Hide # Exploratory Analysis of the Data Set: library(ggplot2) library(dplyr) Attaching package: 'dplyr' The following objects are masked from 'package:stats': filter, lag The following objects are masked from 'package:base': intersect, setdiff, setequal, union Hide library(tidyr) setwd("/Users/faizshaikh/Downloads/archive") Warning: The working directory was changed to /Users/faizshaikh/Downloads/archive inside a notebook chunk. The wo rking directory will be reset when the chunk is finished running. Use the knitr root.dir option in the setup chun k to change the working directory for notebook chunks. Hide heart_data <- read.csv("heart.csv")</pre> # Viewing the structure of the data set: str(heart_data) 'data.frame': 303 obs. of 14 variables: \$ age : int 63 37 41 56 57 57 56 44 52 57 ... \$ sex : int 1 1 0 1 0 1 0 1 1 1 ... : int 3 2 1 1 0 0 1 1 2 2 ... \$ trtbps : int 145 130 130 120 120 140 140 120 172 150 ... \$ chol : int 233 250 204 236 354 192 294 263 199 168 ... \$ fbs : int 1 0 0 0 0 0 0 1 0 ... \$ restecg : int 0 1 0 1 1 1 0 1 1 1 ... \$ thalachh: int 150 187 172 178 163 148 153 173 162 174 ... \$ exng : int 0 0 0 0 1 0 0 0 0 ... \$ oldpeak : num 2.3 3.5 1.4 0.8 0.6 0.4 1.3 0 0.5 1.6 ... \$ slp : int 0 0 2 2 2 1 1 2 2 2 ... \$ caa : int 0 0 0 0 0 0 0 0 0 ... \$ thall : int 1 2 2 2 2 1 2 3 3 2 ... \$ output : int 1 1 1 1 1 1 1 1 1 1 ... Hide # Summary statistics of the numeric variables: summary(heart_data[c("age", "trtbps", "chol", "thalachh")]) chol trtbps thalachh age Min. :29.00 Min. : 94.0 Min. :126.0 Min. :71.0 1st Qu.:47.50 1st Qu.:120.0 1st Qu.:211.0 1st Qu.:133.5 Median :55.00 Median :130.0 Median :240.0 Median :153.0 Mean :54.37 Mean :131.6 Mean :246.3 Mean :149.6 3rd Qu.:61.00 3rd Qu.:140.0 3rd Qu.:274.5 3rd Qu.:166.0 Max. :77.00 Max. :200.0 Max. :564.0 Max. :202.0 Hide # Creating frequency table for the categorical variables: table(heart_data\$Sex) Hide table(heart_data\$exang) Hide table(heart_data\$ca) 0 1 2 3 4 175 65 38 20 5 Hide table(heart_data\$cp) 0 1 2 3 143 50 87 23 Hide table(heart_data\$fbs) 0 1 258 45 Hide table(heart_data\$rest_ecg) Hide table(heart_data\$target) Hide # Plotting distributions of the numeric variables: numeric_vars <- c("age", "trtbps", "chol", "thalachh")</pre> numeric_plots <- lapply(numeric_vars, function(var) {</pre> ggplot(heart_data, aes_string(x = var)) + geom_histogram(binwidth = 5, fill = "blue", color = "black") + labs(title = paste("Distribution of", var), x = var, y = "Frequency") Warning: `aes_string()` was deprecated in ggplot2 3.0.0. Please use tidy evaluation idioms with `aes()`. See also `vignette("ggplot2-in-packages")` for more information. Hide output_print <- ggplot(heart_data, aes(x = as.factor(output))) +</pre> geom_bar(fill = "blue") + scale_y_continuous(labels = scales::comma) + labs(title = "Distribution of Output", x = "Output", y = "Frequency") print(output_print) Distribution of Output 150 **-**Frequency - 001 50 **-**0 -Output Hide gender <- ggplot(heart_data, aes(x = as.factor(sex))) +</pre> geom_bar(fill = "blue") + scale_y_continuous(labels = scales::comma) + labs(title = "Distribution of Gender", x = "Gender", y = "Frequency")print(gender) Distribution of Gender 200 -150 **-**Frequency 001 50 **-**0 -Gender Hide caa_print <- ggplot(heart_data, aes(x = as.factor(caa))) +</pre> geom_bar(fill = "blue") + scale_y_continuous(labels = scales::comma) + labs(title = "Distribution of Major Vessels", x = "Major Vessels", y = "Frequency") print(caa_print) Distribution of Major Vessels 150 **-**Frequency on the second of the 50 **-**Major Vessels Hide cp_print <- ggplot(heart_data, aes(x = as.factor(cp))) +</pre> geom_bar(fill = "blue") + scale_y_continuous(labels = scales::comma) + labs(title = "Distribution of Chest Pain", x = "Chest Pain Level", y = "Frequency") print(cp_print) Distribution of Chest Pain 150 **-**100 -Frequency 50 **-**0 -0 Chest Pain Level Hide # Printing the plots for (plot in numeric_plots) { print(plot) Distribution of age 60 **-**Frequency 20 **-**50 70 40 60 30 age Distribution of trtbps 40 -Frequency 20 **-**120 150 180 90 trtbps Distribution of chol 15 **-**10 -Frequency 5 **-**300 400 200 500 chol Distribution of thalachh 30 **-**Frequency of the second of the 10 -100 200 150 thalachh Hide library(rpart) library(caret) library(e1071) setwd("/Users/faizshaikh/Downloads/archive") Warning: The working directory was changed to /Users/faizshaikh/Downloads/archive inside a notebook chunk. The wo rking directory will be reset when the chunk is finished running. Use the knitr root.dir option in the setup chun k to change the working directory for notebook chunks. Hide heart_data <- read.csv("heart.csv")</pre> set.seed(100) # Factorizing the categorical variables. heart_data\$sex <- as.factor(heart_data\$sex)</pre> heart_data\$cp <- as.factor(heart_data\$cp)</pre> heart_data\$fbs <- as.factor(heart_data\$fbs)</pre> heart_data\$restecg <- as.factor(heart_data\$restecg)</pre> heart_data\$exng <- as.factor(heart_data\$exng)</pre> heart_data\$slp <- as.factor(heart_data\$slp)</pre> heart_data\$caa <- as.factor(heart_data\$caa)</pre> heart_data\$thall <- as.factor(heart_data\$thall)</pre> # Building the decision tree model using RPart() model <- rpart(output ~ cp + fbs + restecg + exng + slp + caa + thall, data = heart_data, method = "class", cp =</pre> 0.000015, minsplit = 10, minbucket = 2) predictions <- predict(model, newdata = heart_data, type = "class")</pre> # Creating the confusion matrix for RPart: print("Confusion Matrix for R-Part:") [1] "Confusion Matrix for R-Part:" Hide cm <- table(heart_data\$output, predictions)</pre> print(cm) predictions 0 1 0 117 21 1 14 151 Hide # Building the Naïve-Bayes model: nb_model <- naiveBayes(output ~ cp + fbs + restecg + exng + slp + caa + thall, data = heart_data)</pre> nb_predictions <- predict(nb_model, newdata = heart_data)</pre> # Creating the confusion matrix for Naïve Bayes: nb_cm <- table(heart_data\$output, nb_predictions)</pre> print("Confusion Matrix for Naïve Bayes:") [1] "Confusion Matrix for Naïve Bayes:" print(nb_cm) nb_predictions 0 1 0 114 24 1 16 149

Calculating accuracy for Naïve-Bayes (based upon confusion matrix results):

Calculating accuracy for RPart (based upon confusion matrix results):

Hide

nb_accuracy <- sum(diag(nb_cm)) / sum(nb_cm)</pre>

rpart_accuracy <- sum(diag(cm)) / sum(cm)</pre>

[1] "Accuracy for rpart: 0.884488448844885"

print(paste("Accuracy for Naïve Bayes:", nb_accuracy))

[1] "Accuracy for Naïve Bayes: 0.867986798679868"

print(paste("Accuracy for rpart:", rpart_accuracy))

Code **▼**

Heart Data Set Analysis