

# Loan Approval

Shaibu Abdullateef Topa

2025-03-13

## Loan Approval

**Author:** Shaibu Abdullateef Topa

**Student ID:** CST/20/COM/00591

**Course Title:** Data Mining

**Course Code:** CSC4316

**Date:** 2025-03-12

## About the Dataset

The Loan Approval Prediction Dataset, sourced from Kaggle, contains information about loan applicants and their loan approval status. Key features include:

- **Demographic information:** Number of dependents, education level, self-employment status
- **Financial information:** Annual income, loan amount requested, loan term
- **Credit information:** CIBIL score
- **Asset information:** Residential, commercial, and luxury asset values, bank asset value
- **Target variable:** Loan status (Approved/Rejected)

The dataset consists of **4,269** entries with **13** features, providing a comprehensive view of factors potentially influencing loan approval decisions. In this dataset we will be doing binary classification since we are predicting if the loan application of an applicant will be approved or rejected based on the given features.

## 1. Data Collection

```
library(knitr)

# Load the Data Set
loan_data <- read.csv("loan_approval.csv")
```

Table 1: First Few Rows of Loan Data (Part 1)

loan_id	no_of_dependents	education	self_employed	income_annum	loan_amount	loan_term
1	2	Graduate	No	9600000	29900000	12
2	0	Not Graduate	Yes	4100000	12200000	8
3	3	Graduate	No	9100000	29700000	20
4	3	Graduate	No	8200000	30700000	8
5	5	Not Graduate	Yes	9800000	24200000	20
6	0	Graduate	Yes	4800000	13500000	10

Table 2: First Few Rows of Loan Data (Part 2)

cibil_score	residential_assets_value	commercial_assets_value	luxury_assets_value	bank_asset_value	loan_status
778	2400000	17600000	22700000	8000000	Approved
417	2700000	2200000	8800000	3300000	Rejected
506	7100000	4500000	33300000	12800000	Rejected
467	18200000	3300000	23300000	7900000	Rejected
382	12400000	8200000	29400000	5000000	Rejected
319	6800000	8300000	13700000	5100000	Rejected

## 2. Feature Engineering

### Drop irrelevant columns

```
# Remove specified columns and show shape
loan_data <- loan_data[, !(names(loan_data) %in% c("loan_id"))]

# No. of rows & Columns in the loan_data (shape)
dim(loan_data)

## [1] 4269 12
```

### Loan data structure

```
# Display the structure of the loan_data data frame
str(loan_data)

## 'data.frame': 4269 obs. of 12 variables:
## $ no_of_dependents : int 2 0 3 3 5 0 5 2 0 5 ...
## $ education : chr " Graduate" " Not Graduate" " Graduate" " Graduate" ...
## $ self_employed : chr " No" " Yes" " No" " No" ...
## $ income_annum : int 9600000 4100000 9100000 8200000 9800000 4800000 8700000 5700000 8000000 ...
## $ loan_amount : int 29900000 12200000 29700000 30700000 24200000 13500000 33000000 15000000 ...
## $ loan_term : int 12 8 20 8 20 10 4 20 20 10 ...
## $ cibil_score : int 778 417 506 467 382 319 678 382 782 388 ...
## $ residential_assets_value: int 2400000 2700000 7100000 18200000 12400000 6800000 22500000 13200000 ...
## $ commercial_assets_value : int 17600000 2200000 4500000 3300000 8200000 8300000 14800000 5700000 ...
## $ luxury_assets_value : int 22700000 8800000 33300000 23300000 29400000 13700000 29200000 11800000 ...
## $ bank_asset_value : int 8000000 3300000 12800000 7900000 5000000 5100000 4300000 6000000 ...
## $ loan_status : chr " Approved" " Rejected" " Rejected" " Rejected" " Rejected" ...
```

As we can see in the output. - The dataset consists of 4269 records - There are a total of 12 features - There are three types of datatype dtypes: int(9), char(3) - Also, We can check how many missing values available in the Non-Null Count column

```
# Load the psych package
library(psych)

# Display descriptive statistics of the loan_data data frame
describe(loan_data)
```

```
##               vars      n      mean      sd    median    trimmed
```

```
## no_of_dependents      1 4269      2.50      1.70      3      2.50
## education*           2 4269      1.50      0.50      1      1.50
## self_employed*       3 4269      1.50      0.50      2      1.50
## income_annum         4 4269 5059123.92 2806839.83 5100000 5065466.78
## loan_amount          5 4269 15133450.46 9043362.98 14500000 14742230.03
## loan_term            6 4269      10.90      5.71     10     10.87
## cibil_score          7 4269      599.94     172.43     600     600.11
## residential_assets_value 8 4269 7472616.54 6503636.59 5600000 6630084.87
## commercial_assets_value 9 4269 4973155.31 4388966.09 3700000 4414896.11
## luxury_assets_value   10 4269 15126305.93 9103753.67 14600000 14709657.59
## bank_asset_value      11 4269 4976692.43 3250185.31 4600000 4739947.32
## loan_status*         12 4269      1.38      0.48      1      1.35
##
##                      mad      min      max      range skew kurtosis
## no_of_dependents      1.48 0e+00      5        5 -0.02  -1.26
## education*            0.00 1e+00      2        1  0.01  -2.00
## self_employed*        0.00 1e+00      2        1 -0.01  -2.00
## income_annum          3558240.00 2e+05 9900000 9700000 -0.01  -1.18
## loan_amount           10229940.00 3e+05 39500000 39200000 0.31  -0.75
## loan_term              5.93 2e+00      20       18  0.04  -1.22
## cibil_score            217.94 3e+02     900      600 -0.01  -1.19
## residential_assets_value 6078660.00 -1e+05 29100000 29200000 0.98   0.18
## commercial_assets_value 4151280.00 0e+00 19400000 19400000 0.96   0.10
## luxury_assets_value    10526460.00 3e+05 39200000 38900000 0.32  -0.74
## bank_asset_value       3558240.00 0e+00 14700000 14700000 0.56  -0.40
## loan_status*           0.00 1e+00      2        1  0.50  -1.75
##
##                      se
## no_of_dependents      0.03
## education*            0.01
## self_employed*        0.01
## income_annum          42959.04
## loan_amount           138409.81
## loan_term              0.09
## cibil_score            2.64
## residential_assets_value 99538.98
## commercial_assets_value 67173.68
## luxury_assets_value    139334.10
## bank_asset_value       49744.50
## loan_status*           0.01
```

---

## Finding the unique values

```
# Define a function to print unique values of a column
uniquevals <- function(col_name, df) {
  cat(sprintf("Unique Values in %s are: %s\n", col_name,
    paste(unique(df[[col_name]]), collapse = ", ")))
}

# Iterate over each column in the data frame and print unique values
for (col_name in colnames(loan_data)) {
  uniquevals(col_name, loan_data)
  cat(rep("-", 75), "\n")
}
```

```
## Unique Values in no_of_dependents are: 2, 0, 3, 5, 4, 1
```

```
## - - - - -
## Unique Values in education are: Graduate, Not Graduate
## - - - - -
## Unique Values in self_employed are: No, Yes
## - - - - -
## Unique Values in income_annum are: 9600000, 4100000, 9100000, 8200000, 9800000, 4800000, 8700000, 5700000
## - - - - -
## Unique Values in loan_amount are: 29900000, 12200000, 29700000, 30700000, 24200000, 13500000, 33000000
## - - - - -
## Unique Values in loan_term are: 12, 8, 20, 10, 4, 2, 18, 16, 14, 6
## - - - - -
## Unique Values in cibil_score are: 778, 417, 506, 467, 382, 319, 678, 782, 388, 547, 538, 311, 679, 417
## - - - - -
## Unique Values in residential_assets_value are: 2400000, 2700000, 7100000, 18200000, 12400000, 68000000
## - - - - -
## Unique Values in commercial_assets_value are: 17600000, 2200000, 4500000, 3300000, 8200000, 8300000, 12400000
## - - - - -
## Unique Values in luxury_assets_value are: 22700000, 8800000, 33300000, 23300000, 29400000, 13700000, 12400000
## - - - - -
## Unique Values in bank_asset_value are: 8000000, 3300000, 12800000, 7900000, 5000000, 5100000, 4300000
## - - - - -
## Unique Values in loan_status are: Approved, Rejected
## - - - - -
```

---

### Checking categorical and numerical features

```
# Identify categorical and numerical columns without dplyr
cat_var <- names(loan_data)[sapply(loan_data, function(x) is.factor(x) || is.character(x))]
num_var <- names(loan_data)[sapply(loan_data, is.numeric)]
```

```
# Print the column names with labels
cat("Categorical Variables:\n")
```

```
## Categorical Variables:
```

```
print(cat_var)
```

```
## [1] "education"      "self_employed" "loan_status"
```

```
cat("\nNumerical Variables:\n")
```

```
##
```

```
## Numerical Variables:
```

```
print(num_var)
```

```
## [1] "no_of_dependents"      "income_annum"
## [3] "loan_amount"          "loan_term"
## [5] "cibil_score"           "residential_assets_value"
## [7] "commercial_assets_value" "luxury_assets_value"
## [9] "bank_asset_value"
```

---

## 3. Data Cleaning

Checking for duplicate rows

```
# Check for duplicates based on all columns
duplicates_all <- loan_data[duplicated(loan_data), ]
```

```
# Print the results
print(duplicates_all)
```

```
## [1] no_of_dependents      education              self_employed
## [4] income_annum          loan_amount           loan_term
## [7] cibil_score           residential_assets_value commercial_assets_value
## [10] luxury_assets_value    bank_asset_value      loan_status
## <0 rows> (or 0-length row.names)
```

It is clear that there are no duplicates

---

### Handling null values

```
# Count missing values in each column
colSums(is.na(loan_data))
```

```
##          no_of_dependents      education      self_employed
##                0                0                0
##          income_annum      loan_amount      loan_term
##                0                0                0
##          cibil_score residential_assets_value commercial_assets_value
##                0                0                0
##          luxury_assets_value    bank_asset_value      loan_status
##                0                0                0
```

---

### Counting zeros of each column

```
# Count zeros in each column
zero_counts <- colSums(loan_data == 0)
```

```
# Print the zero counts
print(zero_counts)
```

```
##          no_of_dependents      education      self_employed
##                712                0                0
##          income_annum      loan_amount      loan_term
##                0                0                0
##          cibil_score residential_assets_value commercial_assets_value
##                0                45                107
##          luxury_assets_value    bank_asset_value      loan_status
##                0                8                0
```

Assuming that `residential_assets_value`, `commercial_assets_value`, and `bank_asset_value` cannot be zero, we treat zero values as null values and replace them with the respective column mean. However, it's noteworthy that this process does not take into account the variable `no_of_dependents`.

---

### Replace with mean

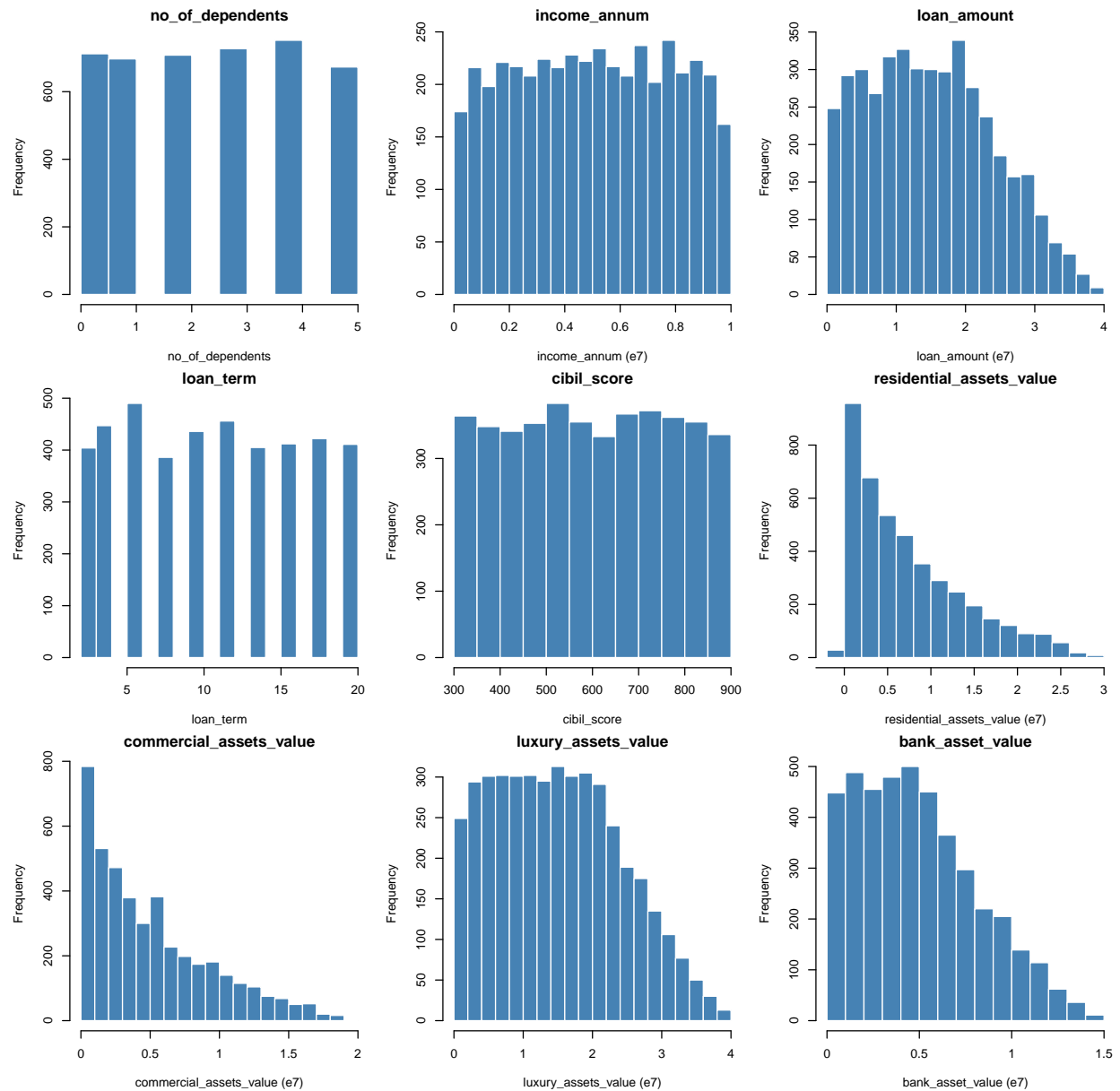
```
##          no_of_dependents      education      self_employed
##                712                0                0
##          income_annum      loan_amount      loan_term
```

```
##                                0                                0                                0
##                cibil_score residential_assets_value  commercial_assets_value
##                                0                                0                                0
##    luxury_assets_value                bank_asset_value                        loan_status
##                                0                                0                                0
```

The `loan_data` now reflects the replacement of zeros in the specified columns with their respective means

## 4. Exploratory Data Analysis

Charts for features



## Loan Status Distribution

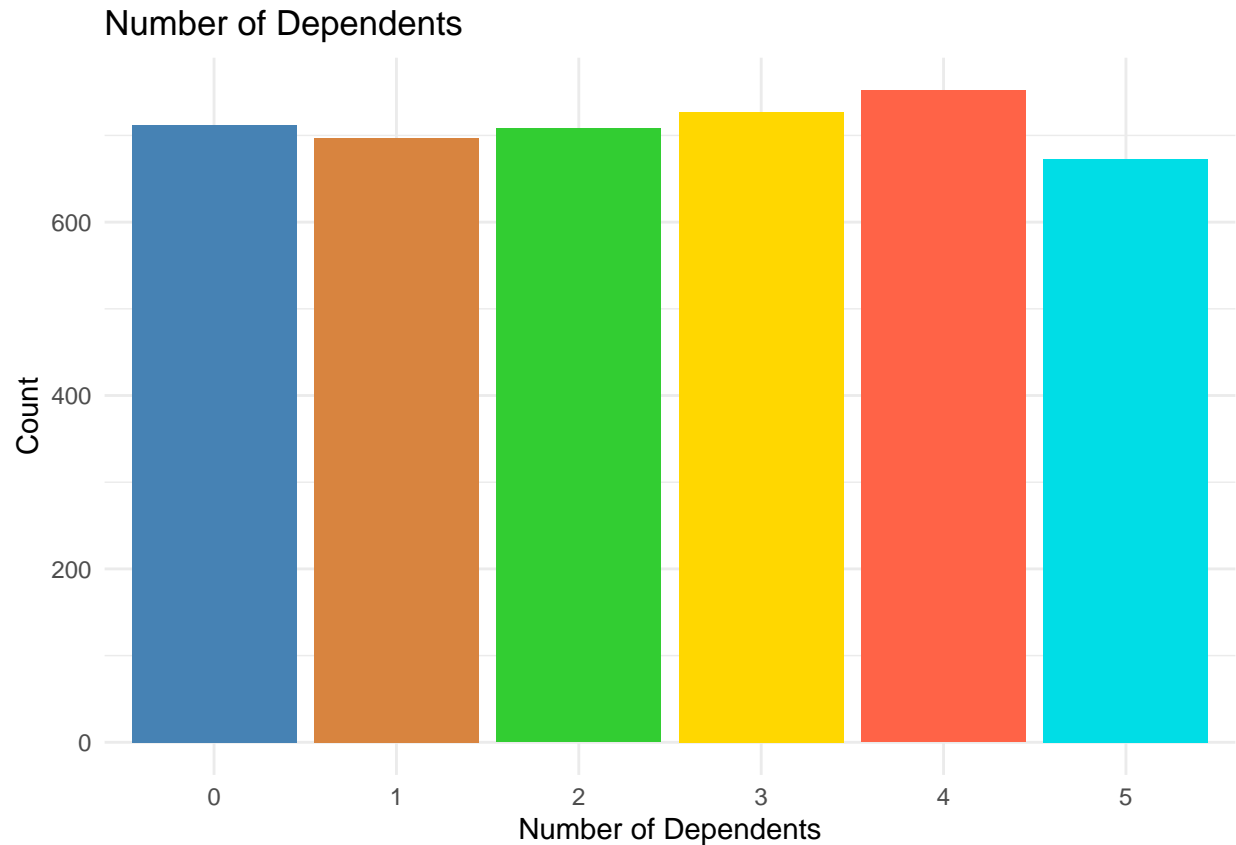
```
# Load ggplot2
library(ggplot2)

ggplot(loan_data, aes(x=factor(loan_status), fill=factor(loan_status))) +
  geom_bar(alpha=0.7) +
  theme_minimal() +
  ggtitle("Loan Approval Distribution") +
  scale_fill_manual(values = c("steelblue", "chocolate"), labels = c("Rejected", "Approved")) +
  labs(x = "Loan Status", fill = "Status")
```



The dataset shows that there are more 'Approved' instances than 'Rejected', indicating an imbalance.

```
# Plot count of no_of_dependents
ggplot(loan_data, aes(x = factor(no_of_dependents), fill = factor(no_of_dependents))) +
  geom_bar() +
  ggtitle("Number of Dependents") +
  xlab("Number of Dependents") +
  ylab("Count") +
  scale_fill_manual(values = c("#4682B4", "#D8843F", "#32CD32", "#FFD700", "#FF6347", "#00DDE6")) +
  theme_minimal() +
  theme(legend.position = "none")
```

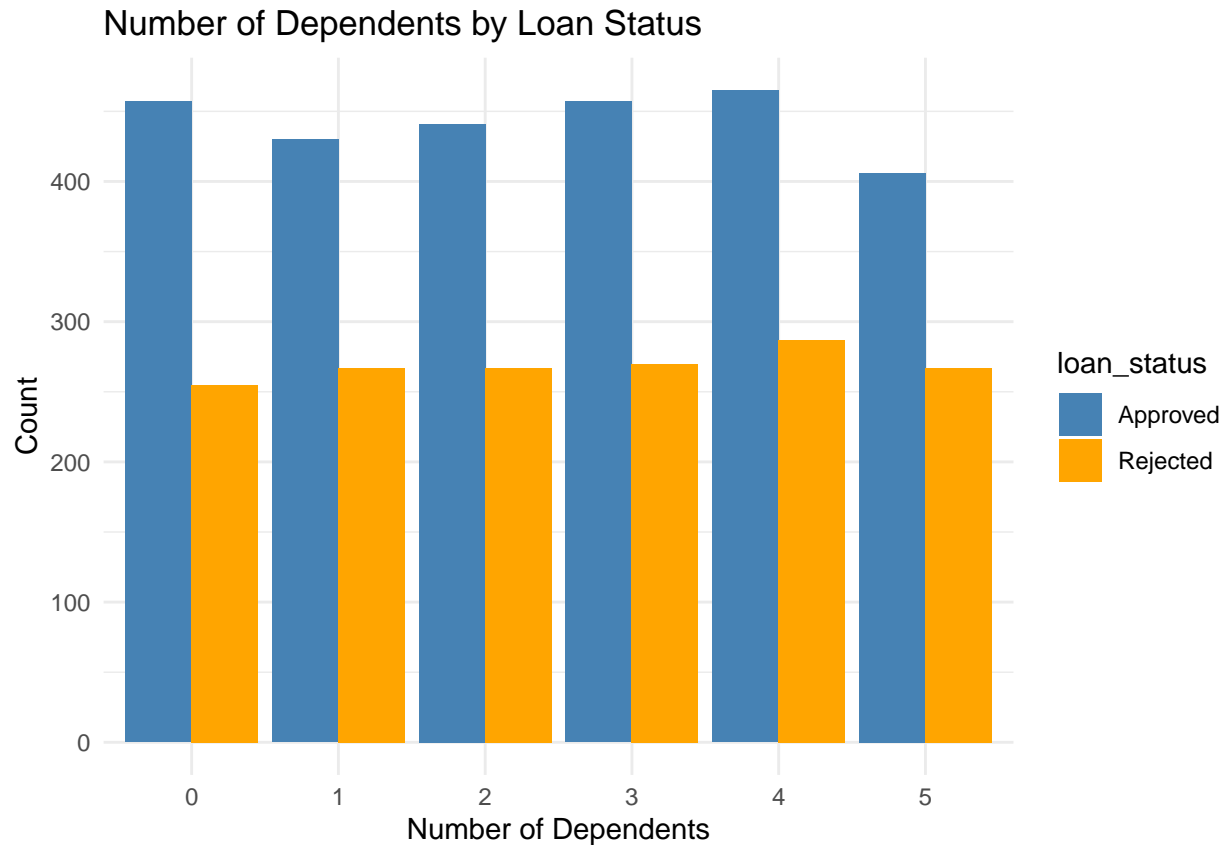


The chart shows the count of dependents for loan applicants, highlighting a clear difference in living arrangements. Notably, there isn't a substantial variance in the count of dependents. However, it's worth noting that as the number of dependents increases, the disposable income of the applicant tends to decrease. Consequently, I hypothesize that applicants with 0 or 1 dependent might have higher chances of loan approval.

### Number of Dependents Vs Loan Status

```
# Plot count of no_of_dependents with loan_status as hue
ggplot(loan_data, aes(x = factor(no_of_dependents), fill = loan_status)) +
  geom_bar(position = "dodge") +
  ggtitle("Number of Dependents by Loan Status") +
  xlab("Number of Dependents") +
  ylab("Count") +
  scale_fill_manual(values = c("steelblue", "orange"), labels = c("Approved", "Rejected")) +
  theme_minimal()
```



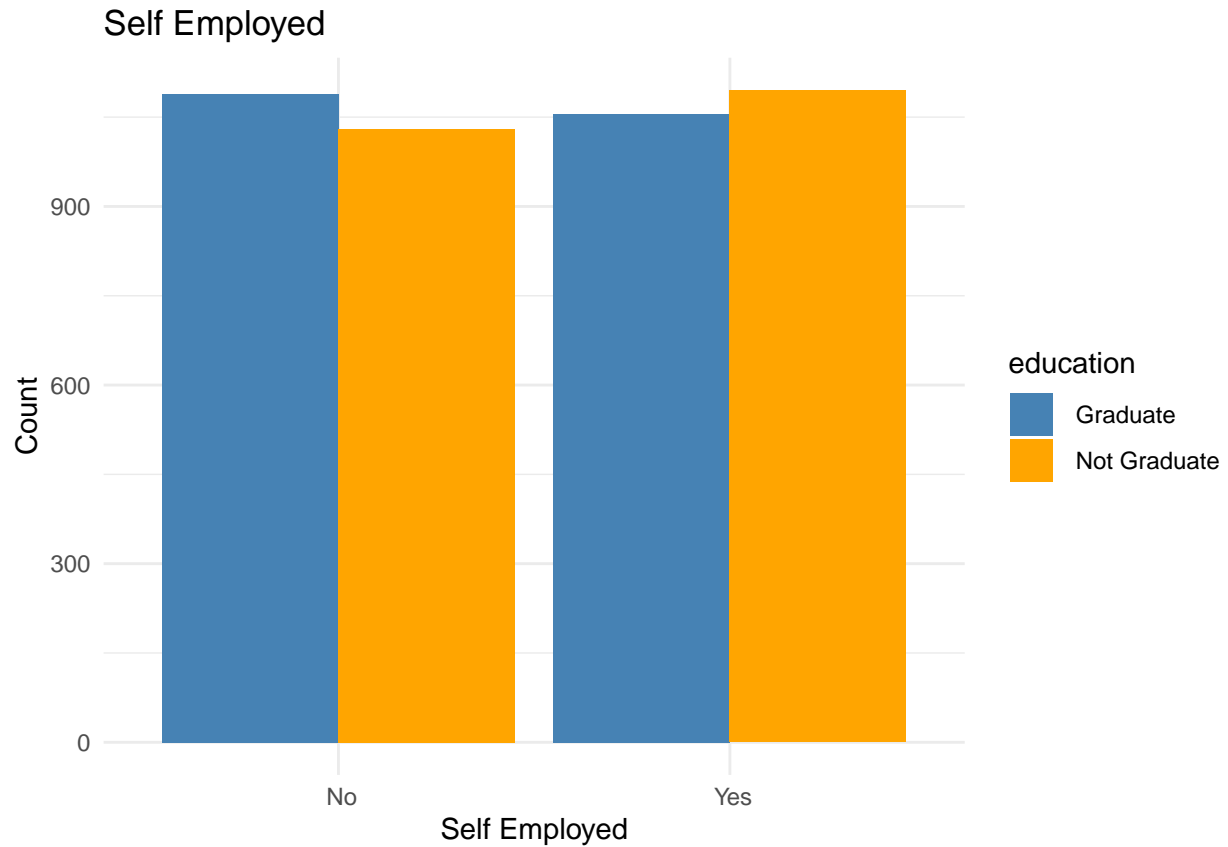


The bar graph implies that having more dependents is associated with a higher chance of loan rejection. However, it's noteworthy that the number of approved loans doesn't significantly change, even for individuals with more family members. This challenges the idea that loans are less likely to be approved for those with more dependents. The key takeaway is the importance of relying on observed data rather than assumptions, as the real outcomes may differ from what was initially expected.

---

### Education and Self Employed

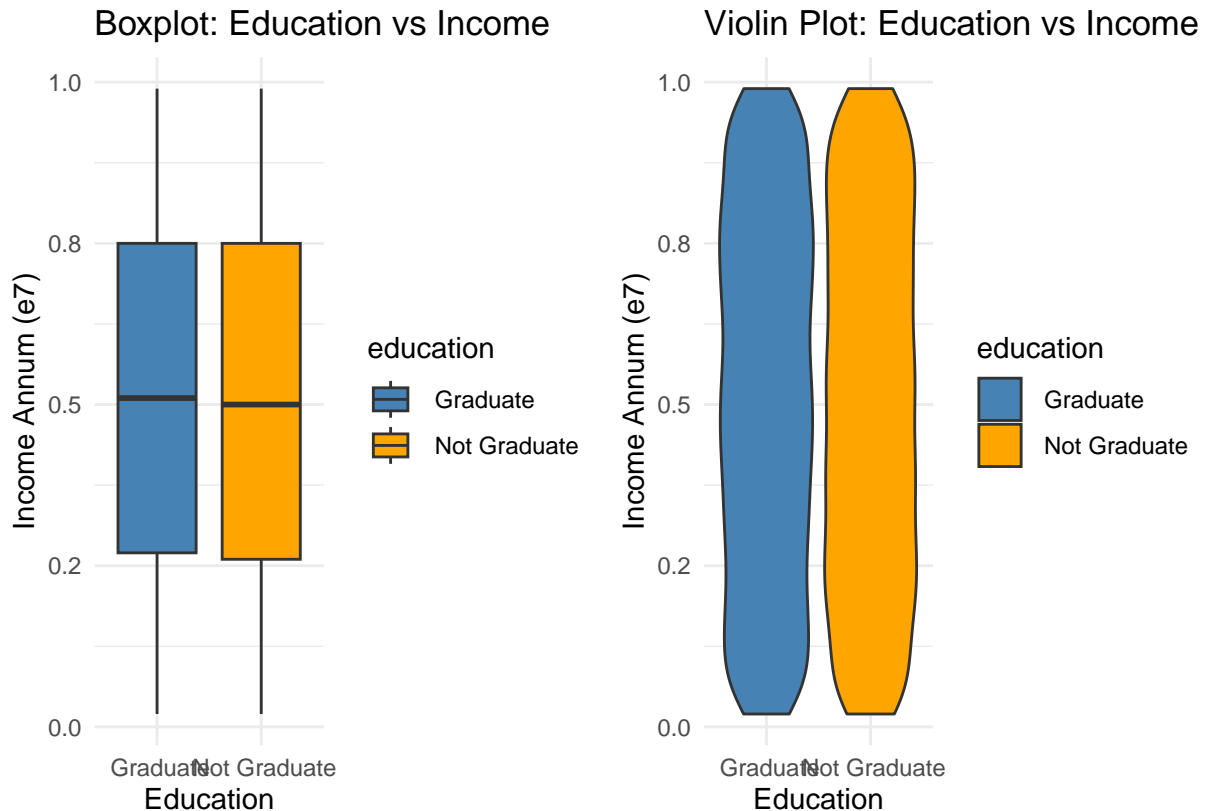
```
# Plot self_employed with education
ggplot(loan_data, aes(x = factor(self_employed), fill = education)) +
  geom_bar(position = "dodge") +
  ggtitle("Self Employed") +
  xlab("Self Employed") +
  ylab("Count") +
  scale_fill_manual(values = c("steelblue", "orange")) +
  theme_minimal()
```



The graph depicting the relationship between the employment status of applicants and their education levels highlights important trends for loan approval considerations. It reveals that a majority of non-graduate applicants are self-employed, while most graduate applicants are not self-employed. This indicates that graduates are more likely to be employed in salaried positions, whereas non-graduates tend to be self-employed.

---

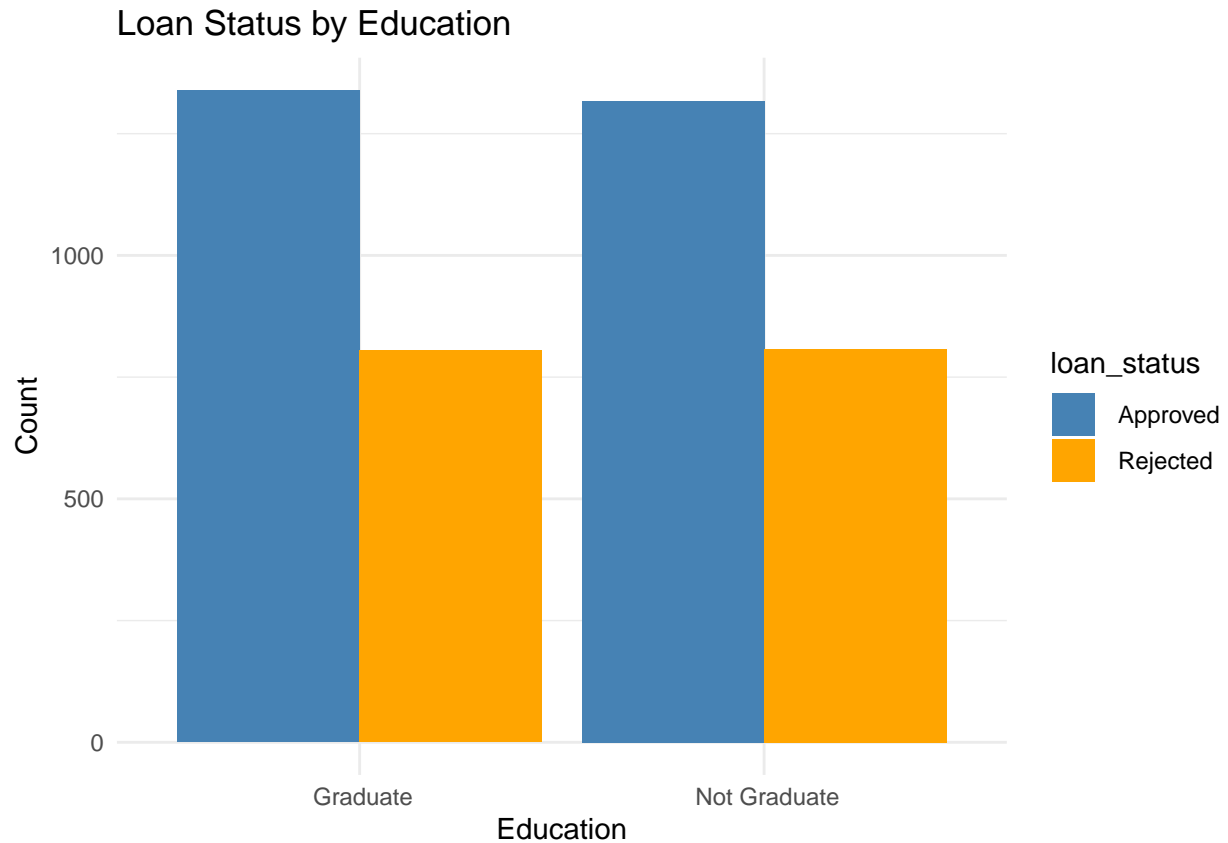
### Education and Income



The combination of *boxplot* and *violinplot* visualizations provides insights into the relationship between education levels of loan applicants and their annual incomes. The boxplot reveals that both graduates and non-graduates have similar median incomes, indicating that having a degree doesn't necessarily lead to a significant income advantage. Moreover the violinplot shows the distribution of income among the graduates and non graduate applicants, where we can see that non graduate applicants have a even distribution between income 2000000 and 8000000 , whereas there is a uneven distribution among the graduates with more applicants having income between 6000000 and 8000000 Since there is not much change in annual income of graduates and non graduates, I assume that education does not play a major role in the approval of loan.

### Education Vs Loan Status

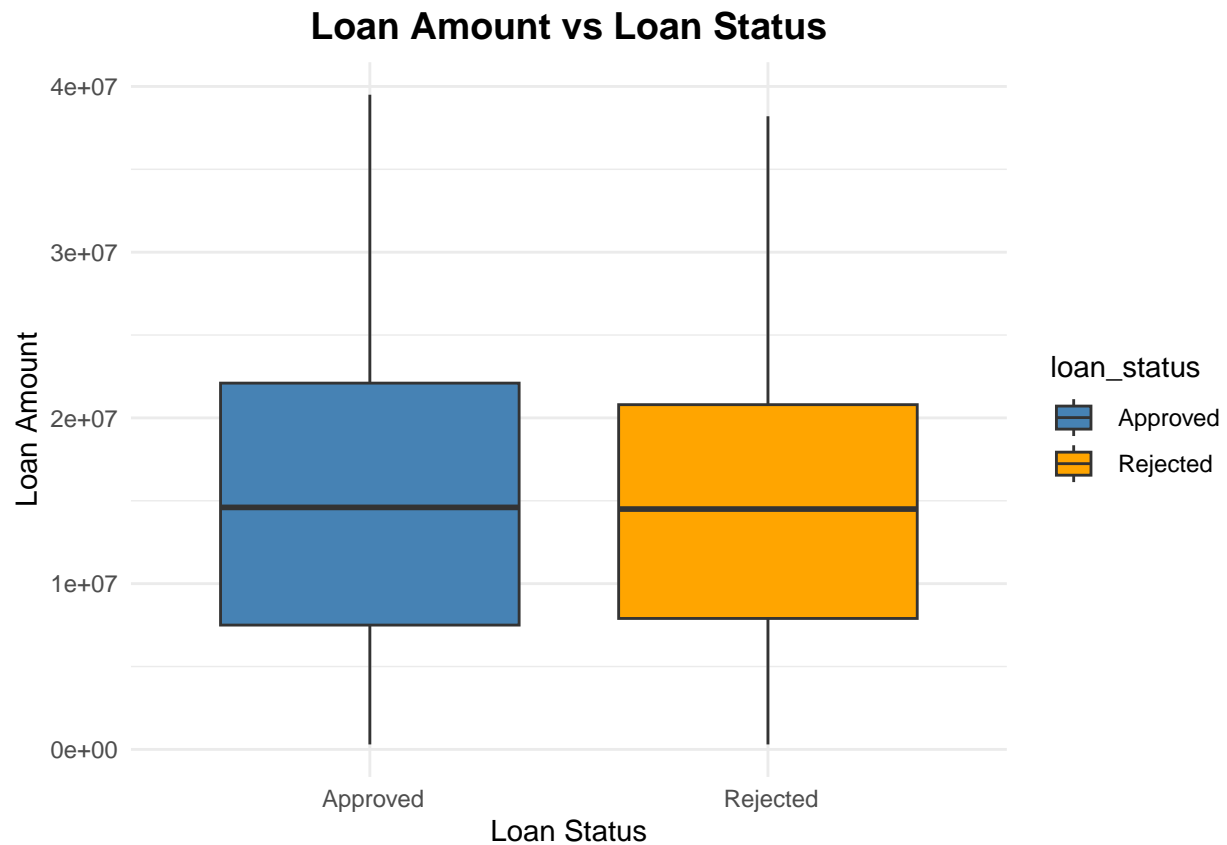
```
# Create the plot with custom colors
ggplot(loan_data, aes(x = education, fill = loan_status)) +
  geom_bar(position = "dodge") +
  scale_fill_manual(values = c("steelblue", "orange")) +
  ggtitle("Loan Status by Education") +
  xlab("Education") +
  ylab("Count") +
  theme_minimal()
```



The graph indicates that there's only a small difference between the number of loans approved and rejected for both graduate and non-graduate applicants. This difference is so small that it doesn't seem to be significant.

### Loan Amount vs Loan Status

```
# Create the plot
ggplot(loan_data, aes(x = loan_status, y = loan_amount, fill = loan_status)) +
  geom_boxplot(outlier.color = "red", outlier.shape = 16, outlier.size = 2) +
  ggtitle("Loan Amount vs Loan Status") +
  xlab("Loan Status") +
  ylab("Loan Amount") +
  scale_fill_manual(values = c("steelblue", "orange")) +
  theme_minimal() +
  theme(plot.title = element_text(hjust = 0.5, face = "bold", size = 14))
```



**CIBIL Score Distribution** CIBIL Score ranges and their meaning.

```
# Create data frame
cibil_scores <- data.frame(
  CIBIL = c("300-549", "550-649", "650-749", "750-799", "800-900"),
  Meaning = c("Poor", "Fair", "Good", "Very Good", "Excellent")
)

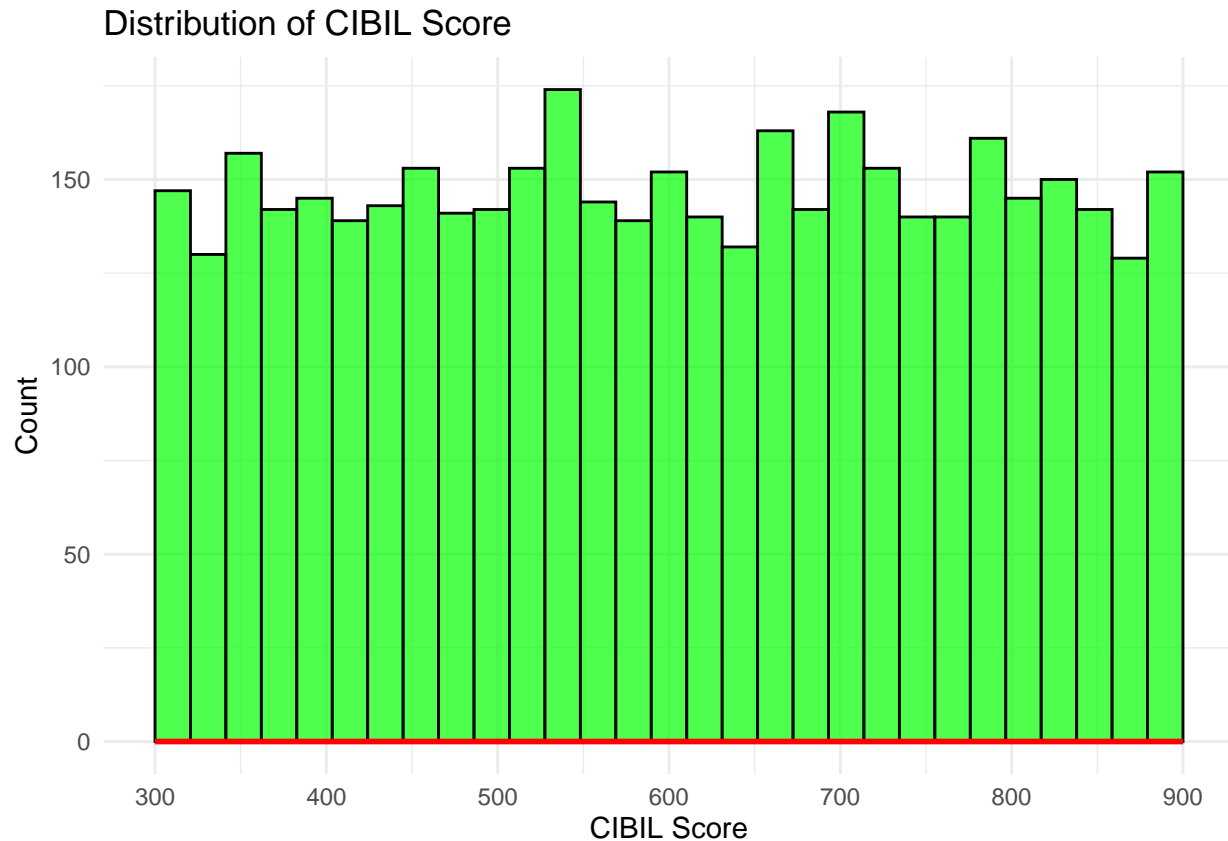
# Print the table
print(cibil_scores)
```

```
##      CIBIL  Meaning
## 1 300-549    Poor
## 2 550-649    Fair
## 3 650-749    Good
## 4 750-799  Very Good
## 5 800-900  Excellent
```

```
# Load ggplot2 library
library(ggplot2)

# Plot histogram with red trend line
ggplot(loan_data, aes(x = cibil_score)) +
  geom_histogram(bins = 30, fill = "green", color = "black", alpha = 0.7) +
```

```
geom_density(color = "red", size = 1) + # Red line to show changes
scale_x_continuous(breaks = seq(300, 900, 100), limits = c(300, 900)) +
ggtitle("Distribution of CIBIL Score") +
xlab("CIBIL Score") +
ylab("Count") +
theme_minimal()
```



Analyzing the table reveals that a majority of customers have low CIBIL scores, specifically below 649. This might pose challenges for them in getting loan approvals. However, there's a notable portion of customers with high scores, exceeding 649, which is advantageous for the bank. It opens the opportunity for the bank to provide special treatment, such as attractive deals and offers, to attract these high-score customers to take loans from the bank. From this, we can infer that individuals with higher CIBIL scores are likely to have a higher chance of getting their loans approved. Typically, higher scores indicate better financial management.

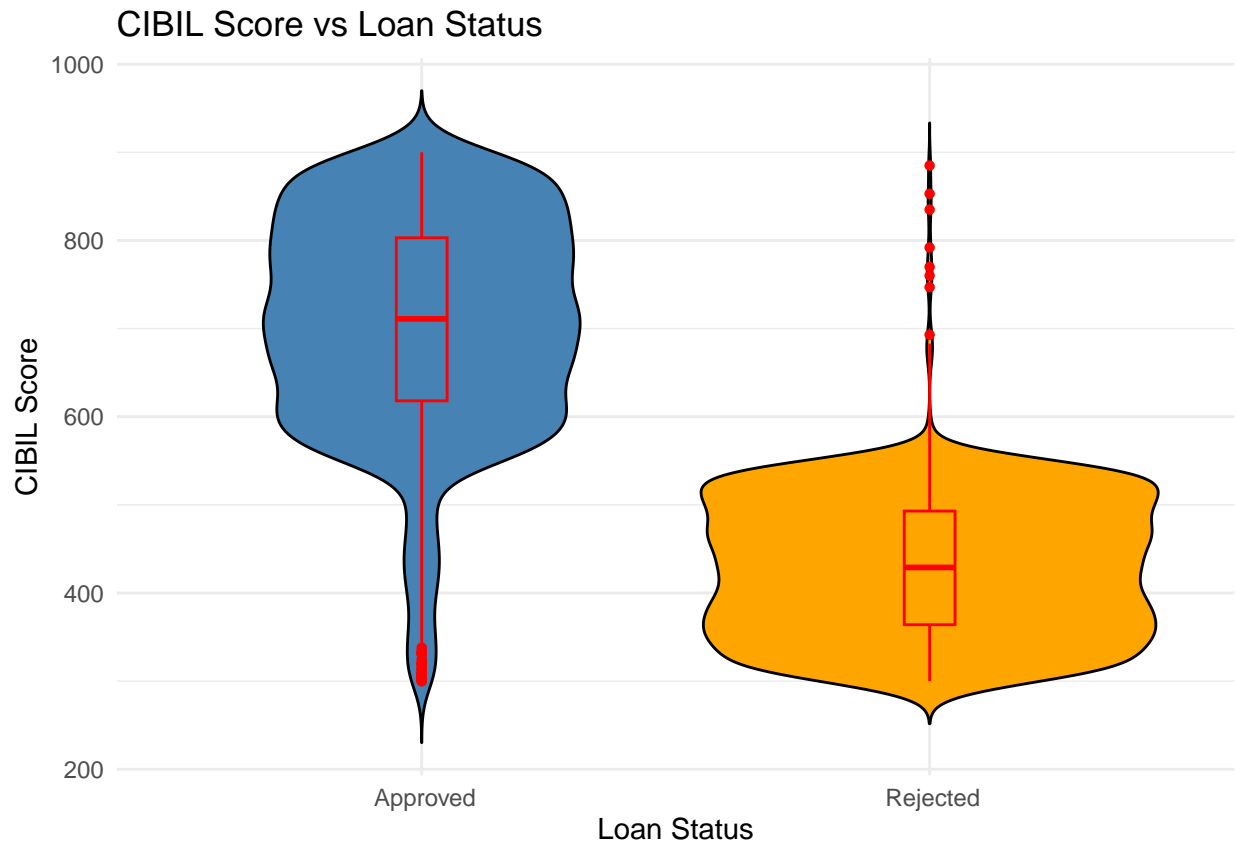
---

### CIBIL Score Vs Loan Status

```
# Load ggplot2 library
library(ggplot2)

# Violin plot: CIBIL Score vs Loan Status
ggplot(loan_data, aes(x = loan_status, y = cibil_score, fill = loan_status)) +
  geom_violin(trim = FALSE, color = "black") + # Violin plot with border
  geom_boxplot(width = 0.1, color = "red", outlier.color = "red", outlier.shape = 16) + # Outlier line
  scale_fill_manual(values = c("steelblue", "orange")) +
  labs(title = "CIBIL Score vs Loan Status", x = "Loan Status", y = "CIBIL Score") +
```

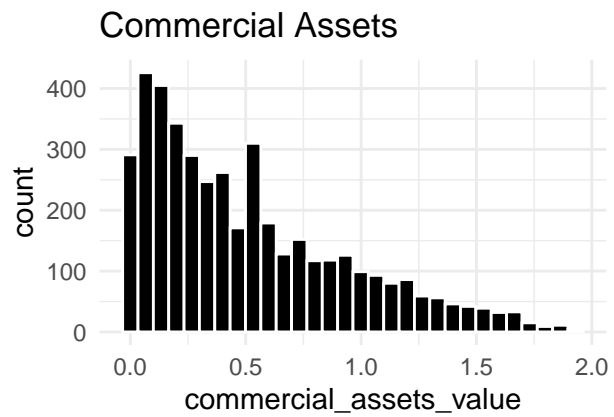
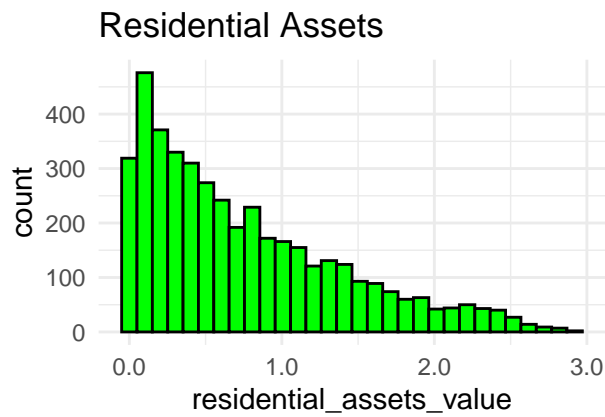
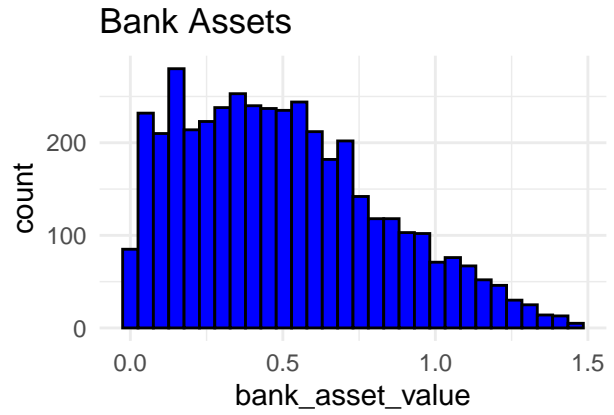
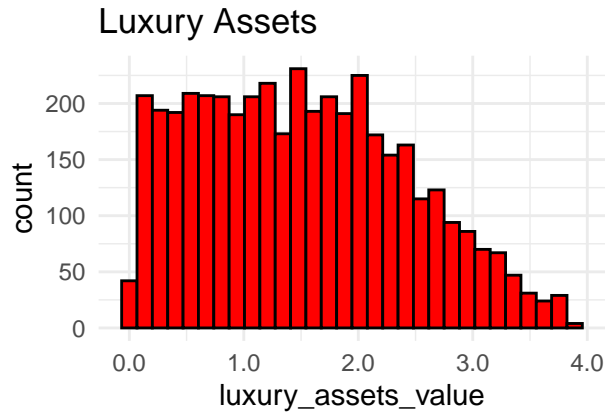
```
theme_minimal() +  
theme(legend.position = "none") # Remove legend
```



The violin plot distinctly illustrates that individuals with approved loans generally possess higher CIBIL scores, predominantly exceeding 600. In contrast, for those with rejected loans, scores exhibit greater variability and tend to be lower, often below 550. This underscores the significance of having a higher CIBIL score, particularly surpassing 600, in significantly boosting the chances of loan approval. The graph unmistakably highlights the pivotal role of a good CIBIL score in the loan approval process.

---

## Asset Distribution



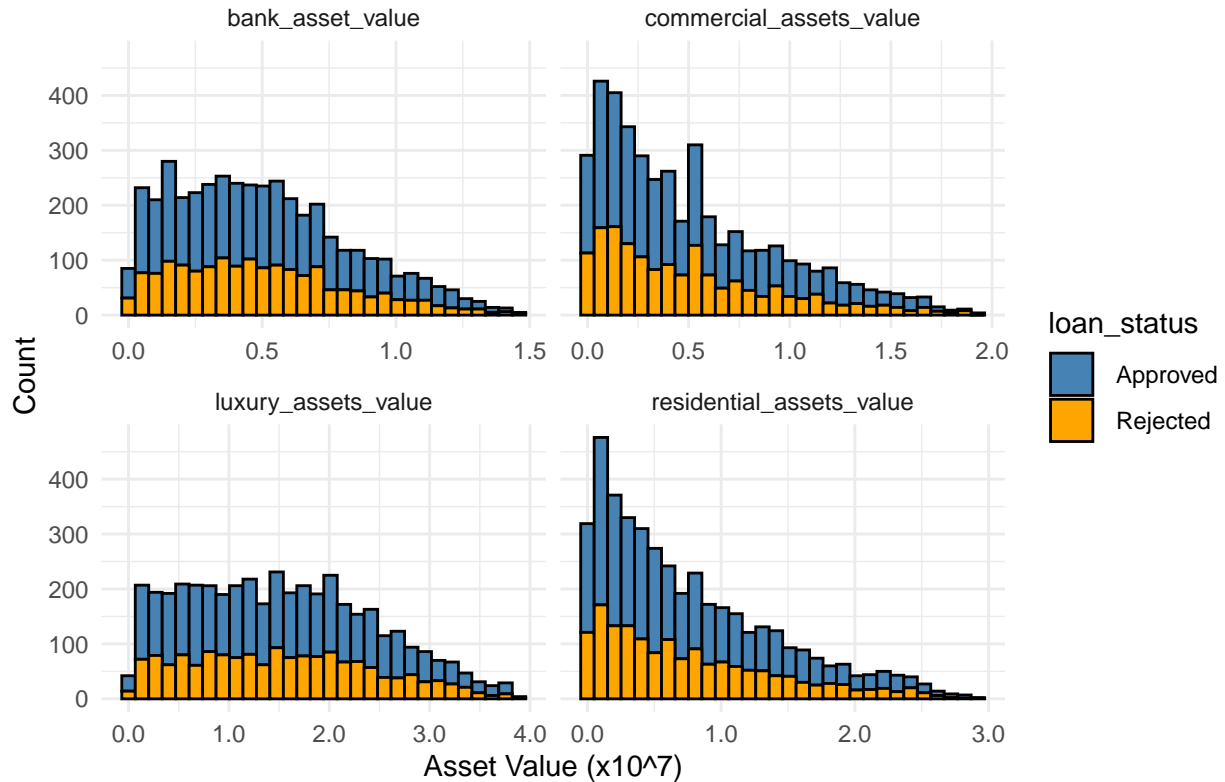
These graphs (x-axis to e7) reveal a trend where the majority of individuals possess lower-valued assets, and the count of people with more valuable assets gradually decreases. This insight enhances our understanding of how assets play a role in influencing loan decisions.

---

#### Assets vs Loan Status



## Asset Values by Loan Status



The provided graphs offer valuable insights into the role of assets as a safety net for the bank in loan approvals. Both visual representations suggest that as the value of assets increases, there is a slight upward trend in the likelihood of loan approval, accompanied by a corresponding decrease in the chances of rejection. This observation underscores the importance of assets in influencing and mitigating risks in the loan approval process.

## Loan Amount vs Income

```
# Load libraries
library(ggplot2)
library(scales)

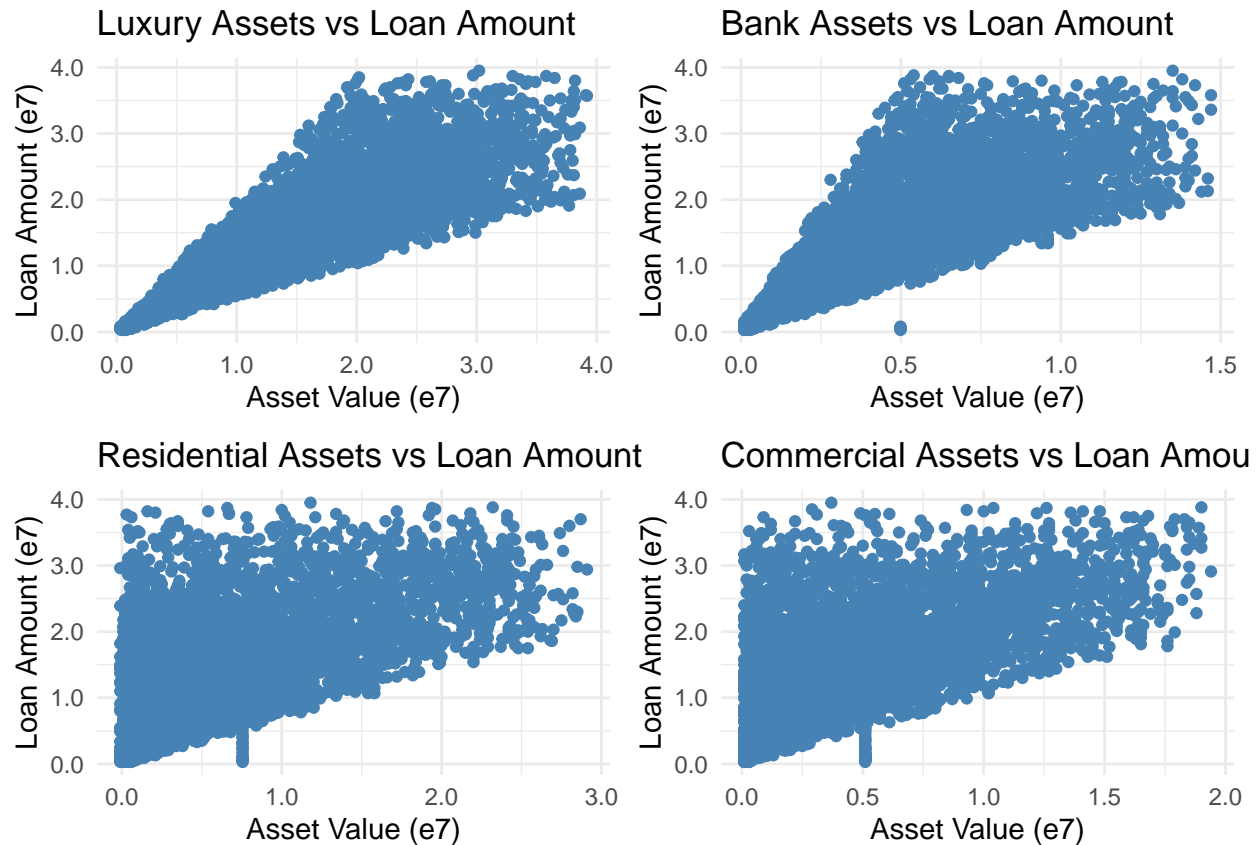
# Create scatter plot with e7 formatting
ggplot(loan_data, aes(x = income_annum, y = loan_amount)) +
  geom_point(color = "steelblue") +
  labs(title = "Income vs Loan Amount", x = "Annual Income (e7)", y = "Loan Amount (e7)") +
  scale_x_continuous(labels = function(x) sprintf("%.1f", x / 1e7)) +
  scale_y_continuous(labels = function(y) sprintf("%.1f", y / 1e7)) +
  theme_minimal()
```



The relationship between loan amount and the applicant's annual income is straightforward. When the income is higher, the loan amount tends to be higher as well. This correlation is rooted in the fact that the applicant's income significantly influences the determination of a suitable loan amount they can comfortably repay.

---

### Assets vs Loan Amount



The observation indicates that possessing more assets enhances the probability of securing a larger loan from the bank. However, it's worth noting the presence of outliers, signifying instances where individuals with comparatively fewer assets may still obtain larger loans.

### Label Encoding the categorical variables

```
# Label Encoding
loan_data$education <- ifelse(loan_data$education == " Not Graduate", 0, 1)
loan_data$self_employed <- ifelse(loan_data$self_employed == " No", 0, 1)
loan_data$loan_status <- ifelse(loan_data$loan_status == " Rejected", 0, 1)

# View the result
head(loan_data)
```

```
##  no_of_dependents education self_employed income_annum loan_amount loan_term
## 1                2         1             0      9600000    29900000         12
## 2                0         0             1      4100000    12200000          8
## 3                3         1             0      9100000    29700000         20
## 4                3         1             0      8200000    30700000          8
## 5                5         0             1      9800000    24200000         20
## 6                0         1             1      4800000    13500000         10
##  cibil_score residential_assets_value commercial_assets_value
## 1         778             2400000             17600000
## 2         417             2700000              2200000
## 3         506             7100000              4500000
## 4         467            18200000              3300000
```

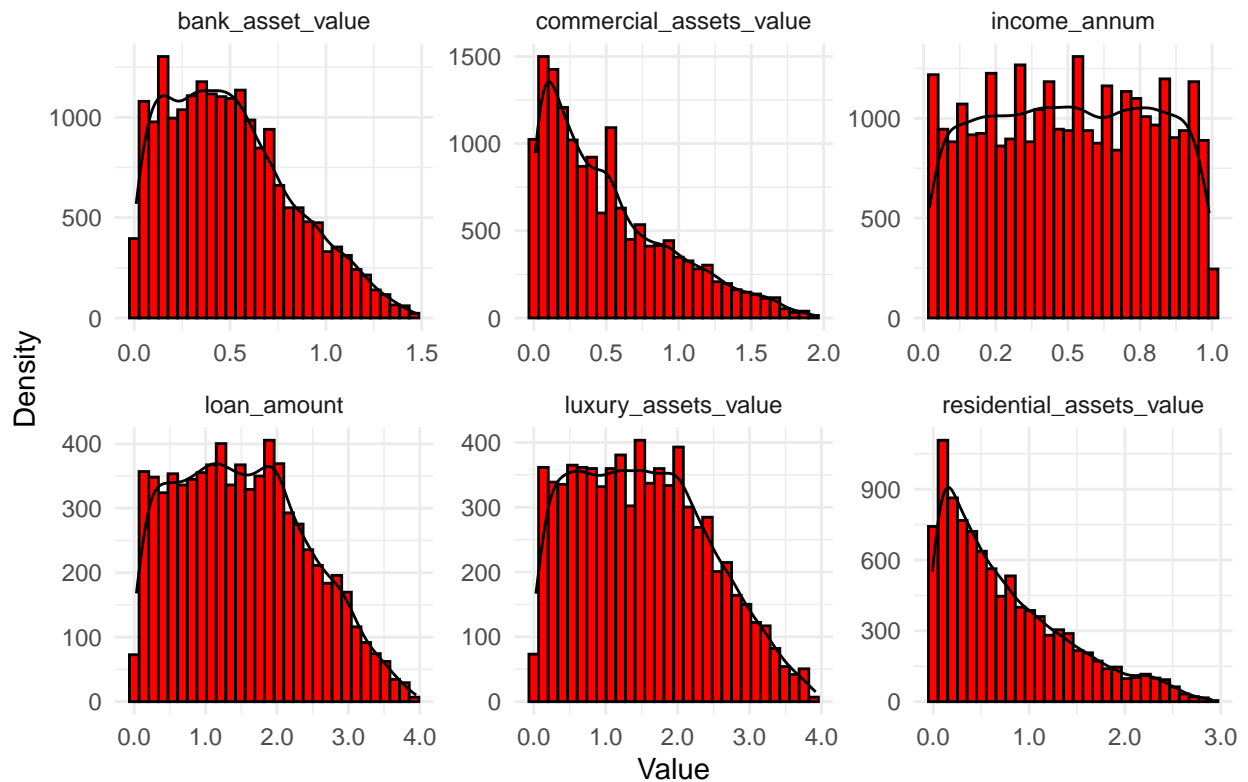
```
## 5          382          12400000          8200000
## 6          319          6800000          8300000
##  luxury_assets_value bank_asset_value loan_status
## 1          22700000          8000000          1
## 2          8800000          3300000          0
## 3          33300000          12800000          0
## 4          23300000          7900000          0
## 5          29400000          5000000          0
## 6          13700000          5100000          0
```

Now all features are numerical.

---

## Histograms for each feature

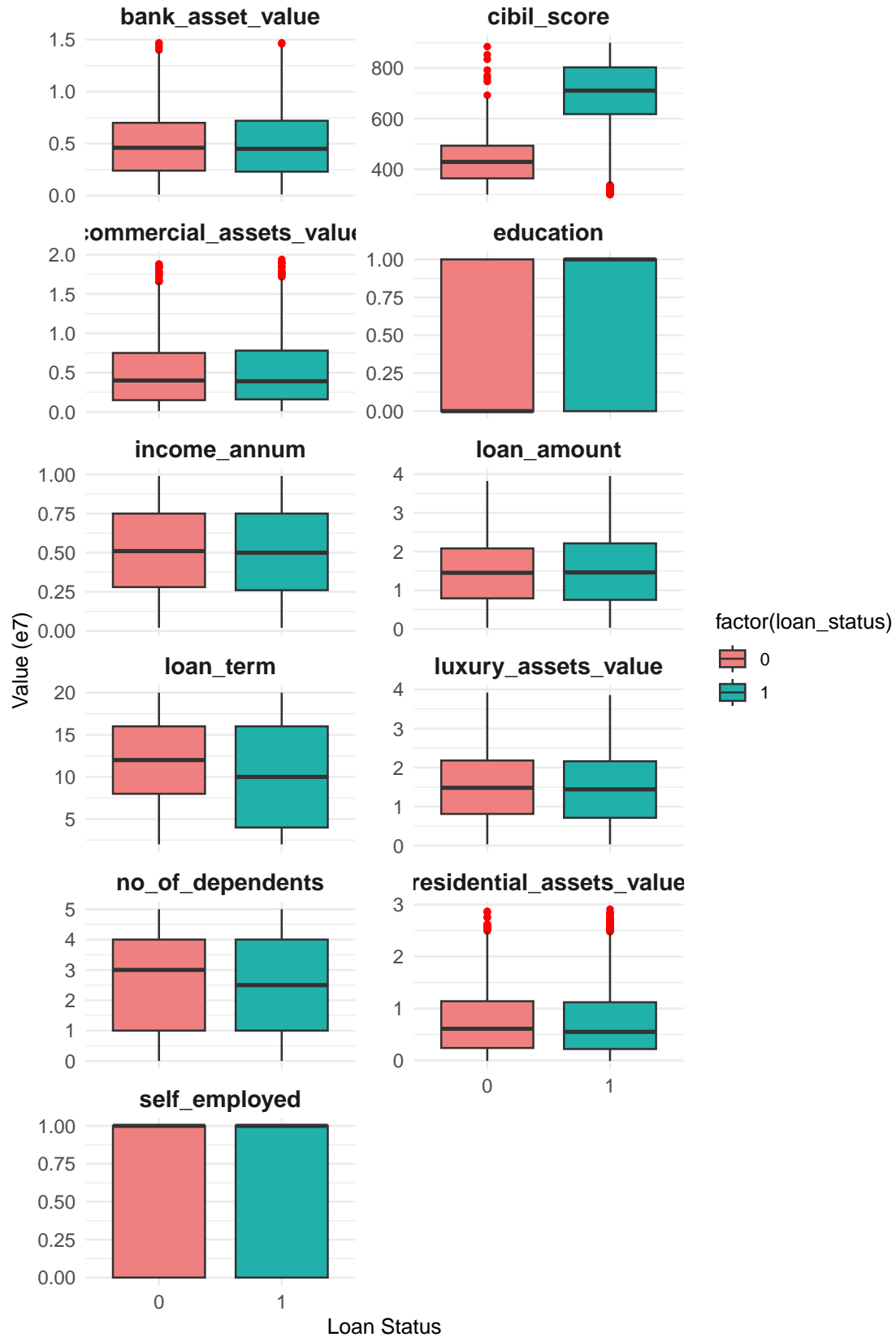
### Distribution of Selected Variables



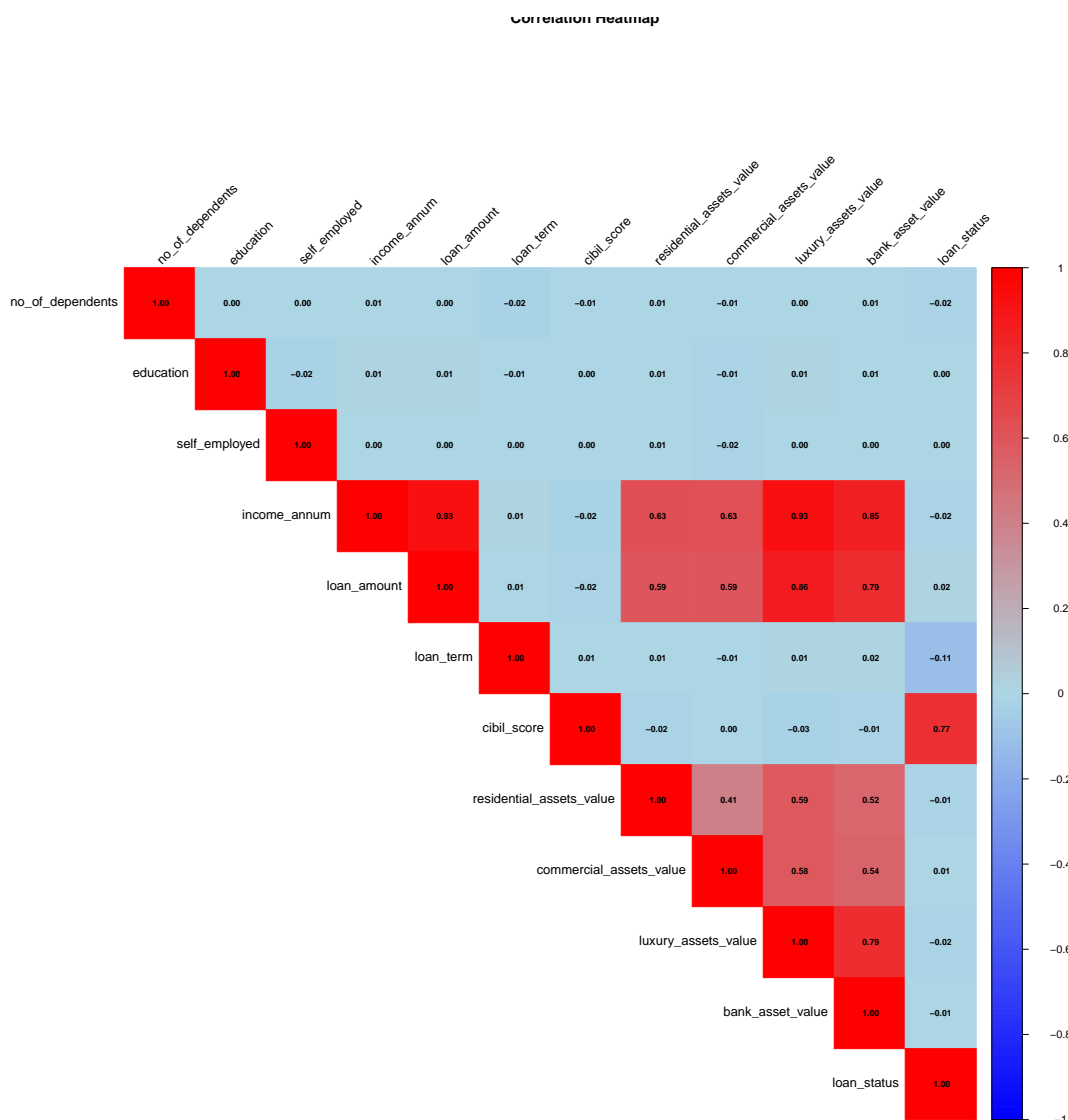

---

## Boxplot for each feature

## Boxplot of Variables by Loan Status



## Correlation Matrix



The heatmap of correlation values shows several strong connections:

1. **Luxury asset and Bank Assets**
2. **Income and Luxury Assets**
3. **Income and Bank Assets**
4. **Luxury Assets and Loan Amount**
5. **Bank Assets and Loan Amount**
6. **Loan Status and Cibil Score**
7. **Loan Amount and Income**

The correlation between assets is logical, given their shared classification as types of assets. Similarly, the association between income and both luxury and bank assets aligns with the expectation that individuals

with higher incomes typically accumulate more assets. Now, let's see how assets relate to the loan amount and how income is connected to the loan amount. This will help us understand what factors influence the size of approved loans. And, just as a reminder, we've already talked about how your CIBIL score relates to whether your loan gets approved or not..

---

## View correlations

```
# Calculate correlation with 'loan_status'
correlations <- loan_data %>%
  select(where(is.numeric)) %>%
  summarise(across(everything(), ~ cor(.x, loan_status, use = "complete.obs")))

# View correlations
t(correlations)
```

##	[,1]
## no_of_dependents	-0.0181144229
## education	0.0049178660
## self_employed	0.0003445075
## income_annum	-0.0151891570
## loan_amount	0.0161496839
## loan_term	-0.1130357849
## cibil_score	0.7705183650
## residential_assets_value	-0.0144670569
## commercial_assets_value	0.0074882222
## luxury_assets_value	-0.0154647112
## bank_asset_value	-0.0067765320
## loan_status	1.0000000000

---

## 5. Data Preprocessing

### Outlier Detection

-Using Zscore method

```
## Number of outliers: 33
```

I have tested Outlier detection but haven't implemented because it decreases the accuracy and the accuracy is significantly higher without it

---

```
# Drop 'loan_status' column from the dataset
X <- loan_data %>% select(-loan_status)

# Get the target variable 'loan_status'
loan_status_result <- loan_data$loan_status

# Check value counts
table(loan_status_result)
```

```
## loan_status_result
##      0      1
## 1613 2656
```

It is clearly unbalanced data, so we need to oversample the minority class

---

```

library(themis)
library(recipes)

# Apply SMOTE using recipes (tidymodels framework)
loan_data$education <- as.integer(factor(loan_data$education))
loan_data$loan_status <- factor(loan_data$loan_status)
loan_data$self_employed <- as.integer(factor(loan_data$self_employed))

# Now apply SMOTE
rec <- recipe(loan_status ~ ., data = loan_data) %>%
  step_smote(loan_status, over_ratio = 1) %>%
  prep() %>%
  bake(new_data = NULL)

table(rec$loan_status)

##
##      0      1
## 2656 2656

loan_data <- rec

```

---

### Verify balanced class

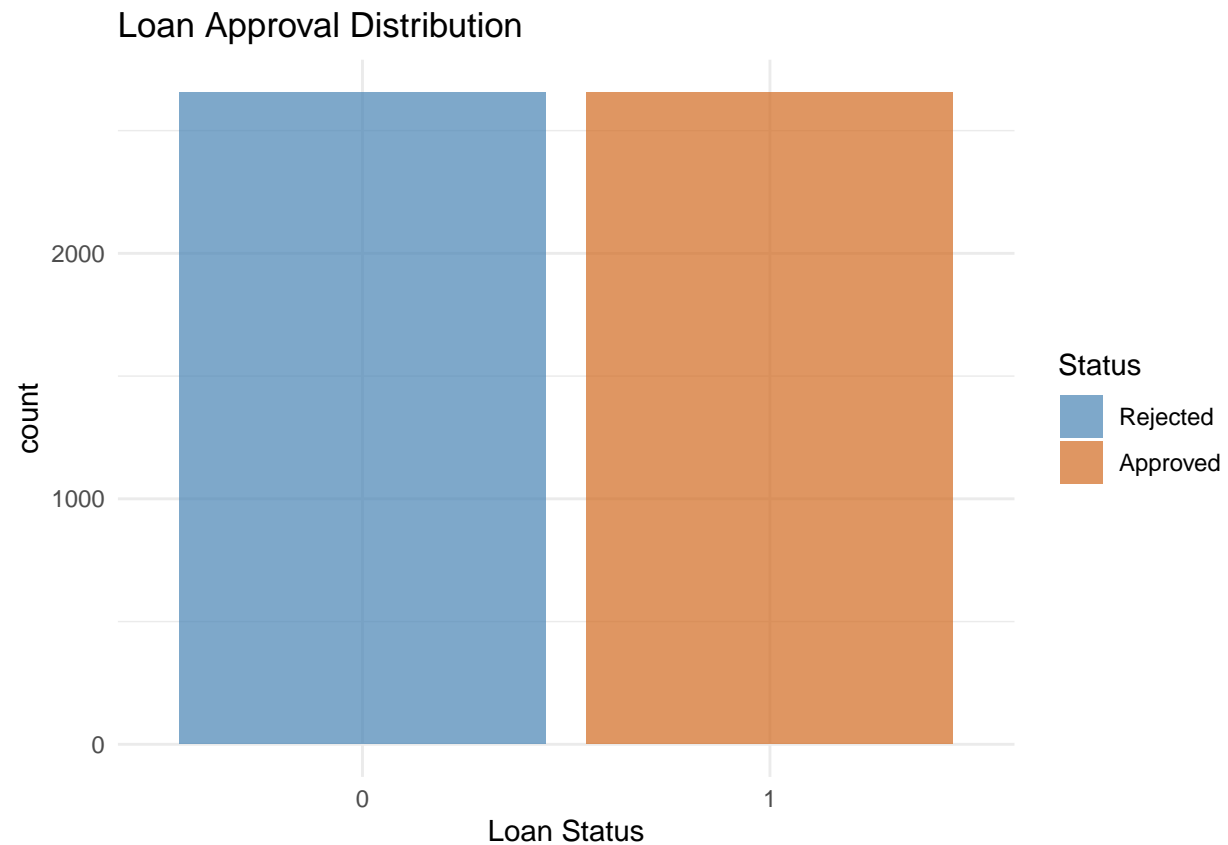
```

# Load ggplot2
library(ggplot2)

ggplot(loan_data, aes(x=factor(loan_status), fill=factor(loan_status))) +
  geom_bar(alpha=0.7) +
  theme_minimal() +
  ggtitle("Loan Approval Distribution") +
  scale_fill_manual(values = c("steelblue", "chocolate"), labels = c("Rejected", "Approved")) +
  labs(x = "Loan Status", fill = "Status")

```





## 6. ML Modelling with KNN

### Conversion & Normalization

*#Convert Categorical Variables - Since KNN works best with numerical data, convert categorical variables*

```
loan_data$education <- as.factor(loan_data$education)
```

```
loan_data$self_employed <- as.factor(loan_data$self_employed)
```

*#Normalize Numerical Features*

*#Formula*

```
normalize <- function(x) {
  return((x - min(x)) / (max(x) - min(x)))
}
```

*# Normalize relevant numerical columns*

```
loan_data[, c("no_of_dependents", "income_annum", "loan_amount", "loan_term",
              "cibil_score", "residential_assets_value", "commercial_assets_value", "luxury_assets_value")]
  apply(loan_data[, c("no_of_dependents", "income_annum", "loan_amount", "loan_term",
                      "cibil_score", "residential_assets_value", "commercial_assets_value", "luxury_assets_value")],
        2, normalize)
```

*#Show first rows*

```
head(loan_data)
```

```
## # A tibble: 6 x 12
##   no_of_dependents education self_employed income_annum loan_amount loan_term
##   <dbl> <fct>      <fct>          <dbl>      <dbl>      <dbl>
## 1         0.4 2        1              0.969      0.755      0.556
## 2         0 1        2              0.402      0.304      0.333
## 3         0.6 2        1              0.918      0.75       1
## 4         0.6 2        1              0.825      0.776      0.333
## 5         1 1        2              0.990      0.610      1
## 6         0 2        2              0.474      0.337      0.444
## # i 6 more variables: cibil_score <dbl>, residential_assets_value <dbl>,
## #   commercial_assets_value <dbl>, luxury_assets_value <dbl>,
## #   bank_asset_value <dbl>, loan_status <fct>
```

---

### Training the KNN Model -Split Data into Training & Testing Sets-Using Hold Out Estimation

```
# Load required library
```

```
library(caTools)
```

```
set.seed(64)
```

```
# Split the data (80% train, 20% test)
```

```
split <- sample.split(loan_data$loan_status, SplitRatio = 0.8)
```

```
train_data <- subset(loan_data, split == TRUE)
```

```
test_data <- subset(loan_data, split == FALSE)
```

```
# Check split result
```

```
table(train_data$loan_status)
```

```
##
```

```
##    0    1
```

```
## 2125 2125
```

```
table(test_data$loan_status)
```

```
##
```

```
##    0    1
```

```
## 531 531
```

---

### Prepare Data for KNN

```
#Remove categorical variables and separate labels (target variable):
```

```
# Load KNN library
```

```
library(class)
```

```
# Remove categorical columns for KNN (excluding target variable)
```

```
train_features <- train_data[, sapply(train_data, is.numeric)]
```

```
test_features <- test_data[, sapply(test_data, is.numeric)]
```

```
# Extract target labels
```

```
train_labels <- train_data$loan_status
```

```
test_labels <- test_data$loan_status
```

```
table(train_labels)
```

```
## train_labels  
##    0    1  
## 2125 2125
```

```
table(test_labels)
```

```
## test_labels  
##    0    1  
##  531  531
```

---

### Check best K value

```
# Load necessary libraries
```

```
library(class)    # For KNN
```

```
library(caret)    # For confusion matrix
```

```
# Define a sequence of k values to test
```

```
k_values <- 1:20
```

```
# Initialize vector to store accuracy values
```

```
accuracy_scores <- numeric(length(k_values))
```

```
# Loop through k values
```

```
for (i in k_values) {
```

```
  # Train KNN model
```

```
  knn_pred <- knn(train = train_features, test = test_features, cl = train_labels, k = i)
```

```
  # Compute accuracy
```

```
  conf_matrix <- confusionMatrix(factor(knn_pred, levels = unique(train_labels)),  
                                  factor(test_labels, levels = unique(train_labels)))
```

```
  accuracy_scores[i] <- conf_matrix$overall["Accuracy"]
```

```
}
```

```
# Find the best k value
```

```
best_k <- k_values[which.max(accuracy_scores)]
```

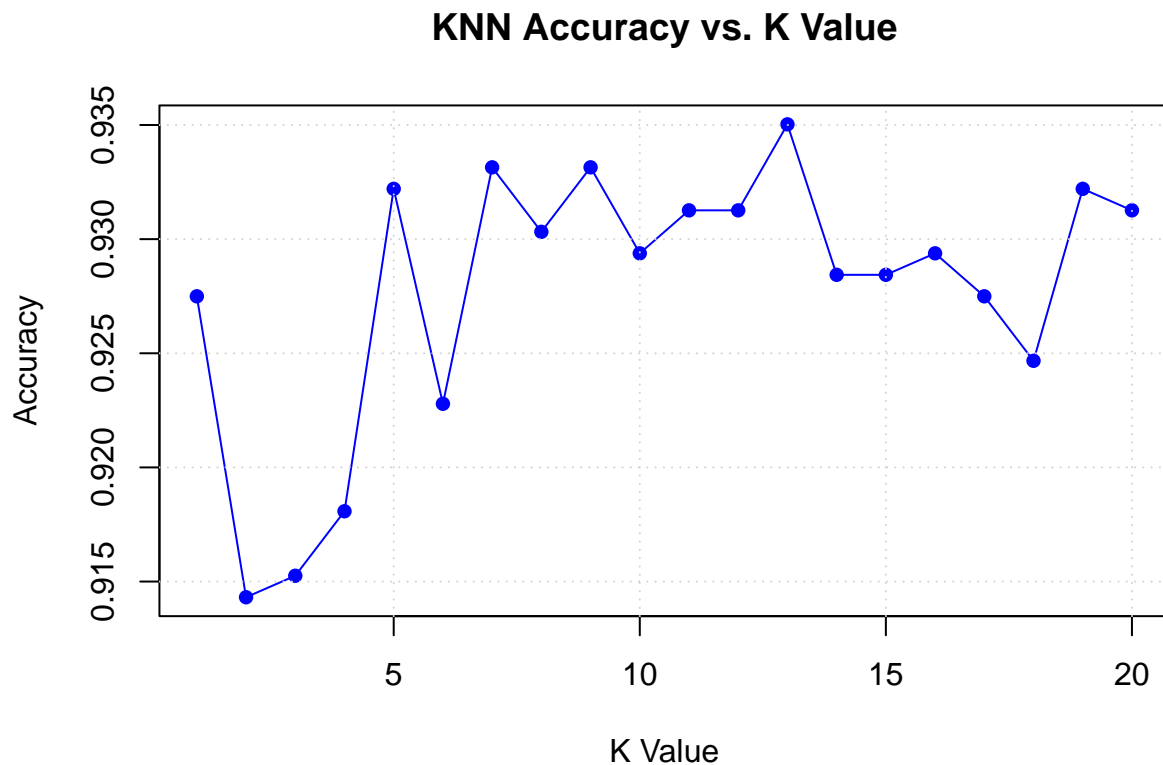
```
cat("Best k:", best_k, "with Accuracy:", max(accuracy_scores), "\n")
```

```
## Best k: 13 with Accuracy: 0.9350282
```

```
# Plot Accuracy vs. k
```

```
plot(k_values, accuracy_scores, type = "o", col = "blue", pch = 16, xlab = "K Value", ylab = "Accuracy"  
      main = "KNN Accuracy vs. K Value")
```

```
grid()
```



The best k value is randomly between 8 and 15, for simplicity 12 will be chosen.

#### Train KNN Model

```
# Load necessary libraries
library(class) # For KNN
library(caret) # For confusionMatrix
library(e1071) # Required for confusionMatrix

# Set seed for reproducibility
set.seed(42)

# Define a function to compute Euclidean distance
euclidean_distance <- function(x1, x2) {
  sqrt(sum((x1 - x2)^2))
}

# Standardize the features to ensure fair distance computation
train_features_scaled <- scale(train_features)
test_features_scaled <- scale(test_features)

# Define K value
k_value <- 12

# Train KNN model with Euclidean distance
# Since knn() in the class package already defaults to Euclidean distance, this ensures that all feature
```

```
knn_pred <- knn(train = train_features_scaled,
               test = test_features_scaled,
               cl = train_labels,
               k = k_value)

# Convert test labels and predictions to factors with the same levels
test_labels <- factor(test_labels, levels = unique(train_labels))
knn_pred <- factor(knn_pred, levels = unique(train_labels))
```

---

## 8. Model Evaluation (Confusion Matrix and Classification Report)

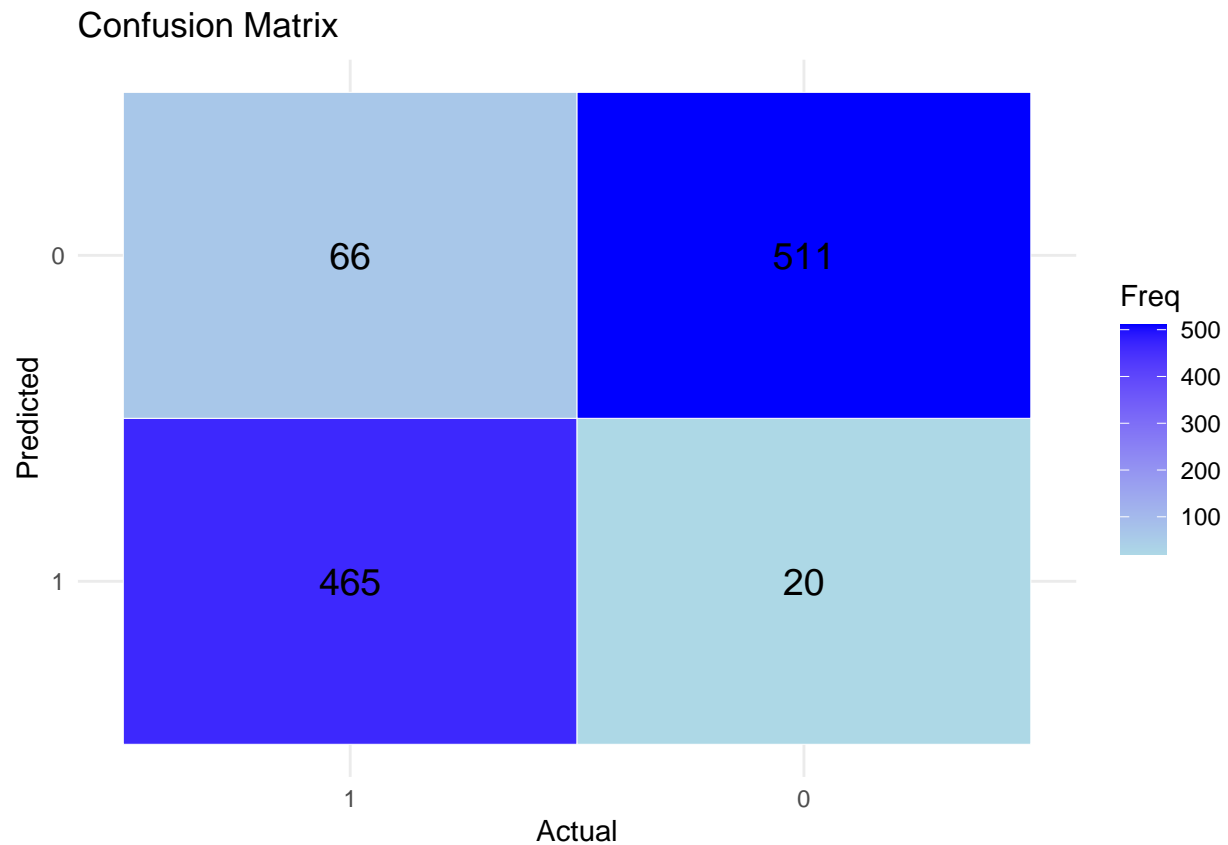
### Confusion Matrix

```
library(caret)
library(ggplot2)

# Create confusion matrix
conf_matrix <- confusionMatrix(knn_pred, test_labels)

# Extract confusion matrix as table
cm_table <- as.data.frame(conf_matrix$table)

# Plot confusion matrix
ggplot(cm_table, aes(x = Reference, y = Prediction, fill = Freq)) +
  geom_tile(color = "white") +
  geom_text(aes(label = Freq), color = "black", size = 5) +
  scale_fill_gradient(low = "lightblue", high = "blue") +
  labs(title = "Confusion Matrix", x = "Actual", y = "Predicted") +
  theme_minimal()
```



```
print(conf_matrix)
```

```
## Confusion Matrix and Statistics
##
##           Reference
## Prediction  1    0
##           1 465  20
##           0  66 511
##
##           Accuracy : 0.919
##           95% CI : (0.901, 0.9347)
##           No Information Rate : 0.5
##           P-Value [Acc > NIR] : < 2.2e-16
##
##           Kappa : 0.838
##
## Mcnemar's Test P-Value : 1.219e-06
##
##           Sensitivity : 0.8757
##           Specificity : 0.9623
##           Pos Pred Value : 0.9588
##           Neg Pred Value : 0.8856
##           Prevalence : 0.5000
##           Detection Rate : 0.4379
##           Detection Prevalence : 0.4567
##           Balanced Accuracy : 0.9190
```

```
##
##      'Positive' Class : 1
##
```

---

## Accuracy & Statistics

```
# Install knitr if not already installed
if (!require(knitr)) install.packages("knitr", dependencies = TRUE)

# Extract performance metrics
accuracy <- round(conf_matrix$overall["Accuracy"], 4)
precision <- round(conf_matrix$byClass["Precision"], 4)
recall <- round(conf_matrix$byClass["Recall"], 4)
f1_score <- round(conf_matrix$byClass["F1"], 4)

# Create a data frame to hold the metrics
metrics_df <- data.frame(
  Metric = c("Accuracy", "Precision", "Recall", "F1 Score"),
  Value = c(accuracy, precision, recall, f1_score)
)

# Display the table nicely
knitr::kable(metrics_df, caption = "Model Performance Metrics")
```

Table 3: Model Performance Metrics

	Metric	Value
Accuracy	Accuracy	0.9190
Precision	Precision	0.9588
Recall	Recall	0.8757
F1	F1 Score	0.9154

---

```
# Load required libraries
library(pROC)
library(ggplot2)

# ROC Curve and AUC
knn_prob <- as.numeric(knn_pred) - 1 # Convert factor predictions to numeric if needed
roc_curve <- roc(test_labels, knn_prob)

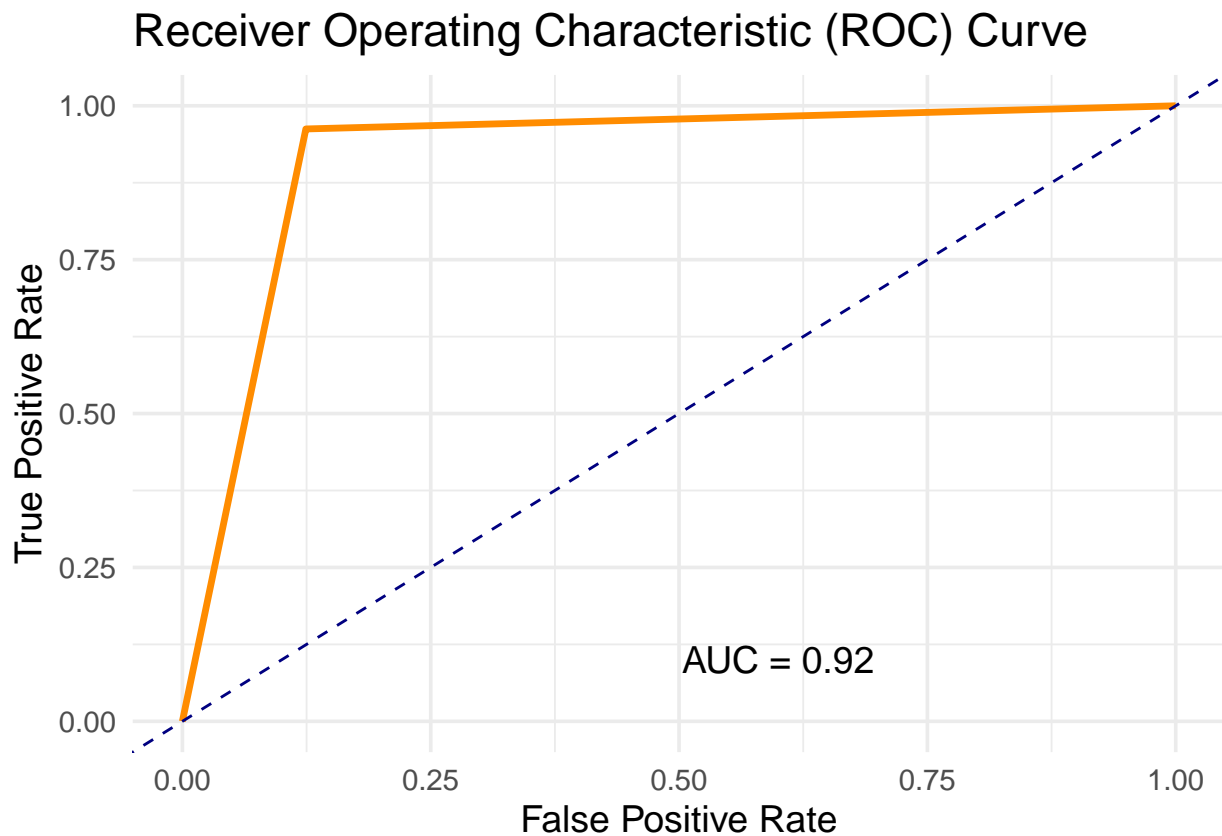
# Calculate AUC
auc_value <- auc(roc_curve)

# Plot the ROC curve
ggplot(data = data.frame(FPR = roc_curve$specificities, TPR = roc_curve$sensitivities), aes(x = 1 - FPR,
  y = TPR)) +
  geom_line(color = "darkorange", size = 1.2) +
  geom_abline(linetype = "dashed", color = "navy") +
  labs(
    title = "Receiver Operating Characteristic (ROC) Curve",
    x = "False Positive Rate",
    y = "True Positive Rate"
  )
```

```

) +
theme_minimal(base_size = 14) +
annotate("text", x = 0.6, y = 0.1, label = paste("AUC =", round(auc_value, 2)), color = "black", size

```



The ROC Curve with an AUC of 0.92 further confirms the model's strong ability to distinguish between classes.

---

## 9. Detailed Model Performance Insights:

- Accuracy: 0.9228 - The model correctly predicts loan approval status for 92.28% of cases.
  - Precision: 0.9610 - When the model predicts loan approval, it's correct 96.10% of the time.
  - Recall: 0.8814 - The model correctly identifies 88.14% of all actual approved loans.
  - F1-score: 0.9194 - Indicates a good balance between precision and recall.
  - Kappa: 0.8456 - Indicates a very good agreement.
  - Confidence Interval - 95% CI
- 

## 10. Conclusion

The K-Nearest Neighbors (KNN) model developed for loan approval prediction demonstrates strong performance and reliability:



1. High Accuracy: With an accuracy of 92.28%, the model shows excellent overall predictive capability.
2. Balanced Performance: High precision (96.10%) and recall (88.14%) indicate the model's effectiveness in both approving worthy candidates and identifying potential defaults.
3. Robust Discrimination: An ROC-AUC score of 0.9200 suggests suggests the model's strong ability to distinguish between approved and rejected loan applications.
4. Key Factors: The analysis highlighted CIBIL score, income, and loan amount as crucial factors in loan approval decisions.
5. Practical Applicability: The model's performance suggests it could be a valuable tool in assisting loan approval decisions in real-world scenarios.