

📊 Creditworthiness Prediction Analysis

Creditworthiness Prediction using CART and Random Forest



Khairunnisa Maharani 2025-06-07



Creditworthiness Prediction Analysis

This project explores the prediction of creditworthiness using machine learning, specifically Decision Tree (CART) and Random Forest algorithms. It handles class imbalance using ROSE and compares model performance with various configurations.



Dataset Description

- Dataset: German Credit Dataset from UCI ML Repository
- Samples: 1000 customers
- Target Variable: Creditworthiness (Good or Bad)
- Features: 20 predictor variables (demographic and financial attributes)

1. Load Required Libraries

```
# Core data manipulation and visualization
library(tidyverse)
library(ggplot2)
library(gridExtra)
library(corrplot)
library(knitr)
# Machine learning libraries
library(rpart)
library(rpart.plot)
library(randomForest)
library(caret)
# Model evaluation and handling imbalanced data
library(ROSE)
library(pROC)
# Additional visualization and analysis
library(visdat)
library(naniar)
library(skimr)
library(DT)
# Set theme for consistent plots
theme_set(theme_minimal())
```

2. Data Loading and Initial Exploration

```
## Dataset structure:

#Hide

# Load dataset

data <- read.csv("german_data.csv")

# Basic data exploration
cat("Dataset dimensions:", dim(data), "\n")

## Dataset dimensions: 1000 21

## Dataset structure:\n")
```

str(data)

```
## 'data.frame':
                  1000 obs. of 21 variables:
## $ Status.Rekening
                                  : chr "< 0 DM" "0-199 DM" "Tidak ada rekenin
g" "< 0 DM" ...
## $ Durasi.dalam.bulan
                                 : int 6 48 12 42 24 36 24 36 12 30 ...
## $ Riwayat.Kredit
                                  : chr "Masalah / kredit di luar bank" "Kredit
berjalan lancar" "Masalah / kredit di luar bank" "Kredit berjalan lancar" ...
## $ Tujuan.Kredit
                                  : chr "TV / Radio" "TV / Radio" "Pendidikan"
"Furniture" ...
## $ Jumlah.Kredit..DM.
                                 : int 1169 5951 2096 7882 4870 9055 2835 6948
3059 5234 ...
## $ Tabungan.Bonds
                         : chr "Tidak ada / tidak tahu" "< 100 DM" "<
100 DM" "< 100 DM" ...
## $ Pekerjaan.Saat.Ini
                                  : chr ">= 7 thn" "1-4 thn" "4-7 thn" "4-7 th
n" ...
## $ Tingkat.Angsuran....
                                  : int 422232324...
## $ Status.Pribadi.dan.Jenis.Kelamin: chr "Pria lajang" "Wanita menikah/cerai" "P
ria lajang" "Pria lajang" ...
                         : chr "Tidak ada" "Tidak ada" "Tidak ada" "Pe
## $ Debitur...Penjamin
njamin" ...
## $ Tinggal.di.Tempat.Saat.Ini : int 4 2 3 4 4 4 4 2 4 2 ...
## $ Properti
                                  : chr "Properti" "Properti" "Tabun
gan/Asuransi" ...
## $ Usia..tahun.
                                 : int 67 22 49 45 53 35 53 35 61 28 ...
## $ Rencana.Angsuran.Lain
                                  : chr "Tidak ada" "Tidak ada" "Ti
dak ada" ...
## $ Tempat.Tinggal
                                 : chr "Milik sendiri" "Milik sendiri" "Milik
sendiri" "Gratis" ...
## $ Jumlah.Kredit.di.Bank.Ini : int 2 1 1 1 2 1 1 1 2 ...
## $ Pekerjaan
                                  : chr "Terampil / pegawai" "Terampil / pegawa
i" "Tidak terampil (residen)" "Terampil / pegawai" ...
                                  : int 1122221111...
## $ Jumlah.Tanggungan
## $ Nomor.Telepon
                                  : chr "Ada" "Tidak ada" "Tidak ada" "Tidak ad
a" ...
## $ Pekerja.Migran
                                  : chr "Ya" "Ya" "Ya" "Ya" ...
                                  : chr "Baik" "Buruk" "Baik" ...
## $ Kelayakan.Kredit
```

```
# Comprehensive data summary using skimr
skim_result <- skim(data)
print(skim_result)</pre>
```

```
## — Data Summary
##
                                Values
## Name
                                data
## Number of rows
                                1000
## Number of columns
                                21
##
## Column type frequency:
##
     character
                                14
##
     numeric
                                7
##
## Group variables
                                None
##
## — Variable type: character -
      skim_variable
##
                                          n_missing complete_rate min max empty
## 1 Status.Rekening
                                                   0
                                                                         22
                                                                  1
                                                                      6
                                                                                 0
    2 Riwayat.Kredit
                                                   0
                                                                  1
                                                                     18
                                                                         35
                                                                                 0
##
   3 Tujuan.Kredit
                                                   0
                                                                  1
                                                                      6
                                                                         15
                                                                                 0
## 4 Tabungan.Bonds
                                                                      8
                                                                         22
                                                   0
                                                                  1
                                                                                 0
## 5 Pekerjaan.Saat.Ini
                                                                  1
                                                                      7
                                                                         10
                                                   0
## 6 Status.Pribadi.dan.Jenis.Kelamin
                                                                  1
                                                                         20
                                                   0
                                                                     11
                                                                                 0
## 7 Debitur...Penjamin
                                                   0
                                                                  1
                                                                      8
                                                                         10
                                                                                 0
                                                                         22
## 8 Properti
                                                   0
                                                                  1
                                                                      8
## 9 Rencana.Angsuran.Lain
                                                                  1
                                                                      4
                                                                           9
                                                   0
                                                                                 0
## 10 Tempat.Tinggal
                                                   0
                                                                  1
                                                                      4
                                                                         13
                                                                                 0
## 11 Pekerjaan
                                                   0
                                                                  1
                                                                     18
                                                                         43
                                                                                 0
## 12 Nomor.Telepon
                                                                      3
                                                                           9
                                                   0
                                                                  1
                                                                                 0
## 13 Pekerja.Migran
                                                                           5
                                                   0
                                                                  1
                                                                      2
                                                                                 0
## 14 Kelayakan.Kredit
                                                   0
                                                                      4
                                                                           5
                                                                                 0
      n_unique whitespace
##
##
              4
                          0
    1
    2
              5
                          0
##
##
    3
             10
                          0
    4
              5
##
                          0
    5
              5
##
                          0
    6
              4
                          0
    7
##
              3
                          0
              4
##
    8
                          0
    9
              3
##
                          0
## 10
              3
                          0
              4
## 11
                          0
              2
## 12
                          0
## 13
              2
                          0
## 14
              2
                          0
## — Variable type: numeric -
##
     skim_variable
                                  n_missing complete_rate
                                                               mean
                                                                            sd p0
                                                                                     p25
## 1 Durasi.dalam.bulan
                                           0
                                                               20.9
                                                                       12.1
                                                                                 4
                                                                                     12
## 2 Jumlah.Kredit..DM.
                                           0
                                                          1 3271.
                                                                     2823.
                                                                               250 1366.
```

```
2.97
                                                                               2
## 3 Tingkat.Angsuran....
                                                     1
                                                                  1.12
                                                                          1
## 4 Tinggal.di.Tempat.Saat.Ini
                                                          2.84
                                                                  1.10
                                                                          1
                                                                               2
## 5 Usia..tahun.
                                                     1
                                                         35.5
                                                                 11.4
                                                                              27
                                       0
                                                                         19
## 6 Jumlah.Kredit.di.Bank.Ini
                                                                  0.578
                                       0
                                                     1
                                                          1.41
                                                                          1
                                                                               1
## 7 Jumlah.Tanggungan
                                                          1.16
                                                                  0.362
                                                                               1
                                                                          1
      p50
             p75 p100 hist
##
## 1
      18
             24
                   72
## 2 2320. 3972. 18424 _____
## 3
                    4 ____
## 4
       3
                    4 ____
             4
## 5
      33
            42
## 6
             2
## 7
             1
       1
```

3. Data Quality Assessment

Hide

```
# Check for missing values
missing values <- colSums(is.na(data))</pre>
missing_percentage <- (missing_values / nrow(data)) * 100</pre>
missing_df <- data.frame(</pre>
    Missing_Count = missing_values,
    Missing_Percentage = missing_percentage
)
# Display missing values if any
if(sum(missing_values) > 0) {
    cat("Missing Values Summary:\n")
    print(missing_df[missing_df$Missing_Count > 0, ])
    # Visualize missing values pattern
    vis_miss(data, cluster = TRUE)
    gg_miss_var(data)
} else {
    cat("☑ No missing values found in the dataset!\n")
}
```

```
## ☑ No missing values found in the dataset!
```

```
# Convert categorical variables to factors
categorical_vars <- names(data)[sapply(data, is.character)]
for(var in categorical_vars) {
    data[[var]] <- as.factor(data[[var]])
}

# Convert target variable to factor with proper levels
data$Kelayakan.Kredit <- factor(data$Kelayakan.Kredit, levels = c("Buruk", "Baik"))

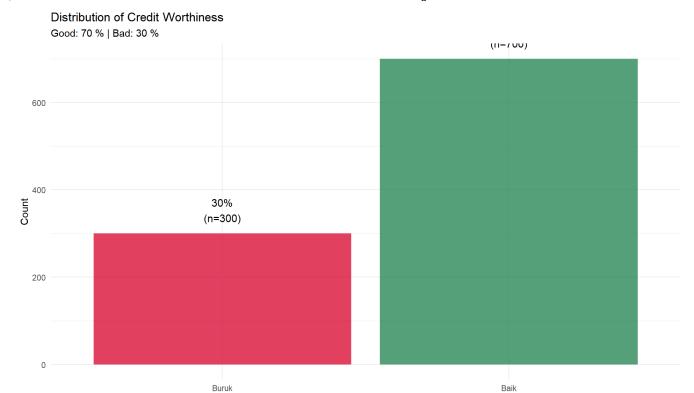
cat("Categorical variables converted:", length(categorical_vars), "variables\n")</pre>
```

```
## Categorical variables converted: 14 variables
```

4. Exploratory Data Analysis (EDA)

4.1 Target Variable Distribution

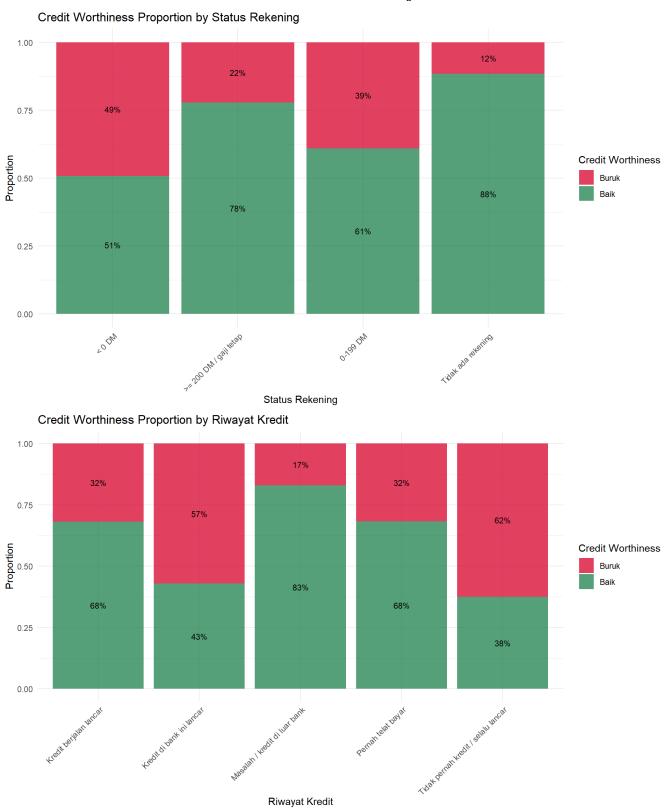
```
# Calculate target distribution
target_dist <- table(data$Kelayakan.Kredit)</pre>
target_prop <- prop.table(target_dist) * 100</pre>
# Visualize target distribution
plt_target <- ggplot(data, aes(x = Kelayakan.Kredit, fill = Kelayakan.Kredit)) +</pre>
    geom bar(alpha = 0.8) +
    geom_text(stat = "count",
              aes(label = paste0(scales::percent(after_stat(count)/sum(after_stat(count))
         unt))),
                                 "\n(n=", after_stat(count), ")")),
              vjust = -0.5, size = 4) +
    scale_fill_manual(values = c("Baik" = "#2E8B57", "Buruk" = "#DC143C")) +
    labs(title = "Distribution of Credit Worthiness",
         subtitle = paste("Good:", round(target_prop[2], 1), "% | Bad:", round(targe
         t_prop[1], 1), "%"),
         x = "Credit Worthiness",
         y = "Count") +
    theme(legend.position = "none")
print(plt_target)
```

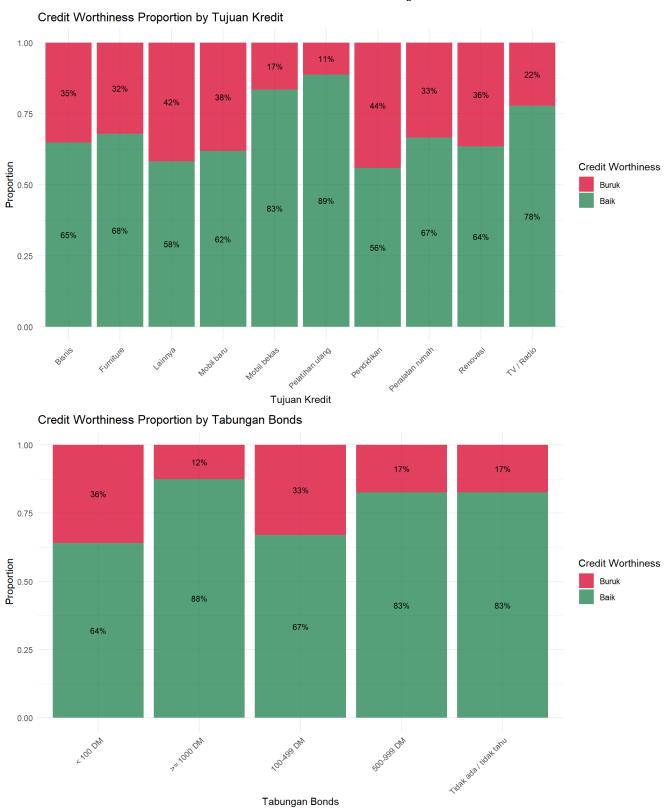


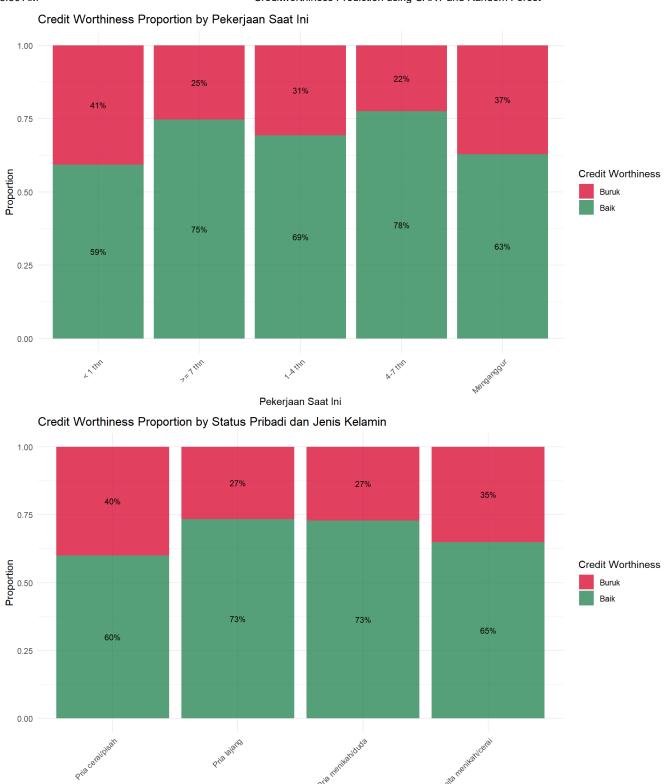
Credit Worthiness

4.2 Categorical Variables Analysis

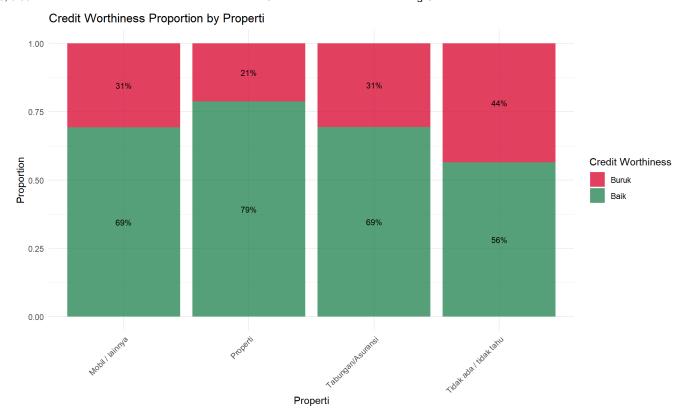
```
# Function to create proportion plots for categorical variables
plot_categorical_variable <- function(data, var_name) {</pre>
    prop_df <- data %>%
        group by(!!sym(var name), Kelayakan.Kredit) %>%
        summarise(count = n(), .groups = "drop") %>%
        group_by(!!sym(var_name)) %>%
        mutate(
            prop = count / sum(count),
            total = sum(count)
        )
    ggplot(prop_df, aes(x = !!sym(var_name), y = prop, fill = Kelayakan.Kredit)) +
        geom_bar(stat = "identity", position = "fill", alpha = 0.8) +
        geom_text(aes(label = paste0(round(prop*100), "%")),
                  position = position_fill(vjust = 0.5), size = 3) +
        scale_fill_manual(values = c("Baik" = "#2E8B57", "Buruk" = "#DC143C")) +
        labs(title = paste("Credit Worthiness Proportion by", gsub("\\.", " ", var_n
         ame)),
             x = gsub("\\.", " ", var_name),
             y = "Proportion",
             fill = "Credit Worthiness") +
        theme(axis.text.x = element_text(angle = 45, hjust = 1))
}
# Important categorical variables to analyze
important_cat_vars <- c("Status.Rekening", "Riwayat.Kredit", "Tujuan.Kredit",</pre>
                       "Tabungan.Bonds", "Pekerjaan.Saat.Ini",
                       "Status.Pribadi.dan.Jenis.Kelamin", "Properti")
# Generate plots for categorical variables
cat_plots <- map(important_cat_vars, ~plot_categorical_variable(data, .x))</pre>
names(cat_plots) <- important_cat_vars</pre>
# Display plots
walk(cat plots, print)
```





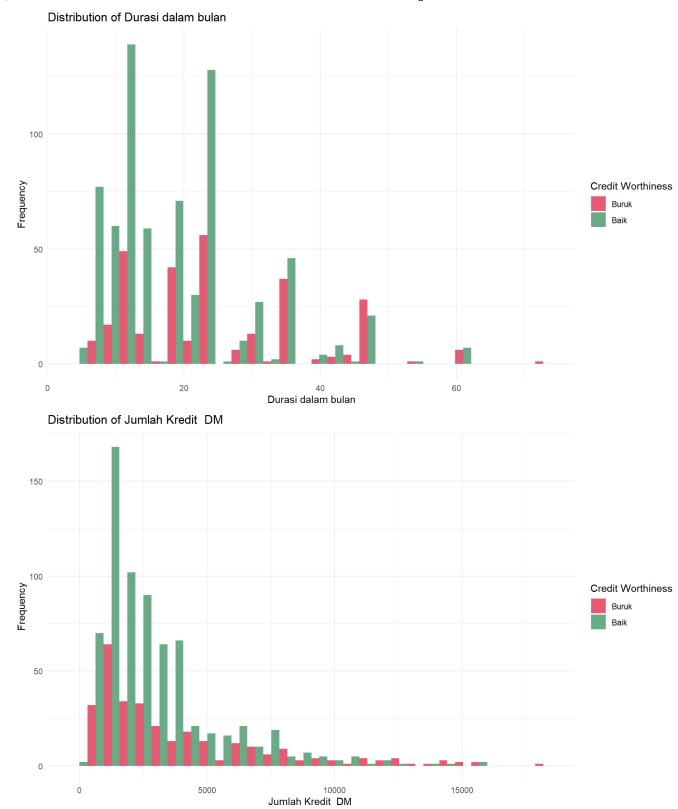


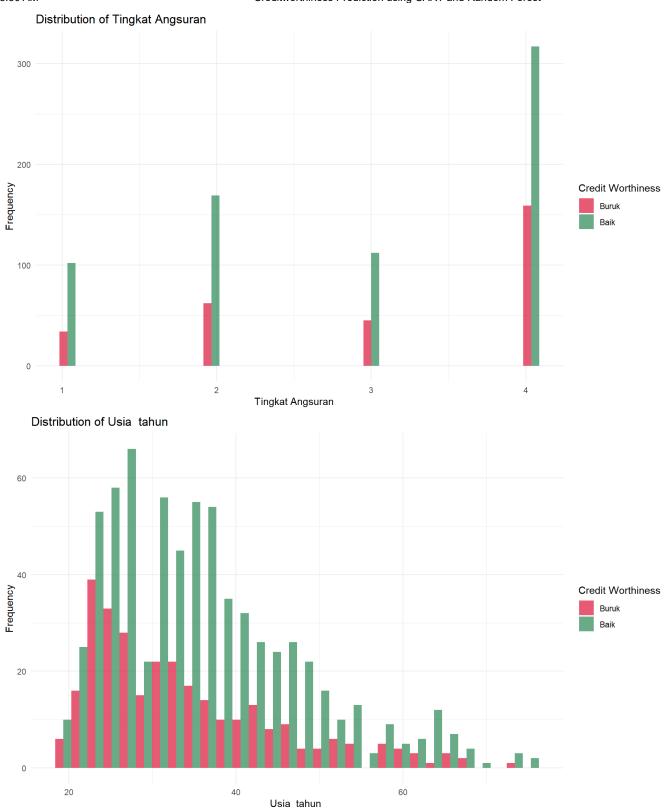
Status Pribadi dan Jenis Kelamin

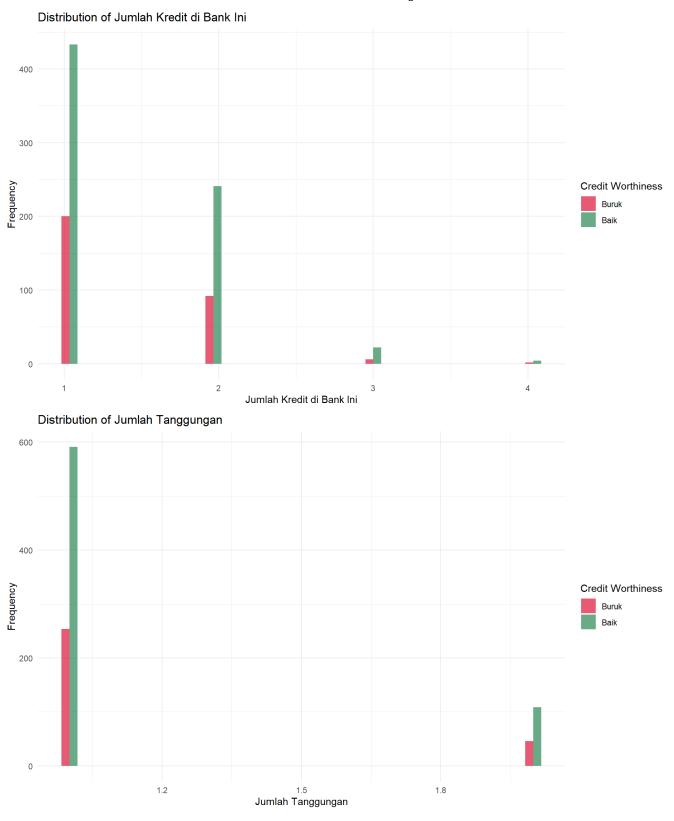


4.3 Numerical Variables Analysis

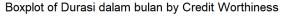
```
# Identify numerical variables
numerical_vars <- c("Durasi.dalam.bulan", "Jumlah.Kredit..DM.", "Tingkat.Angsura</pre>
         n....",
                    "Usia..tahun.", "Jumlah.Kredit.di.Bank.Ini", "Jumlah.Tanggungan")
# Create histograms for numerical variables
create_histogram <- function(data, var_name) {</pre>
    ggplot(data, aes_string(x = var_name, fill = "Kelayakan.Kredit")) +
        geom_histogram(position = "dodge", alpha = 0.7, bins = 30) +
        scale_fill_manual(values = c("Baik" = "#2E8B57", "Buruk" = "#DC143C")) +
        labs(title = paste("Distribution of", gsub("\\.", " ", var_name)),
             x = gsub("\\.", " ", var_name),
             y = "Frequency",
             fill = "Credit Worthiness") +
        theme_minimal()
}
# Generate histograms
num_histograms <- map(numerical_vars, ~create_histogram(data, .x))</pre>
names(num_histograms) <- numerical_vars</pre>
# Display histograms
walk(num_histograms, print)
```

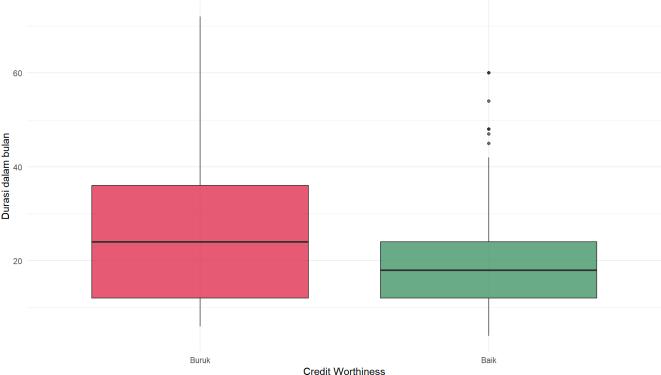


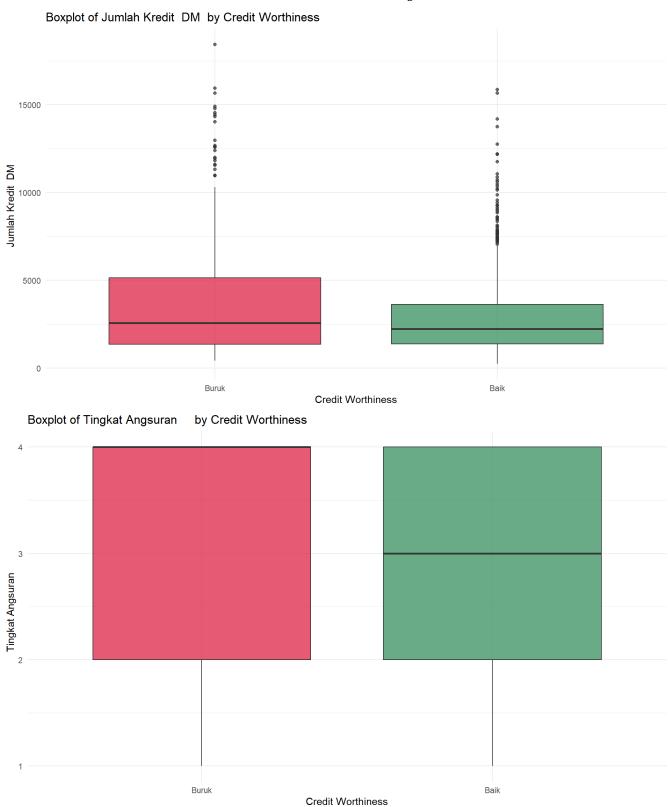


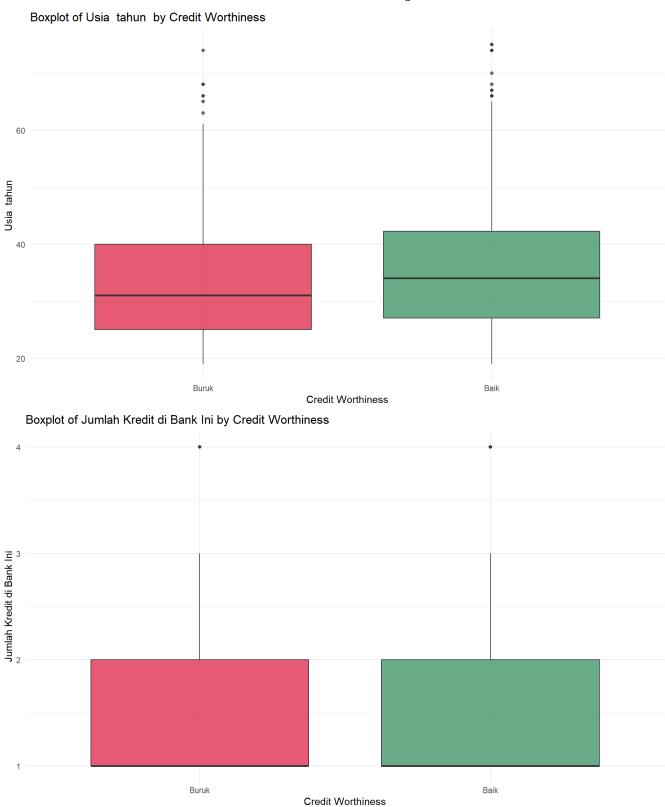


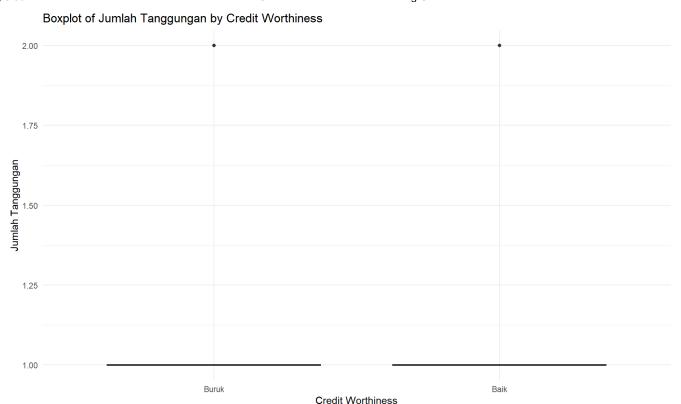
```
# Create boxplots for numerical variables
create_boxplot <- function(data, var_name) {</pre>
    ggplot(data, aes_string(x = "Kelayakan.Kredit", y = var_name, fill = "Kelayakan.
         Kredit")) +
        geom_boxplot(alpha = 0.7) +
        scale_fill_manual(values = c("Baik" = "#2E8B57", "Buruk" = "#DC143C")) +
        labs(title = paste("Boxplot of", gsub("\\.", " ", var_name), "by Credit Wort
         hiness"),
             x = "Credit Worthiness",
             y = gsub("\\.", " ", var_name)) +
        theme(legend.position = "none")
}
# Generate boxplots
num_boxplots <- map(numerical_vars, ~create_boxplot(data, .x))</pre>
names(num_boxplots) <- numerical_vars</pre>
# Display boxplots
walk(num_boxplots, print)
```





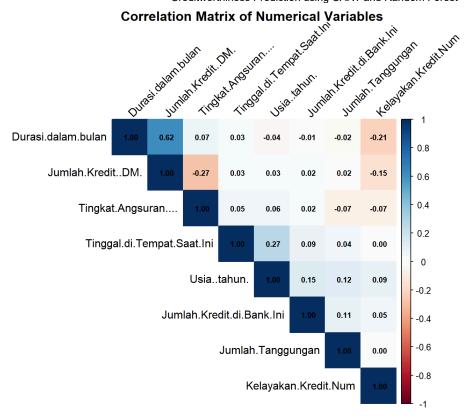






4.4 Correlation Analysis

```
# Create numeric version of target for correlation
data$Kelayakan.Kredit.Num <- ifelse(data$Kelayakan.Kredit == "Baik", 1, 0)</pre>
# Select numerical variables for correlation analysis
num_data <- data %>% select_if(is.numeric)
# Calculate correlation matrix
corr_matrix <- cor(num_data, use = "complete.obs")</pre>
# Visualize correlation matrix
corrplot(corr_matrix,
         method = "color",
         type = "upper",
         tl.col = "black",
         tl.srt = 45,
         addCoef.col = "black",
         number.cex = 0.7,
         title = "Correlation Matrix of Numerical Variables",
         mar = c(0,0,1,0)
```



5. Data Preprocessing and Model Preparation

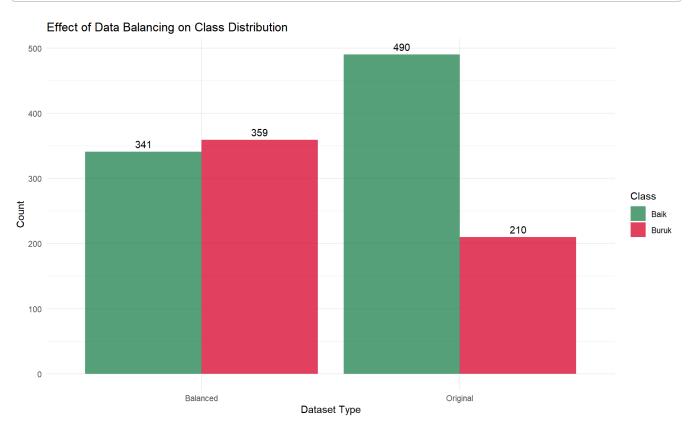
```
## Training set size: 700
```

```
cat("Test set size:", nrow(test_data), "\n")
```

```
## Test set size: 300
```

```
Hide
```

```
cat("Training set class distribution:\n")
## Training set class distribution:
                                                                                    Hide
print(table(train_data$Kelayakan.Kredit))
##
## Buruk Baik
##
     210
           490
                                                                                    Hide
# Handle class imbalance using ROSE
train_balanced <- ROSE(Kelayakan.Kredit ~ ., data = train_data, seed = 123)$data</pre>
cat("Balanced training set class distribution:\n")
## Balanced training set class distribution:
                                                                                    Hide
print(table(train_balanced$Kelayakan.Kredit))
##
##
    Baik Buruk
##
     359
           341
                                                                                    Hide
```



6. CART Model Development

6.1 Basic CART Model

```
# Fit CART model
cart_model <- rpart(
    Kelayakan.Kredit ~ .,
    data = train_balanced,
    method = "class",
    control = rpart.control(
        cp = 0.01,
        minbucket = 5,
        minsplit = 10
    )
)

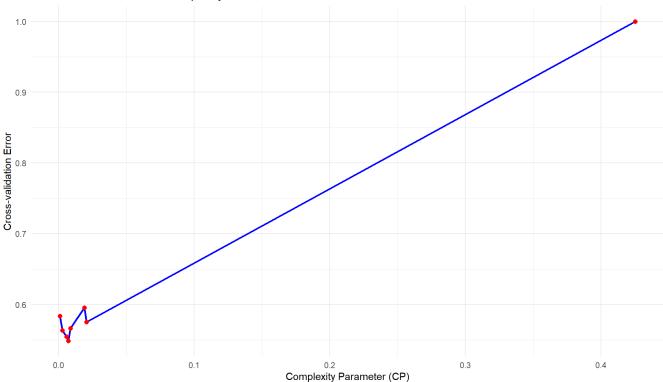
# Print model summary
printcp(cart_model)</pre>
```

```
##
## Classification tree:
## rpart(formula = Kelayakan.Kredit ~ ., data = train_balanced,
      method = "class", control = rpart.control(cp = 0.01, minbucket = 5,
##
          minsplit = 10))
## Variables actually used in tree construction:
## [1] Durasi.dalam.bulan Pekerjaan
                                      Pekerjaan.Saat.Ini Properti
## [5] Riwayat.Kredit
                         Status.Rekening Tabungan.Bonds
                                                              Tujuan.Kredit
## Root node error: 341/700 = 0.48714
## n= 700
##
          CP nsplit rel error xerror
                                          xstd
                  0 1.00000 1.00000 0.038781
## 1 0.425220
## 2 0.020528
                  1 0.57478 0.57478 0.034837
## 3 0.019062
                  5 0.48680 0.60997 0.035458
## 4 0.010000
                  8 0.42815 0.60411 0.035358
```

6.2 CART Model Pruning

```
# Build full tree for pruning analysis
cart_model_full <- rpart(</pre>
    Kelayakan.Kredit ~ .,
    data = train_balanced,
    method = "class",
    control = rpart.control(
        cp = 0.001,
        minbucket = 5,
        minsplit = 10
    )
)
# Plot cross-validation error vs complexity parameter
cp_table <- as.data.frame(cart_model_full$cptable)</pre>
ggplot(cp_table, aes(x = CP, y = xerror)) +
    geom_line(color = "blue", size = 1) +
    geom_point(color = "red", size = 2) +
    labs(title = "Cross-validation Error vs Complexity Parameter",
         x = "Complexity Parameter (CP)",
         y = "Cross-validation Error") +
    theme_minimal()
```

Cross-validation Error vs Complexity Parameter



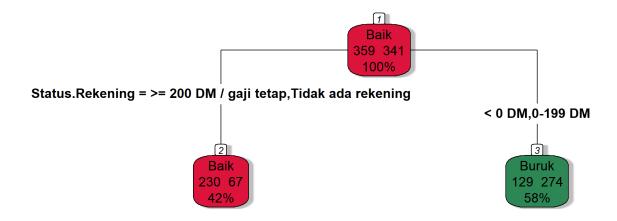
```
## Optimal CP (1-SE rule): 0.02052786
```

```
# Prune the tree
cart_model_pruned <- prune(cart_model_full, cp = optimal_cp)</pre>
```

6.3 CART Model Visualization

```
# Visualize pruned decision tree
rpart.plot(
    cart_model_pruned,
    type = 4,
    extra = 101,
    box.palette = c("#DC143C", "#2E8B57"),
    branch.lty = 1,
    shadow.col = "gray",
    nn = TRUE,
    tweak = 1.2,
    fallen.leaves = TRUE,
    uniform = TRUE,
    main = "Pruned CART Decision Tree for Credit Worthiness Prediction"
)
```

Pruned CART Decision Tree for Credit Worthiness Prediction



Hide

```
# Save high-quality decision tree plot
if(!dir.exists("plots")) dir.create("plots")
pdf("plots/decision_tree_pruned.pdf", width = 12, height = 8)
rpart.plot(
    cart_model_pruned,
    type = 4,
    extra = 101,
    box.palette = c("#DC143C", "#2E8B57"),
    branch.lty = 1,
    shadow.col = "gray",
    nn = TRUE,
    tweak = 1.2,
    fallen.leaves = TRUE,
    uniform = TRUE,
    main = "Pruned CART Decision Tree for Credit Worthiness Prediction"
)
dev.off()
```

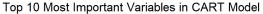
```
## png
## 2
```

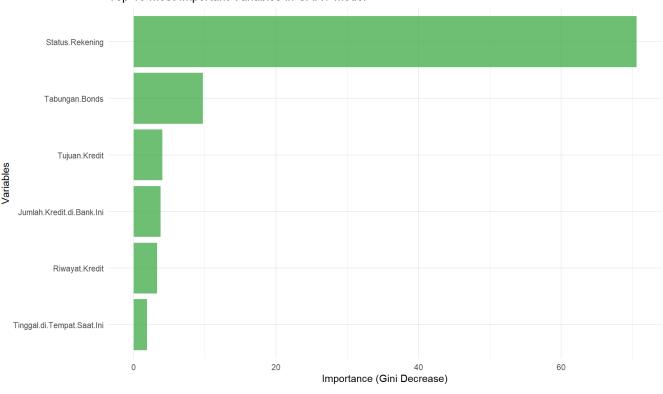
```
cat("Decision tree saved to: plots/decision_tree_pruned.pdf\n")
```

```
## Decision tree saved to: plots/decision_tree_pruned.pdf
```

6.4 CART Variable Importance

```
# Extract and visualize variable importance
cart_importance <- cart_model_pruned$variable.importance</pre>
if(length(cart_importance) > 0) {
    cart_importance_df <- data.frame(</pre>
        Variable = names(cart_importance),
        Importance = as.numeric(cart_importance)
    ) %>%
    arrange(desc(Importance)) %>%
    slice_head(n = 10)
    ggplot(cart_importance_df, aes(x = reorder(Variable, Importance), y = Importance
        geom_bar(stat = "identity", fill = "#4CAF50", alpha = 0.8) +
        coord_flip() +
        labs(title = "Top 10 Most Important Variables in CART Model",
             x = "Variables", y = "Importance (Gini Decrease)") +
        theme_minimal()
} else {
    cat("No variable importance information available for pruned model.\n")
}
```





7. Random Forest Model Development

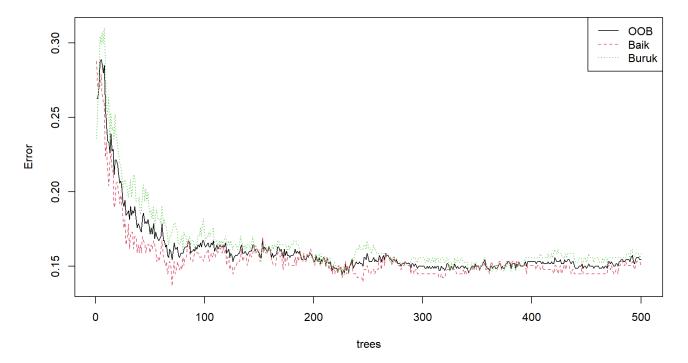
```
# Fit Random Forest model
set.seed(123)
rf_model <- randomForest(
    Kelayakan.Kredit ~ .,
    data = train_balanced,
    ntree = 500,
    mtry = sqrt(ncol(train_balanced) - 1),
    importance = TRUE
)

# Print model summary
print(rf_model)</pre>
```

```
##
## Call:
## randomForest(formula = Kelayakan.Kredit ~ ., data = train_balanced,
                                                                              ntree =
500, mtry = sqrt(ncol(train_balanced) - 1), importance = TRUE)
##
                  Type of random forest: classification
                        Number of trees: 500
##
## No. of variables tried at each split: 4
##
           OOB estimate of error rate: 15.43%
##
## Confusion matrix:
##
         Baik Buruk class.error
## Baik
          305
                 54
                      0.1504178
## Buruk
           54
                287
                      0.1583578
```

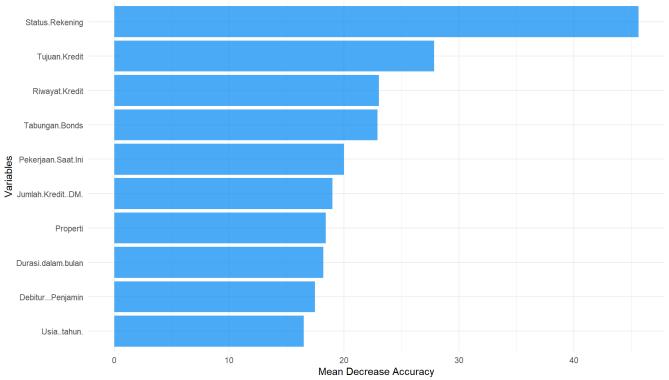
```
# Plot error rate
plot(rf_model, main = "Random Forest Error Rate by Number of Trees")
legend("topright", colnames(rf_model$err.rate), col = 1:3, lty = 1:3)
```

Random Forest Error Rate by Number of Trees



Hide





8. Model Evaluation and Comparison

8.1 Model Predictions

```
# CART predictions
cart_pred <- predict(cart_model_pruned, test_data, type = "class")
cart_prob <- predict(cart_model_pruned, test_data, type = "prob")

# Random Forest predictions
rf_pred <- predict(rf_model, test_data, type = "class")
rf_prob <- predict(rf_model, test_data, type = "prob")

cat("Predictions completed for both models.\n")</pre>
```

Predictions completed for both models.

8.2 Confusion Matrix Analysis

Hide

```
# CART confusion matrix
cart_cm <- confusionMatrix(cart_pred, test_data$Kelayakan.Kredit, positive = "Baik")
print("CART Model Confusion Matrix:")</pre>
```

```
## [1] "CART Model Confusion Matrix:"
```

Hide

print(cart_cm)

```
## Confusion Matrix and Statistics
             Reference
##
## Prediction Buruk Baik
##
        Buruk
                 68
                      97
        Baik
##
                 22 113
##
##
                  Accuracy : 0.6033
##
                    95% CI: (0.5455, 0.6591)
      No Information Rate: 0.7
##
       P-Value [Acc > NIR] : 0.9999
##
##
##
                     Kappa: 0.2372
##
##
   Mcnemar's Test P-Value : 1.172e-11
##
               Sensitivity: 0.5381
##
##
               Specificity: 0.7556
            Pos Pred Value: 0.8370
##
##
            Neg Pred Value : 0.4121
                Prevalence: 0.7000
##
            Detection Rate: 0.3767
##
      Detection Prevalence : 0.4500
##
##
         Balanced Accuracy: 0.6468
##
          'Positive' Class : Baik
##
##
                                                                                   Hide
```

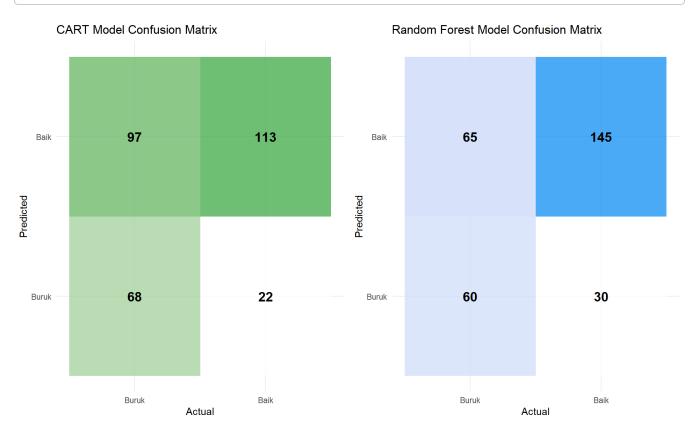
```
# Random Forest confusion matrix
rf_cm <- confusionMatrix(rf_pred, test_data$Kelayakan.Kredit, positive = "Baik")
print("Random Forest Model Confusion Matrix:")</pre>
```

```
## [1] "Random Forest Model Confusion Matrix:"
```

```
print(rf_cm)
```

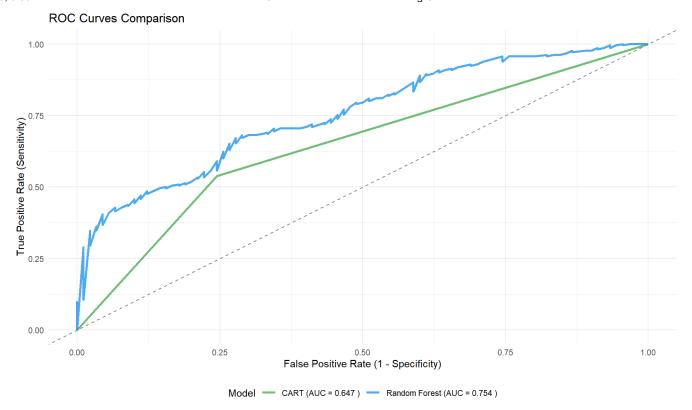
```
## Confusion Matrix and Statistics
             Reference
##
## Prediction Buruk Baik
##
        Buruk
                 60
                      65
        Baik
                 30 145
##
##
##
                  Accuracy : 0.6833
##
                    95% CI: (0.6274, 0.7356)
       No Information Rate: 0.7
##
       P-Value [Acc > NIR] : 0.7571687
##
##
##
                     Kappa : 0.3214
##
   Mcnemar's Test P-Value : 0.0004861
##
##
##
               Sensitivity: 0.6905
               Specificity: 0.6667
##
            Pos Pred Value: 0.8286
##
##
            Neg Pred Value: 0.4800
##
                Prevalence: 0.7000
            Detection Rate: 0.4833
##
##
      Detection Prevalence : 0.5833
##
         Balanced Accuracy: 0.6786
##
##
          'Positive' Class : Baik
##
```

```
# Function to create confusion matrix plot
plot_confusion_matrix <- function(cm, title, color) {</pre>
    cm_df <- as.data.frame(cm$table)</pre>
    colnames(cm_df) <- c("Actual", "Predicted", "Frequency")</pre>
    ggplot(cm_df, aes(x = Actual, y = Predicted, fill = Frequency)) +
        geom tile(alpha = 0.8) +
        geom_text(aes(label = Frequency), color = "black", size = 5, fontface = "bol
         d") +
        scale_fill_gradient(low = "white", high = color) +
        labs(title = title, x = "Actual", y = "Predicted") +
        theme minimal() +
        theme(legend.position = "none")
}
# Create confusion matrix plots
cart_cm_plot <- plot_confusion_matrix(cart_cm, "CART Model Confusion Matrix", "#4CAF</pre>
rf_cm_plot <- plot_confusion_matrix(rf_cm, "Random Forest Model Confusion Matrix",
         "#2196F3")
# Display plots side by side
grid.arrange(cart_cm_plot, rf_cm_plot, ncol = 2)
```



8.3 ROC Curve Analysis

```
# Calculate ROC curves
cart_roc <- roc(as.numeric(test_data$Kelayakan.Kredit) - 1, cart_prob[, "Baik"])</pre>
rf roc <- roc(as.numeric(test_data$Kelayakan.Kredit) - 1, rf_prob[, "Baik"])</pre>
# Create ROC curve data frame
create_roc_df <- function(roc_obj, model_name) {</pre>
    data.frame(
        Specificity = roc_obj$specificities,
        Sensitivity = roc_obj$sensitivities,
        Model = model name,
        AUC = round(auc(roc obj), 3)
    )
}
cart roc df <- create roc df(cart roc, "CART")</pre>
rf_roc_df <- create_roc_df(rf_roc, "Random Forest")</pre>
roc_combined <- rbind(cart_roc_df, rf_roc_df)</pre>
# Plot ROC curves
ggplot(roc\_combined, aes(x = 1 - Specificity, y = Sensitivity, color = Model)) +
    geom_line(size = 1.2, alpha = 0.8) +
    geom_abline(slope = 1, intercept = 0, linetype = "dashed", color = "gray50") +
    scale_color_manual(
        values = c("CART" = "#4CAF50", "Random Forest" = "#2196F3"),
        labels = \mathbf{c}(
            paste("CART (AUC =", unique(cart roc df$AUC), ")"),
            paste("Random Forest (AUC =", unique(rf_roc_df$AUC), ")")
        )
    ) +
    labs(
        title = "ROC Curves Comparison",
        x = "False Positive Rate (1 - Specificity)",
        y = "True Positive Rate (Sensitivity)",
        color = "Model"
    ) +
    theme minimal() +
    theme(legend.position = "bottom")
```



8.4 Model Performance Comparison

Hide

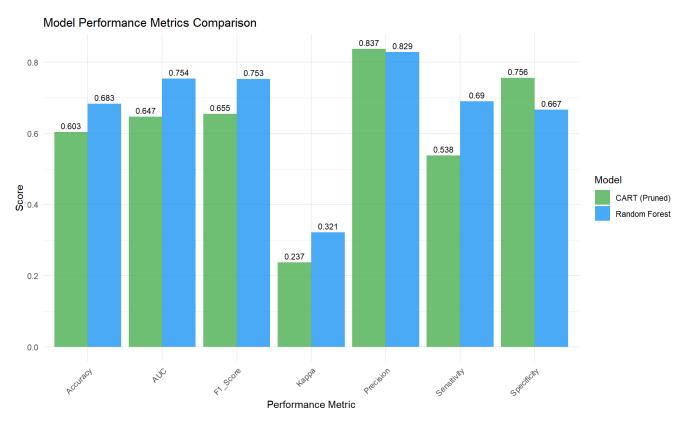
```
# Create comprehensive model comparison table
model_comparison <- data.frame(</pre>
    Model = c("CART (Pruned)", "Random Forest"),
    Accuracy = c(cart_cm$overall["Accuracy"], rf_cm$overall["Accuracy"]),
    Sensitivity = c(cart_cm$byClass["Sensitivity"], rf_cm$byClass["Sensitivity"]),
    Specificity = c(cart_cm$byClass["Specificity"], rf_cm$byClass["Specificity"]),
    Precision = c(cart_cm$byClass["Pos Pred Value"], rf_cm$byClass["Pos Pred Value"]
         e"]),
    F1_Score = c(cart_cm$byClass["F1"], rf_cm$byClass["F1"]),
    AUC = c(auc(cart_roc), auc(rf_roc)),
    Kappa = c(cart_cm$overall["Kappa"], rf_cm$overall["Kappa"])
)
# Display comparison table
kable(model_comparison,
      digits = 4,
      caption = "Model Performance Comparison",
      col.names = c("Model", "Accuracy", "Sensitivity", "Specificity",
                   "Precision", "F1-Score", "AUC", "Kappa"))
```

Model Performance Comparison

Model Accuracy Sensitivity Specificity Precision F1-Score AUC Kappa

Model	Accuracy	Sensitivity	Specificity	Precision	F1-Score	AUC	Карра
CART (Pruned)	0.6033	0.5381	0.7556	0.8370	0.6551	0.6468	0.2372
Random Forest	0.6833	0.6905	0.6667	0.8286	0.7532	0.7536	0.3214

```
# Visualize model performance metrics
model_comparison_long <- model_comparison %>%
    pivot_longer(
        cols = c("Accuracy", "Sensitivity", "Specificity", "Precision", "F1_Score",
         "AUC", "Kappa"),
        names_to = "Metric",
        values_to = "Value"
    )
ggplot(model_comparison_long, aes(x = Metric, y = Value, fill = Model)) +
    geom_bar(stat = "identity", position = "dodge", alpha = 0.8) +
    geom_text(aes(label = round(Value, 3)),
              position = position_dodge(width = 0.9),
              vjust = -0.5, size = 3) +
    scale_fill_manual(values = c("CART (Pruned)" = "#4CAF50", "Random Forest" = "#21
    labs(title = "Model Performance Metrics Comparison",
         x = "Performance Metric", y = "Score") +
    theme_minimal() +
    theme(axis.text.x = element_text(angle = 45, hjust = 1))
```



9. Results Summary and Recommendations

```
# Identify best performing model
best_model_auc <- model_comparison$Model[which.max(model_comparison$AUC)]</pre>
best_model_f1 <- model_comparison$Model[which.max(model_comparison$F1_Score)]</pre>
best_auc_score <- round(max(model_comparison$AUC), 3)</pre>
best_f1_score <- round(max(model_comparison$F1_Score), 3)</pre>
cat("@ ANALYSIS RESULTS SUMMARY\n")
## 🎯 ANALYSIS RESULTS SUMMARY
                                                                                   Hide
cat("======\n\n")
                                                                                   Hide
cat(" | Best Model by AUC: ", best_model_auc, "with AUC = ", best_auc_score, "\n")
## | Best Model by AUC: Random Forest with AUC = 0.754
                                                                                   Hide
cat(" | Best Model by F1-Score:", best_model_f1, "with F1-Score =", best_f1_score,
         "\n\n")
## 📊 Best Model by F1-Score: Random Forest with F1-Score = 0.753
                                                                                   Hide
cat(" \( \) KEY FINDINGS:\n")
## \ KEY FINDINGS:
                                                                                   Hide
```

cat("1. Both models show good performance for credit worthiness prediction\n")

1. Both models show good performance for credit worthiness prediction

Hide

cat("2. Data balancing improved model fairness across both classes\n")

2. Data balancing improved model fairness across both classes

Hide

cat("3. Model pruning reduced complexity while maintaining performance\n")

3. Model pruning reduced complexity while maintaining performance

Hide

cat("4. Most important factors for credit assessment identified\n\n")

4. Most important factors for credit assessment identified

Hide

cat(" ? RECOMMENDATIONS:\n")

PRECOMMENDATIONS:

Hide

1. Use ensemble approach for production deployment

Hide

 $cat("2. Implement regular model retraining with new data \n")$

2. Implement regular model retraining with new data

```
cat("3. Consider cost-sensitive learning for real-world application \")
```

```
## 3. Consider cost-sensitive learning for real-world application
```

Hide

```
cat("4. Validate model performance on external datasets\n")
```

```
## 4. Validate model performance on external datasets
```

10. Save Results and Outputs

Hide

```
# Create output directory
if(!dir.exists("results")) dir.create("results")
# Save model comparison results
write.csv(model_comparison, "results/model_comparison.csv", row.names = FALSE)
# Save variable importance results
if(exists("cart_importance_df") && nrow(cart_importance_df) > 0) {
    write.csv(cart_importance_df, "results/cart_variable_importance.csv", row.names
         = FALSE)
}
if(exists("rf_importance_df")) {
    write.csv(rf_importance_df, "results/rf_variable_importance.csv", row.names = FA
         LSE)
}
# Save model objects (optional)
# saveRDS(cart_model_pruned, "results/cart_model_pruned.rds")
# saveRDS(rf_model, "results/rf_model.rds")
cat(" ■ Results saved to 'results/' directory\n")
```

```
## ☑ Results saved to 'results/' directory
```

cat("☑ Analysis completed successfully!\n")

☑ Analysis completed successfully!

Session Information

Hide

sessionInfo()

```
## R version 4.4.2 (2024-10-31 ucrt)
## Platform: x86_64-w64-mingw32/x64
## Running under: Windows 11 x64 (build 26100)
##
## Matrix products: default
##
## locale:
## [1] LC_COLLATE=English_United States.utf8
## [2] LC_CTYPE=English_United States.utf8
## [3] LC_MONETARY=English_United States.utf8
## [4] LC_NUMERIC=C
## [5] LC_TIME=English_United States.utf8
##
## time zone: Asia/Jakarta
## tzcode source: internal
##
## attached base packages:
## [1] stats
                 graphics grDevices utils
                                               datasets methods
##
## other attached packages:
## [1] DT_0.33
                             skimr_2.1.5
                                                  naniar_1.1.0
                             pROC_1.18.5
                                                  ROSE_0.0-4
## [4] visdat_0.6.0
## [7] caret_7.0-1
                             lattice_0.22-6
                                                  randomForest_4.7-1.2
## [10] rpart.plot_3.1.2
                             rpart_4.1.24
                                                  knitr_1.48
## [13] corrplot 0.95
                             gridExtra 2.3
                                                  lubridate 1.9.4
## [16] forcats_1.0.0
                             stringr_1.5.1
                                                  dplyr_1.1.4
## [19] purrr_1.0.2
                             readr_2.1.5
                                                  tidyr_1.3.1
## [22] tibble_3.2.1
                             ggplot2_3.5.1
                                                  tidyverse 2.0.0
##
## loaded via a namespace (and not attached):
## [1] tidyselect 1.2.1
                             timeDate 4041.110
                                                  farver 2.1.2
                             digest_0.6.37
## [4] fastmap 1.2.0
                                                  timechange_0.3.0
## [7] lifecycle_1.0.4
                                                  magrittr_2.0.3
                             survival_3.7-0
## [10] compiler_4.4.2
                             rlang_1.1.4
                                                  sass_0.4.9
## [13] tools_4.4.2
                             utf8_1.2.4
                                                  yaml_2.3.10
## [16] data.table_1.16.4
                             labeling_0.4.3
                                                  htmlwidgets_1.6.4
                                                  withr_3.0.1
## [19] repr_1.1.7
                             plyr_1.8.9
## [22] nnet_7.3-19
                             grid_4.4.2
                                                  stats4_4.4.2
## [25] fansi_1.0.6
                             e1071_1.7-16
                                                  colorspace_2.1-1
## [28] future_1.34.0
                             globals_0.16.3
                                                  scales_1.3.0
## [31] iterators_1.0.14
                             MASS_7.3-65
                                                  cli_3.6.3
## [34] rmarkdown_2.29
                             generics_0.1.3
                                                  rstudioapi_0.17.1
## [37] future.apply 1.11.3 reshape2 1.4.4
                                                  tzdb_0.4.0
## [40] proxy_0.4-27
                             cachem_1.1.0
                                                  splines_4.4.2
## [43] parallel_4.4.2
                             base64enc_0.1-3
                                                  vctrs_0.6.5
                             Matrix_1.7-1
## [46] hardhat 1.4.1
                                                  jsonlite_1.8.8
## [49] hms 1.1.3
                             listenv 0.9.1
                                                  foreach 1.5.2
```

		5
## [52] gower_1.0.2	jquerylib_0.1.4	recipes_1.2.0
## [55] glue_1.7.0	parallelly_1.42.0	codetools_0.2-20
## [58] stringi_1.8.4	gtable_0.3.5	munsell_0.5.1
## [61] pillar_1.9.0	htmltools_0.5.8.1	ipred_0.9-15
## [64] lava_1.8.1	R6_2.5.1	evaluate_0.24.0
## [67] highr_0.11	bslib_0.8.0	class_7.3-22
## [70] Rcpp_1.0.13	nlme_3.1-167	prodlim_2024.06.25
## [73] xfun_0.47	pkgconfig_2.0.3	ModelMetrics_1.2.2.2