

```
1 # Code converted to pdf using https://tarikjaber.github.io/Code-to-PDF/
2 # Word converted all comments into headings :(
3
4 # B01003963, B00778496, B01006361, B01009661
5
6 # IMPORTING AND LOADING LIBRARIES
7 packages <- c("tidyverse", "ggplot2", "dplyr", "DataExplorer", "caret",
8 "corrplot", "pROC", "randomForest")
9 install_if_missing <- function(p) {
10   if (!require(p, character.only = TRUE)) {
11     install.packages(p, dependencies = TRUE)
12     library(p, character.only = TRUE)
13   }
14 }
15 invisible(sapply(packages, install_if_missing))
16
17
18 # SECTION ONE: Loading the dataset
19
20 # Load the dataset (pleasenmake sure cardio_train.csv is in your working
21 # directory)
22 data <- read.csv("data/cardio_train.csv", sep = ";")
23
24 head(data)
25 str(data)
26 summary(data)
27 print("Data Loading Completed.")
28
29
30 # SECTION TWO: Data Cleansing
31
32 # Load the dataset (make sure cardio_train.csv is in your working directory)
33 df <- read.csv("data/cardio_train.csv", sep = ";")
34
35
36 # View basic information about the dataset
37 print("Structure of dataset:")
38 str(df)
39
40 print("Summary statistics:")
41 summary(df)
42
43 print("Missing values per column:")
44 print(colSums(is.na(df)))
45
46 print("Column names:")
47 print(names(df))
```

```
48
49 print("First few rows of dataset:")
50 head(df)
51
52 # Basic Exploratory Data Analysis, more in eda.R
53 # Distribution of target variable (cardio)
54 ggplot(df, aes(x = factor(cardio))) +
55   geom_bar(fill = "steelblue") +
56   labs(title = "Distribution of Cardiovascular Disease", x = "Cardio (0=No,
1=Yes)", y = "Count")
57
58 # Correlation matrix plot
59 # Notable correlations:
60 #   gender&height, height&weight, gender&smoke, chol&gluc, age&cardio,
weight&cardio, chol&cardio
61 numeric_df <- df[sapply(df, is.numeric)]
62 cor_matrix <- cor(numeric_df)
63 corrplot(cor_matrix, method = "color", type = "upper", tl.cex = 0.8)
64
65 # Data Cleaning
66 # Removing outliers for height, weight, and blood pressure
67 df <- df %>% filter(height > 100 & height < 250)
68 df <- df %>% filter(weight > 30 & weight < 200)
69 df <- df %>% filter(ap_hi > 80 & ap_hi < 250)
70 df <- df %>% filter(ap_lo > 40 & ap_lo < 200)
71
72 # Remove duplicate rows
73 df <- df[!duplicated(df), ]
74
75 # Add new column for age in years (Dataset uses days) and round down to the
nearest whole number
76 df <- df%>%
77   mutate(age_years = floor(age / 365))
78
79 # Feature Engineering
80 # Create BMI variable
81 df$bmi <- df$weight / ((df$height / 100)^2)
82
83 df <- df %>%
84   mutate(bmi = weight / ((height / 100) ^ 2)) # height is in cm
85
86 # Convert categorical variables to factors
87 df$gender <- as.factor(df$gender)
88 df$cholesterol <- as.factor(df$cholesterol)
89 df$gluc <- as.factor(df$gluc)
90 df$smoke <- as.factor(df$smoke)
91 df$alco <- as.factor(df$alco)
92 df$active <- as.factor(df$active)
93 df$cardio <- as.factor(df$cardio)
94
95 # Save the cleaned dataset
```

```
96 write.csv(df, "data/cardio_cleaned.csv", row.names = FALSE)
97
98 print("Data exploration, cleaning, and preprocessing completed. Cleaned
dataset saved as cardio_cleaned.csv")
99
100
101
102
103 # SECTION THREE: Exploratory Data Analysis - THIS SHOULD MAYBE BE THREE
104 # I need to sort this out a little bit, not sure if should use data or df
here
105 # Maybe should rename variables to reflect status, like initial, cleaned,
model_ready
106 # Should be saving plots, not just displaying them
107 #library(scales)
108
109 # Add new column for age in years (Dataset uses days)
110 data <- data%>%
111   mutate(age_years = age / 365)
112
113 # Age Distribution
114 gg_age <- ggplot(df, aes(x = age_years)) +
115   geom_histogram(bins = 30, fill = "steelblue") +
116   xlab("Age (Years)") +
117   ggtitle("Age Distribution of Patients")
118 gg_age
119
120 # Cholesterol Levels (categorical)
121 gg_chol <- ggplot(df, aes(x = cholesterol, fill = cholesterol)) +
122   geom_bar() +
123   xlab("Cholesterol Level") +
124   ggtitle("Cholesterol Levels in Dataset") +
125   theme(legend.position = "none")
126 gg_chol
127
128
129 # Comparing Gender Distribution
130 # Important for ethical considerations, need to consider how this data
effects the results and if it's a fair distribution
131 # Since there's a far greater number (roughly 65%)of female patients we need
to consider why in the report
132 ggplot(data, aes(x=factor(gender))) +
133   geom_bar(fill= "steelblue") +
134   xlab("Gender (1 = Women, 2 = Men)") +
135   ggtitle("Gender Distribution of Patients")
136
137 # Comparing gender distribution against CVD
138 # Very slightly more men have CVD proportionally than women but statistically
insignificant
139 ggplot(data, aes(x = factor(gender), fill = factor(cardio))) +
140   geom_bar(position = "fill") +
```

```
141   scale_y_continuous(labels = scales::percent) +
142   xlab("Gender (1 = Women, 2 = Men)") +
143   ylab("Proportion of Patients") +
144   ggtitle("Gender by Cardiovascular Disease Status") +
145   labs(fill = "CVD (0 = No, 1 = Yes)")
146
147 # Comparing gender distribution against Age
148 # Accounts for difference in number of men and women
149 ggplot(data, aes(x = age_years, fill = factor(gender))) +
150   geom_density(alpha = 0.4, adjust = 1.5) + # Adjusting for a smoother curve
151   xlab("Age (Years)") +
152   ylab("Density") +
153   ggtitle("Normalised Age Distribution against Gender") +
154   labs(fill = "Gender (1 = Women, 2 = Men)")
155
156 # Comparing gender distribution against height
157 ggplot(data, aes(x = height, fill = factor(gender))) +
158   geom_density(alpha = 0.4, adjust = 1.5) + # Adjusting for a smoother curve
159   xlab("Height (cm)") +
160   ylab("Density") +
161   ggtitle("Normalised Height Distribution against Gender") +
162   labs(fill = "Gender (1 = Women, 2 = Men)")
163
164 # Comparing gender distribution against weight
165 ggplot(data, aes(x = weight, fill = factor(gender))) +
166   geom_density(alpha = 0.4, adjust = 1.5) + # Adjusting for a smoother curve
167   xlab("Weight (Kg)") +
168   ylab("Density") +
169   ggtitle("Normalised Weight Distribution against Gender") +
170   labs(fill = "Gender (1 = Women, 2 = Men)")
171
172 # Comparing smoking by gender
173 # Smoking is far more prevalent in men
174 ggplot(data, aes(x = factor(gender), fill = factor(smoke))) +
175   geom_bar(position = "fill") +
176   scale_y_continuous(labels = scales::percent) +
177   xlab("Gender (1 = Women, 2 = Men)") +
178   ylab("Proportion of Patients") +
179   ggtitle("Smoking against Gender") +
180   labs(fill = "Smoking (0 = No, 1 = Yes)")
181
182 # Comparing Alcohol consumption by gender
183 # Alcohol consumption is also more prevalent in men, but less so than smoking
184 ggplot(data, aes(x = factor(gender), fill = factor(alco))) +
185   geom_bar(position = "fill") +
186   scale_y_continuous(labels = scales::percent) +
187   xlab("Gender (1 = Women, 2 = Men)") +
188   ylab("Proportion of Patients") +
189   ggtitle("Alcohol Consumption against Gender") +
190   labs(fill = "Alcohol (0 = No, 1 = Yes)")
191
```

```
192 # Comparing cholesterol levels against gender
193 # Women in the population have a generally higher cholesterol level than men
194 ggplot(data, aes(x = factor(gender), fill = factor(cholesterol))) +
195   geom_bar(position = "fill") +
196   scale_y_continuous(labels = scales::percent) +
197   xlab("Gender (1 = Women, 2 = Men)") +
198   ylab("Proportion of Patients") +
199   ggtitle("Cholesterol Levels against Gender") +
200   labs(fill = "Cholesterol Level (1 = Normal, 2 = Above Normal, 3 = Well
Above Normal)")
201
202 # Comparing Systolic and Diastolic blood pressure against gender
203 # Not entirely happy about how these look, might be better way to represent
204 # However women have a greater max blood pressure than men
205 ggplot(data, aes(x = factor(gender), y = ap_hi, fill = factor(gender))) +
206   geom_boxplot() +
207   xlab("Gender (1 = Women, 2 = Men)") +
208   ylab("Systolic Blood Pressure") +
209   ggtitle("Systolic Blood Pressure against Gender")
210
211 ggplot(data, aes(x = factor(gender), y = ap_lo, fill = factor(gender))) +
212   geom_boxplot() +
213   xlab("Gender (1 = Women, 2 = Men)") +
214   ylab("Diastolic Blood Pressure") +
215   ggtitle("Diastolic Blood Pressure against Gender")
216
217 # Comparing Glucose levels between men and women
218 # Women have slightly higher glucose levels
219 ggplot(data, aes(x = factor(gender), fill = factor(gluc))) +
220   geom_bar(position = "fill") +
221   scale_y_continuous(labels = scales::percent) +
222   xlab("Gender (1 = Women, 2 = Men)") +
223   ylab("Proportion of Patients") +
224   ggtitle("Glucose Levels against Gender") +
225   labs(fill = "Glucose Level (1 = Normal, 2 = Above Normal, 3 = Well Above
Normal)")
226
227 # Comparing activity by gender
228 # Activity levels are very similar
229 ggplot(data, aes(x = factor(gender), fill = factor(active))) +
230   geom_bar(position = "fill") +
231   scale_y_continuous(labels = scales::percent) +
232   xlab("Gender (1 = Women, 2 = Men)") +
233   ylab("Proportion of Patients") +
234   ggtitle("Activity against Gender") +
235   labs(fill = "Activity (0 = Not Active, 1 = Active)")
236
237 # Visualize age distribution
238 ggplot(data, aes(x = age_years)) +
239   geom_histogram(bins = 30) +
240   xlab("Age (Years)") +
```

```
241 ggtitle("Age Distribution Of Patients")
242
243 # Visualize distribution of cholesterol levels
244 ggplot(data, aes(x = factor(cholesterol))) +
245   geom_bar() +
246   xlab("Cholesterol Level (1: Normal, 2: Above Normal, 3: Well Above
Normal)") +
247   ggtitle("Cholesterol Level Distribution")
248
249 # Cardiovascular disease class balance
250 gg_cardio <- ggplot(df, aes(x = factor(cardio))) +
251   geom_bar(fill = "darkred") +
252   xlab("Cardiovascular Disease (0 = No, 1 = Yes)") +
253   ggtitle("Distribution of Cardiovascular Disease Cases")
254 gg_cardio
255
256 # Age vs Disease status (boxplot)
257 gg_age_cardio <- ggplot(df, aes(x = factor(cardio), y = age_years, fill =
factor(cardio))) +
258   geom_boxplot() +
259   xlab("Cardiovascular Disease (0 = No, 1 = Yes)") +
260   ylab("Age (Years)") +
261   ggtitle("Age Distribution by Cardiovascular Disease") +
262   theme(legend.position = "none")
263 gg_age_cardio
264
265 # Smoking vs. Disease (stacked proportion)
266
267 gg_smoke <- ggplot(df, aes(x = factor(smoke), fill = factor(cardio))) +
268   geom_bar(position = "fill") +
269   scale_y_continuous(labels = scales::percent) +
270   xlab("Smoking (0 = No, 1 = Yes)") +
271   ylab("Proportion of Patients") +
272   ggtitle("Smoking Habits of Cardiovascular Disease Status") +
273   labs(fill = "CVD (0 = No, 1 = Yes)")
274 gg_smoke
275
276 # Correlation matrix for numeric variables
277 numeric_vars <- df%>%select_if(is.numeric)
278 cor_matrix <- cor(numeric_vars)
279 corrplot(cor_matrix, method = "color", type = "upper", tl.cex = 0.7)
280
281 # BMI distribution by disease
282 gg_bmi <- ggplot(df, aes(x = bmi, fill = factor(cardio))) +
283   geom_histogram(position = "identity", alpha = 0.5, bins = 30) +
284   xlab("BMI") +
285   ggtitle("BMI Distribution by Cardiovascular Disease") +
286   labs(fill = "CVD (0 = No, 1 = Yes)")
287 gg_bmi
288
289 # Relationship between alcohol status intake and cardiovascular disease
```

```
290 gg_alcohol <- ggplot(data, aes(x = factor(alco), fill = factor(cardio))) +
291   geom_bar(position = "fill") +
292   scale_y_continuous(labels = scales::percent) +
293   xlab("Alcohol Consumption (0 = No, 1 = Yes)") +
294   ylab("Proportion of Patients") +
295   ggtitle("Alcohol Consumption by Cardiovascular Disease Status") +
296   labs(fill = "CVD (0 = No, 1 = Yes)")
297
298 # Display gg_alcohol
299 gg_alcohol
300
301 # Required variables are treated correctly
302 df$gluc <- factor(df$gluc, levels = c(1, 2, 3),
303                  labels = c("Normal", "Above Normal", "Well Above Normal"))
304 df$active <- factor(df$active, levels = c(0, 1), labels = c("Inactive",
305 "Active"))
306 df$cardio <- factor(df$cardio, levels = c(0, 1), labels = c("No", "Yes"))
307
308 # Height Distribution
309 gg_height <- ggplot(df, aes(x = height)) +
310   geom_histogram(bins = 30, fill = "seagreen") +
311   xlab("Height (cm)") +
312   ggtitle("Height Distribution of Patients")
313 gg_height
314
315 # Weight Distribution
316 gg_weight <- ggplot(df, aes(x = weight)) +
317   geom_histogram(bins = 30, fill = "purple") +
318   xlab("Weight (kg)") +
319   ggtitle("Weight Distribution of Patients")
320 gg_weight
321
322 # Glucose Levels by Disease Status
323 gg_gluc <- ggplot(df, aes(x = gluc, fill = cardio)) +
324   geom_bar(position = "fill") +
325   scale_y_continuous(labels = scales::percent) +
326   xlab("Glucose Level") +
327   ylab("Proportion of Patients") +
328   ggtitle("Glucose Levels by Cardiovascular Disease Status") +
329   labs(fill = "CVD")
330 gg_gluc
331
332 # Physical Activity by Disease Status
333 gg_active <- ggplot(df, aes(x = active, fill = cardio)) +
334   geom_bar(position = "fill") +
335   scale_y_continuous(labels = scales::percent) +
336   xlab("Physical Activity") +
337   ylab("Proportion of Patients") +
338   ggtitle("Physical Activity by Cardiovascular Disease Status") +
339   labs(fill = "CVD")
```

```
340 gg_active
341
342
343 # SECTION FOUR: Feature Engineering
344
345 # Load cleaned data
346 df <- read.csv("data/cardio_cleaned.csv")
347
348 # Drop unnecessary columns
349 # ID as not relevant to our research
350 # COMMENTING THIS OUT TEMPORARILY, NEED TO CHECK SOMETHING, GETTING ERROR,
    NOT SHOWING UP IN COLUMN NAMES
351 df <- df %>%
352   select(-id)
353
354 # Scale numeric features
355 # So easier to compare
356 numeric_cols <- c("age_years", "height", "weight", "ap_hi", "ap_lo", "bmi")
357 df[numeric_cols] <- scale(df[numeric_cols])
358
359 # Changes chol & gluc to meaningful labels
360 df$cholesterol <- factor(df$cholesterol, levels = c(1,2,3),
361   labels = c("normal", "above_normal",
    "well_above_normal"))
362 df$gluc <- factor(df$gluc, levels = c(1,2,3),
363   labels = c("normal", "above_normal", "well_above_normal"))
364
365 # Convert chol & gluc to dummy variables, one-hot encoding
366 df <- cbind(df, model.matrix(~ cholesterol + gluc - 1, data = df))
367 df <- df %>% select(-cholesterol, -gluc)
368
369 # Convert cardio values to factors
370 # Changing this temporarily, model will need this as 0 or 1
371 df$cardio <- factor(df$cardio, levels = c(0,1), labels = c("No", "Yes"))
372 #df$cardio <- factor(df$cardio, levels = c(0,1))
373
374
375 # Save
376 write.csv(df, "data/cardio_model_ready.csv")
377
378
379 # SECTION FIVE: Data Splitting
380
381 # Using the feature engineered dataset
382 # Selecting the whole dataset
383 df <- read.csv("data/cardio_model_ready.csv")
384
385 # Setting a seed so can be performed again
386 set.seed(1)
387
388 # Creating the partition
```



```
389 # 70% training, 30% testing
390 partition <- createDataPartition(df$cardio, p = 0.7, list=FALSE)
391
392 # Allocating training and testing data
393 trainingData <- df[partition, ]
394 testData <- df[-partition, ]
395
396 # Checking proportion of cardio data is roughly even
397 table(trainingData$cardio)
398 table(testData$cardio)
399
400 # Saving for use in model
401 write.csv(trainingData, "data/trainingData.csv", row.names = FALSE)
402 write.csv(testData, "data/testData.csv", row.names = FALSE)
403
404 print("Data splitting completed. Datasets saved as testData.csv and
trainingData.csv")
405
406 # SECTION SIX: Creation and evaluation of a simple logistic regression model
407 # We need to consider why we've used this and alternative things
408 # Do we want to see if we can predict it based on only one category like age
or BMI?
409 # This model has an accuracy of roughly 0.7279 which is pretty decent!
410 # This is adapted from the lab with the logistic regression tutorial
411 # cholesterol well above normal coming up as NA
412
413 log_model <- glm(cardio ~., data= trainingData,
family=binomial(link="logit"))
414
415 # Make Predictions
416 log_predictions <- predict(log_model, testData, type = "response")
417 log_predicted_classes <- ifelse(log_predictions > 0.5, 1, 0)
418
419
420 # Begin evaluation of model using a confusion matrix
421 # Big issues here, data and reference should be factors with the same levels
422 # Fix here, but if have time need to go back and check
423 log_conf_matrix <- confusionMatrix(
424   factor(log_predictions, levels = levels(factor(testData$cardio))),
425   factor(testData$cardio, levels = levels(factor(testData$cardio)))
426 )
427
428 print(log_conf_matrix)
429 # Get F1 from the confusion matrix
430 log_conf_matrix$byClass["F1"]
431
432
433 # Printing a summary of the model here
434 summary(log_model)
435
436
```

```
437 # Get feature importance using p-values
438 p_values <- model_summary$coefficients[-1, 4]
439 importance_df_p <- data.frame(
440   Feature = names(p_values),
441   P_Value = p_values,
442   Importance = -log10(p_values)
443 )
444 importance_df_p <- importance_df_p[order(-importance_df_p$Importance),]
445
446 most_significant_feature <- importance_df_p[1, "Feature"]
447 most_significant_pvalue <- importance_df_p[1, "P_Value"]
448
449 # Print paste seems to work for formatting strings
450 print(paste("Most statistically significant feature:",
451            most_significant_feature,
452            "with p-value:", most_significant_pvalue))
453
454 # computing odds ratios and confidence intervals
455 exp(coef(log_model))
456 exp(cbind(OR = coef(log_model), confint(log_model)))
457
458 # AUC ROC - Sensitivity vs specificity
459 roc_object <- roc( testData$cardio, log_predictions)
460 rocCurve <- ggroc(roc_object) + ggtitle("ROC Curve for Logistic Regression
461 Model")
462 ggsave("results/logistic/roc.png", plot = rocCurve)
463
464 # AUC = 0.7897, the closer the auc is to 1, the better the model
465 # Close to 1: Good at distinguishing between positive and negative classes
466 # Close to 0.5: Performs no better than random guessing
467 # https://www.geeksforgeeks.org/plotting-roc-curve-in-r-programming/
468 auc(roc_object)
469
470 # Trying out logistic curves
471 # Original go results in a very spiky plot, this makes things readable
472 # The shaded section shows the confidence interval, automatically 95%
473 ageLC <- ggplot(testData, aes(x = age_years, y = log_predictions)) +
474   geom_smooth(method = "loess", color = "steelblue") +
475   labs(title = "Smoothed Logistic Regression Curve for age",
476        x = "Age in years",
477        y = "Probability of CVD")
478 ggsave("results/logistic/cvdByAgeLC.png", plot = ageLC)
479
480 bmiLC <- ggplot(testData, aes(x = bmi, y = log_predictions)) +
481   geom_smooth(method = "loess", color = "steelblue") +
482   labs(title = "Logistic Regression Curve for BMI",
483        x = "BMI",
484        y = "Predicted Probability of CVD")
485 ggsave("results/logistic/cvdByBMILC.png", plot = bmiLC)
```

```
486 ap_hiLC ← ggplot(testData, aes(x = ap_hi, y = log_predictions)) +
487   geom_smooth(method = "loess", color = "steelblue") +
488   labs(title = "Logistic Regression Curve for ap_hi",
489     x = "ap_hi",
490     y = "Predicted Probability of CVD")
491 ggsave("results/logistic/cvdByap_hiLC.png", plot = ap_hiLC)
492
493 ap_loLC ← ggplot(testData, aes(x = ap_lo, y = log_predictions)) +
494   geom_smooth(method = "loess", color = "steelblue") +
495   labs(title = "Logistic Regression Curve for ap_lo",
496     x = "ap_lo",
497     y = "Predicted Probability of CVD")
498 ggsave("results/logistic/cvdByap_loLC.png", plot = ap_loLC)
499
500
501 # HEYO! Gender is NOT a continuous variable, however we have it represented
502 # as 1 or 2.
503 # This plot shows us that this model predicts that men (2) are more likely to
504 # have CVD,
505 # however it presents it as a continuous variable. The question is do we
506 # leave it as this,
507 # or do we find a new type of plot that represents binary variables better?
508 # I really like the look of this graph, it shows things clearly
509 genderPlotLC ← ggplot(testData, aes(x = gender, y = log_predictions)) +
510   geom_smooth(method = "loess", color = "steelblue") +
511   labs(title = "Logistic Regression Curve for Gender",
512     x = "Gender (1 = Woman, 2 = Man)",
513     y = "Predicted Probability of CVD")
514 ggsave("results/logistic/cvdByGenderLC.png", plot = genderPlotLC)
515
516 # Not as nice graph, but technically represents things better
517 genderPlot ← ggplot(testData, aes(x = gender, y = log_predictions)) +
518   stat_summary(fun = mean, geom = "bar", position = position_dodge(),
519     fill="steelblue") +
520   labs(title = "Probability of CVD by Gender",
521     x = "Gender (1 = Woman, 2 = Man)",
522     y = "Predicted Probability of CVD")
523 ggsave("results/logistic/cvdByGender.png", plot = genderPlot)
524
525
526 smokePlot ← ggplot(testData, aes(x = smoke, y = log_predictions)) +
527   stat_summary(fun = mean, geom = "bar", position = position_dodge(),
528     fill="steelblue") +
529   labs(title = "Probability of CVD by Smoking",
530     x = "Smoking (0 = Does not smoke, 1 = Does smoke)",
531     y = "Predicted Probability of CVD")
532 ggsave("results/logistic/cvdBySmoke.png", plot = smokePlot)
533
534
535 alcoPlot ← ggplot(testData, aes(x = smoke, y = log_predictions)) +
536   stat_summary(fun = mean, geom = "bar", position = position_dodge(),
```

```
fill="steelblue") +
532   labs(title = "Probability of CVD by Alcohol",
533         x = "Alcohol (0 = Does not drink, 1 = Does drink)",
534         y = "Predicted Probability of CVD")
535 ggsave("results/logistic/cvdByAlco.png", plot = alcoPlot)
536
537
538 activePlot ← ggplot(testData, aes(x = active, y = log_predictions)) +
539   stat_summary(fun = mean, geom = "bar", position = position_dodge(),
fill="steelblue") +
540   labs(title = "Probability of CVD by Activity",
541         x = "Activity (0 = Inactive, 1 = Active)",
542         y = "Predicted Probability of CVD")
543 ggsave("results/logistic/cvdByActivity.png", plot = activePlot)
544
545
546 # I wanted to explore the possibility of smoking interfering with the gender
predictions
547 # This doesn't seem to be the case, this shows the probability of men and
women who smoke and don't
548 # The model reckons we should all take up smoking
549 genderSmoke ← ggplot(testData, aes(x = interaction(gender, smoke), y =
log_predictions, fill = interaction(gender, smoke))) +
550   stat_summary(fun = mean, geom = "bar", position = position_dodge()) +
551   labs(title = "Probability of CVD by Gender and Smoking",
552         x = "Group",
553         y = "Mean Probability of CVD")
554 ggsave("results/logistic/cvdByGenderSmoking.png", plot = genderSmoke)
555
556 # Similarly, the odds of having CVD and drinking alcohol are lower across the
board
557 genderAlco ← ggplot(testData, aes(x = interaction(gender, alco), y =
log_predictions, fill = interaction(gender, alco))) +
558   stat_summary(fun = mean, geom = "bar", position = position_dodge()) +
559   labs(title = "Probability of CVD by Gender and Alcohol",
560         x = "Group",
561         y = "Mean Probability of CVD")
562 ggsave("results/logistic/cvdByGenderAlcohol.png", plot = genderAlco)
563
564 # Saving the model
565 # We might want to move this up before the evaluation, not sure it really
matters
566 saveRDS(log_model, "results/models/cardio_logistic_model.rds")
567
568 print("Logistic Regression model created and saved as
cardio_logistic_model.rds")
569
570 # SECTION SEVEN: Creation and Evaluation of a K-Nearest Neighbors (KNN) Model
571 # We're using KNN to classify whether a person has cardiovascular disease
based on all available features.
572 # This approach looks at the 'k' nearest patients and makes a prediction
```

```
based on what class most of them belong to.
573 # NEEDS ROC, AUC AND SUMMARY, PROB CONF INTERVALS
574
575 # Train the KNN model using 10-fold cross-validation
576 # We scale the data to ensure fairness in distance calculations
577 # The issue here is because cardio is being treated as numeric instead of
    categorical
578 # Fixed this
579 trainingData$cardio ← as.factor(trainingData$cardio)
580 testData$cardio ← as.factor(testData$cardio)
581
582 knn_model ← train(cardio ~ ., data = trainingData, method = "knn", preProcess
    = c("center", "scale"), trControl = trainControl(method = "cv", number = 10))
583
584 # Make predictions on the test set
585 knn_predictions ← predict(knn_model, testData)
586 knn_probs ← predict(knn_model, testData, type = "prob")[,2] #
    Probabilities for class 1
587
588 # Performs best when using 9 nearest neighbours
589 print(knn_model)
590
591 # FIXED THIS
592 # Evaluate the model using a confusion matrix
593 knn_conf_matrix ← confusionMatrix(
594   factor(knn_predictions, levels = levels(testData$cardio)),
595   factor(testData$cardio, levels = levels(testData$cardio))
596 )
597
598 # Display the results
599 print("Confusion Matrix for KNN Model:")
600 print(knn_conf_matrix)
601
602
603 # Printing a summary of the model here
604 summary(knn_model)
605
606
607 # AUC ROC - Sensitivity vs specificity
608 # 0.7629
609 knn_roc_object ← roc( testData$cardio, knn_probs)
610 knn_rocCurve ← ggroc(knn_roc_object) + ggtitle("ROC Curve for K-NN Model")
611 ggsave("results/knn/roc.png", plot = rocCurve)
612
613 # AUC
614 auc(knn_roc_object)
615
616 knn_conf_matrix$byClass["F1"]
617
618
619 # Visualize predicted CVD status across BMI
```

```
620 knn_bmi_plot <- ggplot(testData, aes(x = bmi, fill = knn_predictions)) +
621   geom_density(alpha = 0.5) +
622   labs(title = "KNN: Predicted CVD Probability by BMI",
623        x = "BMI",
624        fill = "Predicted CVD")
625 ggsave("results/knn/knnByBMI.png", plot = knn_bmi_plot)
626
627 # Visualize predicted CVD by gender
628 knn_gender_plot <- ggplot(testData, aes(x = gender, fill = knn_predictions))
+ geom_bar(position = "fill") +
629   scale_y_continuous(labels = scales::percent) +
630   labs(title = "KNN: Gender Distribution by Predicted CVD",
631        x = "Gender (1 = Woman, 2 = Man)",
632        y = "Proportion of Patients",
633        fill = "Predicted CVD")
634 ggsave("results/knn/knnByGender.png", plot = knn_gender_plot)
635
636 # Get variable importance
637 importance <- varImp(knn_model)
638 importance_df$Feature <- rownames(importance_df)
639 importance_df <- importance_df[order(-importance_df[,1]),]
640 most_important_feature <- importance_df[1, "Feature"]
641 most_important_value <- importance_df[1, 1]
642 print(paste("Most important feature:", most_important_feature,
643            "with importance value:", most_important_value))
644
645
646
647 # Save the trained model for later evaluation
648 saveRDS(knn_model, "results/models/cardio_knn_model.rds")
649 print("KNN model trained and saved as cardio_knn_model.rds")
650
651
652
653 # SECTION EIGHT: Creation and Evaluation of a Random Forest Model
654
655
656 # Ensure cardio is treated as factor for classification
657 trainingData$cardio <- as.factor(trainingData$cardio)
658 testData$cardio <- as.factor(testData$cardio)
659
660 # Train the Random Forest model
661 rf_model <- randomForest(cardio ~ ., data = trainingData, ntree = 100,
+ importance = TRUE)
662
663 # Save the trained model
664 saveRDS(rf_model, "results/models/cardio_randomforest_model.rds")
665 print(" Random Forest model saved as cardio_randomforest_model.rds")
666
667 # --- Predictions ---
668 rf_predictions <- predict(rf_model, testData, type = "class") #
```

```
Predicted classes
669 rf_probs <- predict(rf_model, testData, type = "prob")[,2] #
    Probabilities for class 1
670
671 # --- Confusion Matrix ---
672 rf_conf_matrix <- confusionMatrix(
673   factor(rf_predictions, levels = levels(testData$cardio)),
674   factor(testData$cardio, levels = levels(testData$cardio))
675 )
676 print("Confusion Matrix for Random Forest Model:")
677 print(rf_conf_matrix)
678
679 # Get F1
680 rf_conf_matrix$byClass["F1"]
681
682
683 # --- ROC & AUC ---
684 rf_roc <- roc(testData$cardio, rf_probs)
685 rf_auc <- auc(rf_roc)
686 print(paste("AUC for Random Forest:", rf_auc))
687
688 # Save ROC plot
689 rf_roc_plot <- ggroc(rf_roc) + ggtitle("ROC Curve for Random Forest Model")
690 ggsave("results/randomforest/roc.png", plot = rf_roc_plot)
691
692 # --- Variable Importance ---
693 varImpPlot(rf_model, main = "Random Forest - Variable Importance")
694 write.csv(importance(rf_model), "results/randomforest/
variable_importance.csv")
695
696 # --- Probability Summary with Confidence Intervals ---
697 # Bootstrap confidence interval for mean predicted probabilities
698 rf_ci <- t.test(rf_probs ~ testData$cardio)$conf.int
699 print(" 95% Confidence Interval for Predicted Probabilities (CVD vs No
CVD):")
700 print(rf_ci)
701
702 # Save ROC stats and confidence interval summary
703 rf_eval <- data.frame(
704   Accuracy = rf_conf_matrix$overall["Accuracy"],
705   Kappa = rf_conf_matrix$overall["Kappa"],
706   Sensitivity = rf_conf_matrix$byClass["Sensitivity"],
707   Specificity = rf_conf_matrix$byClass["Specificity"],
708   AUC = as.numeric(rf_auc),
709   CI_Lower = rf_ci[1],
710   CI_Upper = rf_ci[2]
711 )
712
713 write.csv(rf_eval, "results/randomforest/evaluation_summary.csv", row.names =
FALSE)
714 print("Evaluation summary for Random Forest saved to results/randomforest/
```

```
evaluation_summary.csv")
715
716 # Print summary table to console
717 print(rf_eval)
718
719
720 # SECTION NINE: Any further eval & model comparison
721
722 # Reckon it's useful to combine all rocs into one plot
723 roc_list <- list(roc_object, knn_roc_object, rf_roc)
724
725 # Create the plot for all ROC curves
726 # 1= Log, 2= knn, 4= rf
727 all_rocs_plot <- ggroc(roc_list) +
728   ggtitle("The ROC curves for Models")
729 ggsave("results/all_rocs_plot.png", plot =all_rocs_plot)
730
```