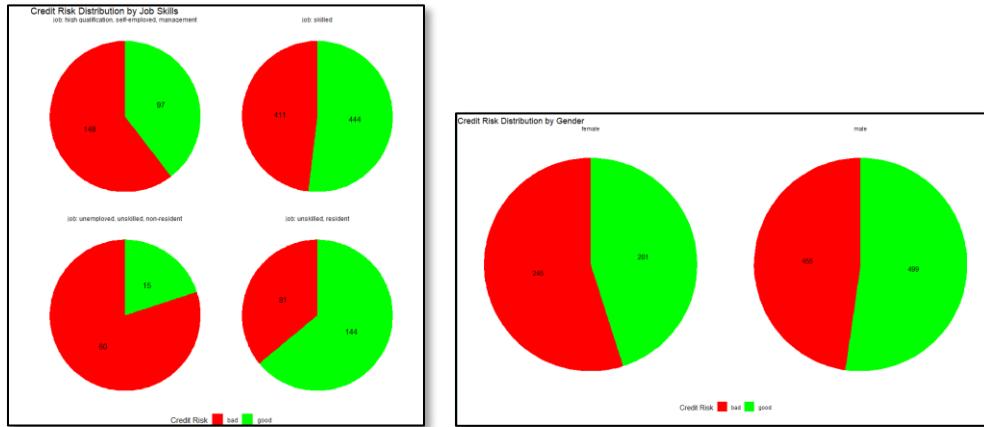


Figure 28 & 29: Pie Chart Distribution Based on Job Skills / Pie Chart Distribution Based on Gender

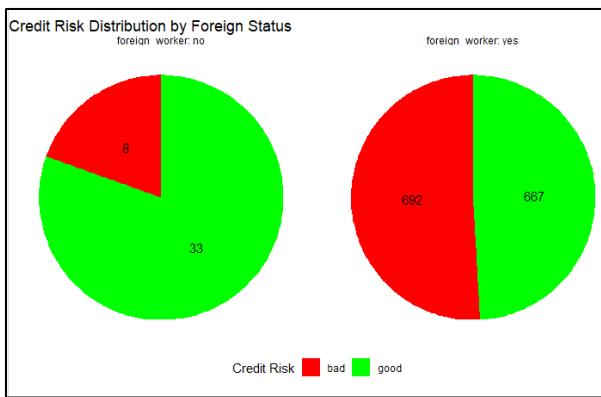


Figure 7: Pie Chart Distribution Based on Foreign Worker Status

The analysis of credit risk distribution reveals significant effects of both gender, foreign status and job skills on creditworthiness. Individuals who got higher qualifications or managerial roles show a higher proportion of bad credit risk compared to good, and skilled workers able to maintain balanced distribution with a little toward good credit risk. Unskilled and unemployed individuals, particularly non-residents, show the highest proportion of bad credit risk, which indicates that they might have financial vulnerability. Considering on gender factors, males got a more balanced credit risk profile, with a higher number of good credit cases compared to females, who show a greater proportion of bad credit risk. The combined effects suggest that skilled males are more likely to have favourable credit outcomes, while females, particularly in unskilled or unemployed categories, face greater credit challenges. These trends highlight the need for targeted interventions, such as job skill development and financial literacy programs, especially for vulnerable groups.

### **3.1.6 Additional Features**

- Tibble functions

```
# Display the tibble
personal_status_frequency_tibble
status_class_counts <- as.data.frame(table(ADATA$personal_status, ADATA$class))
colnames(status_class_counts) <- c("personal_status", "class", "count")
```

- Chi-squared test

```
> chi_table <- table(ADATA$personal_status, ADATA$class)
>
> cctest <- chisq.test(chi_table)
>
> cctest

Pearson's chi-squared test

data: chi_table
X-squared = 81.799, df = 3, p-value < 2.2e-16
```

## 3.2 Objective 3-2: To Analyse The Effect Of Loan Related Factors On Credit Risk Classification – Tham Wing Hein (TP067080)

### 3.2.1 Is there any relationship between loan purpose and credit risk classification?

```
> #Objective 2: To analyze the effect of loan-related factors on credit risk classification
> #analysis 2-1: Is there any relationship between loan purpose and credit risk classification? (THAM WING HEIN TP067080)
> #view frequency of values for 'purpose' and 'class' separately
> table(AData$purpose)

  business domestic appliance      education furniture/equipment      new car       other
    122        36          36           89           268          344          70
  radio/tv      repairs     retraining      used car
    313         35          18          105

> table(AData$class)

  bad  good
    700   700

> #create contingency table for loan purpose and credit risk classification
> PurposeClass_table <- table(AData$purpose, AData$class)
> print(PurposeClass_table)

  business  domestic appliance  education furniture/equipment  new car       other
  bad        59        63          28          8           122        36          36
  good       63        7           21         14           313        35          35
  radio/tv      repairs     retraining      used car
  95        218          10          8           19         145          18          105

> prop.table(PurposeClass_table,1)

                                bad      good
  business             0.4836666  0.5163934
  domestic appliance   0.7777778  0.2222222
  education            0.6853933  0.3146067
  furniture/equipment  0.5410448  0.4589552
  new car              0.5784884  0.4215116
  other                0.9000000  0.1000000
  radio/tv             0.3035144  0.6964856
  repairs               0.6000000  0.4000000
  retraining            0.5555556  0.4444444
  used car             0.1809524  0.8190476
```

Figure 8: Output

```
#generate bar plot of the contingency table
ggplot(AData, aes(x = purpose, fill = class)) +
  geom_bar(position = "dodge") +
  labs(title = "Loan Purpose vs Credit Risk Classification", x = "Loan Purpose", y = "Count") +
  theme(axis.text.x = element_text(angle = 45, hjust = 1)) +
  scale_fill_manual(values = c("pink", "lightgreen"), name = "Credit Risk") +
  geom_text(stat = "count", aes(label = ..count..), position = position_dodge(width = 0.8), vjust = -0.5)
```

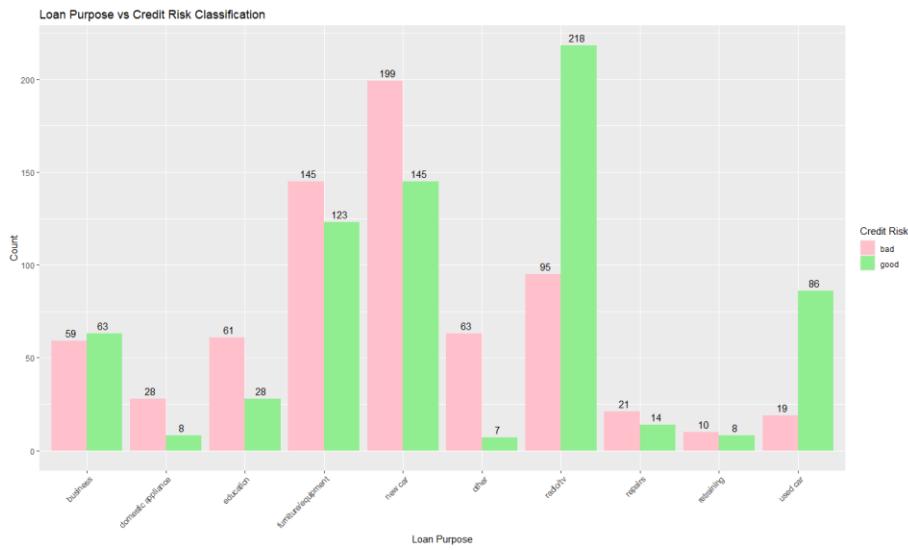


Figure 32 &amp; 33: Bar Graph &amp; its Code Snippet

In an analysis of the credit risk classification dataset comprising 1400 individuals, it was found that there are 700 people, or 50% of the total customers are having bad credit class. Among the 1400 individuals, the top 3 most common purposes of taking the credit are for “New Car” (344 cases), “Radio/TV” (313 cases) and “Furniture/Equipment” (268 cases) while the least common purposes include “Retraining” at 18 cases and “Repairs” at 35 cases.

Based on the contingency table and the bar plot generated, they revealed varying levels of risk associated with different loan purposes. The highest credit risk among the purposes is “Other” which 90% of it is classified as Bad, this marks the riskiest loan purpose in the dataset.

“Domestic Appliance” comes along and becomes the second highest credit risk at around 77.78% classified as Bad, indicating a high likelihood of credit risk. On the other hand, “Used Car” scores 81.90% Good classification making it the least risky credit while “Radio/TV” contains 69.65% Good classification, suggesting a stronger association with good credit risk as well. At moderate risk categories, we found out that “Furniture/Equipment” is having 54.10% Bad classification showing a relatively balanced risk level. “Business” on the other hand indicates 51.64% Good classification showing a moderate risk.

```
> #run Chi-squared Test of Independence
> PurposeClass_chisqr <- chisq.test(PurposeClass_table)
> print(PurposeClass_chisqr)

Pearson's Chi-squared test

data: PurposeClass_table
X-squared = 171.27, df = 9, p-value < 2.2e-16
```

*Figure 34: Chi-squared relationship*

The Chi-squared test result reveals Chi-squared statistic of 171.27 with 9 degree of freedom and a p-value significantly less than 0.05 (at 2.2e-16 or 0.0000000000000022). This result indicates that the loan purpose is significantly associated with credit risk classification. In other words, it means that there is strong evidence of a relationship between loan purpose and credit risk classification in the dataset. For example, some certain loan purposes might be more likely to have a “bad” classification, while others might be correlate with a “good” classification.

### 3.2.2 Does the duration of the loan (installment duration) influence the credit risk classification?

```
analysis 2-2: Does the duration of the loan (installment duration) influence the credit risk classification?
str(AData) #check structure of data
unique(AData$class) #check unique values in 'class' column
summary(AData$duration) #statistics of 'duration' column

#generate density plot of the loan duration across credit risk classifications
ggplot(AData, aes(x = duration, fill = class)) +
  geom_density(alpha = 0.5) +
  labs(title = "Density of Loan Duration by Credit Risk Classification",
       x = "Loan Duration (months)",
       y = "Density") +
  scale_fill_manual(values = c("bad" , "good"),name = "Credit_Risk")
```

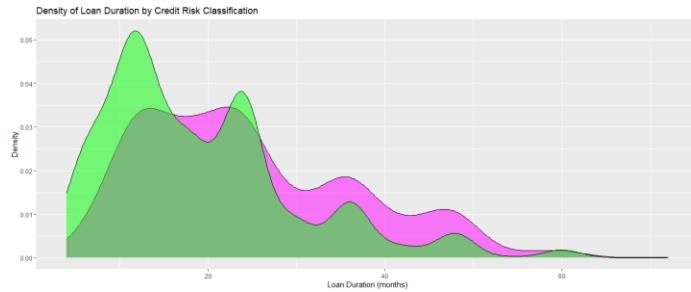


Figure 35 & 36: Density plot and its Code Snippet

|    | DurationClass_table <- table(AData\$duration, AData\$class) |  |
|----|---|--|
|    | > print(DurationClass_table)                                |  |
|    |   |  |
|    | bad good  |  |
| 4  | 8 6   |  |
| 5  | 0 1   |  |
| 6  | 12 66   |  |
| 7  | 4 5   |  |
| 8  | 3 6   |  |
| 10 | 35  |  |
| 12 | 25  |  |
| 11 | 10 9  |  |
| 12 | 88 130  |  |
| 13 | 10 4  |  |
| 14 | 15 3  |  |
| 15 | 24 52   |  |
| 16 | 13 1  |  |
| 17 | 0 0   |  |
| 18 | 56 71   |  |
| 19 | 16 0  |  |
| 20 | 12 7  |  |
| 21 | 27 21   |  |
| 22 | 12 2  |  |
| 23 | 10 0  |  |
| 24 | 93 128  |  |
| 25 | 0 0   |  |
| 26 | 5 1   |  |
| 27 | 9 8   |  |
| 28 | 9 2   |  |
| 29 | 10 0  |  |
| 30 | 17 27   |  |
| 31 | 4 0   |  |
| 32 | 7 0   |  |
| 33 | 8 2   |  |
| 34 | 5 0   |  |
| 35 | 11 0  |  |
| 36 | 56 46   |  |
| 37 | 2 0   |  |
| 38 | 4 0   |  |
| 39 | 5 4   |  |
| 40 | 9 0   |  |
| 41 | 6 0   |  |
| 42 | 4 8   |  |
| 43 | 7 0   |  |
| 44 | 4 0   |  |
| 45 | 6 1   |  |
| 46 | 6 0   |  |
| 47 | 4 1   |  |
| 48 | 36 20   |  |
| 49 | 1 0   |  |
| 50 | 1 0   |  |
| 52 | 2 0   |  |
| 54 | 1 1   |  |
| 55 | 1 0   |  |
| 57 | 2 0   |  |
| 59 | 1 0   |  |
| 60 | 6 7   |  |
| 72 | 1 0   |  |

Figure 37 & 38: Output Data

Based on the density plot and data generated, they revealed that the loan durations of 12 months have the highest count for both Good (at 130 records) and Bad (at 88 records) classification. This suggests that short-term loans are common across both categories, but Good is slightly more than Bad classifications at this duration. While at longer durations of 36 to 60 months, there are more Bad classifications (56) than Good (46) at 36 months, then this trend continues at 48 months with 36 Bad and 20 Good classifications. When it reaches beyond 60 months, although loans are rare, but the majority are classified as bad credit risks. Shorter durations which are below 12 months on the other hand have very few Bad classifications and significantly more good classifications. For instance, at 6 months, there are 12 Bad and a whopping of 66 Good classifications, showing a strong tendency for short-term loans to be associated with good credit risks.

### 3.2.3 Are loan amounts a strong predictor for credit risk classification?

```
#analysis 2-3: Are loan amounts a strong predictor for credit risk classification?
#summarize credit amount by credit risk classification
summary(AData$credit_amount)
summary(AData$class)
str(AData$credit_amount)
str(AData$class)

Class_summary <- AData %>%
  group_by(class) %>%
  summarize(CreditAmount_mean = mean(credit_amount,na.rm = TRUE),
            CreditAmount_median = median(credit_amount,na.rm = TRUE),
            CreditAmount_sd = sd(credit_amount,na.rm = TRUE))
print(Class_summary)

> #analysis 2-3: Are loan amounts a strong predictor for credit risk classification?
> #summarize credit amount by credit risk classification
> summary(AData$credit_amount)
>   Min. 1st Qu. Median Mean 3rd Qu. Max.
>    759     1255    2326    3383    4226   18424
> summary(ADatasclass)
>   length   Class      Mode
>    1400 character character
> str(ADatasclass)
> num [1:1400] 1 2 3 4 5 6 7 8 9 10 11 12 13 14 15 16 17 18 19 20 21 22 23 24 25 26 27 28 29 30 31 32 33 34 35 36 37 38 39 40 41 42 43 44 45 46 47 48 49 50 51 52 53 54 55 56 57 58 59 60 61 62 63 64 65 66 67 68 69 70 71 72 73 74 75 76 77 78 79 80 81 82 83 84 85 86 87 88 89 80 91 92 93 94 95 96 97 98 99 100 101 102 103 104 105 106 107 108 109 110 111 112 113 114 115 116 117 118 119 110 120 121 122 123 124 125 126 127 128 129 120 130 131 132 133 134 135 136 137 138 139 130 140 141 142 143 144 145 146 147 148 149 140 150 151 152 153 154 155 156 157 158 159 150 160 161 162 163 164 165 166 167 168 169 160 170 171 172 173 174 175 176 177 178 179 170 180 181 182 183 184 185 186 187 188 189 180 190 191 192 193 194 195 196 197 198 199 190 200 201 202 203 204 205 206 207 208 209 200 210 211 212 213 214 215 216 217 218 219 210 220 221 222 223 224 225 226 227 228 229 220 230 231 232 233 234 235 236 237 238 239 230 240 241 242 243 244 245 246 247 248 249 240 250 251 252 253 254 255 256 257 258 259 250 260 261 262 263 264 265 266 267 268 269 260 270 271 272 273 274 275 276 277 278 279 270 280 281 282 283 284 285 286 287 288 289 280 290 291 292 293 294 295 296 297 298 299 290 300 301 302 303 304 305 306 307 308 309 300 310 311 312 313 314 315 316 317 318 319 310 320 321 322 323 324 325 326 327 328 329 320 330 331 332 333 334 335 336 337 338 339 330 340 341 342 343 344 345 346 347 348 349 340 350 351 352 353 354 355 356 357 358 359 350 360 361 362 363 364 365 366 367 368 369 360 370 371 372 373 374 375 376 377 378 379 370 380 381 382 383 384 385 386 387 388 389 380 390 391 392 393 394 395 396 397 398 399 390 400 401 402 403 404 405 406 407 408 409 400 410 411 412 413 414 415 416 417 418 419 410 420 421 422 423 424 425 426 427 428 429 420 430 431 432 433 434 435 436 437 438 439 430 440 441 442 443 444 445 446 447 448 449 440 450 451 452 453 454 455 456 457 458 459 450 460 461 462 463 464 465 466 467 468 469 460 470 471 472 473 474 475 476 477 478 479 470 480 481 482 483 484 485 486 487 488 489 480 490 491 492 493 494 495 496 497 498 499 490 500 501 502 503 504 505 506 507 508 509 500 510 511 512 513 514 515 516 517 518 519 510 520 521 522 523 524 525 526 527 528 529 520 530 531 532 533 534 535 536 537 538 539 530 540 541 542 543 544 545 546 547 548 549 540 550 551 552 553 554 555 556 557 558 559 550 560 561 562 563 564 565 566 567 568 569 560 570 571 572 573 574 575 576 577 578 579 570 580 581 582 583 584 585 586 587 588 589 580 590 591 592 593 594 595 596 597 598 599 590 600 601 602 603 604 605 606 607 608 609 600 610 611 612 613 614 615 616 617 618 619 610 620 621 622 623 624 625 626 627 628 629 620 630 631 632 633 634 635 636 637 638 639 630 640 641 642 643 644 645 646 647 648 649 640 650 651 652 653 654 655 656 657 658 659 650 660 661 662 663 664 665 666 667 668 669 660 670 671 672 673 674 675 676 677 678 679 670 680 681 682 683 684 685 686 687 688 689 680 690 691 692 693 694 695 696 697 698 699 690 700 701 702 703 704 705 706 707 708 709 700 710 711 712 713 714 715 716 717 718 719 710 720 721 722 723 724 725 726 727 728 729 720 730 731 732 733 734 735 736 737 738 739 730 740 741 742 743 744 745 746 747 748 749 740 750 751 752 753 754 755 756 757 758 759 750 760 761 762 763 764 765 766 767 768 769 760 770 771 772 773 774 775 776 777 778 779 770 780 781 782 783 784 785 786 787 788 789 780 790 791 792 793 794 795 796 797 798 799 790 800 801 802 803 804 805 806 807 808 809 800 810 811 812 813 814 815 816 817 818 819 810 820 821 822 823 824 825 826 827 828 829 820 830 831 832 833 834 835 836 837 838 839 830 840 841 842 843 844 845 846 847 848 849 840 850 851 852 853 854 855 856 857 858 859 850 860 861 862 863 864 865 866 867 868 869 860 870 871 872 873 874 875 876 877 878 879 870 880 881 882 883 884 885 886 887 888 889 880 890 891 892 893 894 895 896 897 898 899 890 900 901 902 903 904 905 906 907 908 909 900 910 911 912 913 914 915 916 917 918 919 910 920 921 922 923 924 925 926 927 928 929 920 930 931 932 933 934 935 936 937 938 939 930 940 941 942 943 944 945 946 947 948 949 940 950 951 952 953 954 955 956 957 958 959 950 960 961 962 963 964 965 966 967 968 969 960 970 971 972 973 974 975 976 977 978 979 970 980 981 982 983 984 985 986 987 988 989 980 990 991 992 993 994 995 996 997 998 999 990 1000
```

Figure 39 & 40: Summary and Statistics

```
#generate box plot of the loan amount vs credit risk classification
ggplot(AData, aes(x = class, y = credit_amount, fill = class)) +
  geom_boxplot() +
  labs(title = "Loan Amount vs Credit Risk Classification",
       x = "Credit Risk Classification",
       y = "Loan Amount") +
  scale_fill_manual(values = c("pink","lightgreen"),name = "Credit Risk")
```

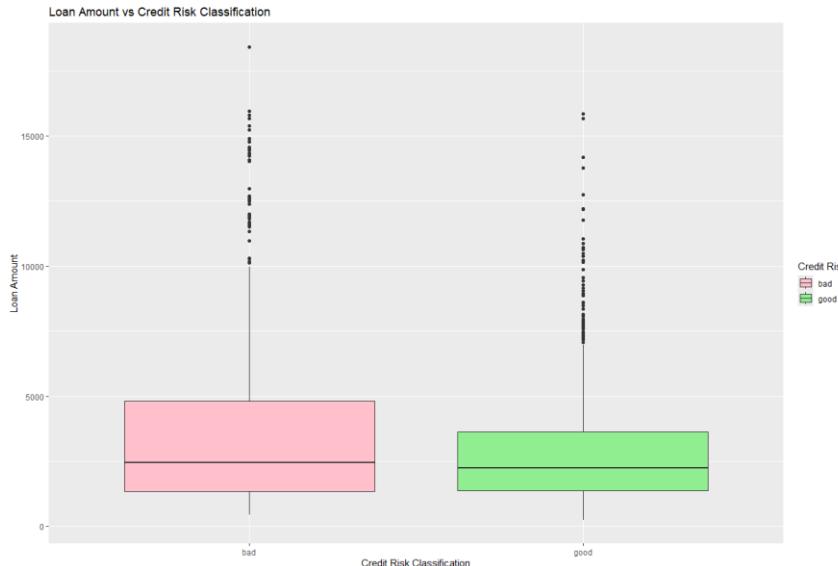


Figure 41 & 42: Boxplot and its Code Snippet

Based on the statistics and boxplot generated, we found out that individuals classified as “Bad” credit risk tend to have slightly higher loan amounts compared to individuals classified as “Good” credit risk. This can be proved with the median loan amount which Bad credit risk is 2450 while Good credit risk is slightly lower at 2244. Besides that, the loan amounts for “Bad” credit risk group show greater variability at 3405 comparing with “Good” credit risk group at 2401. This indicates that while “Bad” borrowers generally have higher loan amounts, there is

more inconsistency in the loan amounts within this group. On the other hand, we can see from the boxplot that both groups have outliers, but the “Bad” group have more extreme values which means higher loan amounts. These outliers indicate a tendency for some “Bad” credit risk borrowers to take out unusually large loans.

```
> #Run Welch Two Sample t-test (Extra Feature)
> CreditAmountClass_ttest <- t.test(credit_amount ~ class, data = AData)
> print(CreditAmountClass_ttest)

Welch Two Sample t-test

data: credit_amount by class
t = 5.0472, df = 1256.4, p-value = 5.144e-07
alternative hypothesis: true difference in means between group bad and group good is not equal to 0
95 percent confidence interval:
 485.9122 1103.8680
sample estimates:
mean in group bad mean in group good
 3780.347      2985.457
```

Figure 43: T-test results

The T-test result reveals T-value of 5.0472 with 1256.4 degree of freedom and a p-value of 5.144e-7 or 0.0000005144 which is much smaller than common significant level of 0.05. This result indicates that the loan amounts is statistically significant difference between the credit risk classification. This means that the T-test confirms that loan amounts are significantly higher for “Bad” credit risk borrowers compared to “Good” one. However, this variable alone might not be able to fully explain credit risk classification due to the overlap in the distributions as seen in the boxplot.

### 3.2.4 Does employment status interacts with loan duration to influence credit risk classification?

```
#analysis 2-4: Does employment status interacts with loan duration to influence credit risk classification?
#convert employment status and credit risk classification columns to factors
AData$employment <- as.factor(AData$employment)
AData$class <- as.factor(AData$class)
AData$duration <- as.numeric(AData$duration)
str(AData) #check the structure

#create and output summarized dataset for interaction plotting
summary2_4 <- AData %>%
  group_by(employment, duration, class) %>%
  summarize(count = n()) %>%
  ungroup()
print(summary2_4, n = 267)
```

Figure 44: Code snippet

```
#create interaction plot of employment status and loan duration by credit risk classification (Extra Feature)
ggplot(summary2_4, aes(x = duration, y = count, color = employment, group = employment)) +
  geom_line(size = 1) +
  geom_point(size = 2) +
  facet_wrap(~ class) +
  labs(
    title = "Interaction of Employment Status and Loan Duration on Credit Risk Classification",
    x = "Loan Duration (Months)",
    y = "Count of Classification",
    color = "Employment Status"
  )
```

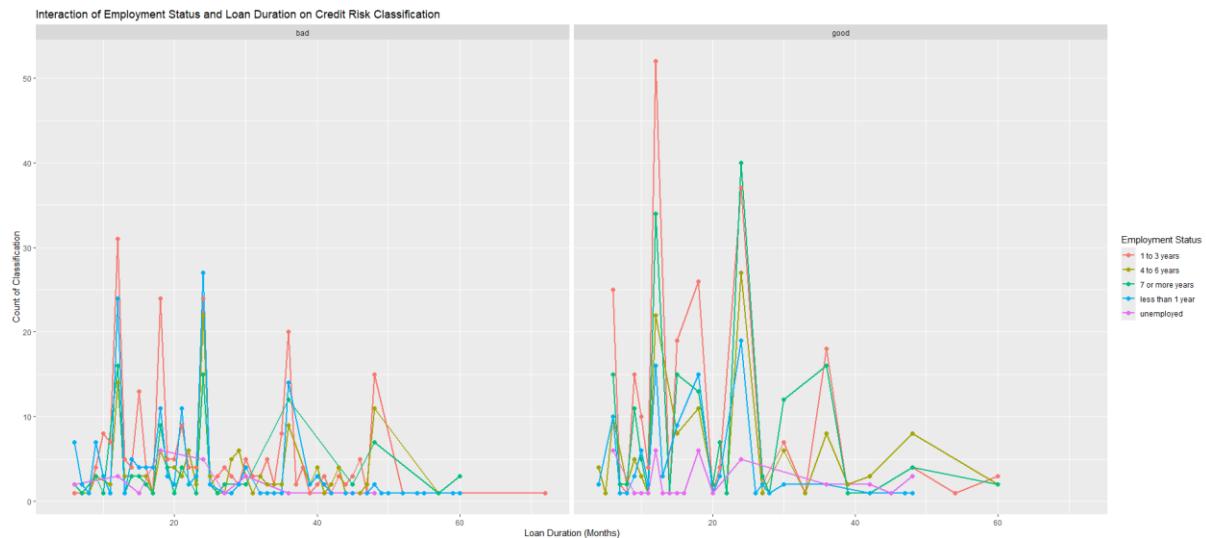


Figure 45 & 46: Interaction Plot and its Code Snippet (Extra Feature)

According to the interaction plot generated, we can identify that unemployed individuals have a significantly higher bad credit rate comparing with others who are employed, although there is some good credit, but it's much less frequent. For people in the “1 to 3 years” employment group, good credit risk is relatively more common across different loan durations, although there are instances where bad credit risk appears especially for longer durations at around 12 to 24 months.

Meanwhile based on the loan duration's role in credit risk, shorter loan durations like 6 to 12 months seems to have more balanced or better credit outcomes across most employment status. Then as loan duration increases, the tendency for bad credit increases in certain employment group, especially for the "unemployed" and those with "1 to 3 years" of employment.

Based on the visualization, good credit tends to be more prevalent in individuals with longer employment histories like "7 or more years" especially in loan duration of 12-24 months. Bad credit classification on the other hand is more frequent among unemployed individuals across all loan durations, indicating that employment status plays a major role in an individual's financial stability and ability to manage debt.

```
> #Run Logistic Regression Model with interaction term
> model2_4 <- glm(class ~ employment * duration, data = AData, family = "binomial")
> summary(model2_4)

Call:
glm(formula = class ~ employment * duration, family = "binomial",
     data = AData)

Coefficients:
                                         Estimate Std. Error z value Pr(>|z|)
(Intercept)                         1.0380192  0.2022975  5.131 2.88e-07 ***
employment4 to 6 years              -0.2049345  0.3390466 -0.604  0.5455
employment7 or more years           0.2913978  0.3361316  0.867  0.3860
employmentless than 1 year          -0.5038148  0.3497896 -1.440  0.1498
employmentunemployed                -0.3379752  0.5696981 -0.593  0.5530
duration                            -0.0530114  0.0086444 -6.132 8.65e-10 ***
employment4 to 6 years:duration    0.0140047  0.0134849  1.039  0.2990
employment7 or more years:duration 0.0219562  0.0136659  1.607  0.1081
employmentless than 1 year:duration -0.0001525  0.0161580 -0.009  0.9925
employmentunemployed:duration      0.0447443  0.0237436  1.884  0.0595 .
---
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

(Dispersion parameter for binomial family taken to be 1)

Null deviance: 1940.8  on 1399  degrees of freedom
Residual deviance: 1803.3  on 1390  degrees of freedom
AIC: 1823.3

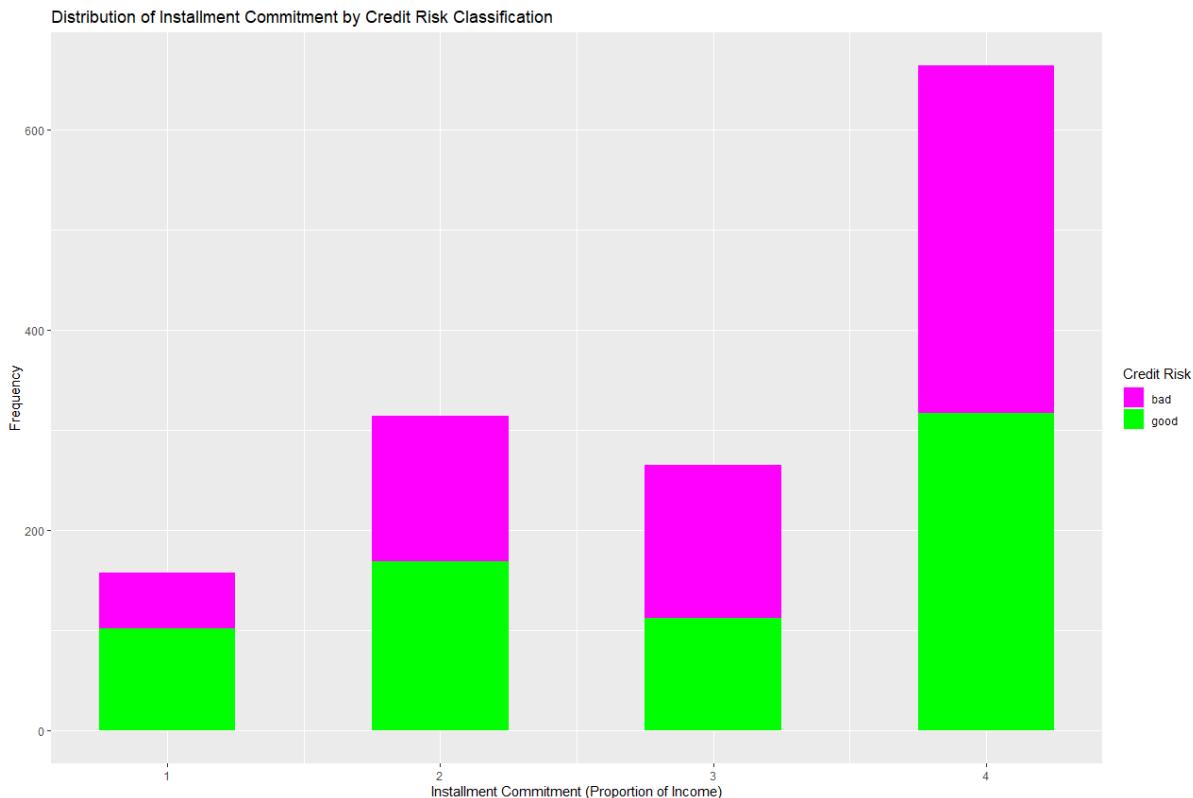
Number of Fisher Scoring iterations: 4
```

Figure 47: Logistic Regression Model

The Logistic Regression model shows an intercept of 1.038, meaning that when all factors are at their baseline, the likelihood of being classified as bad credit is higher. Employment status doesn't strongly affect the chances of being bad credit, except for those with "less than 1 year" of employment. Most employment categories on the other hand don't have a significant impact when compared to the "1 to 3 years" group. However, loan duration is a key factor. The negative coefficient for loan duration (-0.053) indicates that longer loan durations are linked to a lower chance of being classified as bad credit, suggesting that longer loans are associated with lower credit risk.

### **3.2.5 Does the proportion of income committed to loan payments (installment commitment) influence the likelihood of higher credit risk classification?**

```
> summary(AData$installment_commitment) #statistics of 'installment_commitment' column
  Min. 1st Qu. Median Mean 3rd Qu. Max.
  1.000   2.000   3.000  3.026   4.000   4.000
>
> #summarize installment commitment and credit risk classification
> summary2_5 <- ADATA %>%
+   count(installment_commitment, class)
> print(summary2_5)
  installment_commitment class n
1                      1   bad 55
2                      1  good 102
3                      2   bad 145
4                      2  good 169
5                      3   bad 153
6                      3  good 112
7                      4   bad 347
8                      4  good 317
```

*Figure 48: Summary**Figure 49: Histogram and its Code Snippet*

Based on the output and histogram generated, as instalment commitment increases, the number of applicants classified as “bad credit risk” increases, particularly at higher level of commitment (3 and 4). This shows that the individuals are using more of their income to loan payment, which lead to the high-risk classification. On the other hand, the number of “good credit risk”

applicants generally decrease as instalment commitment increases, which declines at higher commitment levels. This can be seen at the moderate commitment levels (1 and 2) has a higher proportion of good credit risk especially at lowest commitment level (1).

### **3.2.6 Conclusion of Objective 3-2**

The analysis showed how different loan-related factors (like loan purpose, duration, amount, and income commitment) affect credit risk classification.

- **Analysis 1** showed that while personal loans were more likely to be associated with bad credit risk, the loan purpose alone didn't significantly affect credit risk.
- **Analysis 2** found that longer loan durations were linked to higher chances of being classified as bad credit risk, but this effect wasn't as strong for shorter loans.
- **Analysis 3** looked at all factors together and found that applicants with high income commitment and longer loan durations were more likely to be high-risk. However, only a small percentage of applicants fit this description.
- **Analysis 5** showed that applicants who committed more of their income to loan repayments were more likely to be classified as bad credit risk, especially those with high commitments.

As conclusion, the findings suggest that loan-related factors alone don't significantly determine credit risk. Only a small group of applicants with high loan amounts, long durations, and high commitment were considered high-risk. Other factors, like employment status, income, and financial behavior, likely play a bigger role in determining credit risk.

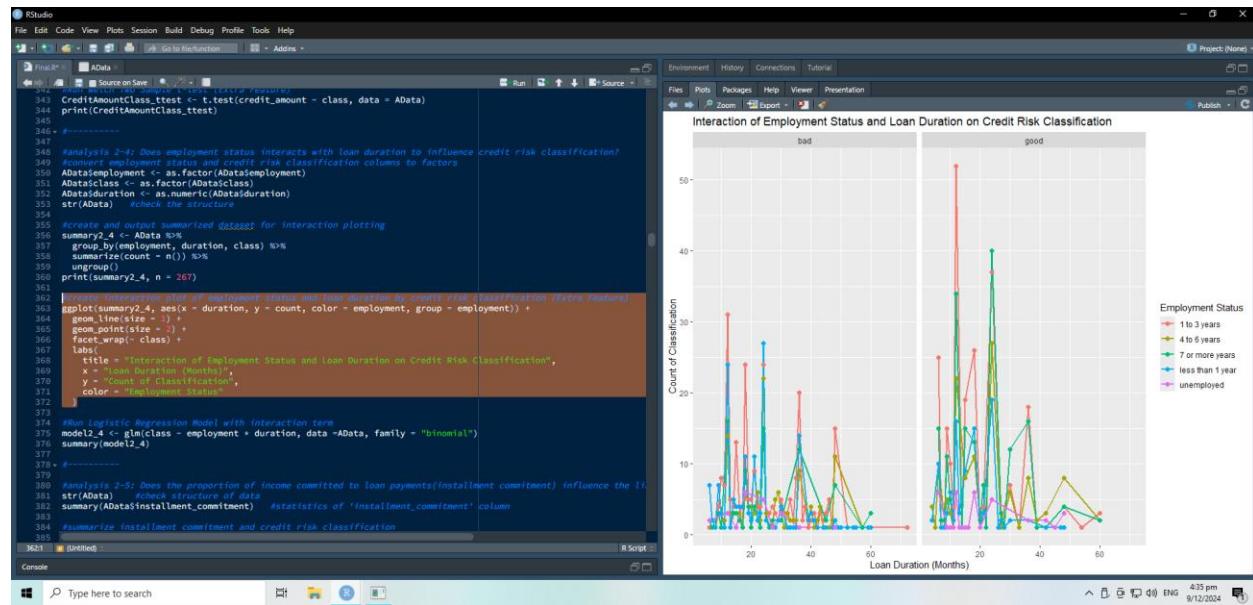
### 3.2.7 Additional Features

#### Welch Two Sample T-test

```
#Run Welch Two Sample t-test (Extra Feature)
CreditAmountClass_ttest <- t.test(credit_amount ~ class, data = AData)
print(CreditAmountClass_ttest)
```

This can be found at analysis 3. It was being used to statistically assess whether there is a significant difference in the mean loan amounts between the two credit risk groups (“good” and “bad”). It is used because it accounts for differences in variance between the two groups.

#### Interaction Plot



This can be found at analysis 4. An interaction plot is a helpful visualization tool when we want to examine whether the effect of one predictor (independent variable) on the dependent variable changes depending on the level of another predictor. In the context of my analysis, the interaction plot is used to explore whether the effect of loan duration on credit risk classification (good vs. bad) depends on the employment status of the individual.

### 3.3 Objective 3-3: To Assess the Role of Financial Commitments on Credit Risk Classification. Eleanor Permata Fry (TP072606)

#### 3.3.1 Is There A Relationship Between The Instalment Commitments And Credit Risk Classification?

##### Step 1: Data Overview

```
> # Step 1: Data Overview
> # Check Structure of Data
> str(Adata)
'data': 1400 obs. of  22 variables:
 $ checking_status : chr  "<0" "0-X<200" "no checking" "<0" ...
 $ duration        : num  6.48 12.42 24.36 24.36 12.30 ...
 $ credit_history  : chr  "critical account" "existing paid" "critical accounts" "existing paid" ...
 $ number_credits  : num  1 1 1 1 1 1 1 1 1 1 ...
 $ credit_amount   : num  1169 5931 2096 7882 4870 ...
 $ savings_status  : chr  "no known savings" "<100" "<100" ...
 $ employment      : chr  "unemployed" "<40k" "<40k" ...
 $ installment_commitment: num  4.2 2.2 3.2 3.2 2.4 ...
 $ personal_status  : chr  "male single" "female divorced/divorced/married" "male single" "male single" ...
 $ other_parties    : chr  "none" "none" "none" "guarantor" ...
 $ resided_since   : num  2.3 4.4 4.4 4.2 ...
 $ property_magnitude: chr  "real estate" "real estate" "real estate" "life insurance" ...
 $ age             : num  67.22 49.45 53.35 53.35 61.28 ...
 $ other_payment_plans: chr  "own" "own" "own" "stores" ...
 $ housing          : chr  "own" "own" "own" "for free" ...
 $ existing_credits: num  2.1 1.2 3.2 1.1 2.2 ...
 $ sex              : num  1 1 2 2 2 2 1 1 1 1 ...
 $ num_dependents  : num  1 1 2 2 2 2 1 1 1 1 ...
 $ own_telephone    : chr  "yes" "none" "none" "none" ...
 $ foreign_worker   : chr  "yes" "yes" "yes" "yes" ...
 $ class            : chr  "good" "good" "good" "good" ...
 $ income_level     : chr  "no known savings" "<100" "<100" "<100" ...
```

Figure 50: Data Structure

##### Step 2: Descriptive Statistics

```
> # Step 2: Descriptive Statistics
> # Calculate Mean, Median, and Standard Deviation of Installment Commitments by Credit Risk Class
> installment_summary <- Adata %>%
+   group_by(class) %>%
+   summarise(
+     avg_installment = mean(installment_commitment, na.rm = TRUE),
+     median_installment = median(installment_commitment, na.rm = TRUE),
+     sd_installment = sd(installment_commitment, na.rm = TRUE)
+   )
> print(installment_summary)
# A tibble: 2 × 4
  class avg_installment median_installment sd_installment
  <chr>           <dbl>                <dbl>           <dbl>
1 bad              3.13                 3            1.00
2 good             2.92                 3            1.13
```

Figure 51: Summarise

The descriptive statistics above provides a summary of the central tendency (mean and median) and variability (standard deviation) of instalment commitments for each credit risk class.

##### Step 3: Data Visualization

```
# Step 2: Data Visualization
# 1.1 Density Plot: Counts of Credit Risk by Other Payment Plans
library(ggplot2) # Library for Density Plot

ggplot(Adata, aes(x = installment_commitment, fill = class)) +
  geom_density(alpha = 0.5, color = "black") + # Add Border
  scale_fill_manual(values = c("good" = "#00A0A0", "bad" = "#FF0000")) +
  labs(title = "Density Plot of Installment Commitments by Credit Risk Classification",
       x = "Number of Installment Commitments",
       y = "Density",
       fill = "Credit Risk Class") +
  theme_minimal(base_size = 12) + # Increase Font Size
  theme(
    plot.title = element_text(hjust = 0.5, size = 13, face = "bold"),
    panel.grid.major = element_line(color = "#F0F0F0"), # Add Grid Lines
    panel.grid.minor = element_blank()
  )
```

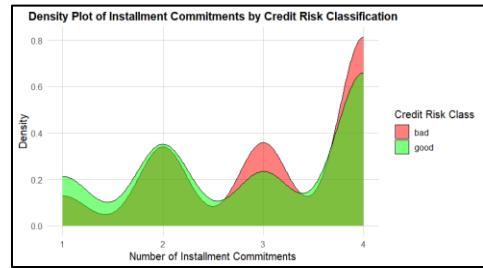


Figure 52 &amp; 53: Density Plot

```
# 1.2 Density Ridgeline Plot: Comparing Installment Between 'Good' and 'Bad' Credit Risk Classes
library(ggridges) # Library for Ridgeline Plot

ggridges(Data, aes(x = installment_commitment, y = class, fill = class)) +
  geom_density_ridges(alpha = 0.6, color = "white") +
  scale_fill_manual(values = c("good" = "#00B050", "bad" = "#E63333")) +
  labs(
    title = "Ridgeline Plot of Installment Commitments by Credit Risk Classification",
    x = "Number of Instalment Commitments",
    y = "Credit Risk Class",
    fill = "Credit Risk Class"
  ) +
  theme_minimal(base_size = 12) +
  theme(
    plot.title = element_text(hjust = 0.5, size = 13, face = "bold"),
    panel.grid.major = element_line(color = "#F0F0F0"),
    panel.grid.minor = element_blank()
  )
```

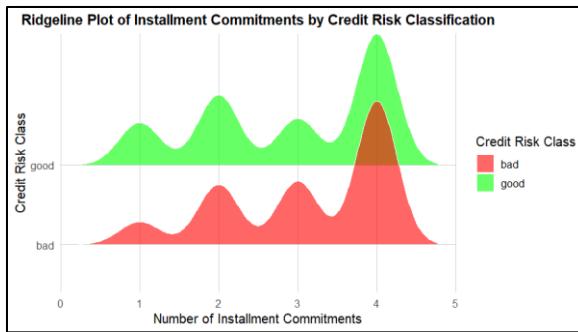


Figure 54 & 55: Ridgeline Plot

From the two figures above, both the density and ridgeline plot shares identical shapes. The distribution of “good” credit risk class peaks around the 2 instalment commitments, while the “bad” credit risk class shows peak around 3-4 instalment commitments, with a wider spread towards higher values than the “good” class. This indicates that the “good” credit risks have typically less instalment commitments, while the “bad” has more, suggesting higher financial exposure.

#### Step 4: Statistical Test

```
> # Step 4: Statistical Test
> # T-Test (Evaluate Differences in Mean Instalment Commitments Between Credit Risk Classes)
> good_class <- AData %>%
+   filter(class == "good") %>%
+   pull(installment_commitment)
>
> bad_class <- AData %>%
+   filter(class == "bad") %>%
+   pull(installment_commitment)
>
> t_test_class_result <- t.test(good_class, bad_class, var.equal = FALSE)
> print(t_test_class_result)

Welch Two Sample t-test

data: good_class and bad_class
t = -3.7096, df = 1378.4, p-value = 0.0002158
alternative hypothesis: true difference in means is not equal to 0
95 percent confidence interval:
-0.32323390 -0.09962324
sample estimates:
mean of x mean of y
2.920000 3.131429
```

Figure 56: Welch's T-Test

The statistical test was done using the Welch's t-test, to compare instalment commitments between two credit risk classes. The results showed a significant difference, with a t-statistic value of -3.71, indicating a strong difference. The p-value was extremely low, indicating the difference was unlikely due to coincidence. The "good" credit class had a mean value of 2.92, while the "bad" credit class had a mean value of 3.13. This suggests that individuals with "bad" credit risk tend to have higher average instalment commitments.

## **Conclusion**

Based on the analysis done above, we can conclude that **there is a significantly strong relationship between instalment commitments and credit risk classification**. The “bad” credit risk has a statistically significantly higher mean of instalment commitment than the “good” credit risk class. The higher number of instalment commitments in the “bad” class could suggest an increase in financial strain. On the contrary, the “good” credit class has lesser commitments, suggesting better financial responsibility especially in managing loans.

### 3.3.2 Does Income Level Predict Credit Risk Classification When Considered Alongside Instalment Commitment?

#### Step 1: Data Visualization

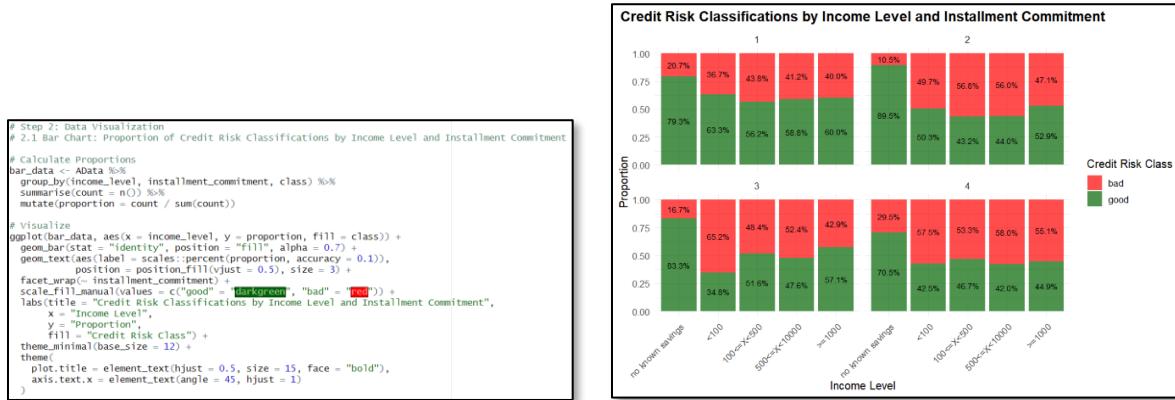


Figure 57 &amp; 58: Stacked Bar Chart

Before the visualization, the ‘*savings\_status*’ column is used as a representative for the ‘*income\_level*’, suggesting more savings means a higher income. From the stacked bar chart above, we can analyse that individuals with lower income levels (<100) have a higher proportion of “bad” credit risk classification, especially when the instalment commitment increases. For higher income levels (>=1000), the proportion of “bad” credit risk tends to be lower, even with the increase of instalment commitments. Suggesting that, a high income helps individuals to better manage their instalment commitments, despite the amount of isntalmetns.

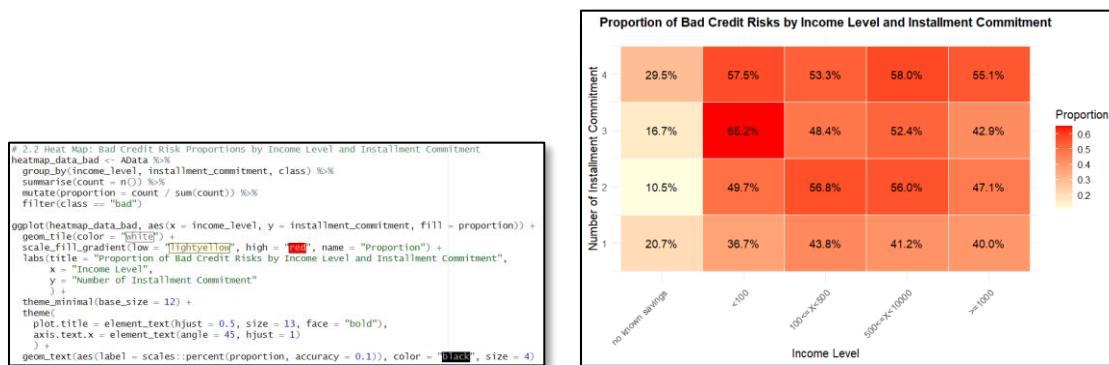


Figure 59 &amp; 60: “Bad” Heatmap Snippet

In the heatmap above, a large proportion of “bad” credit risk class can be seen in the low income (<100) and 3 instalment commitments (65.2%). Individuals with no known savings / income and 2 instalments commitments has the lowest proportion of “bad” credit risks (10.5%).

The share of “bad” credit risks increases with the amount of instalment commitments. Meanwhile, the increase of income level indicates a decrease in “bad” credit risks for any given instalments.

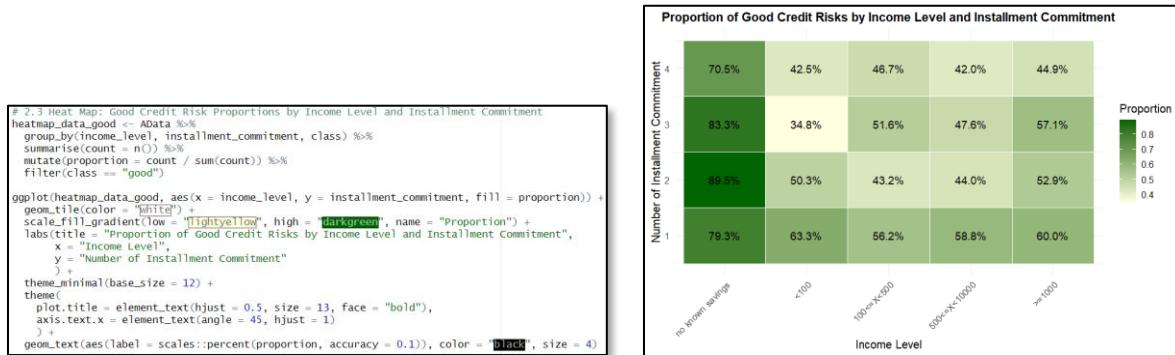


Figure 61 & 62: “Good” Heatmap Snippet

In the heatmap above, a large proportion of “good” credit risk class can be seen in the low income (<100) and with one instalment commitments (63.3%). Individuals with no know income, regardless of amount of instalments, has the highest share of “good” credit risks. Similar to the “bad” credit risk heatmap, the more the instalment commitment an individual has, despite the income level, the lower the “good” credit risks become, suggesting a notice of financial burden.

## Step 2: Statistical Tests

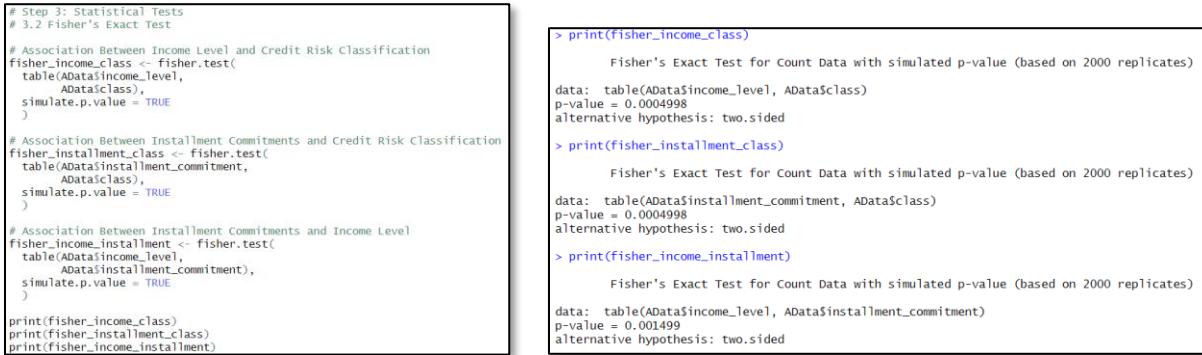


Figure 63 & 64: Fisher's Exact Test

A Fisher’s Exact Test was used for the statistical analysis between income level, instalment commitments, and credit risk classification. Based on all three’s result, a p-value of 0.0005 was received from the fisher’s test between income level, instalment commitments, and credit risk classification. Meanwhile, a p-value of 0.0015 was received between income level and instalment commitments. All three p-value are below the value of 0.005, which indicate that all three relationships are statistically significant to one another.

## **Conclusion**

Based on the analysis done above, we can conclude that **income level does predict credit risk classification when considered alongside instalment commitments.** Instalment commitments has a big impact on credit risk, but if alone, it doesn't give a full picture. The descriptive statistics, visualizations and statistical tests show how important income level is in determining creditworthiness. Income levels provide essential context, as they show a person's ability to manage financial commitments, without it can lead to overestimating or underestimating credit risk classification.

### 3.3.3 How Does The Presence Of the Financial Commitments (e.g., Other Payment Plans) Influence Credit Risk Classification?

#### Step 1: Data Visualization

```
# 3.3 Objective 3: To Assess The Role of Financial Commitments on Credit Risk Classification
# 3.3.3 How Does The Presence Of the Financial Commitments (e.g., Other Payment Plans) Influence Credit Risk Classification?

# Step 1: Data Visualization
# 1.1 Bar Chart: Counts of Credit Risk by Other Payment Plans
ggplot(data = adata %>%
  group_by(other_payment_plans, class) %>%
  summarise(count = n()) %>%
  ungroup() %>%
  group_by(other_payment_plans) %>%
  mutate(proportion = count / sum(count)) * 100) # Calculate Number of Individuals as Percentage
  aes(x = other_payment_plans, y = "Number of individuals (Counts)", fill = "Credit Risk Class")
  geom_bar(stat = "identity", color = "#F0A0A0", width = 0.7) +
  geom_text(stat = "count", aes(label = ..count..), position = position_dodge(0.9), vjust = -0.5, size = 4) +
  scale_fill_manual(values = c("#FF9999", "#99CC99")) +
  labs(
    title = "Counts of Credit Risk by Other Payment Plans",
    x = "Other Payment Plans",
    y = "Number of individuals (Counts)",
    fill = "Credit Risk Class"
  ) +
  theme(
    plot.title = element_text(hjust = 0.5, size = 13, face = "bold"),
    panel.grid.major = element_line(color = "#CCCCCC"),
    panel.grid.minor = element_line(),
    axis.text.x = element_text(angle = 45, hjust = 1) # Rotate X-Axis
  )
  
```

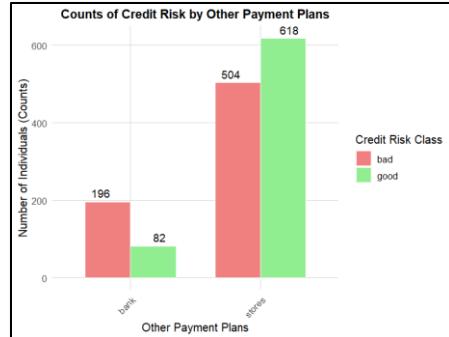


Figure 65 &amp; 66: Bar Chart Snippet

The bar chart shows that a significant number of individuals prefer store payment plans, with 1122 individuals in total. Store-based plans are primarily used by individuals with "good" credit risks (618), while bank-based plans are more popular with "bad" credit risks (196). This preference is due to easier accessibility, flexible terms, and less demanding qualifying conditions. Bank payment plans are more likely to be linked to "bad" credit risks, suggesting that individuals with financial struggles may face more difficulties under these plans.

```
# 3.3 Pie chart: Proportion of Credit Risk Classes (Good and Bad) by Other Payment Plans
pie_chart_data <- Adata %>
  group_by(other_payment_plans, class) %>
  summarise(count = n()) %>#
  ungroup() %>#
  group_by(other_payment_plans) %>#
  mutate(proportion = count / sum(count)) * 100) # Calculate Number of Individuals as Percentage
  aes(x = "", y = "proportion", fill = "class")
  geom_bar(stat = "identity", width = 1, color = "#CCCCCC") +
  facet_wrap(~other_payment_plans, ncol = 2) + # Facet by Other Payment Plans
  scale_fill_manual(values = c("#FF9999", "#99CC99")) +
  labs(
    title = "Proportion of Credit Risk Classes by Other Payment Plans",
    fill = "Credit Risk Class"
  ) +
  theme_void(base_size = 12) +
  geom_text(aes(label = paste0(round(proportion, 1), "%")),
            position = position_stack(vjust = 0.5),
            size = 4)
  
```

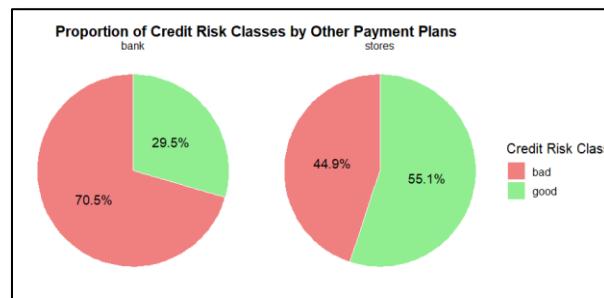


Figure 67 &amp; 68: Pie Chart Snippet

From the pie charts above, we can evaluate that bank payment plans are mostly associated by individuals with "bad" credit risks (70.5%), while store payment plans are evenly split though slightly skewed to the "good" credit risks (55.1%). The result of these pie charts complements the bar chart as it emphasizes the differences in risk class distribution between the two payment types.

## **Step 2: Statistical Test**

```
> # Step 2: Statistical Test
> # Fisher's Exact Test (Association Between Other Payment Plans and Credit Risk Classification)
> fisher_plans_class <- fisher.test(
+   table(Adata$other_payment_plans,
+         AData$class)
+ )
> print(fisher_plans_class)

Fisher's Exact Test for Count Data
data: table(Adata$other_payment_plans, AData$class)
p-value = 1.87e-14
alternative hypothesis: true odds ratio is not equal to 1
95 percent confidence interval:
2.190406 3.941332
sample estimates:
odds ratio
2.928662
```

*Figure 69: Fisher's Exact Test*

The Fisher's Test revealed a statistically significant relationship between payment plan type and credit risk classification. The p-value is extremely low, confirming the significant relationship. The odds ratio indicates a strong association between bank payment plans and a higher chance of "bad" credit risk classification compared to store payment plans.

## **Conclusion**

The analysis reveals a significant relationship between payment plan types and credit risk classification. Store payment plans are more suited to "good" credit risk individuals, indicating lower financial risk. However, a significant number of these individuals have "bad" credit, suggesting they may be financially secure but burdened due to their accessibility and less-strict repayment rules. Bank payment plans are more suited to "bad" credit risk individuals, making them more difficult to use. This relationship could help financial institutions analyse risks associated with specific payment plans and adjust their credit regulations accordingly.

### 3.3.4 What Demographic Factors (e.g., Age, Marital Status) Moderate The Relationship Between Instalment Commitments and Credit Risk Classification?

#### Step1 1: Data Visualization

```
# 3.3.3 Objective 3: To Assess the Role of Financial Commitments on Credit Risk Classification
# 3.3.4 What Demographic Factors (e.g., Age, Marital Status) Moderate The Relationship Between Instalment Commitments and Credit Risk Classification?

# Step 1: Data Visualization
# 1.1 Jittered Scatter Plot: Visualize Age by Instalment Commitments and Credit Risk
ggplot(Adat, aes(x = personal.age, y = installment.commitment, color = class)) +
  geom_jitter(alpha = 0.6, width = 0.2, height = 0.2) +
  labs(
    title = "Relationship between Age, Instalment Commitments, and Credit Risk Classification",
    x = "Age (years)",
    y = "Number of Instalment Commitments",
    color = "Credit Risk Class"
  ) +
  theme_minimal() +
  theme(
    plot.title = element_text(hjust = 0.5, size = 13, face = "bold"),
    panel.grid.major = element_line(color = "#D9E1F2"),
    panel.grid.minor = element_line(color = "#D9E1F2")
  )
scale_color_manual(name = "Credit Risk Class", values = c("red", "darkgreen"))


```

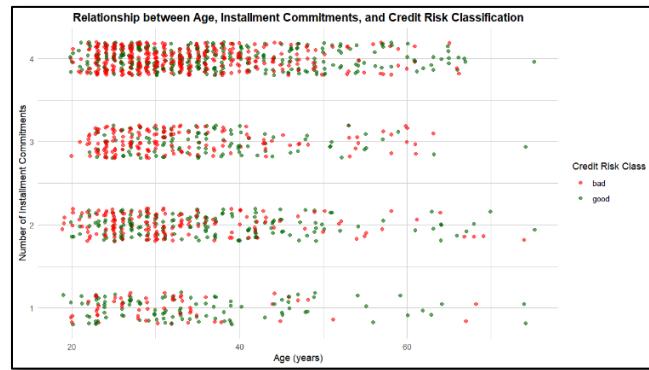


Figure 70 & 71: Jittered Scatter Plot Snippet

The jittered plot above shows that aged 20-40 of individuals with high instalments fall under "bad" credit risk, indicating financial vulnerabilities. This highlights the need for financial literacy and supervision of high-commitment loans. Those above 40 have less commitments and are more stable.

```
# 2.1 Dot Plot: Visualize Marital Status by Instalment Commitments and Credit Risk
ggplot(Adat, aes(x = personal.status, y = installment.commitment, color = class)) +
  geom_jitter(width = 0.2, alpha = 0.6, height = 0.2) +
  labs(
    title = "Relationship between Marital Status, Instalment Commitments, and Credit Risk Classification by Marital Status",
    x = "Marital Status",
    y = "Number of Instalment Commitments",
    color = "Credit Risk Class"
  ) +
  theme_minimal(base_size = 12) +
  theme(
    plot.title = element_text(hjust = 0.5, size = 13, face = "bold"),
    panel.grid.major = element_line(color = "#D9E1F2"),
    panel.grid.minor = element_line(color = "#D9E1F2")
  )
scale_color_manual(name = "Credit Risk Class", values = c("red", "darkgreen"))


```

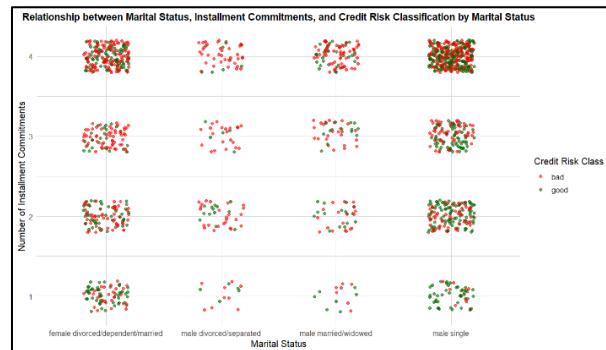


Figure 72 & 73: Dot Plot Snippet

The dot plot shows that males who are divorced or separated are more susceptible to "bad" credit risk classifications, while females exhibit better credit behaviour at lower commitments, and single males display balanced credit behaviour, but "bad" risks spike at the fourth instalment.

## Step 2: Statistical Test

```
# Step 2: Statistical Test
# ANOVA Test
# Convert Numerical Data to Factor
abata %>%
  mutate(class = as.factor(class),
        personal_status = as.factor(personal_status))

# Specify Formula for ANOVA
formula_str <- "instalment_commitment ~ class + personal_status + age"

# Perform ANOVA
anova_result <- aov(as.formula(formula_str), data = abata)

# Summarize the ANOVA Results
summary(anova_result)

# Revert Variable Back to Original Class
abata$class <- as.character(abata$class)
abata$personal_status <- as.character(abata$personal_status)

# Check
str(abata)
```

```
> # Summarize the ANOVA Results
> summary(anova_result)

Df Sum Sq Mean Sq F value    Pr(>F)
class          1   15.6   15.646  14.020 0.000188 ***
personal_status 3   32.0   10.673   9.564 2.98e-06 ***
age             1     1.7     1.717   1.538 0.215069
Residuals      1394 1555.7   1.116
---
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
```

Figure 74 & 75: ANOVA Test

The ANOVA test revealed that class significantly influences instalment commitments, while marital/personal status has a low impact. Age does not directly moderate the relationship between instalment commitments and credit risk classification, with a p-value of 0.215069 over 0.05. These findings suggest that age does not directly influence instalment commitments.

## Conclusion

The study reveals that **personal/marital status moderates the relationship** between instalment commitments and credit risk classification. Groups like divorced/separated males face the highest risk, while single males, females, and married men vary in risk patterns. **Age doesn't directly moderate** this relationship but can influence risk classifications if it interacts with instalment commitments and marital status. Younger individuals are more vulnerable to "bad" credit risks with higher commitments, while older individuals show better financial stability.

### 3.3.5 How Does The Duration Of Residence Impact The Relationship Between Financial Commitments And Credit Risk Classification?

#### Step 1: Data Visualization

```
# 3.3 Objective 3: To Assess The Role of Financial Commitments on Credit Risk Classification
# 3.3.5 How Does The Duration of Residence Impact The Relationship Between Financial Commitments And Credit Risk Classification?
library(dplyr)

# Step 1: Data Visualization
# 1.1 Title Plot: Relationship between Residence Duration, Installment Commitments, and Credit Risk Class
# Generate Title Plot
title_data <- Adata %>
  group_by(residence_since, class) %>%
  summarise(mean_commitment = mean(installment_commitment, na.rm = TRUE)) %>%
  ungroup()

# Generate Title Plot
ggplot(title_data, aes(x = factor(residence_since), y = class, fill = mean_commitment)) +
  geom_tile(color = "#E0FFFF", size = 0.5) + # color = "#E0FFFF", size = 0.5
  scale_fill_gradient_low("lightblue", high = "darkblue", name = "Mean Commitment") +
  labs(title = "Title Plot of Mean Financial Commitments by Residence Duration and Credit Risk Class",
       x = "Duration of Residence (years)",
       y = "Credit Risk Class") +
  theme_minimal(base_size = 12) +
  plot.title = element_text(hjust = 0.5, size = 13, face = "bold")
  point_text(aes(label = round(mean_commitment, 1), color = "#00FFFF", size = 4) # Mean Commitment Labels
```

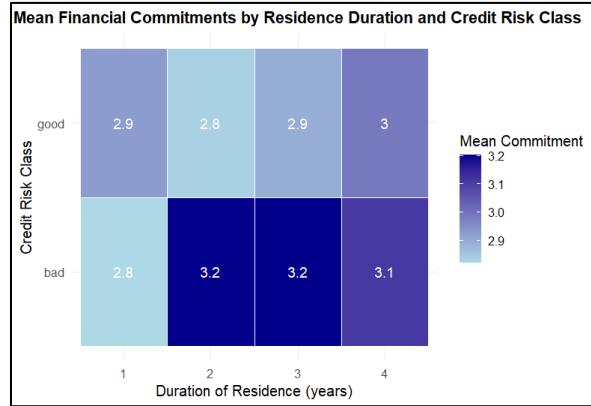


Figure 76 & 77: Tile Plot Snippet

The tile plot above summarizes the mean instalment commitments between residence duration and credit risk class. The tile plot shows that individuals with “bad” credit risk tends to have a higher mean of instalment commitments, specifically with the residence duration of 2 or 3 years.

```
# 1.2 Violin Plot: Distribution of Instalment Commitments by Residence Duration and Credit Risk Class
ggplot(Adata, aes(x = factor(residence_since), y = installment_commitment, fill = class)) +
  geom_violin(trim = FALSE, alpha = 0.6) +
  geom_boxplot(outlier.colour = NA, outlier.size = 0.5) +
  position_dodge(0.9) +
  scale_fill_manual(values = c("good" = "#1f77b4", "bad" = "#ff0000")) +
  labs(title = "Violin Plot of Financial Commitments by Residence Duration and Credit Risk Class",
       x = "Duration of Residence (years)",
       y = "Number of Instalment Commitment",
       fill = "Credit Risk Class") +
  theme_minimal(base_size = 12) +
  plot.title = element_text(hjust = 0.5, size = 13, face = "bold")
```

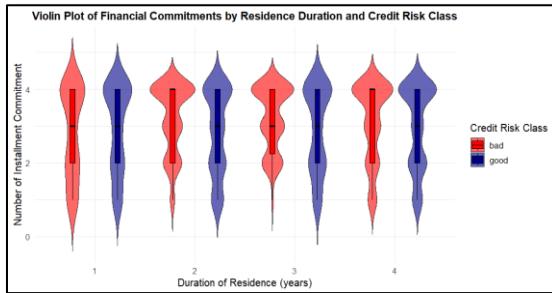


Figure 78 & 79: Violin Plot Snippet

In the violin plot above, individuals that falls under the “bad” credit risk class tend to have a wider and higher range of instalment commitments than “good” credit risk individuals throughout all residence durations years. A prominent difference can be seen clearly in year 2 and 3 of the residence duration. A nearly identical shape of violin can be seen in duration 1 and 4 years.

## **Step 2: Statistical Test**

```
> # Step 2: Statistical Test  
> # Kruskal-Wallis Test ( Compare Distributions Among Multiple Groups )  
> kruskal_result <- kruskal.test(  
+   instalment_commitment ~ interaction(residence_since, class),  
+   data = AData  
+ )  
> print(kruskal_result)  
  
Kruskal-Wallis rank sum test  
  
data: instalment_commitment by interaction(residence_since, class)  
Kruskal-Wallis chi-squared = 18.126, df = 7, p-value = 0.01141
```

*Figure 80: Kruskal-Wallis Test*

The Kruskal-Wallis's test was used to compare the distributions of residence duration, instalment commitments, and credit risk classification across multiple groups. The results showed a significant difference ( $p\text{-value} = 0.01141$ ) at a 5% significance level, indicating that the combination of residence duration and credit risk classification significantly influences the distribution of instalment commitments.

## **Conclusion**

The study reveals that individuals with short residence durations maintain small instalment commitments, regardless of their credit risk class. Residence duration significantly influences instalment commitments, with the "bad" class showing higher commitments at 2 and 3 years. In the 4 years duration category, commitments slightly decrease, suggesting individuals in the "bad" credit risk class are trying to stabilize their instalment commitment amount. This analysis can help financial institutions design targeted interventions based on residence duration, making it a crucial factor in credit risk classification analysis.

### **3.3.6 Additional Features**

From the analysis done above, several additional features were used and added in order to support the analysis process. Here are some of the additional features used:

#### **3.3.6.1 Fisher's Exact Test**

```
# Step 2: Statistical Test
> # Fisher's Exact Test (Association Between Other Payment Plans and Credit Risk Classification)
> fisher_plans_class <- fisher.test(
+   table(AData$other_payment_plans,
+         AData$class),
+   )
> print(fisher_plans_class)

Fisher's Exact Test for Count Data

data: table(AData$other_payment_plans, AData$class)
p-value = 1.87e-14
alternative hypothesis: true odds ratio is not equal to 1
95 percent confidence interval:
0.941332 3.941332
sample estimates:
odds ratio
2.928862
```

A Fisher's Exact test was used to determine if there is a significant relationship between two categorical variables, especially for small sample size. In the figure above, is used to determine the relationship between other payment plans and credit risk classes as they both are categorical (chr) variables.

#### **3.3.6.2 Kruskal-Wallis Test**

```
# Step 2: Statistical Test
> # Kruskal-Wallis Test ( Compare Distributions Among Multiple Groups )
> kruskal_result <- kruskal.test(
+   installment_commitment ~ interaction(residence_since, class),
+   data = AData
+   )
> print(kruskal_result)

Kruskal-Wallis rank sum test

data: installment_commitment by interaction(residence_since, class)
Kruskal-wallis chi-squared = 18.126, df = 7, p-value = 0.01141
```

A Kruskal-Wallis test was used for the statistical test in the fifth analysis. This test is a non-parametric test for determining whether two or more groups' medians are from the same distribution. In the figure above, the test was used to compare the instalment commitments across different credit risk groups while also taking residence duration as one of the influences.

#### **3.3.6.3 ANOVA (Analysis of Variance) Test**

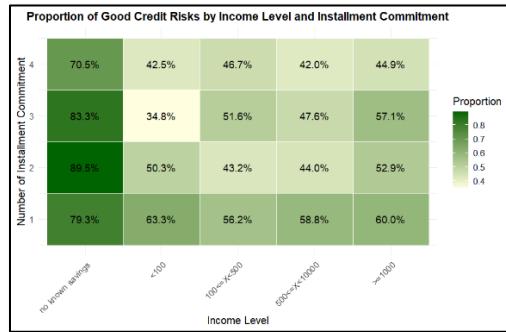
```
> # Summarize the ANOVA Results
> summary(anova_result)

Df Sum Sq Mean Sq F value    Pr(>F)
class          1   15.6   15.646  14.020 0.000188 ***
personal_status 3   32.0   10.673   9.564 2.98e-06 ***
age            1   1.7    1.717    1.538 0.215069
Residuals     1394 1555.7   1.116
---
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
```

An ANOVA test was used to compare the means of three or more groups, assuming identical variances and normal distribution. To avoid incorrect results, variables must be converted

to categorical grouping variables. The test was used to assess differences between demographic factors like age and personal status across credit risk classes.

### 3.3.6.4 Heatmap



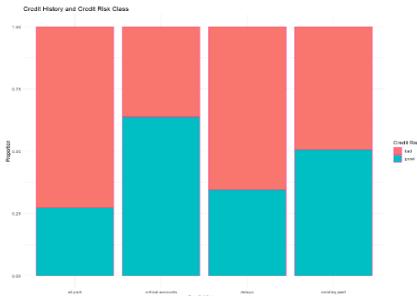
A heatmap is used to serve as a graphical representation of data in which the intensity of each hue represents a distinct value. A heatmap is often used to visualize correlations, associations, or patterns in complicated datasets. The key feature of a heatmap is the hue indicator where the value is differentiated by the intensity of that hue.

### 3.4 Objective 4: To explore the impact of customer credit history on credit risk classification. - CHONG ZHI YUE (TP067869)

3.4.1 Is there a significant relationship between the history of past loans and current credit risk classification?

```
# > # 1.1 Is there a relationship between credit history and credit risk classification?
# > # Pearson's Chi-squared Test
# > # H0: There is no relationship between credit history and credit risk classification
# > # H1: There is a relationship between credit history and credit risk classification
# > # Chi-Squared Test Statistic: Chi-Squared = 73.082, df = 3, p-value = 9.335e-16
# > # P-value is very small, so we reject the null hypothesis
# > # Therefore, there is a strong predictor of credit risk classification!
# > # We can also use a bar chart to visualize the proportions of credit risk classes
# > # Based on credit history
# > # Bar Chart
# > # Credit History vs Credit Risk
```

*Figure 81 & 82: Analysis 4.1*



*Figure 83: Create Bar Chart*

Based on the chi-squared test of independence between the credit history and class variables in the contingency table, it shows that a p-value of 9.335e – 16 greatly below the significance level of 0.05. As a result, there is a highly significant relationship between credit history and credit class, we can reject the null hypothesis.

Categories for all paid and existing paid have a higher proportion of individuals with good credit risk; for critical accounts and delays have a higher proportion of individuals with bad credit risk. Customers who have always paid their bills on time and have a history of paying off existing accounts are good credit risk. Those with critical accounts, signifying overdue or missed payments, are strongly likely to be high-risk borrowers.

```
> #4.4.2 Is credit history a strong predictor of credit risk classification?
> #Conversion
> AData$Class <- as.factor(AData$Class)
> #Create Logistic Regression Model
> logistic_risk<-glm(class ~ credit_history, data = AData, family = binomial)
> summary(logistic_risk)

Call:
glm(formula = class ~ credit_history, family = binomial, data = AData)

Coefficients:
Estimate Std. Error z value Pr(>|z|)
(Intercept) -0.9808    0.1954 -5.019 5.20e-07 ***
credit_historycritical accounts 1.0466   0.2268  6.378 1.71e-09 ***
credit_historydelays 0.3390   0.2523  1.344  0.179
credit_historyexisting paid  1.0062   0.2093  4.807 1.53e-06 ***
---
Signif. codes: 0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

(Dispersion parameter for binomial family taken to be 1)

Null deviance: 1940.8 on 1390 degrees of freedom
Residual deviance: 1866.0 on 1396 degrees of freedom
AIC: 1874

Number of Fisher Scoring iterations: 4
```

*Figure 84: Testing Credit History Is a Strong Predictor or Not*

The p-values for critical accounts, delays, and existing paid are all relevant to the 0.01 level. This shows that we can be certain that these credit history categories do affect the probability of being in the bad credit. When credit history information for a new customer is provided, the logistic model can be estimated the probability of the customer being a high-risk borrower. The probability can be used to inform decisions like approving or rejecting a loan application.

Low p-values and high coefficients show that critical accounts and already paid accounts have a large and meaningful impact on predicting credit risk. Although delays have some effect, they are not statistically significant, which suggests that they may not be a reliable indicator of credit risk on their own.

```
> #Predicted Probabilities
> AData$predicted_risk <- predict(logistic_risk, type = "response")
>
> #Calculate Mean Predicted Probabilities for Each Credit History Category
> mean_predicted_risk <- AData %>%
+ group_by(credit_history) %>%
+ summarise(mean_predicted_risk = mean(predicted_risk, na.rm = TRUE))
>
> #Bar Plot
> ggplot(mean_predicted_risk, aes(x = credit_history, y = mean_predicted_risk, fill = credit_history)) +
+ geom_bar(stat = "identity", color = "black") +
+ labs(title = "Mean Predicted Probability of Credit Risk by Credit History",
+ x = "Credit History",
+ y = "Mean Predicted Probability") +
+ theme_minimal() +
+ theme(axis.text.x = element_text(angle = 45, hjust = 1)) +
+ guides(fill = "none")
```

*Figure 85: Calculation of Mean Predicted Probabilities & Create Bar Plot*

| predicted_risk |
|----------------|
| 0.6377953      |
| 0.5063114      |
| 0.6377953      |
| 0.5063114      |
| 0.3448276      |
| 0.5063114      |
| 0.5063114      |
| 0.6377953      |
| 0.5063114      |
| 0.6377953      |
| 0.5063114      |
| 0.5063114      |
| 0.6377953      |
| 0.5063114      |
| 0.5063114      |
| 0.6377953      |
| 0.5063114      |
| 0.6377953      |
| 0.5063114      |
| 0.5063114      |
| 0.6377953      |

Figure 86: New Column for Predicted Probabilities of Credit Risk

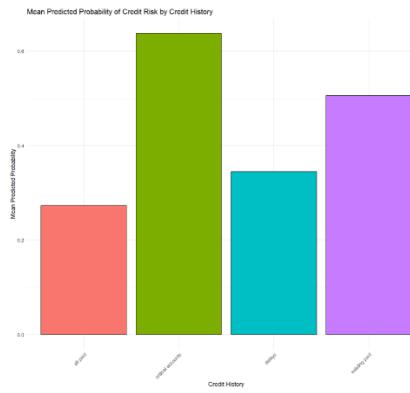


Figure 87: Bar Plot

From the bar chart above, the critical accounts categories have the highest predicted probability of credit risk. All paid categories have the lowest predicted probability, indicating that individuals with no prior credit issues are the least likely to default. They would require consideration of additional factors that could affect reputation such as the duration of payment history.

### 3.4.2 Does the number of previous loan defaults predict the likelihood of being classified as a high-risk customer?

```
#analysis 4-2: Does the number of previous loan defaults predict the likelihood of being classified as a high-risk customer?
#conversion
#AnatRisk <- as.factor(AnatRisk)
#Logistic regression Model
logit <- glm(formula = class ~ credit_history, data = AnatRisk, family = binomial)
summary(logit[,defaultrs])
predicted_probabilities
AnatRisk$predicted_risk_defaults <- predict(logit,defaults, type = "response")
#Calculate mean predicted Probabilities for each number of Defaults
mean_prob_defaults <- AnatRisk %>%
  group_by(defaultrs) %>%
  summarise(mean_predicted_risk = mean(predicted_risk_defaults, na.rm = TRUE))

#Line Plot for Predicted Probabilities
ggplot(mean_prob_defaults, aes(x = credit_history, y = mean_predicted_risk, group = 1)) +
  geom_point(color = "#E69138", size = 20) +
  geom_point(color = "#E69138", size = 20) +
  labs(title = "Mean Predicted Probability of High Risk by Credit History",
       subtitle = "Number of Previous Defaults",
       y = "Mean Predicted Probability") +
  theme_minimal()

#Violin Plot
ggplot(AnatRisk, aes(x = as.factor(credit_history), y = predicted.risk_defaults, fill = as.factor(credit_history))) +
  geom_violin(trim = FALSE) +
  facet_wrap(~ credit_history, nrow = 2, scales = "free_y") +
  labs(title = "Predicted Probability by Credit History",
       subtitle = "Number of Previous Defaults",
       y = "Predicted Probability") +
  theme_minimal()
```

Figure 88: Analysis 4.2

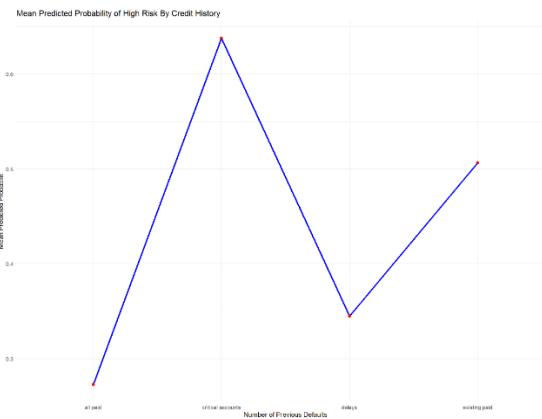
*Figure 89: New Column for Predicted Risk Defaults*

```

> #Predicted Probabilities
> Adata$predicted_risk_defaults <- predict(logistic_defaults, type = "response")
>
> #calculate Mean Predicted Probabilities for Each number of Defaults
+ mean_prob_defaults <- Adata %>%
+   group_by(credit_history) %>%
+   summarise(mean_predicted_risk = mean(predicted_risk_defaults, na.rm = TRUE))
> #Line Plot for Predicted Probabilities
+ ggplot(mean_prob_defaults, aes(x = credit_history, y = mean_predicted_risk, group = 1)) +
+   geom_line(color = "blue", linewidth = 1) +
+   geom_point(color = "red", size = 2) +
+   labs(title = "Mean Predicted Probability of High Risk By Credit History",
+        x = "Number of Previous Defaults",
+        y = "Mean Predicted Probability") +
+   theme_minimal()

```

*Figure 90: Calculation of Mean Predicted Probabilities & Create Point Plot*



*Figure 91: Point Plot*

For the line plot graph, this shows the average predicted probability of high risk for each credit history group. Individuals with a history of critical accounts have the highest mean predicted probability of credit risk. In contrast, all paid accounts have the lowest risk among the other 3 categories.

```
> #Violin Plot
> ggplot(data_aes, aes(x = as.factor(credit_history), y = predicted_risk_defaults, fill = as.factor(credit_history))) +
  + geom_violin(trim = FALSE) +
  + labs(title = "Violin Plot of Predicted Probability by Credit History",
        subtitle = "Number of Previous Defaults",
        x = "Credit History" +
        y = "Predicted Probability") +
  + theme_minimal()
```

*Figure 92: Create Violin Plot*

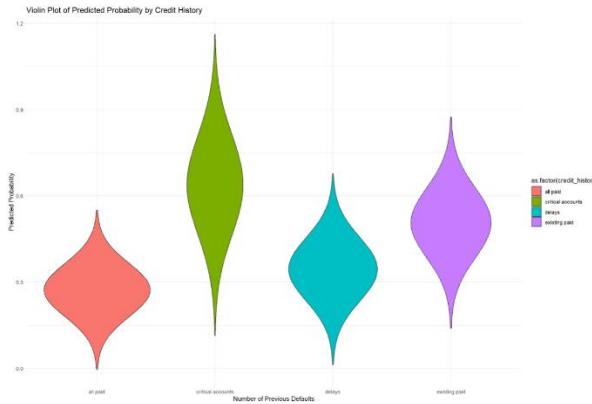


Figure 93: Violin Plot

In the violin plot, it shows that all paid group has a wider spread, with many individuals are predicted to have a low risk of bad credit and high predicted probabilities. The critical accounts have a narrow range centred around a more concentrated distribution of predicted probability, which means that they have higher predicted probabilities of bad risk.

### 3.4.3 How does the length of credit history impact the credit risk classification?

```
> #4.3 Length of Credit History and Credit Risk Classification
> #Conversion
> AData$duration <- as.numeric(as.character(AData$duration))
> #Convert Duration into Broader Categorical Group
> AData$duration <- cut(AData$duration,
+                         breaks = c(-Inf, 12, 24, 36, Inf), # Adjust ranges as needed
+                         labels = c("0-12", "13-24", "25-36", "36+"),
+                         right = TRUE)
> #Create Tables for Duration and Credit class
> duration_class_table <- table(AData$duration, AData$class)
> duration_class_table
```

| Duration | 0   | 1   |
|----------|-----|-----|
| 0-12     | 146 | 283 |
| 13-24    | 295 | 289 |
| 25-36    | 150 | 86  |
| 36+      | 109 | 42  |

```
> #chi-Square Test
> duration_chisq <- chisq.test(duration_class_table)
> duration_chisq
```

Pearson's Chi-squared test  
data: duration\_class\_table  
X-squared = 90.897, df = 3, p-value < 2.2e-16

Figure 94: Length of Credit History Impact Credit Risk Classification

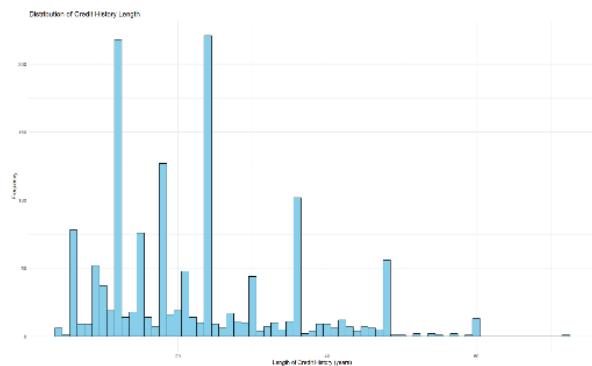


Figure 95: Histogram

The histogram clearly shows that most people in this dataset have a credit history between 20 and 25 years. On the histogram, we can observe a long tail, showing that some individuals have credit histories that are 40 years or more lengthy. A modest right-skewedness in the distribution shows that there are more people with shorter credit histories than those with longer ones.

```
> #Boxplot
> ggplot(AData, aes(x = class, y = duration, fill = class)) +
+   geom_boxplot() +
+   labs(title = "Credit History Length and Credit Risk Classification",
+       x = "Credit Class",
+       y = "Duration (Credit History Length)") +
+   theme_minimal()
```

Figure 96: Creating Boxplot

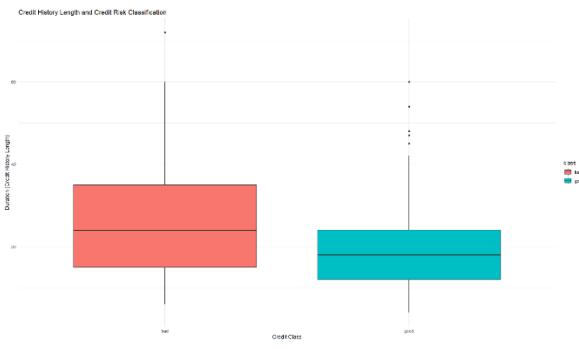


Figure 97: Boxplot

In comparison to people with good credit risk, those with bad credit risk typically have shorter median credit history lengths. Longer credit history is frequently linked to greater reliability. The boxplots for both good and negative credit risks show that the lengths of credit histories in each group range widely. The boxplot for bad credit risks is broader, showing that there is greater variance in the lengths of credit histories among those who have bad credit.

```
> #Convert duration into broader categorical group
> AData$duration <- cut(AData$duration,
+                         breaks = c(-Inf, 9, 24, 36, Inf), labels = TRUE,
+                         right = TRUE)
> #Create tables for duration and credit class
> duration_class_table <- table(AData$duration, AData$class)
> duration_class_table
```

| duration  | bad | good |
|-----------|-----|------|
| (-Inf, 9] | 283 | 283  |
| [9, 24]   | 295 | 295  |
| [24, 36]  | 86  | 86   |
| [36, Inf] | 42  | 309  |

```
> fisher <- fisher.test(duration_class_table, simulate.p.value = TRUE)
> fisher$fisher_exact.p <- fisher.test(duration_class_table, simulate.p.value = TRUE)
> fisher$fisher_exact.p
Fisher's Exact Test for Count Data with simulated p-value (based on 2000 replicates)

data: duration_class_table
p.value = 0.0004998
alternative hypothesis: two.sided

> chisq <- chisq.test(duration_class_table, simulate.p.value = TRUE)
> chisq$chisq
```

| duration  | bad | good |
|-----------|-----|------|
| (-Inf, 9] | 283 | 283  |
| [9, 24]   | 295 | 295  |
| [24, 36]  | 86  | 86   |
| [36, Inf] | 42  | 309  |

```

Pearson's Chi-squared test with simulated p-value (based on 2000 replicates)

data: duration_class_table
X-squared = 88.807, df = NA, p-value = 0.0004998
```

Figure 98: Conversion, New Table, Fisher's Exact Test & Chi-Squared Test

After converting the variable into different categories, we would need to create a contingency table which can show the frequency counts of each combination of duration category and credit risk class. When sample sizes are limited, the fisher test is utilized to determine if two categorical variables are independent. In the chi-squared test, it reveals a strong relationship between credit risk classification and credit history duration, with both tests producing a p-value of 0.004998.

```
> #Conversion  
> AData$duration <- as.factor(AData$duration)  
> AData$class <- as.factor(AData$class)  
> #Bar Chart  
> ggplot(AData, aes(x = duration, fill = class)) +  
+   geom_bar(position = "dodge") +  
+   labs(title = "Credit History Length and Credit Risk",  
+       x = "Duration",  
+       y = "Count",  
+       fill = "Credit Class") +  
+   theme_minimal()
```

Figure 99: Conversion & Create Bar Chart

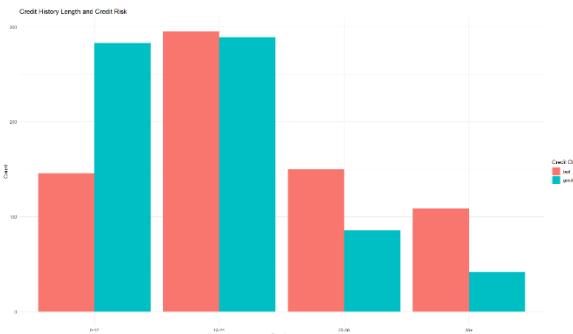


Figure 100: Bar Chart

This bar chart represents the distribution of good and bad credit classes across different duration categories. The result shows that the longer the credit histories, the lower the credit risk that might be associated with. From the graph for the category of 25 – 36 months and 36+ months, the proportion of bad classes decrease, the longer credit histories will associate with lower risk as the duration increases.

### *Additional Features – Showing Percentage in Graph*

```
> #Calculate Percentages for Duration and Class
> duration_class_df <- as.data.frame(duration_class_table)
> colnames(duration_class_df) <- c("Duration", "Class", "Count")
>
> duration_class_df <- duration_class_df %>%
+   group_by(Duration) %>%
+   mutate(Percentage = Count / sum(Count) * 100)
> #Bar Chart with Percentages
> ggplot(duration_class_df, aes(x = Duration, y = Count, fill = Class)) +
+   geom_bar(stat = "identity") +
+   geom_text(aes(label = paste0(round(Percentage, 1), "%")),
+             position = position_stack(vjust = 0.5)) +
+   labs(title = "Credit History Length and Risk Classification",
+        x = "Credit History Length (Duration)",
+        y = "Count",
+        fill = "Credit Class") +
+   theme_minimal()
```

Figure 101: Calculation Percentages & Create Bar Chart with Percentages

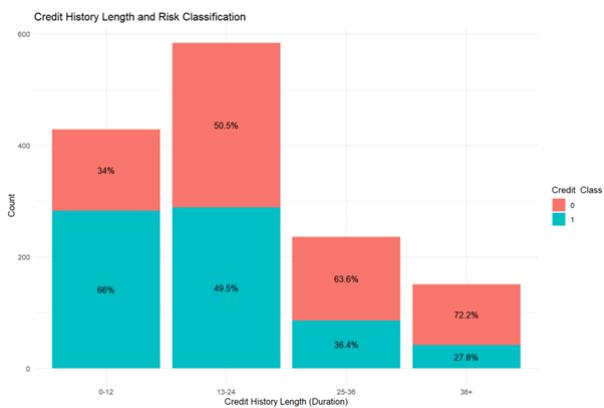


Figure 102: Bar Chart with Percentages

Based on the graph, it shows that:

- For 0-12 months duration, most individuals are classified as high-risk, making up 66% of the total for becoming risky.
- For 36+ months duration, the longest credit histories have a higher percentage of non-risky loans, with 72.2% of individuals classified as “Class 0” (non-risky).

3.4.4 What are the external factors that interact with past credit history (e.g., income, employment status) to influence credit risk classification?

```

> Analysis 4.4: What are the external factors that interact with past credit history (e.g., income, employment status) to influence credit risk classification?
> m4.4 <- train(logit ~ ., data = Abaca, family = "binomial")
> m4.4
Call: train(logit ~ ., data = Abaca, family = "binomial")
Parameters: maxDepth = 4, ntree = 500
Number of Trees: 500
Error measures: misclassification, 0.001
Training time: 0.13 sec
Test time: 0.00 sec
Number of Fisher Scoring Iterations: 4

> summary(m4.4)

Logistic regression model with Interaction between Credit History and Employment
History, employ, interaction ~ g1/(class ~ credit.history + employment ~ credit.history*employment, data = Abaca, family = binomial)
summary(m4.4)$logit$interaction

g1/(class ~ credit.history + employment ~ credit.history*employment,
family = "binomial", data = Abaca)

coefficients:
Estimate Std. Error z value Pr(>|z|)
(Intercept) 0.10798 0.10017 1.871 0.06135
credit.history*accounts 0.00000 0.00000 0.000 1.00000
credit.history*keys 0.07983 0.38577 0.204 0.81612
employment*keys -0.00000 0.00000 0.000 1.00000
employment*4 to years -0.78700 0.86219 -1.806 0.19417
employed*or new year -0.14431 0.57021 -0.251 0.80222
employed*unemployed 0.12584 0.71670 -0.176 0.72146
employment*6 years -0.17514 0.81760 -0.136 0.89411
employment*12 years 0.00000 0.00000 0.000 1.00000
employment*6 years*employment*6 years 0.44996 0.77560 0.176 0.84218
employment*6 years*employment*12 years 0.00000 0.00000 0.000 1.00000
credit.history*accounts*employment*6 years 1.12156 0.81543 1.382 0.08969
credit.history*keys*employment*6 years 0.00000 0.00000 0.000 1.00000
credit.history*accounts*employment*12 years 1.12297 0.61323 1.842 0.06552
credit.history*keys*employment*12 years 0.00000 0.00000 0.000 1.00000
credit.history*keys*employment*6 years -0.15938 0.76924 -0.201 0.81514
credit.history*keys*employment*12 years -0.09940 0.36048 1.098 0.08919
employment*keys*employment*6 years 0.61352 0.25243 0.490 0.62423
employment*keys*employment*12 years 0.72024 0.39964 0.723 0.45957

Signif. codes: 0 '***' 0.001 '**' 0.01 '*' 0.1 ' ' 1

Optimization parameter for firthr: family = "binomial"
Number of Fisher Scoring Iterations: 13
> m4.4
Call: train(logit ~ ., data = Abaca, family = "binomial")
Number of Trees: 500
Error measures: misclassification, 0.001
Training time: 0.13 sec
Test time: 0.00 sec
Number of Fisher Scoring Iterations: 4

> summary(m4.4)

Logistic regression model with Interaction between Credit History and Employment
History, employ, interaction ~ g1/(class ~ credit.history + employment ~ credit.history*employment, data = Abaca, family = binomial)
summary(m4.4)$logit$interaction

g1/(class ~ credit.history + employment ~ credit.history*employment,
family = "binomial", data = Abaca)

coefficients:
Estimate Std. Error z value Pr(>|z|)
(Intercept) 0.10798 0.10017 1.871 0.06135
credit.history*accounts 0.00000 0.00000 0.000 1.00000
credit.history*keys 0.07983 0.38577 0.204 0.81612
employment*keys -0.00000 0.00000 0.000 1.00000
employment*4 to years -0.78700 0.86219 -1.806 0.19417
employed*or new year -0.14431 0.57021 -0.251 0.80222
employed*unemployed 0.12584 0.71670 -0.176 0.72146
employment*6 years -0.17514 0.81760 -0.136 0.89411
employment*12 years 0.00000 0.00000 0.000 1.00000
employment*6 years*employment*6 years 0.44996 0.77560 0.176 0.84218
employment*6 years*employment*12 years 0.00000 0.00000 0.000 1.00000
credit.history*accounts*employment*6 years 1.12156 0.81543 1.382 0.08969
credit.history*keys*employment*6 years 0.00000 0.00000 0.000 1.00000
credit.history*accounts*employment*12 years 1.12297 0.61323 1.842 0.06552
credit.history*keys*employment*12 years 0.00000 0.00000 0.000 1.00000
credit.history*keys*employment*6 years -0.15938 0.76924 -0.201 0.81514
credit.history*keys*employment*12 years -0.09940 0.36048 1.098 0.08919
employment*keys*employment*6 years 0.61352 0.25243 0.490 0.62423
employment*keys*employment*12 years 0.72024 0.39964 0.723 0.45957
```

*Figure 103: Conversion & Analysis 4.4.I*

The probability of default compared to the reference category to the reference category, which is most likely to all paid, is displayed by the coefficients for the different credit history categories. For instance, if employment status maintains the same, a positive coefficient for critical accounts indicates that those with critical accounts are more likely to default than those with all paid accounts.

For credit history, those who have always made their payments on time are less likely to have bad credit than those who have a history of critical accounts and paid accounts. Individuals with credit histories of 4 to 6 years or less than a year are less likely to have bad credit compared to those with a credit history of 7 or more years. While unemployment might slightly decrease the odds of bad credit, this effect is not statistically significant.

```
# Predict Probabilities using Interaction Node?
# Abxapredicted_prob <- predict(history_employ_interaction, type = "response")
# Summary of Predicted Probabilities
# interaction_summary <- Abxapredicted_prob %>%
#   group_by(credit_history, employment) %>%
#   summarise(mean_prob = mean(predicted_prob))
# summarise() has grouped output by 'credit_history'. You can override using the '.groups' argument
# A tibble: 20 × 3
# Groups: credit_history [4]
#   credit_history employment    mean_prob
#   <fct>        <fct>           <dbl>
# 1 all paid      1 to 3 years   0.362
# 2 all paid      4 to 6 years   0.208
# 3 all paid      7 or more years 0.286
# 4 all paid      less than 1 year 0.176
# 5 all paid      unemployed   0.333
# 6 critical accounts 1 to 3 years 0.609
# 7 critical accounts 4 to 6 years 0.005
# 8 critical accounts 7 or more years 0.788
# 9 critical accounts less than 1 year 0.397
# 10 critical accounts unemployed 0.789
# 11 delays       1 to 3 years   0.380
# 12 delays       4 to 6 years   0.400
# 13 delays       7 or more years 0.414
# 14 delays       less than 1 year 0.11
# 15 existing paid 1 to 3 years   0.500
# 16 existing paid 4 to 6 years   0.466
# 17 existing paid 7 or more years 0.465
# 18 existing paid less than 1 year 0.469
# 19 existing paid unemployed   0.613
```

*Figure 104: Interaction Model for Predicted Probabilities*

| predicted prob |
|----------------|
| 0.1764706      |
| 0.2857143      |
| 0.3617021      |
| 0.3617021      |
| 0.2857143      |
| 0.3617021      |
| 0.3617021      |
| 0.2083333      |
| 0.1764706      |
| 0.2857143      |
| 0.2083333      |
| 0.2857143      |
| 0.2857143      |
| 0.3333333      |
| 0.3617021      |
| 0.2857143      |
| 0.3617021      |

Figure 105: New Column for Predicted Probabilities

Customers with a history of critical accounts are at the highest risk of bad credit, regardless of their employment status. Those with a history of delays also face a significant risk, albeit lower than the critical accounts.

Customers with existing paid accounts have a lower risk, and those with all paid accounts have the lowest risk of bad credit. Unemployed individuals and those with less than a year of employment are at a higher risk of bad credit. On the other hand, those with 4 to 6 years or 7 or more years of employment are at a lower risk. Therefore, credit risk is greatly influenced by the mix of work position and credit history.

```
> #Bar Plot between Credit History and Employment Status on Predicted Probabilities
> ggpplot(interaction_summary, aes(x = credit_history, y = mean_prob, fill = employment)) +
+   geom_bar(stat = "identity", position = "dodge") +
+   geom_text(aes(label = round(mean_prob, 2)),
+             position_dodge(width = 0.9), vjust = -0.5) +
+   labs(title = "Interaction Effect between Credit History and Employment on Credit Risk",
+        x = "Credit History",
+        y = "Mean Predicted Probability",
+        fill = "Employment") +
+   theme_minimal()
```

Figure 106: Create Bar Plot

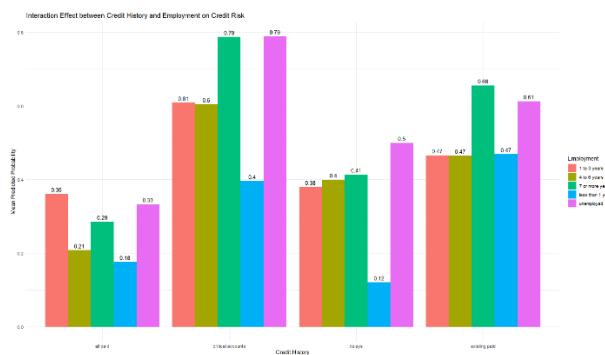


Figure 107: Bar Plot

From the observation on the graph above, critical accounts category has the highest mean predicted mean probabilities across all employment statuses in credit history effect. A less risky profile is suggested by the typically lower expected probabilities for people having a history of existing paid accounts. Within each credit history category, the bars' different heights show the interaction impact. For instance, the impact of having critical accounts is greater among those without jobs than for those with higher employment opportunities.

```
> #4.4.2
> #Interaction between Age Group and Credit class
> Adatas <- Adatas %>
+   mutate(age_group = case_when(
+     age < 30 ~ "0",
+     age > 30 & age <= 50 ~ "30-50",
+     age > 50 ~ "50 more"
+   ))
> #Create Tables for age group and credit class
> age_class_table <- table(Adatas$age_group, Adatas$class)
> age_class_table

      0    1
< 30 356 263
30-50 292 355
50 more 52 82
> #Chi-Square Test
> age_chisq <- chisq.test(age_class_table)
> age_chisq

Pearson's Chi-squared test
data: age_class_table
X-squared = 26.823, df = 2, p-value = 1.498e-06
```

Figure 108: Analysis 4.4.2



Figure 109: New Column for Age Group

Since the p-value of 1.498e-06 is smaller than the significance level of 0.05, we can reject the null hypothesis and assumed that there is a statistically significant association between age group and credit class.

```
> #Bar Plot between Age Group and credit class
> ggplot(Adata, aes(x = age_group, fill = class)) +
+   geom_bar(position = "dodge") +
+   labs(title = "Age Group and Credit Risk Class",
+        x = "Age Group",
+        y = "Count") +
+   theme_minimal()
```

Figure 110: Create Bar Plot

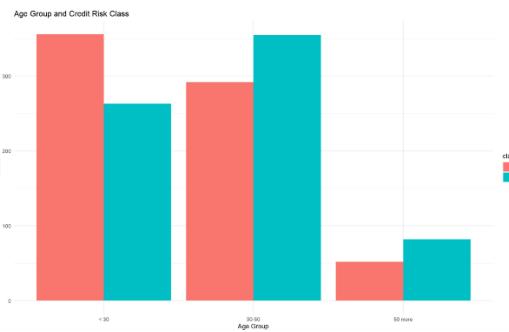


Figure 111: Bar Chart

Most of the individuals in the data set are between 30 and 50 years old, with a smaller number being younger than 30 and an even smaller number being older than 50. While older people, especially those who are over 50 often have a stronger credit history, young people, especially those under 30, are more likely to have bad credit.

```
> #4.4.3
> #Relationship between saving status and Employment
> # Grouping savings_status
> Adatasavings_status <- recode(Adatasavings_status,
+                                 "<100" = "<$500",
+                                 ">100<=x<500" = "<$500",
+                                 ">500<x<10000" = ">$500",
+                                 ">x>10000" = ">=$500",
+                                 "no known savings" = "no known savings")
> # Grouping employment
> Adatasemployment <- recode(Adatasemployment,
+                                "1" = "short_term",
+                                "<x<x<x" = "mid_term",
+                                ">x<x<x" = "long_term",
+                                ">x" = "long-term",
+                                "unemployed" = "unemployed")
```

Figure 112: Grouping

Since the p-value in the Chi-Square test is smaller, which is 0.001013, we can reject the null hypothesis. Therefore, the distribution of saving statuses is not related to employment status.

```
> #Bar Plot between Savings Status and Employment
> ggplot(Adata, aes(x = savings_status, fill = employment)) +
+   geom_bar(position = "dodge") +
+   labs(title = "Savings Status and Employment",
+        x = "Savings Status",
+        y = "Count") +
+   theme_minimal()
```

Figure 113: Create Bar Chart

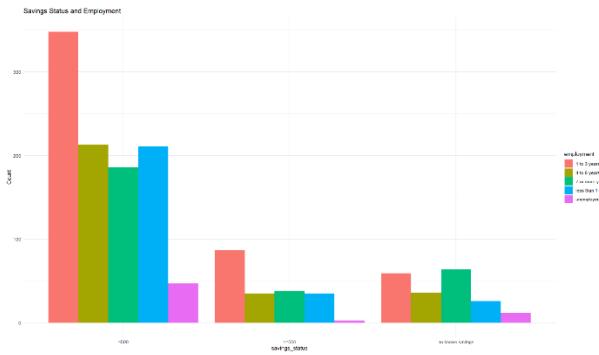


Figure 114: Bar Chart

Each bar represents a specific saving status category, such as less than 500, 500 or more, and no known savings. Those who save less, which are less than 500, often have shorter job histories, which might mean they are less financially stable or earn less money. In contrast, those who save more might have longer job histories, which are 500 or more than 500, showing that they are more financially secure and earn more money.

```
#4.4.4
#Boxplot between Credit Amount and Credit class
ggplot(data, aes(x = class, y = credit_amount, fill = class)) +
  geom_boxplot() +
  labs(title = "Credit Amount by Credit Risk Class",
       x = "Credit Class",
       y = "Credit Amount") +
  theme_minimal()
```

Figure 115: Create Boxplot

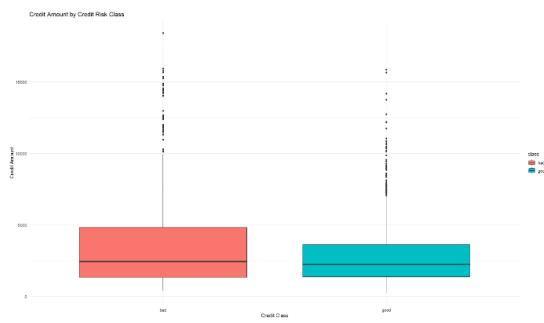


Figure 116: Boxplot

Class 0 has a larger median credit amount, which probably indicates good credit, compared to class 1, which probably indicates bad credit. This suggests that people who represent lesser credit risk are typically granted more credit. The distribution of credit amounts within each class is shown in the boxplots. The box displays the middle 50% of the data, or the interquartile range (IQR). The IQR is shown by the box and includes the middle 50% of the data.

### *Additional Features – Density Plot*

```
> #Density plot for Credit Amount
> ggplot(adata, aes(x = credit_amount)) +
+   geom_density(alpha = 0.3) +
+   labs(title = "Density Plot for Credit Amount by credit class",
+        x = "Credit Amount",
+        y = "Density") +
+   theme_minimal()
```

Figure 117: Creating Density Plot

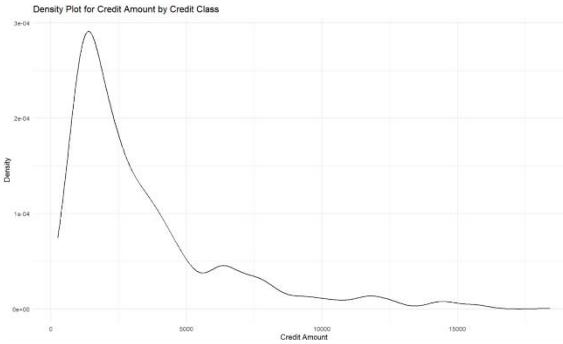


Figure 118: Density Plot

From the graph above, the density plot shows that several credit amounts fall within this range with a peak at a credit amount of around 2000 dollars. With a few big credit amounts setting the tail to the right, the distribution is skewed to the right. The mode, or the most typical credit amount is represented by the peak of the curve.

3.4.5 Is there any relationship between the purpose of having the loan and the credit risk classification?

*Figure 119: Analysis 4.5*

```

</conversion>
<AbcClass> <= class (BusinessAppliance)
  <create> <= factory (BusinessAppliance)
  <createTables> for Purpose of loan and credit class
  <purposes>_list<sub>i</sub> <--> (Table (AbcAppliancePurpose), AbcTable (list))
  <purposes>_list<sub>i</sub>
    business
    domestic
    appliance
    education
    furniture/equipment
    new car
    other
    radio/tv
    repairs
    restraining
    used car
    used equipment
    used furniture
    used radio/tv
    used repairs
    used restraining
    used used car

```

*Figure 120: Conversion & Creating Table*

To find developments, patterns, and possible risk factors, the analysis uses contingency tables and the conversion of categorical data to factors. The contingency tables show the frequency of each combination of purpose and credit class.

```

>fisher.test_toHandle_Larger_Tables
>fisher.test_pur_class$sim <- fisher.test(purpose_class_table, simulate.p.value = TRUE)
>fisher.test_pur_class$sim

    Fisher's Exact Test for Count Data with simulated p-value (based on 2000 replicates)

data: purpose_class_table
p-value = 0.0004998
alternative hypothesis: two.sided

>chiSquareTest
>purpose_class_chisq_sim <- chisq.test(purpose_class_table, simulate.p.value = TRUE)
>purpose_class_chisq_sim

    Pearson's chi-squared test with simulated p-value (based on 2000 replicates)

data: purpose_class_table
X-squared = 12600, df = NA, p-value = 0.0004998

```

Figure 121: Fisher's Exact Test & Chi-Squared Test

```
> #calculate percentages  
> purpose_class_df <- as.data.frame(purpose_class_table)  
> colnames(purpose_class_df) <- c("Purpose", "Class", "Count")  
> purpose_class_df <- purpose_class_df %>%  
+   group_by(Purpose) %>%  
+   mutate(Percentage = Count / sum(Count) * 100)
```

*Figure 122: Calculation Percentages*

```

> # Bar Chart
> ggplot(purpose_class_df, aes(x = Purpose, y = count, fill = class))
+   geom_bar(stat = "identity")
+   geom_text(aes(label = paste0(round(Percentage, 1), "%")),
+             position = position_stack(vjust = 0.5), size = 3.5) +
+   labs(title = "Loan Purpose and Risk Classification",
+        x = "Loan Purpose",
+        y = "Count",
+        subtitle = "Credit Risk Class") +
+   theme_minimal()

```

*Figure 123: Creating Bar Chart*

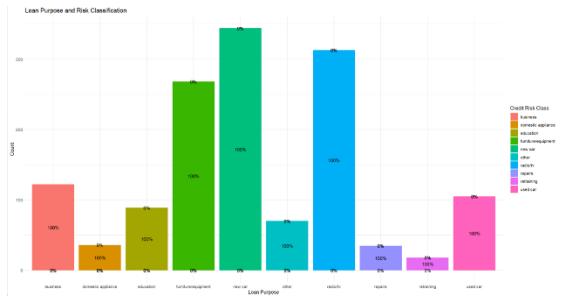


Figure 124: Bar Chart with Percentages

Based on the graph above, some loan purposes, such as new car and furniture equipment have a significantly higher number of loans compared to others loan purposes like retraining and repairs. For most loan purposes, most loans are categorized as good, and it means that a lot of loan purposes generally have a low default rate. Although the general pattern being positive, some loan purposes like repairs, retraining and used car have a comparatively greater proportion of bad loans, which linked to higher credit risk.

### **3.4.6 Is there any relationship between the categorical variables (e.g., credit history, purpose, and employment status)?**

*Figure 125: Analysis 4.6*

*Figure 126: Cross-Tabulation*

People who have good credit histories, like they always pay their bills on time, often apply for loans to buy cars or furniture. In contrast, people with bad credit histories, like always having overdue payments or not making payments, might apply for loans for riskier things like starting a business or retraining. For employment status, people with good credit histories often have stable jobs, while those with bad credit histories may be unemployed or have short job histories.

```

> Fisher's Exact Test to Handle Larger Tables
> Fisher.test purpose_sim <- Fisher.test(credit_hist$purpose_table, simulate.p.value = TRUE)
> Fisher.test purpose_sim

Fisher's Exact Test for Count Data with simulated p-value (based on 2000 replicates)

data: credit_hist$purpose_table
p-value = 0.0004989
alternative hypothesis: two.sided

> chisq.test for Credit History and Purpose
> chisq.test(credit_hist$purpose_sim) <- chisq.test(credit_hist$purpose_table, simulate.p.value = TRUE)
> chisq.test(credit_hist$purpose_sim)

Pearson's chi-squared test with simulated p-value (based on 2000 replicates)

data: credit_hist$purpose_table
X-squared = 108.24, df = 9, p-value = 0.0004988

> chisq.test for Credit History and employment status
> chisq.test(credit_hist$employment) <- chisq.test(credit_hist$employment_table)
> chisq.test(credit_hist$employment)

Pearson's chi-squared test

data: credit_hist$employment_table
X-squared = 35.815, df = 12, p-value = 0.0003469

```

*Figure 127: Fisher's Exact Test, Chi-Squared Test*

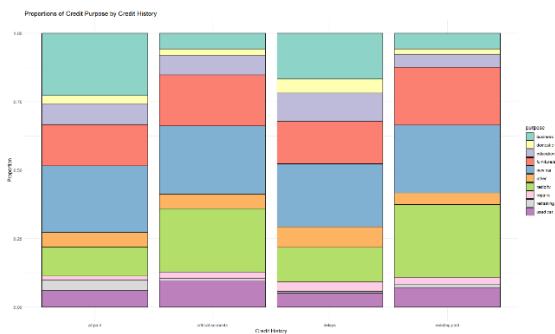
The statistical test that was conducted included the fisher test and chi-squared test. In the fisher test, the p-value is 0.0004998, meaning that there is a statistically significant association between credit history and purpose. In the chi-squared test that is used to analyze the relationship between credit history and purpose, the p-value is 0.0004998, which is the same as the Fisher's Exact Test and reinforces the conclusion that they do have strong relationships with each other.

On the other hand, the p-value in the Chi-squared Test for analyzing the relationship between credit history and employment status, the p-value is 0.0003469.

## Additional Features – Stacked Bar Plot

```
> #Stacked Bar Plot for Credit History and Purpose
> ggplot(data, aes(x = credit_history, fill = purpose)) +
+   geom_bar(position = "fill", color = "black") +
+   labs(title = "Proportions of Credit Purpose by Credit History",
+        x = "Credit History",
+        y = "Proportion") +
+   scale_fill_brewer(palette = "set3") +
+   theme_minimal()
```

*Figure 128: Creating Stacked Bar Plot*



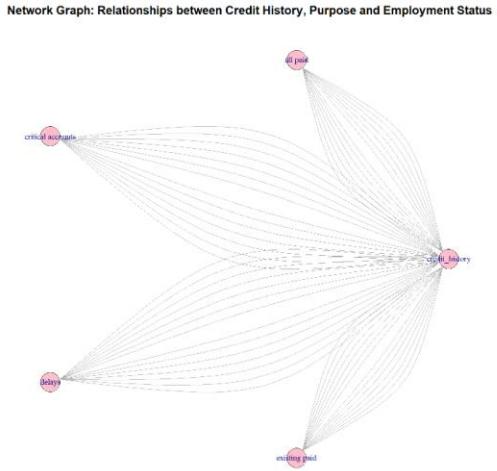
*Figure 129: Stacked Bar Plot*

In this diagram, it represents the distribution of different credit purposes across various credit history categories. For example, those with all paid credit history seems to have a higher proportion of loans for new car and radio/TV purposes compared to other groups. In contrast, individuals with delays in their credit history have a larger proportion of loans for repairs and retraining. Good credit holders, such as those who have consistently made their bill payments on time, are able to apply for loans for a variety of uses, including the purchase of furniture or cars. However, individuals with poor credit histories, such as those who have late payments, may only be eligible for loans for uses, and the conditions of these loans may be less advantageous.

## Additional Features – Network Graph

```
> #Merging Both Tables Into A List To Create A Bipartite Graph
> edges <- rbind(
+   cbind(credit_history, propcredit_history_purpose_table), rnames(credit_history_purpose_table), as.vector(credit_history_purpose_table)),
+   cbind(credit_history, propcredit_history_employment_table), rnames(credit_history_employment_table), as.vector(credit_history_employment_table))
+ )
> #Create the graph
> credit_network <- graph_from_edgelist(edges, 1:2), directed = FALSE)
> plot(credit_network)
> plot(credit_network,
+       vertex.size = 10,
+       vertex.label.cex = 0.5,
+       layout = layout_as_star,
+       edge.width = 1*(credit_network)$weight / 10,
+       main = "Network Graph: Relationships between credit history, Purpose and Employment Status")
```

*Figure 130: Creating Network Graph*



*Figure 131: Network Graph*

The graph clearly shows pink circles representing nodes, which correspond to categorical variables. The size of the nodes may reflect its relative frequency of connection to other nodes. The gray line, which is the edges representing or associations between the categories.

The graph highlights central nodes like existing paid, all paid which are heavily connected and suggests that they function as key hubs or important points in the network. In contrast, the peripheral nodes like critical accounts and delays are less connected as they may show lower interactions with other categories and highlighting the challenges faced by individuals with bad credit histories.

## **4.0 Conclusion**

### **4.1 Results**

In the analysis, all the main factors from the initial hypothesis were individually evaluated before being analysed collectively.

**Demographic Factors:** Analysis 3.1.1 shows that an individual's marital status affect their credit risk classification. From the visualization, the female group shares nearly an equal amount in both classes, although it is skewed to the "bad" risks more (54.9%). Meanwhile for males, those who are married or has separated can be seen with more "bad" risks (75.8% for separated males). While those males who are single, majority are still in "good" credit risks (61.5%). In analysis 3.1.2, job skills also do play a role though it's unlikely compared to other demographic factors as it shows little differences in both classes regardless of the individuals' skills. As for analysis 3.1.3 and 3.1.4, indicates the influence of being a foreign work, which does have a significance where foreign worker is mostly classified as "bad" risk, meanwhile gender has little potential with males being more classified as "good" risks.

**Loan-Related Factors:** Analysis 3.2.1 shows that the loan purpose does affect the credit risk classification with half of the individuals under the "bad" credit risks. The highest "bad" credit risk ratio on loan purpose is spent on vehicles (199 "bad", 145 "good"), furniture (145 "bad", 123 "good"), other purposes (63 "bad", 7 "good"). Another loan-related factor that significantly impacts credit risks is the duration of the loan. From analysis 3.2.2, it shows that individuals with a loan duration around 12 to 36 months has a fair share of both credit classes, above 36 months individuals with "bad" credit risks starts to increase, and above 60 months gather a majority of "bad" than "good" credit risks individuals.

**Instalment Commitments:** Analysis 3.3.1 shows from the Welch's t-test the "bad" credit risk class has a much higher mean value (3.13), than the "good" class (2.92) which indicates that the "bad" credit risks have a higher average in instalment commitments, the low p-value also supports this conclusion. The visualizations also add up to the fact that individuals with a high number of instalment commitments (3 and 4) tends to fall into the "bad" credit risk class, suggesting financial strain within this category.

**Credit History:** Analysis 3.4.1 shows a chi-squared test done; a low p-value was received indicating the conclusion of the significance relationship. The visualization also shows that “good” credit risk individuals tends to have paid their credit beforehand, while those who have made a delay in their payment are more classified with “bad” risks (>50%).

When we analyse all the four main factors above, we can agree that the hypotheses are falls under acceptance as all these factors plays a strong role in determining credit risk classification. Despite some variables being not as strong as others (e.g., age is not a strong demographic factor to determine, but personal status is a strong factor) the entirety of the main factor still helps to determine the credit risk. For the first hypothesis, demographic factor plays a significant role similar to loan term and credit amount, although job status does not affect directly, it still helps to support the hypothesis if combined with the other factor.

For the second hypothesis, we can agree that this hypothesis falls under acceptance as higher credit amount do indicate “bad” credit risks proven from the analysis 3.2.1. Meanwhile, loan duration also proves to have an influence, where the longer a loan duration is, the more individuals is classified as “bad” credit risk indicated from analysis 3.3.2.

## **4.2 Recommendations**

Based on the results from multiple analysis, our group recommends that banks should adopt a data-driven strategy to improve credit risk management by using customer data to create specific credit classification models. Real-time monitoring systems and advanced analytics techniques like AI can help identify complex patterns in customer data. Customized financial products and services can be provided by dividing the client base based on risk profiles, allowing banks to identify inconspicuous patterns. To reduce risks, banks should set higher standards for lending loans with larger amounts and longer durations and adjust policies for high-risk customers. A risk-based approach to lending is recommended, considering the variables that affect credit risk.

## **4.3 Limitations and Future Direction**

This assignment faced limitations in its credit risk classification variables, including a limited range of relevant factors and time constraints. To improve accuracy and generalization, the analysis could have included a larger range of variables and more processing methods.

Additionally, a more in-depth examination of the associations between factors and their impact on credit risk could have improved the analysis. Despite these limitations, our group is satisfied with the analysis and plans to improve in future studies. The analysis provides valuable insights into data analysis and the importance of considering factors in credit risk classification.