



Faculty of Computing

Group NO 06

Heart Disease Dataset Analysis

Table of Contents

1. Introduction	3
3. Using Library	3
3. Data Preprocessing/Cleaning	4
3.1 Handling Missing Data.....	4
3.2 Converting Categorical Variables	4
3.3 Splitting the Data.....	5
4. Important Insights and Visualizations	5
4.1 Visualizations of Key Attributes	5
4.2 Box Plot for Age by Heart Disease Status	7
4.3 Prediction Decision Tree	8
5. Classification Models	9
5.1 Logistic Regression.....	9
5.2 Decision Tree.....	9
5.3 Random Forest	10
5.4 Evaluate the models on the test set.....	10
5.5 Model Accuracies	11
5.5 Confusion Matrices.....	11
6. Clustering Models.....	12
6.1 K-means Clustering	12
6.2 Determining Optimal Number of Clusters.....	13
6.3 Visualizing Clusters	14
6.4 Summary Statistics for Clusters	15
6.5 Scatter plot for clustering results.....	16
7. Heart Disease Analysis UI	17
Learning Development	17
8. Conclusion.....	26
9. Contribution.....	26
10. Reference.....	27

1. Introduction

Cardiovascular disease is the leading cause of death worldwide, and early diagnosis and accurate risk identification are critical for effective treatment and prevention. "Heart Disease" is a very good dataset to identify heart disease risks based on patient characteristics. This dataset includes various types of medical data such as age, gender, type of chest pain, resting blood pressure, cholesterol level, fasting blood sugar, resting ECG results, peak heart rate, exercise-induced angina, old peak, slop, fluoroscopy-stained vessels. The parameters include thalassemia, and the target variable indicating the presence or absence of heart disease and the objective of this analysis is to effectively pre-process the data, gain significant insights and visualize and develop robust classification models for heart disease prediction. Additionally, clustering techniques will be used to uncover hidden patterns in the data set. Using these methods, we aim to improve our understanding of risk factors associated with heart disease and improve the accuracy of predictions, as well as ultimately contribute to better health outcomes for patients.

3. Using Library

Library	Description
readr	readr is used to read and write data in R. It provides faster and more user-friendly functions for data import, offering support for CSV, TSV, and other formats through functions like <code>read_csv()</code> and <code>write_csv()</code> .
ggplot2	ggplot2 is a powerful visualization package based on the grammar of graphics. It allows for the creation of complex and aesthetic graphics by layering components such as points, lines, and text.
rpart.plot	rpart.plot is an extension of the rpart package used for recursive partitioning and regression trees. It provides functions to visualize decision trees created by rpart, making them easier to interpret.
caret	The caret package (short for Classification and Regression Training) is a comprehensive framework for building predictive models in R. It offers a unified interface to various machine learning algorithms, streamlining the process of model training, tuning, and evaluation.
pROC	pROC is designed for visualizing and analyzing receiver operating characteristic (ROC) curves. It is used for binary classification tasks to evaluate model performance and determine optimal classification thresholds.
tidyverse	tidyverse is a collection of R packages designed for data science, including ggplot2, dplyr, tidyr, and others that share a common design philosophy. It allows for seamless integration of data manipulation, visualization, and analysis functions.
cluster	The cluster package provides functions for cluster analysis, which involves grouping similar data points together. It includes various clustering algorithms and tools to visualize clustering results, useful for exploratory data analysis and pattern recognition.
factoextra	factoextra facilitates the visualization and interpretation of multivariate data analysis results. It is used to visualize clustering results, principal component analysis (PCA), and other multivariate techniques, enhancing the understanding of complex data sets.

Cleaning for data preprocessing is an essential step in preparing the data set for analysis and modeling.

First, it is important to identify missing values in the data set. Depending on the amount of missing data, different strategies can be applied. If the number of missing values by deletion is minimal, the affected rows or columns can be discarded.

```
26 # Identity missing values
27 missing <- colSums(is.na(data))
28 print(missing)
29
30 # Remove rows with missing values
31 data_clean <- na.omit(data)
32 print(data_clean)
```

```
> # Identify missing values
> missing <- colSums(is.na(data))
> print(missing)
```

	age	sex	chest_pain_type
resting_blood_pressure	0	0	0
cholesterol	0	0	0
fasting_blood_sugar	0	0	0
rest_ecg	0	0	0
Max_heart_rate	0	0	0
exercise_induced_angina	0	0	0
oldpeak	0	0	0
slope	0	0	0
vessels_colored_by_fluoroscopy	0	0	0
thalassemia	0	0	0
target	0	0	0

```
> # Remove rows with missing values
> data_clean <- na.omit(data)
> print(data_clean)
```

	age	sex	chest_pain_type	resting_blood_pressure	cholesterol	fasting_blood_sugar		
1	52	Male	Typical	angina	125	212	Lower	than 120 mg/ml
2	53	Male	Typical	angina	140	203	Greater	than 120 mg/ml
3	70	Male	Typical	angina	145	174	Lower	than 120 mg/ml
4	61	Male	Typical	angina	148	203	Lower	than 120 mg/ml
5	62	Female	Typical	angina	138	294	Greater	than 120 mg/ml
6	58	Female	Typical	angina	100	248	Lower	than 120 mg/ml
7	58	Male	Typical	angina	114	318	Lower	than 120 mg/ml
8	55	Male	Typical	angina	160	289	Lower	than 120 mg/ml
9	45	Male	Typical	angina	130	240	Lower	than 120 mg/ml

To ensure that our models correctly handle categorical data, we convert relevant columns (gender, type of chest pain, and others) into factors. This is important for correct model interpretation and accuracy.

```
34 # Preprocessing the data
35 # Convert categorical variables to factors
36 data$ssex <- as.factor(data$ssex)
37 data$chest_pain_type <- as.factor(data$chest_pain_type)
38 data$fasting_blood_sugar <- as.factor(data$fasting_blood_sugar)
39 data$rest_ecg <- as.factor(data$rest_ecg)
40 data$exercise_induced_angina <- as.factor(data$exercise_induced_angina)
41 data$slope <- as.factor(data$slope)
42 data$vessels_colored_by_fluoroscopy <- as.factor(data$vessels_colored_by_fluoroscopy)
43 data$thalassemia <- as.factor(data$thalassemia)
44 data$target <- as.factor(data$target)
```

```
> # Preprocessing the data
> # Convert categorical variables to factors
> data$sex <- as.factor(data$sex)
> data$chest_pain_type <- as.factor(data$chest_pain_type)
> data$fasting_blood_sugar <- as.factor(data$fasting_blood_sugar)
> data$rest_ecg <- as.factor(data$rest_ecg)
> data$exercise_induced_angina <- as.factor(data$exercise_induced_angina)
> data$slope <- as.factor(data$slope)
> data$vessels_colored_by_flouoropsy <- as.factor(data$vessels_colored_by_flouoropsy)
> data$thalassemia <- as.factor(data$thalassemia)
> data$target <- as.factor(data$target)
```

3.3 Splitting the Data

We split the dataset into training and testing sets, 80% of the data is used for training and 20% is reserved for testing. This partitioning allows us to evaluate the performance of our models on unseen data as well as their generalization.

```
# Split the dataset into training
set.seed(123)
train_index <- createDataPartition(data$target, p = 0.8, list = FALSE)
train_data <- data[train_index, ]
test_data <- data[-train_index, ]

# Display the dimensions of the training and testing sets
dim(train_data)
dim(test_data)
```

```
> set.seed(123)
> train_index <- createDataPartition(data$target, p = 0.8, list = FALSE)
> train_data <- data[train_index, ]
> test_data <- data[-train_index, ]
> # Display the dimensions of the training and testing sets
> dim(train_data)
[1] 821 14
> dim(test_data)
[1] 204 14
```

4. Important Insights and Visualizations

4.1 Visualizations of Key Attributes

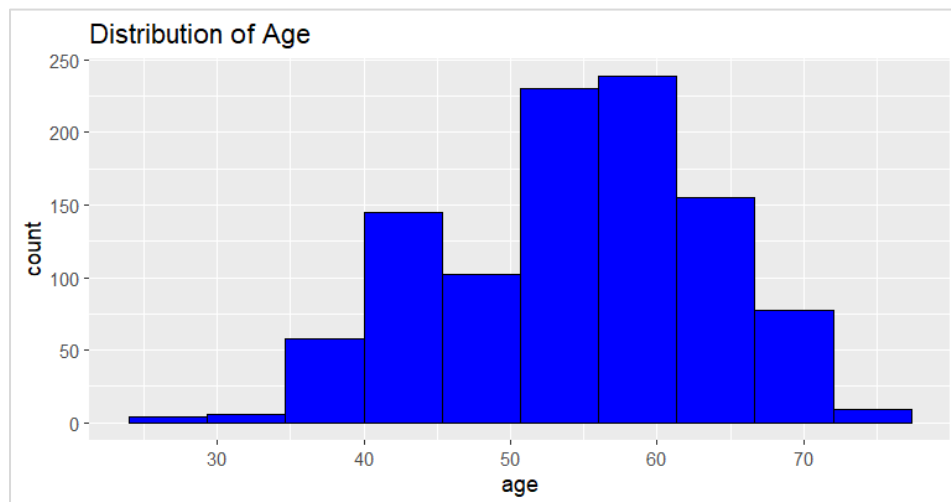
Visualizing the key properties of the data set, as well as helping to understand the distribution, relationships and potential patterns within the data.

Age Distribution

Histograms can be used to visualize the age distribution of the data set, helping us understand the age range and possible age-related trends.

```
age_plot <- ggplot(data_clean, aes(x = age)) +
  geom_histogram(bins = 10, fill = 'blue', color = 'black') +
  ggtitle('Distribution of Age')
print(age_plot)
```

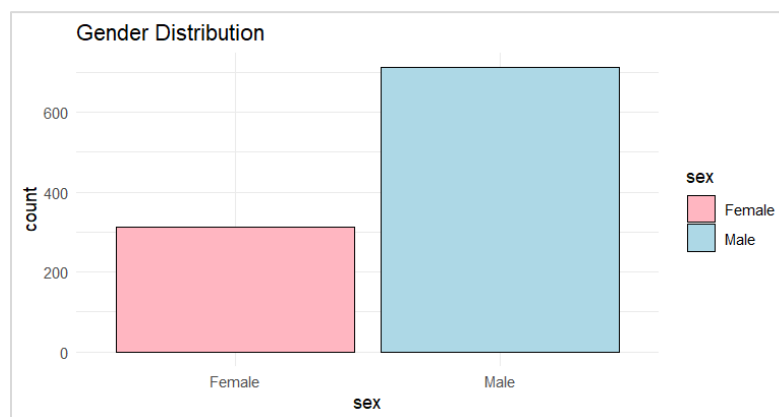
```
> age_plot <- ggplot(data_clean, aes(x = age)) +
+   geom_histogram(bins = 10, fill = 'blue', color = 'black') +
+   ggtitle('Distribution of Age')
> print(age_plot)
```



Gender Distribution

Gender distribution analysis visualizes the ratio of males to females in relation to heart disease status using a bar graph. Males and females are also colored differently, with light blue representing males and light pink representing females. Columns were grouped by gender and heart disease status, with men generally having a higher risk of heart disease than women. The graph shows that the data set includes both genders, with men showing a higher prevalence of heart disease, and should be targeted in male-specific health care strategies.

```
gender_plot <- ggplot(data_clean, aes(x = sex, fill = sex)) +  
  geom_bar(color = 'black') +  
  ggtitle('Gender Distribution') +  
  scale_fill_manual(values = c("Female" = "lightpink", "Male" = "lightblue")) +  
  theme_minimal()  
print(gender_plot)  
  
> gender_plot <- ggplot(data_clean, aes(x = sex, fill = sex)) +  
  + geom_bar(color = 'black') +  
  + ggtitle('Gender Distribution') +  
  + scale_fill_manual(values = c("Female" = "lightpink", "Male" = "lightblue")) +  
  + theme_minimal()  
> print(gender_plot)
```

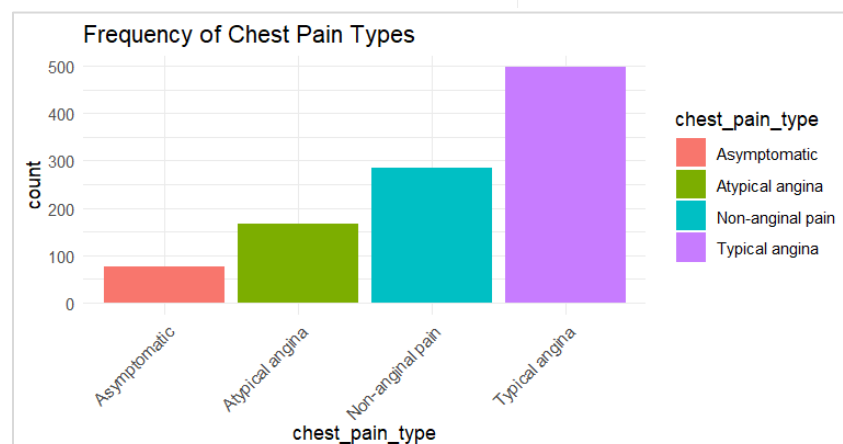


Chest Pain Type

Frequency

We used this bar chart to show the frequency of different chest pains. This can provide insight into the prevalence of different types of chest pain among patients.

```
cp_plot <- ggplot(data_clean, aes(x = chest_pain_type, fill = chest_pain_type)) +  
  geom_bar() +  
  ggtitle('Frequency of Chest Pain Types') +  
  theme_minimal() +  
  theme(axis.text.x = element_text(angle = 45, hjust = 1))  
print(cp_plot)  
  
> cp_plot <- ggplot(data_clean, aes(x = chest_pain_type, fill = chest_pain_type)) +  
  + geom_bar() +  
  + ggtitle('Frequency of Chest Pain Types') +  
  + theme_minimal() +  
  + theme(axis.text.x = element_text(angle = 45, hjust = 1))  
> print(cp_plot)
```

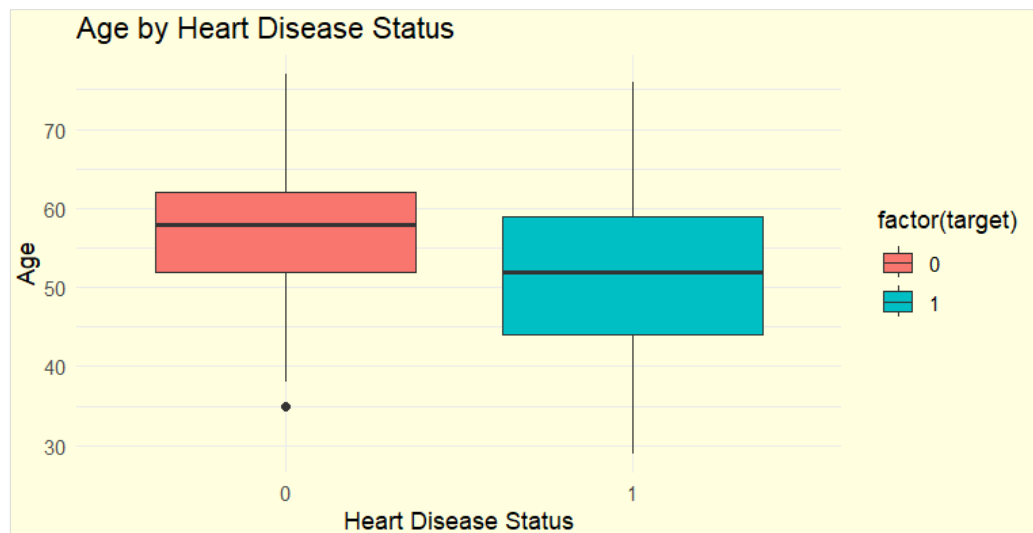


4.2 Box Plot for Age by Heart Disease Status

Box plots visualize the relationship between age and heart disease status, providing insight into how age varies across different heart disease categories. The plot compares the age distribution for people with and without heart disease. The x-axis represents heart disease status, with 0 indicating no heart disease and 1 indicating heart disease. The y-axis shows age in years. There are two box plots, one for each heart condition. The box area for people without heart disease (0) has a median age of about 60 years with an interquartile range of 55 to 65 years and shows a relatively even age distribution within this group. Boxes for people with heart disease (1) The median age in the plot is about 55 years, with an interquartile range of 45 to 60 years. The box for people without heart disease is slightly skewed to the right, suggesting that there may be a few older people with heart disease, and the presence of a single outlier in the "no heart disease" group warrants further investigation. Suggesting a link, the risk of heart disease increases with age.

```
# Box Plot for Age |
ggplot(data_clean, aes(x = factor(target), y = age, fill = factor(target))) +
  geom_boxplot() +
  labs(title = 'Age by Heart Disease Status', x = 'Heart Disease Status', y = 'Age') +
  theme_minimal() +
  theme(plot.background = element_rect(fill = 'lightyellow'))

> ggplot(data_clean, aes(x = factor(target), y = age, fill = factor(target))) +
+   geom_boxplot() +
+   labs(title = 'Age by Heart Disease Status', x = 'Heart Disease Status', y = 'Age') +
+   theme_minimal() +
+   theme(plot.background = element_rect(fill = 'lightyellow'))
```



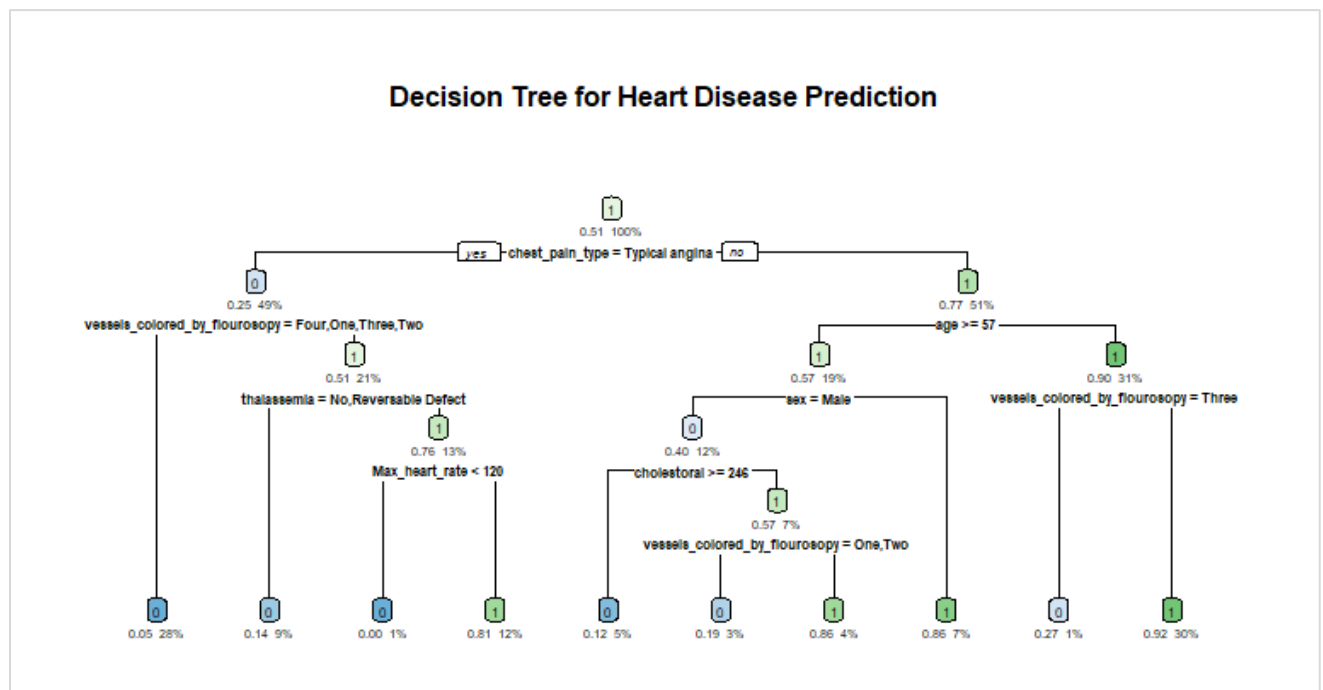
4.3 Prediction Decision Tree

We used this decision tree to estimate the probability of heart disease based on patient attributes. It works by making a series of decisions based on the values of these attributes, dividing the data into smaller subsets starting from the root node and moving down the tree. The bottom most nodes, called leaf nodes, represent the final predictor of either **0 (no heart disease)** or **1 (heart disease)**, with numbers within leaf nodes showing the percentage of cases in that subset belonging to each class. The decision tree algorithm selects the best attribute to split the data at each node, aiming to create as clean subsets as possible. Key attributes in this decision tree include type of *chest pain*, *age*, *thalassemia*, *sex*, *cholesterol*, *peak heart rate* and number of vessels stained by fluoroscopy.

However, it is important to note that this decision tree is based on a specific data set and cannot be generalized to other populations, so it is essential to evaluate the performance of the model using appropriate metrics and consider other factors before making any clinical decisions.

```
#Prediction Decision Tree
fit <- rpart(target ~ ., data = train_data, method = 'class')
rpart.plot(fit, main = 'Decision Tree for Heart Disease Prediction',
           , extra = 106, under = TRUE, faclen = 0)

fit <- rpart(target ~ ., data = train_data, method = 'class')
rpart.plot(fit, main = 'Decision Tree for Heart Disease Prediction',
           , extra = 106, under = TRUE, faclen = 0)
```



5. Classification Models

5.1 Logistic Regression

We used this logistic regression as a binary outcome prediction algorithm in patients. In the context of predicting heart disease, it helps to estimate the probability of heart attack based on patient characteristics. The 'train()' function in the 'caret' package in R is used to train the logistic regression model. The formula is 'target ~ .'. The target variable 'target' must be predicted using all other variables in the data set. The 'method' parameter is set to 'glm' to indicate that we are using generalized linear models, and 'family' parameter is set to 'binomial'. The dataset 'train_data' is used to train the model. Once the model is trained, it can be used to make predictions on new data by estimating the probability of a patient having heart disease. The predicted probabilities can be constrained to classify patients as having or not having heart disease based on a chosen cutoff value.

```
# Logistic Regression
logistic_model <- train(target ~ ., data = train_data,
                        method = "glm", family = "binomial")
logistic_model <- train(target ~ ., data = train_data,
                        method = "glm", family = "binomial")
```

5.2 Decision Tree

Iterative partitioning of the data based on the feature that provides the most informative benefit creates a tree-like decision structure. In the context of cardiac disease prediction, a decision tree model can be trained using the 'train()' function from the 'caret' package in R. The formula 'target ~ .' should predict the target variable 'target'. Using all other variables in the dataset The 'method' parameter is set to 'rpart' to indicate that we are using the recursive partitioning algorithm. The dataset 'train_data' is used to train the model. Once the model is trained, it can be used to make predictions about new data by traversing the decision tree based on the feature values of the new instance. Predictions are made based on the leaf node reached at the end of the decision tree path. Decision trees have the advantage of being easy to interpret and visualize, making them useful for understanding relationships between patient attributes and heart disease risk.

```
# Decision Tree
decision_tree_model <- train(target ~ .,
                             data = train_data,
                             method = "rpart")
# Decision Tree
decision_tree_model <- train(target ~ .,
                             data = train_data,
                             method = "rpart")
```

5.3 Random Forest

We use the random forest model to improve the accuracy and robustness of predictions, as well as due to a popular ensemble learning method that combines multiple decision trees. A model can be trained. The formula 'target ~ .' can predict the target variable 'target'. The 'method' parameter is set to 'rf' to indicate that we are using the random forest algorithm. The 'train_data' dataset is used to train the model. Once the model is trained, it can be used to make predictions on new data by averaging the predictions from each individual decision tree in the forest. Random forests can handle high-dimensional data, capture complex interactions between features, and perform robustly even when the dataset contains noisy or irrelevant features.

```
# Decision Tree
decision_tree_model <- train(target ~ .,
                             data = train_data,
                             method = "rpart")

# Decision Tree
decision_tree_model <- train(target ~ .,
                             data = train_data,
                             method = "rpart")
```

5.4 Evaluate the models on the test set

To evaluate the performance of logistic regression, decision tree and random forest models on the test set, we use the 'predict()' function from the 'caret' package in R. This function takes the trained model and test data and returns the predicted values. Stores the predicted class labels (0 or 1) for each patient. Random Forest The variable 'random_forest_pred' stores the predicted class labels (0 or 1) for each patient in the test set. These predicted values are then compared with the actual target values of the test set to calculate performance metrics such as precision, recall, and F1-score. This evaluation helps assess the ability of models to generalize to new, unseen data and provides insight into their strengths and weaknesses.

```
# Evaluate the models on the test set
logistic_pred <- predict(logistic_model,
                         newdata = test_data)
decision_tree_pred <- predict(decision_tree_model,
                             newdata = test_data)
random_forest_pred <- predict(random_forest_model,
                             newdata = test_data)

# Evaluate the models on the test set
logistic_pred <- predict(logistic_model,
                         newdata = test_data)
decision_tree_pred <- predict(decision_tree_model,
                             newdata = test_data)
random_forest_pred <- predict(random_forest_model,
                             newdata = test_data)
```

5.5 Model Accuracies

To evaluate the performance of the logistic regression, decision tree and random forest models on the test set, the 'confusionMatrix()' function in the 'caret' package in R is used to calculate various performance metrics based on the predicted values. The 'confusionMatrix()' function returns a list containing metrics such as precision, accuracy, recall and F1-score of the actual target values of the test set. Accuracy values for each model are extracted and stored in a data frame for comparison. The results show that the logistic regression model has an accuracy of 81.63%, the decision tree model has an accuracy of 75.51%, and the random forest model has an accuracy of 83.67% on the test set.

```
# Model Accuracies
logistic_confusion <- confusionMatrix(logistic_pred, test_data$target)
decision_tree_confusion <- confusionMatrix(decision_tree_pred,
                                           test_data$target)
random_forest_confusion <- confusionMatrix(random_forest_pred,
                                           test_data$target)

model_accuracies <- data.frame(
  Model = c("Logistic Regression", "Decision Tree", "Random Forest"),
  Accuracy = c(logistic_confusion$overall["Accuracy"],
               decision_tree_confusion$overall["Accuracy"],
               random_forest_confusion$overall["Accuracy"])
)
print(model_accuracies)
```

```
# Model Accuracies
logistic_confusion <- confusionMatrix(logistic_pred, test_data$target)
decision_tree_confusion <- confusionMatrix(decision_tree_pred,
                                           test_data$target)
random_forest_confusion <- confusionMatrix(random_forest_pred,
                                           test_data$target)

model_accuracies <- data.frame(
  Model = c("Logistic Regression", "Decision Tree", "Random Forest"),
  Accuracy = c(logistic_confusion$overall["Accuracy"],
               decision_tree_confusion$overall["Accuracy"],
               random_forest_confusion$overall["Accuracy"])
)
print(model_accuracies)
```

Model	Accuracy
Logistic Regression	0.8676471
Decision Tree	0.7990196
Random Forest	1.0000000

5.5 Confusion Matrices

To show the performance of logistic regression, decision tree and random forest models on the test set, the 'confusionMatrix()' function in the 'caret' package in R is used to calculate the confusion matrix for logistic regression, decision tree, and random forest model heatmaps are shown. The 'ggplot2' package is used to create plots that show the frequency of each combination of actual and predicted values. Confusion matrices are converted to data frames using 'as.data.frame()', and model names are added to each data frame. The data frames are combined into a single data frame using 'rbind()', and the combined data frame is used to create a heatmap using 'ggplot()'. The 'geom_tile()' function is used to create the heatmap, and the 'facet_wrap()' function is used to separate the plots according to the model. The color scale is set to range from white to blue using 'scale_fill_gradient()' and the plot is labeled "Confusion Matrices for Classification Models" and the axes are labeled "Actual" and "Predicted". The plot is styled using the 'theme_minimal()' function and the 'axis.text.x' element is styled to display the x-axis labels at a 45 degree angle. The resulting plot shows the confusion matrix for each model, with actual values on the x-axis and predicted values on the y-axis, with the frequency of each combination of actual and predicted values represented by the color in the heatmap, dark blue indicating higher frequencies.

```
# Confusion Matrices
logistic_cm <- as.data.frame(logistic_confusion$table)
decision_tree_cm <- as.data.frame(decision_tree_confusion$table)
random_forest_cm <- as.data.frame(random_forest_confusion$table)

logistic_cm$Model <- "Logistic Regression"
decision_tree_cm$Model <- "Decision Tree"
random_forest_cm$Model <- "Random Forest"

combined_cm <- rbind(logistic_cm, decision_tree_cm,
                    random_forest_cm)

ggplot(combined_cm, aes(x = Reference, y = Prediction,
                       fill = Freq)) +
  geom_tile(color = "white") +
  facet_wrap(~ Model) +
  scale_fill_gradient(low = "white", high = "blue") +
  labs(title = "Confusion Matrices for Classification Models",
       x = "Actual", y = "Predicted") +
  theme_minimal() +
  theme(axis.text.x = element_text(angle = 45,
                                   hjust = 1))
```

```
# Confusion Matrices
logistic_cm <- as.data.frame(logistic_confusion$table)
decision_tree_cm <- as.data.frame(decision_tree_confusion$table)
random_forest_cm <- as.data.frame(random_forest_confusion$table)
logistic_cm$Model <- "Logistic Regression"
decision_tree_cm$Model <- "Decision Tree"
random_forest_cm$Model <- "Random Forest"
combined_cm <- rbind(logistic_cm, decision_tree_cm,
                    random_forest_cm)
ggplot(combined_cm, aes(x = Reference, y = Prediction,
                       fill = Freq)) +
  geom_tile(color = "white") +
  facet_wrap(~ Model) +
  scale_fill_gradient(low = "white", high = "blue") +
  labs(title = "Confusion Matrices for Classification Models",
       x = "Actual", y = "Predicted") +
  theme_minimal() +
  theme(axis.text.x = element_text(angle = 45,
                                   hjust = 1))
```



6. Clustering Models

6.1 K-means Clustering

To perform k-means clustering on our data set, we used the `set.seed(123)` function to ensure reproducibility of results by fitting the random number generator to a specific seed. The `'apply()'` function is used to calculate the WSS for We used the `'kmeans()'` function to perform the clustering and calculate the WSS of the `'tot.withinss'` component of the output. The `'select_if(is.numeric)'` function in the `'dplyr'` package from the `'data_clean'` dataset by selecting numeric data. Used to select only numeric columns. Scaling the data The `'scale()'` function is used to scale the data to have zero mean and unit variance. This ensures that all features are on the same scale and improves the performance of the clustering algorithm. The goal of this code is to find the optimal number of clusters (`'k'`) for the dataset by minimizing the WSS. WSS is a measure of the overall variance within the clusters and its minimization helps to identify the best clustering solution. The output of this code is the scale data, which is standardized as the mean 0 and standard deviation of 1 for each feature and thereby clusters all the features. contribute equally to the process. This helped us to select the appropriate number of clusters for this data set.

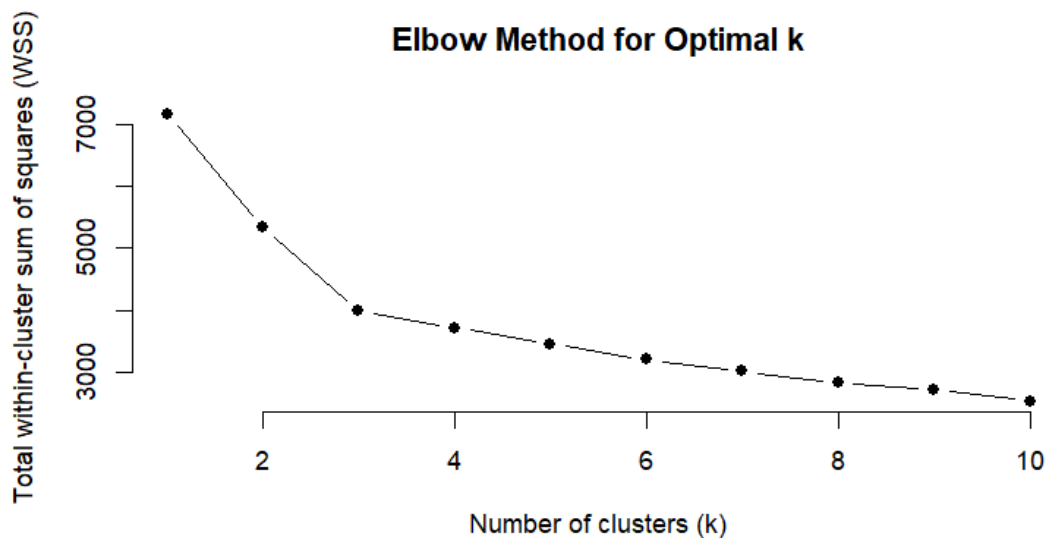
```
# K-means Clustering
set.seed(123)
wss <- sapply(1:10, function(k) {kmeans(scaled_data,
                                         centers = k,
                                         nstart = 25)$tot.withinss})
numeric_data <- data_clean %>% select_if(is.numeric)
scaled_data <- scale(numeric_data)
[[reached getOption("max.print") -- omitted 883 rows ]]
attr(,"scaled:center")
  age resting_blood_pressure cholesterol1 Max_heart_rate oldpeak
54.4341463 131.6117073 246.0000000 149.1141463 1.0715122
target cluster
0.5131707 1.9268293
attr(,"scaled:scale")
  age resting_blood_pressure cholesterol1 Max_heart_rate oldpeak
9.0722902 17.5167180 51.5925102 23.0057237 1.1750533
target cluster
0.5000705 0.7431406
```

6.2 Determining Optimal Number of Clusters

The Elbow method plot helps determine the optimal number of clusters (k) for a data set. It shows the relationship between the number of clusters and the within-cluster sum (WSS). WSS measures the compactness of clusters, with lower values indicating tighter clusters. As the number of clusters increases, the WSS generally decreases. The "elbow" in the plot means that adding more clusters does not significantly decrease the WSS. The optimal number of clusters may be around 3, where the rate of decrease in WSS starts to level off and forms an "elbow" shape.

```
# Determine the optimal number of clusters using the Elbow method
optimal_k <- 3
kmeans_result <- kmeans(scaled_data, centers = optimal_k, nstart = 25)
plot(1:10, wss, type = "b", pch = 19, frame = FALSE,
     xlab = "Number of clusters (k)",
     ylab = "Total within-cluster sum of squares (WSS)",
     main = "Elbow Method for Optimal k")

# Determine the optimal number of clusters using the Elbow method
optimal_k <- 3
kmeans_result <- kmeans(scaled_data, centers = optimal_k, nstart = 25)
plot(1:10, wss, type = "b", pch = 19, frame = FALSE,
     xlab = "Number of clusters (k)",
     ylab = "Total within-cluster sum of squares (WSS)",
     main = "Elbow Method for Optimal k")
```

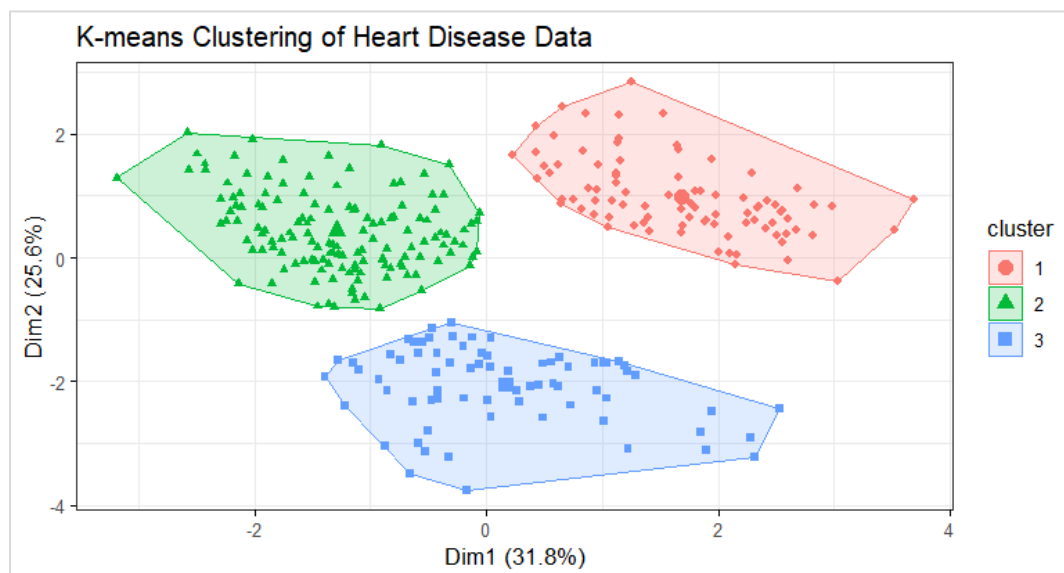


6.3 Visualizing Clusters

The X-axis represents the first principal component Dim1 (31.8%) which accounts for 31.8% of the variance in the data, while the Y-axis Dim2 (25.6%) shows us the second principal component which accounts for 25.6% of the variance in the data. The plot has three distinct clusters represented by different colors and shapes. That is Cluster 1 (red) places the cluster in the upper right quadrant of the plot. Cluster 2 (green) places the cluster in the lower left quadrant of the plot. Cluster 3 (blue) places the cluster in the lower middle portion of the plot. From this we can determine that the clusters appear to be reasonably well separated, and we can conclude that the K-means algorithm has successfully identified distinct groups in the data. The cluster shape exhibits some degree of circularity, indicating that the variables used in cluster analysis are normally distributed. The density of data points within each cluster varies depending on the data density, with cluster 1 (red) being the most populated. This K-means cluster visualization thus suggests that cardiac data can be divided into three meaningful groups based on their characteristics.

```
# Visualizing Clusters
cluster_plot <- fviz_cluster(kmeans_result, data = scaled_data,
                             geom = "point",
                             ellipse.type = "convex",
                             ggtheme = theme_bw()) +
  labs(title = "K-means Clustering of Heart Disease Data")
print(cluster_plot)
```

```
# Visualizing Clusters
cluster_plot <- fviz_cluster(kmeans_result, data = scaled_data,
                             geom = "point",
                             ellipse.type = "convex",
                             ggtheme = theme_bw()) +
  labs(title = "K-means Clustering of Heart Disease Data")
print(cluster_plot)
```



6.4 Summary Statistics for Clusters

Age, cholesterol, resting blood pressure, and maximum heart rate were calculated for each cluster identified in the heart disease dataset. These results show distinct characteristics for each cluster. Cluster 1 had relatively lower mean age, cholesterol, and maximum heart rate compared to the other clusters. Cluster 2 shows the youngest average age and the highest average maximum heart rate. Cluster 3 shows the highest mean values for age, cholesterol and resting blood pressure. These summary statistics give us insight into the differences between the patient groups represented by each cluster.

```
# Summary Statistics for Clusters
cluster_summary <- aggregate(data_clean[, c("age", "cholesterol",
                                             "resting_blood_pressure",
                                             "Max_heart_rate")],
                             by = list(Cluster = data_clean$cluster),
                             FUN = mean)

print(cluster_summary)
```

```
# Summary Statistics for Clusters
cluster_summary <- aggregate(data_clean[, c("age", "cholesterol",
                                             "resting_blood_pressure",
                                             "Max_heart_rate")],
                             by = list(Cluster = data_clean$cluster),
                             FUN = mean)

print(cluster_summary)
```

Cluster	age	cholesterol	resting_blood_pressure	Max_heart_rate
1	56.81115	235.5975	126.9195	130.8576
2	48.61674	230.7665	125.6608	163.0969
3	61.98790	287.4355	148.6169	147.2944

```
# Count the number of observations in each cluster
cluster_counts <- table(data_clean$cluster)
print(cluster_counts)
```

```
> # Count the number of observations in each cluster
> cluster_counts <- table(data_clean$cluster)
> print(cluster_counts)
```

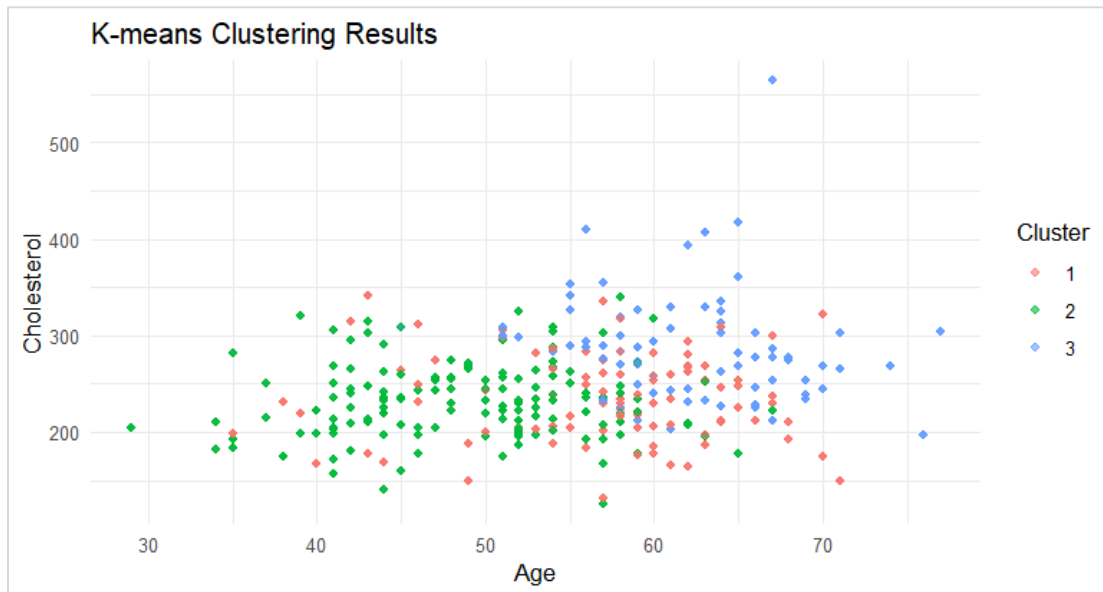
1	2	3
323	454	248

6.5 Scatter plot for clustering results

It appears that individuals in this K-means clustering dataset can be classified into three distinct groups based on their age and cholesterol level. Cluster 1 represents a young population with low cholesterol levels, Cluster 2 a mixed group with a wide range of values, and Cluster 3 an older population with high cholesterol levels. This has the potential to provide valuable insight into potential risk factors for cluster heart disease.

```
# Scatter plot for clustering results
ggplot(data_clean, aes(x = age, y = cholesterol,
                      color = factor(cluster))) +
  geom_point(alpha = 0.6) +
  scale_color_discrete(name = "Cluster") +
  labs(title = "K-means Clustering Results",
       x = "Age",
       y = "Cholesterol") +
  theme_minimal()
```

```
# Scatter plot for clustering results
ggplot(data_clean, aes(x = age, y = cholesterol,
                      color = factor(cluster))) +
  geom_point(alpha = 0.6) +
  scale_color_discrete(name = "Cluster") +
  labs(title = "K-means Clustering Results",
       x = "Age",
       y = "Cholesterol") +
  theme_minimal()
```



7. Heart Disease Analysis UI

Learning Development

We developed this UI over several days to learn more about shiny development and improve our knowledge. We have created a dashboard using two libraries shiny and shiny dashboard. Through this project, we have deepened our understanding of predictive modeling and data visualization, all with enthusiasm and a sense of commitment, often creating our best for remote knowledge. And we have worked hard on this project, gaining valuable insights and skills.

The basic functions of application

The user has the ability to view data summaries and key summaries to understand the structure and content of the cardiology dataset. This allows users to assess each model's accuracy and potential bias.

```
1 library(shiny)
2 library(shinydashboard)
3 library(ggplot2)
4 library(caret)
5 library(rpart)
6 library(rpart.plot)
7 library(dplyr)
8 library(cluster)
9 library(factoextra)
10
11
12 data <- read.csv("C:\\Users\\shard\\Desktop\\Ishara\\MY 2ed year\\End exam\\R\\Final project\\HeartDiseaseTrain-Test.csv")
13
14
15 data$sex <- as.factor(data$sex)
16 data$chest_pain_type <- as.factor(data$chest_pain_type)
17 data$fasting_blood_sugar <- as.factor(data$fasting_blood_sugar)
18 data$rest_ecg <- as.factor(data$rest_ecg)
19 data$exercise_induced_angina <- as.factor(data$exercise_induced_angina)
20 data$slope <- as.factor(data$slope)
21 data$vessels_colored_by_fluoroscopy <- as.factor(data$vessels_colored_by_fluoroscopy)
22 data$thalassemia <- as.factor(data$thalassemia)
23 data$target <- as.factor(data$target)
24
25
26 set.seed(123)
27 train_index <- createDataPartition(data$target, p = 0.8, list = FALSE)
28 train_data <- data[train_index, ]
29 test_data <- data[-train_index, ]
30
31
32 logistic_model <- train(target ~ ., data = train_data, method = "glm", family = "binomial")
33 decision_tree_model <- train(target ~ ., data = train_data, method = "rpart")
34 random_forest_model <- train(target ~ ., data = train_data, method = "rf")
35
36
37 logistic_pred <- predict(logistic_model, newdata = test_data)
38 decision_tree_pred <- predict(decision_tree_model, newdata = test_data)
39 random_forest_pred <- predict(random_forest_model, newdata = test_data)
40
41 logistic_confusion <- confusionMatrix(logistic_pred, test_data$target)
42 decision_tree_confusion <- confusionMatrix(decision_tree_pred, test_data$target)
43 random_forest_confusion <- confusionMatrix(random_forest_pred, test_data$target)
44
45 model_accuracies <- data.frame(
46   Model = c("Logistic Regression", "Decision Tree", "Random Forest"),
47   Accuracy = c(logistic_confusion$overall["Accuracy"], decision_tree_confusion$overall["Accuracy"],
48               random_forest_confusion$overall["Accuracy"])
49 )
50
51
52 numeric_data <- data %>% select_if(is.numeric)
53 scaled_data <- scale(numeric_data)
54 set.seed(123)
55 optimal_k <- 3
56 kmeans_result <- kmeans(scaled_data, centers = optimal_k, nstart = 25)
57
58 data$cluster <- kmeans_result$cluster
```

```

60 ui <- dashboardPage(
61   dashboardHeader(title = "Heart Disease Analysis"),
62   dashboardSidebar(
63     sidebarMenu(
64       menuItem("Data Summary", tabName = "summary", icon = icon("table")),
65       menuItem("Model Accuracy", tabName = "accuracy", icon = icon("chart-bar")),
66       menuItem("Confusion Matrices", tabName = "confusion", icon = icon("th")),
67       menuItem("Clustering", tabName = "clustering", icon = icon("project-diagram")),
68       menuItem("Visualizations", tabName = "visualizations", icon = icon("chart-line")),
69       menuItem("Prediction Tree", tabName = "prediction_tree", icon = icon("tree"))
70     )
71   ),
72   dashboardBody(
73     tabItems(
74       tabItem(tabName = "summary",
75         h2("Data Summary"),
76         verbatimTextOutput("dataSummary"),
77         verbatimTextOutput("dataHead")
78       ),
79       tabItem(tabName = "accuracy",
80         h2("Model Accuracy"),
81         plotOutput("accuracyPlot")
82       ),
83       tabItem(tabName = "confusion",
84         h2("Confusion Matrices"),
85         plotOutput("confusionMatrixPlot")
86       ),
87       tabItem(tabName = "clustering",
88         h2("Clustering"),
89         plotOutput("clusterPlot"),
90         plotOutput("elbowPlot"),
91         verbatimTextOutput("clustersSummary"),
92         verbatimTextOutput("clusterCounts")
93       ),
94       tabItem(tabName = "visualizations",
95         h2("Visualizations"),
96         plotOutput("agePlot"),
97         plotOutput("genderPlot"),
98         plotOutput("cpPlot"),
99         plotOutput("ageBoxPlot")
100      ),
101      tabItem(tabName = "prediction_tree",
102        h2("Prediction Tree"),
103        plotOutput("predictionTreePlot"),
104        verbatimTextOutput("predictionTreesSummary")
105      )
106    )
107  )
108 )

```

```

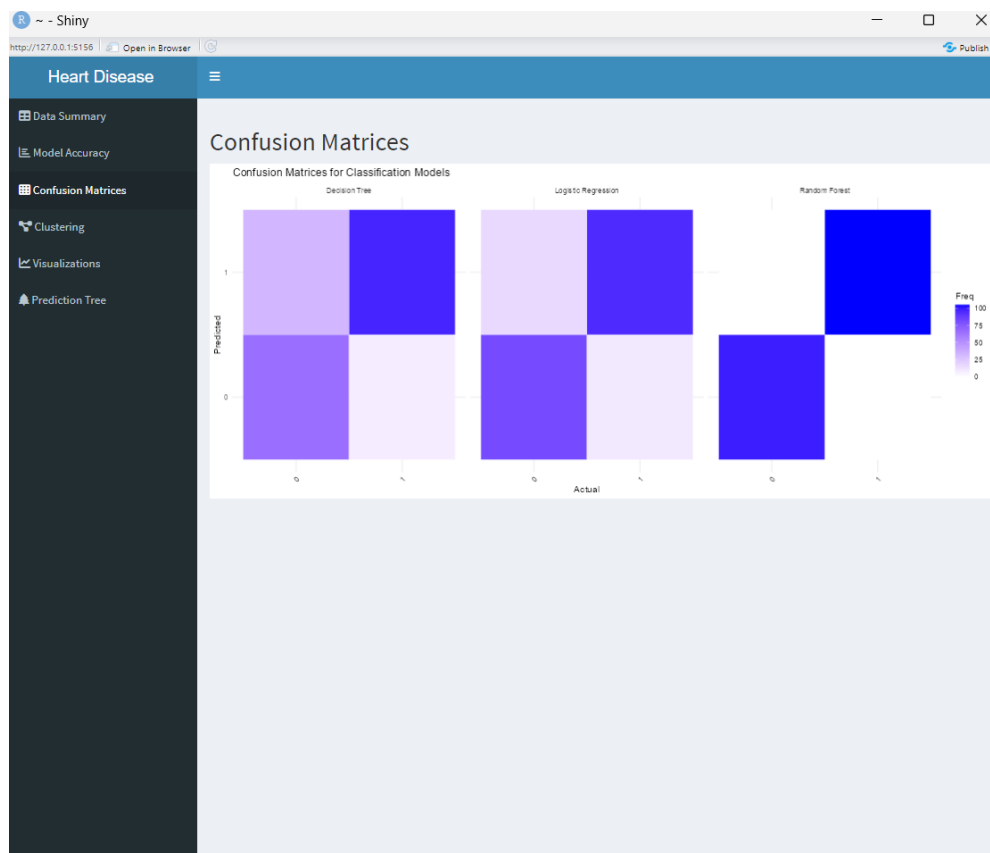
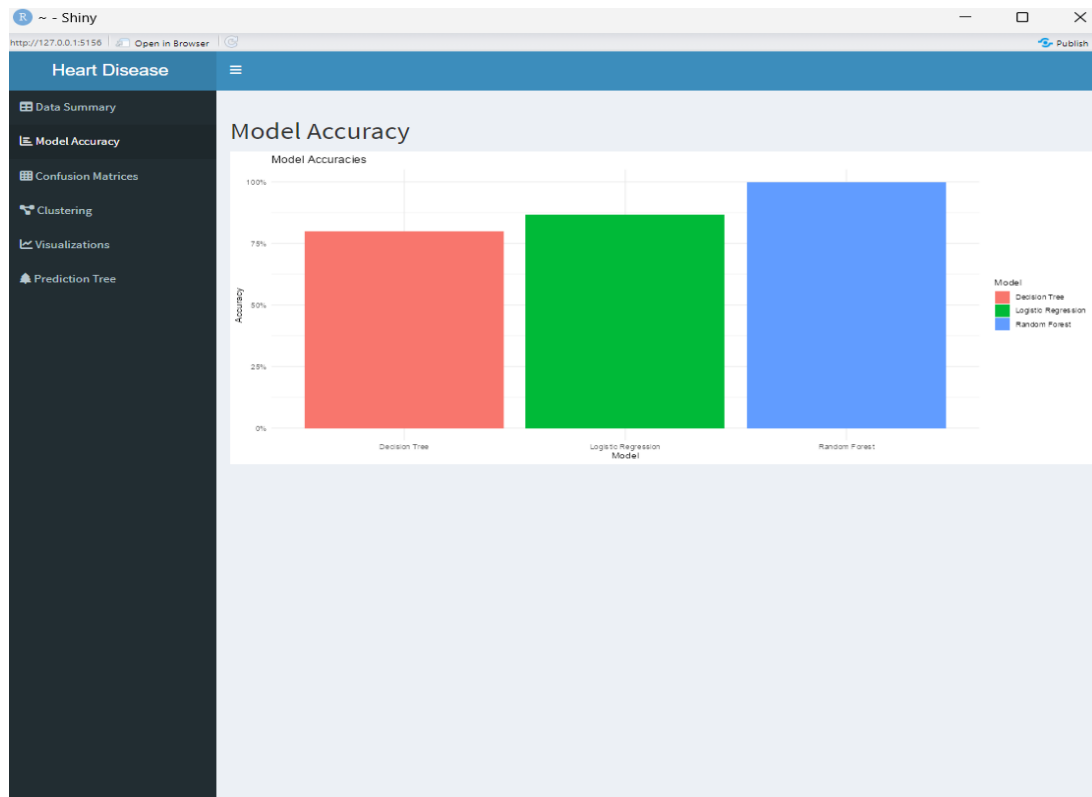
111 server <- function(input, output) {
112   output$dataSummary <- renderPrint({
113     str(data)
114   })
115
116   output$dataHead <- renderPrint({
117     head(data)
118   })
119
120   output$accuracyPlot <- renderPlot({
121     ggplot(model_accuracies, aes(x = Model, y = Accuracy, fill = Model)) +
122       geom_bar(stat = "identity") +
123       theme_minimal() +
124       labs(title = "Model Accuracies", x = "Model", y = "Accuracy") +
125       scale_y_continuous(labels = scales::percent)
126   })
127
128   output$confusionMatrixPlot <- renderPlot({
129     logistic_cm <- as.data.frame(logistic_confusion$table)
130     decision_tree_cm <- as.data.frame(decision_tree_confusion$table)
131     random_forest_cm <- as.data.frame(random_forest_confusion$table)
132     logistic_cm$Model <- "Logistic Regression"
133     decision_tree_cm$Model <- "Decision Tree"
134     random_forest_cm$Model <- "Random Forest"
135     combined_cm <- rbind(logistic_cm, decision_tree_cm, random_forest_cm)
136
137     ggplot(combined_cm, aes(x = Reference, y = Prediction, fill = Freq)) +
138       geom_tile(color = "white") +
139       facet_wrap(~ Model) +
140       scale_fill_gradient(low = "white", high = "blue") +
141       labs(title = "Confusion Matrices for Classification Models", x = "Actual", y = "Predicted") +
142       theme_minimal() +
143       theme(axis.text.x = element_text(angle = 45, hjust = 1))
144   })
145
146   output$clusterPlot <- renderPlot({
147     fviz_cluster(kmeans_result, data = scaled_data,
148                 geom = "point",
149                 ellipse.type = "convex",
150                 ggtheme = theme_bw()) +
151     labs(title = "K-means Clustering of Heart Disease Data")
152   })
153
154   output$elbowPlot <- renderPlot({
155     wss <- sapply(1:10, function(k){kmeans(scaled_data, centers = k, nstart = 25)$tot.withinss})
156     plot(1:10, wss, type = "b", pch = 19, frame = FALSE,
157          xlab = "Number of clusters (k)",
158          ylab = "Total within-cluster sum of squares (WSS)",
159          main = "Elbow Method for optimal k")
160   })

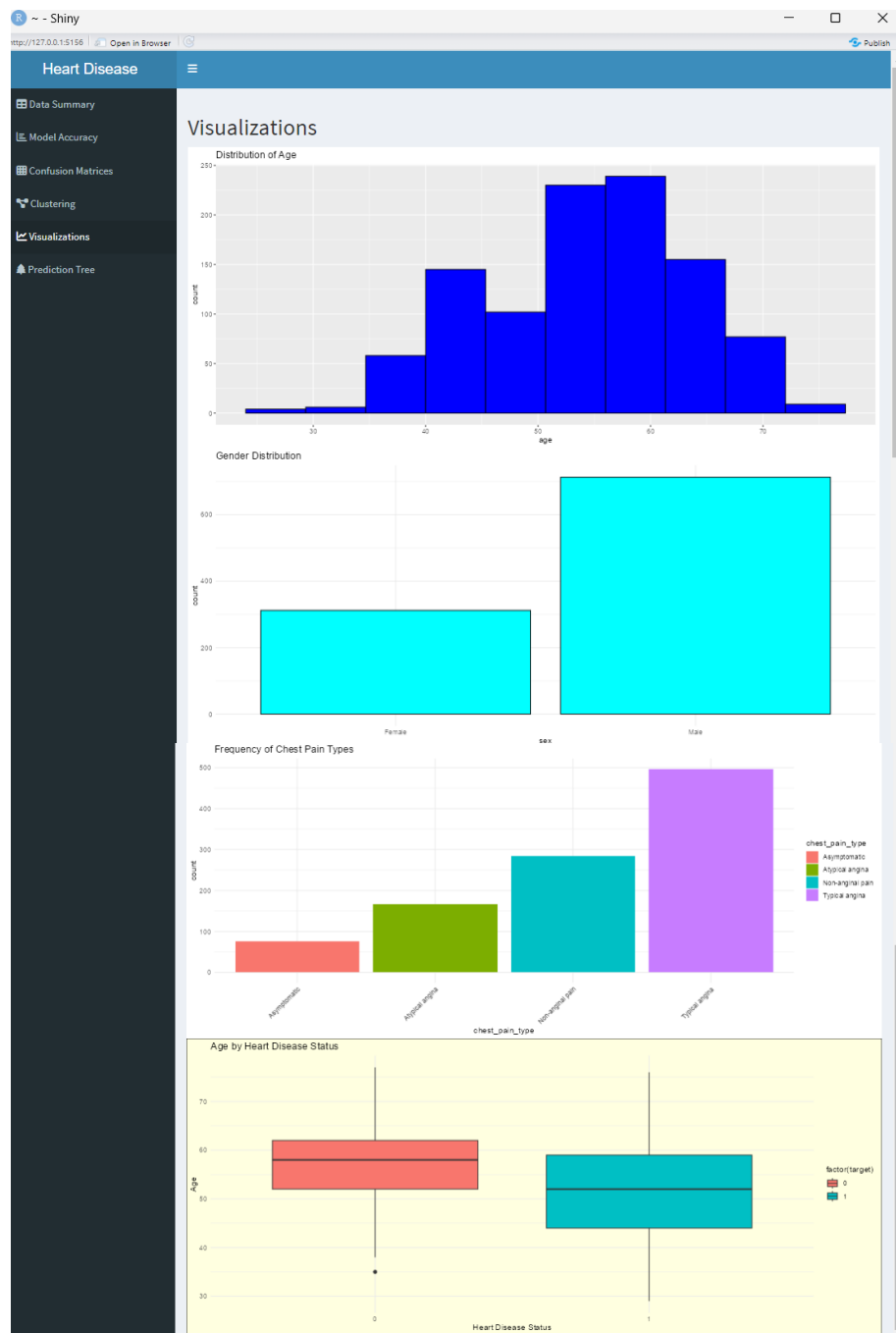
```

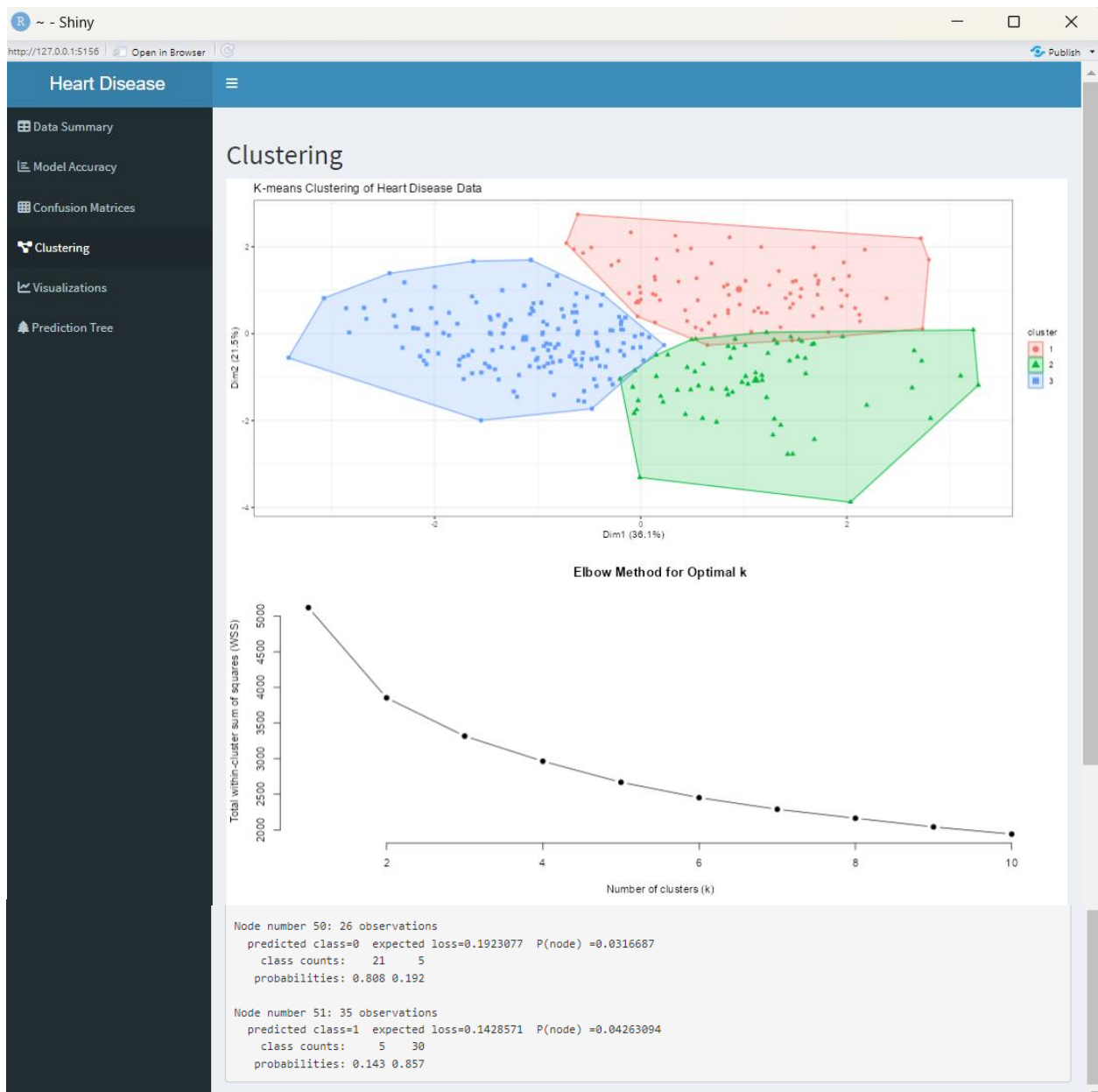
```

162 ~ output$clusterSummary <- renderPrint({
163   cluster_summary <- aggregate(. ~ cluster, data = data, FUN = mean)
164   print(cluster_summary)
165 ~ })
166
167 ~ output$clusterCounts <- renderPrint({
168   cluster_counts <- table(data$cluster)
169   print(cluster_counts)
170 ~ })
171
172 ~ output$agePlot <- renderPlot({
173   ggplot(data, aes(x = age)) + geom_histogram(bins = 10, fill = 'blue', color = 'black') +
174   ggtitle('Distribution of Age')
175 ~ })
176
177 ~ output$genderPlot <- renderPlot({
178   ggplot(data, aes(x = sex)) + geom_bar(fill = 'cyan', color = 'black') +
179   ggtitle('Gender Distribution') +
180   theme_minimal()
181 ~ })
182
183 ~ output$cpPlot <- renderPlot({
184   ggplot(data, aes(x = chest_pain_type, fill = chest_pain_type)) + geom_bar() +
185   ggtitle('Frequency of Chest Pain Types') +
186   theme_minimal() +
187   theme(axis.text.x = element_text(angle = 45, hjust = 1))
188 ~ })
189
190 ~ output$ageBoxPlot <- renderPlot({
191   ggplot(data, aes(x = factor(target), y = age, fill = factor(target))) +
192   geom_boxplot() +
193   labs(title = 'Age by Heart Disease Status', x = 'Heart Disease Status', y = 'Age') +
194   theme_minimal() +
195   theme(plot.background = element_rect(fill = 'lightyellow'))
196 ~ })
197
198 ~ output$predictionTreePlot <- renderPlot({
199   fit <- rpart(target ~ ., data = train_data, method = 'class')
200   rpart.plot(fit, main = 'Heart Disease Prediction', extra = 100, under = TRUE, facLen = 0)
201 ~ })
202
203 ~ output$predictionTreesSummary <- renderPrint({
204   fit <- rpart(target ~ ., data = train_data, method = 'class')
205   summary(fit)
206 ~ })
207 ~ }
208
209
210 shinyApp(ui = ui, server = server)

```





~ - Shiny

http://127.0.0.1:5156

Open in Browser

Publish

Heart Disease

Data Summary

Model Accuracy

Confusion Matrices

Clustering

Visualizations

Prediction Tree

Prediction Tree

Heart Disease Prediction

```

graph TD
    Root((100%)) -->|yes - chest_pain_type = Typical angina| Node1((49%))
    Root -->|no| Node2((51%))
    Node1 -->|vessels_colored_by_fluoroscopy = Four,One,Three,Two| Node1L((0 28%))
    Node1 -->|vessels_colored_by_fluoroscopy = Four,One,Three,Two| Node1R((1 21%))
    Node1R -->|thalassemia = No,Reversible Defect| Node1RL((0 9%))
    Node1R -->|thalassemia = No,Reversible Defect| Node1RR((1 13%))
    Node1RR -->|Max_heart_rate < 120| Node1RRL((0 1%))
    Node1RR -->|Max_heart_rate < 120| Node1RRR((1 12%))
    Node2 -->|age >= 57| Node2L((19%))
    Node2 -->|age >= 57| Node2R((31%))
    Node2L -->|sex = Male| Node2LL((12%))
    Node2L -->|sex = Male| Node2LR((7%))
    Node2LL -->|cholesterol >= 246| Node2LLL((0 5%))
    Node2LL -->|cholesterol >= 246| Node2LLR((1 3%))
    Node2LR -->|vessels_colored_by_fluoroscopy = One,Two| Node2LRL((0 4%))
    Node2LR -->|vessels_colored_by_fluoroscopy = One,Two| Node2LRR((1 7%))
    Node2R -->|vessels_colored_by_fluoroscopy = Three| Node2RL((0 1%))
    Node2R -->|vessels_colored_by_fluoroscopy = Three| Node2RR((1 30%))
  
```

Call:

```
rpart(formula = target ~ ., data = train_data, method = "class")
n= 821
```

CP	nsplit	rel error	xerror	xstd	
1	0.51000	0	1.0000	1.0000	0.03580465
2	0.06625	1	0.4900	0.5650	0.03199493
3	0.02625	3	0.3575	0.3750	0.02768059
4	0.01750	7	0.2425	0.2925	0.02504077
5	0.01250	8	0.2250	0.2875	0.02486106
6	0.01000	9	0.2125	0.2875	0.02486106

Variable importance

Variable	Importance
chest_pain_type	19
Max_heart_rate	13
exercise_induced_angina	8
resting_blood_pressure	2
vessels_colored_by_fluoroscopy	16
oldpeak	9
sex	5
fasting_blood_sugar	1
thalassemia	14
age	9
cholesterol	3
rest_ecg	1

cluster	age	sex	chest_pain_type	resting_blood_pressure	cholesterol	fasting_blood_sugar	rest_ecg
1	59.05017	1.789298	3.608696	126.9130	228.8194	1.882943	2.387960
2	60.24473	1.540084	3.075949	150.8819	292.0633	1.751055	2.320675
3	48.79550	1.713701	2.950920	125.1452	234.1800	1.879346	2.625767

Max_heart_rate	exercise_induced_angina	oldpeak	slope	vessels_colored_by_fluoroscopy	thalassemia	target
125.6856	1.541806	1.9153846	1.923077	3.478261	2.795987	1.240803
140.1983	1.413502	1.3101266	1.679325	4.075949	2.552743	1.388186
163.8834	1.173824	0.4398773	1.394683	4.288344	1.940695	1.740286

1	2	3
299	237	489

8. Conclusion

The analysis of the heart disease dataset has provided significant insights into the various factors associated with heart disease risk. By employing different classification models such as logistic regression, decision tree, and random forest, we were able to predict heart disease with varying degrees of accuracy. The random forest model emerged as the most accurate, followed by logistic regression and decision tree models. Clustering analysis further revealed distinct groups within the data, providing a deeper understanding of how age, cholesterol levels, and other factors cluster among patients. The development and deployment of a user-friendly UI for heart disease analysis also enhanced the accessibility and usability of our findings. Our study underscores the importance of early diagnosis and accurate risk prediction in managing heart disease. By leveraging advanced data analytics techniques, we can contribute to better health outcomes and more targeted interventions for those at risk of heart disease.

9. Contribution

Student Name	Student ID	Contribution
SIMS Wimalasiri (Team leader)	29785	<ul style="list-style-type: none">○ Clustering Models○ Important Insights and Visualizations-Box plot○ Heart Disease Analysis UI
SST Silva	30017	<ul style="list-style-type: none">○ Introduction○ Using libraries○ Data preprocessing and cleaning- Assisted in building UI
KDT Kumara	29883	<ul style="list-style-type: none">○ Important Insights and Visualizations- Prediction Decision Tree- Visualizations of Key Attributes- Assisted in building UI
RMDD Rathnayaka	29618	<ul style="list-style-type: none">○ Classification Models-Model Accuracies-Confusion Matrices-Logistic Regression- Assisted in building UI
MVS Adikari	27613	<ul style="list-style-type: none">○ Classification Models-Valuate the models on the test set-Random Forest-Decision Tree- Assisted in building UI

10. Reference

- [www.kaggle.com. \(n.d.\). *Heart Disease Dataset UCI*. \[online\] Available at: <https://www.kaggle.com/datasets/ketangangal/heart-disease-dataset-uci>.](https://www.kaggle.com/datasets/ketangangal/heart-disease-dataset-uci)
- [rstudio.github.io. \(n.d.\). *Web Application Framework for R — shiny-package*. \[online\] Available at: <https://rstudio.github.io/shiny/reference/shiny-package.html> \[Accessed 27 Jul. 2024\].](https://rstudio.github.io/shiny/reference/shiny-package.html)