**DOUGLAS**COLLEGE

**Diabetes Dataset Analysis**

CSIS4290: Special Topics in Data Analytics

Instructor: Parissa Ahmadi Abkenari

WELCOME

Group Members:

'Aafrin Zahid Memon

Akintunde Akinro

Bruno do Nascimento Beserra

Douglas College respectfully acknowledges that our campuses are located on the unceded traditional and ancestral lands of the Coast Salish Peoples, including the territories of the q̓íc̓əy̓ (Katzie), q̓ʷa:n̓ƛ̓ən̓ (Kwantlen), kʷikʷəƛ̓əm (Kwikwetlem), xʷməθkʷəy̓əm (Musqueam), and qiqéyt (Qayqayt) First Nations.

# Contents

- **Dataset Introduction**
  - Data Visualizations
  - Exploratory Data Analysis
  - Data Cleaning
- **Creating Train, Test datasets**
  - Feature Selection
    - LASSO
    - RIDGE
    - Correlation

- **Models Analysis**
  - Decision Tree
  - Naive Bayes
  - K-Nearest Neighbors
  - Support Vector Machine
  - OLS Regression
  - GLM Regression
  - Bagging
  - XGBoost
  - Random Forest
- **Models Comparison**
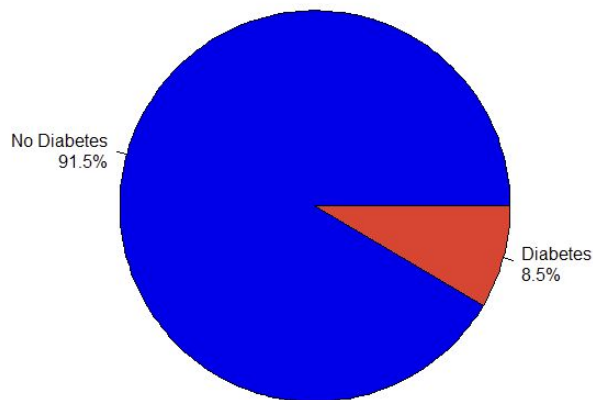- **Research Findings**
- **Conclusion**
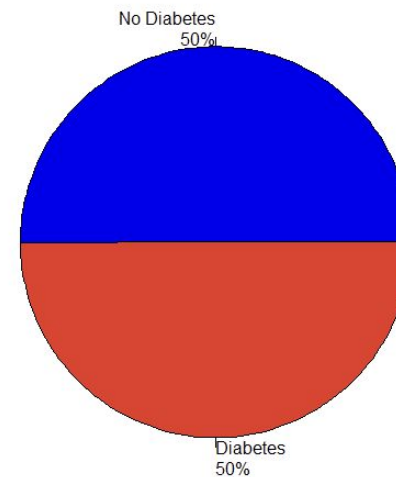
# Dataset Introduction

- Diabetes Dataset

- The Diabetes dataset is a collection of medical and demographic data from patients, along with their diabetes status (positive or negative).

- The data includes features such as age, gender, body mass index (BMI), hypertension, heart disease, smoking history, HbA1c level, and blood glucose level.

- Source: Kaggle

# Data Visualization - Target Balance Comparison

**Diabetes Class Balance - Before**

**Diabetes Class Balance - After**



No Diabetes
91.5%

Diabetes
8.5%

No Diabetes
50%

Diabetes
50%

Size of dataset was 99982 before the balancing process and 17007 after it.
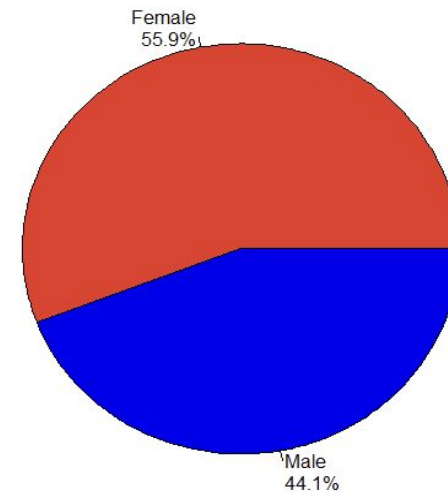
| | |
|---|---|
| ▶ data | 99982 obs. of 16 variables |
| ▶ df | 17007 obs. of 20 variables |

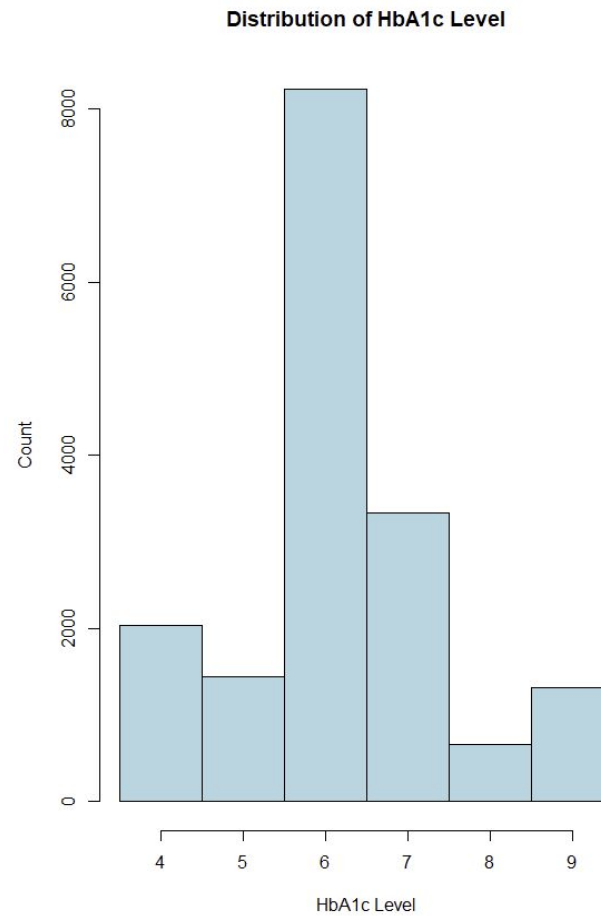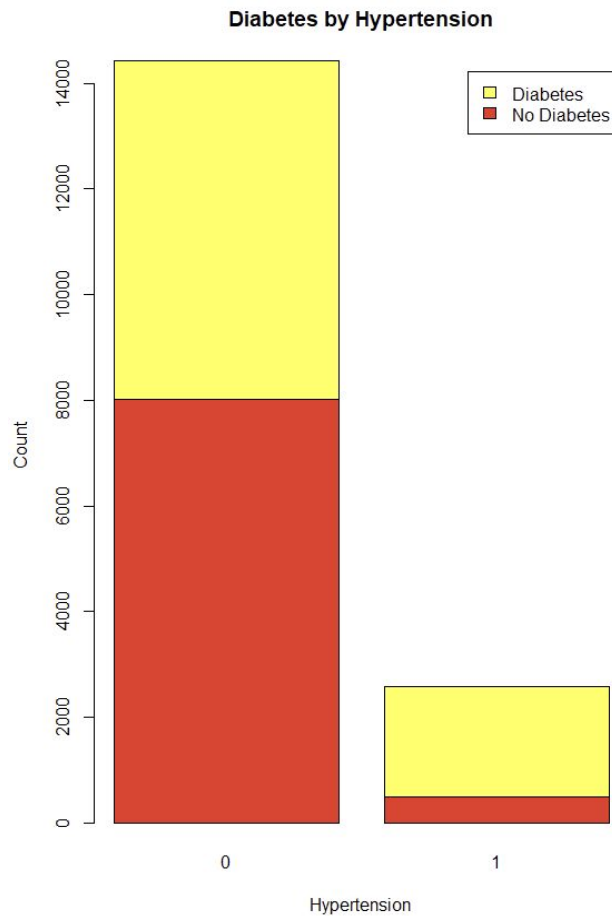# Data Visualization - Age and Gender Proportion

**Age Distribution**

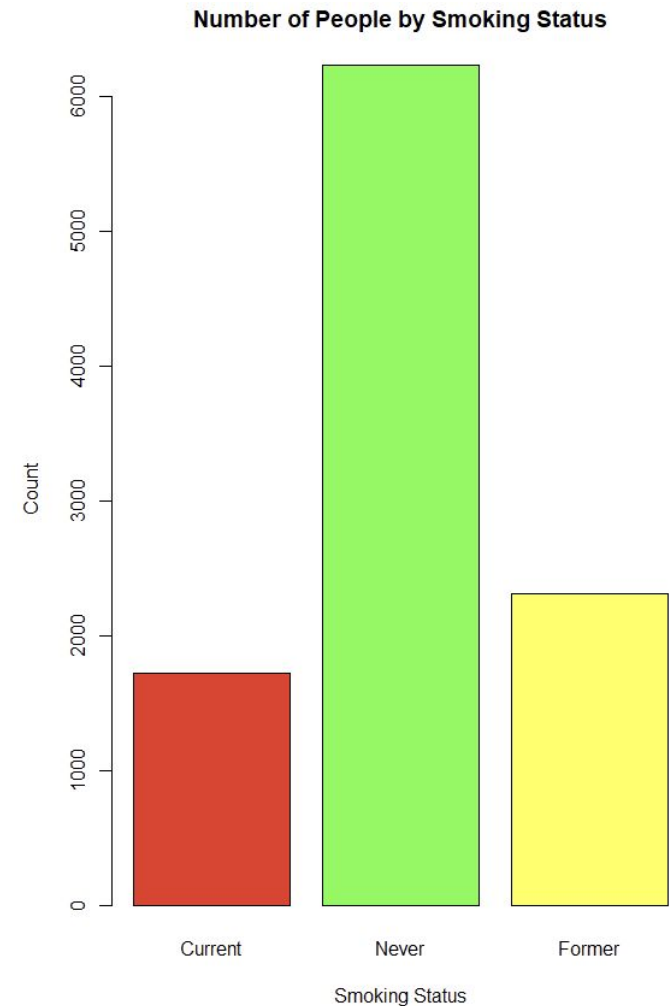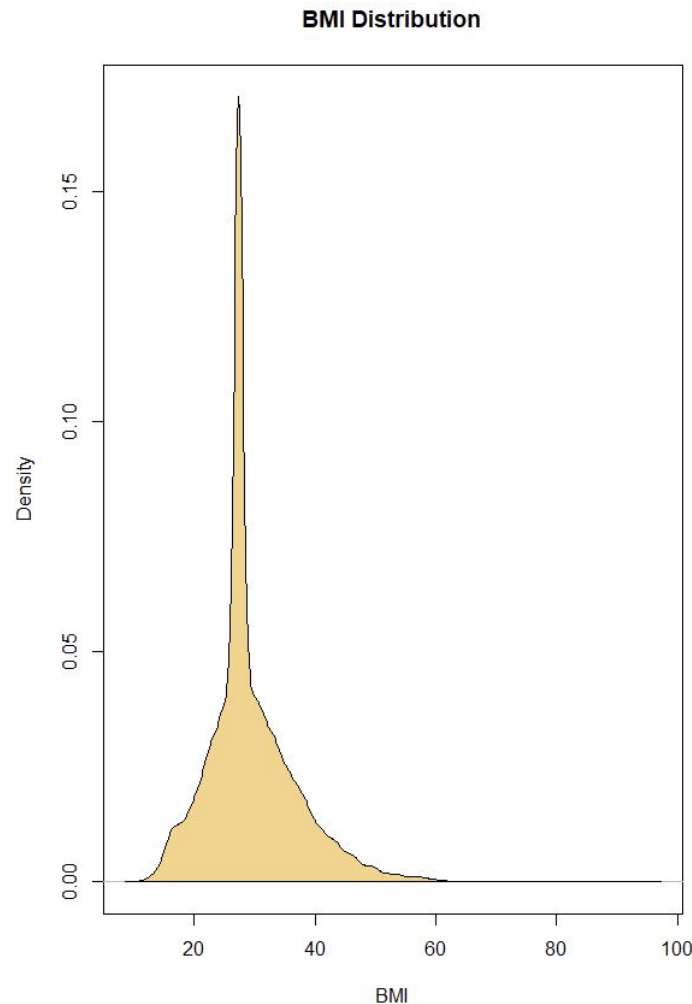**Gender Distribution**
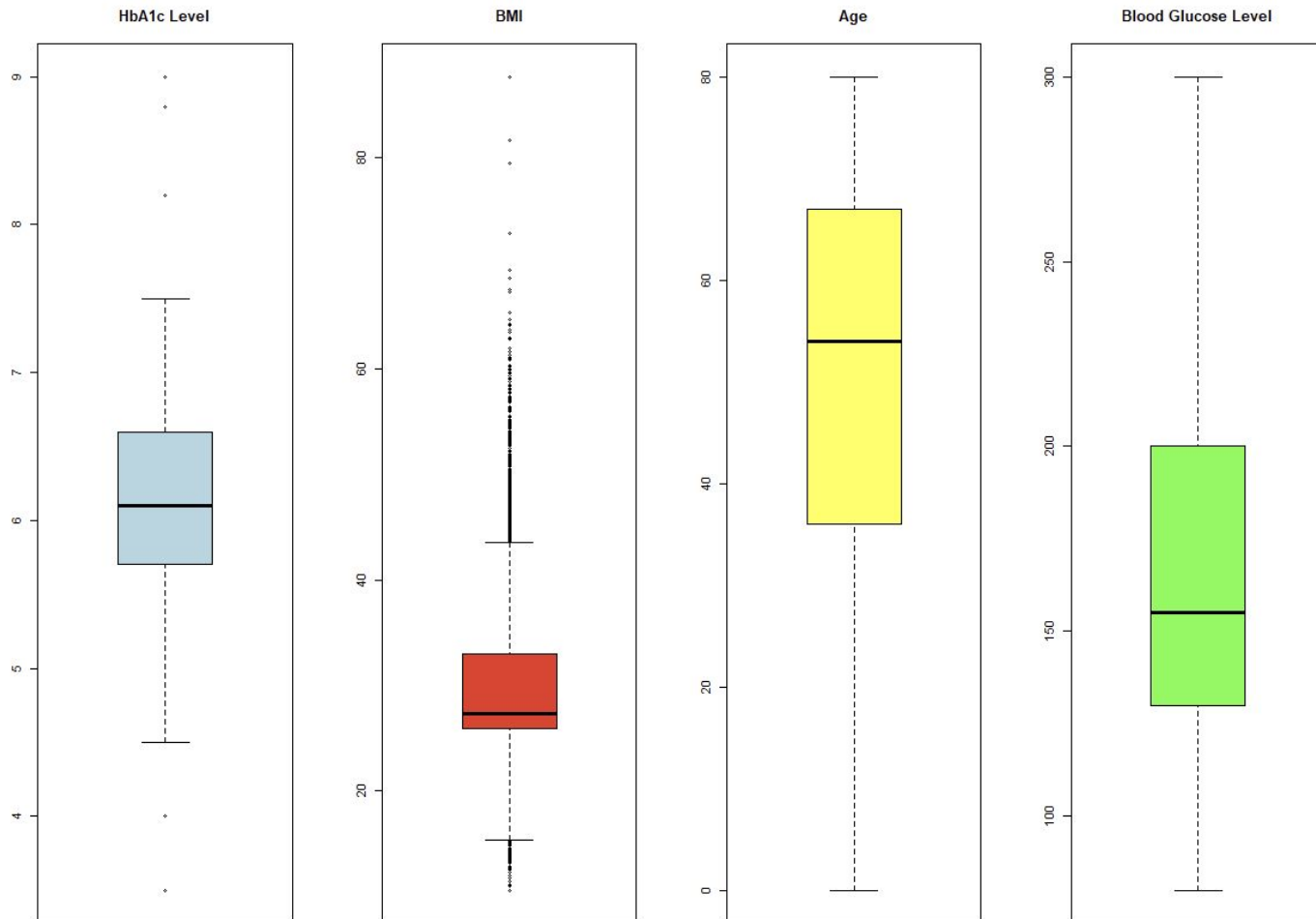
# Data Visualization - Diabetes Risk Factors

# Data Visualization - BMI and Smoking Distribution



**BMI Distribution**

**Number of People by Smoking Status**

# Data Visualization - Box plots

# Data Visualization - Codes

```r
207  par(mfrow = c(1, 2))
208  # Diabetes Class Distribution
209  counts <- table(data$diabetes)
210  percentages <- round(100 * counts / sum(counts), 1)
211  labels <- paste(c("No Diabetes", "Diabetes"), "\n", percentages, "%", sep = "")
212  pie(counts,
213      main = "Diabetes Class Balance - Before",
214      labels = labels,
215      col = c("blue", "red"))
216
217  counts <- table(df$diabetes)
218  percentages <- round(100 * counts / sum(counts), 1)
219  labels <- paste(c("No Diabetes", "Diabetes"), "\n", percentages, "%", sep = "")
220  pie(counts,
221      main = "Diabetes Class Balance - After",
222      labels = labels,
223      col = c("blue", "red"))
224
242  # Bar Plot of Diabetes by Hypertension
243  barplot(table(df$diabetes, df$hypertension),
244          main = "Diabetes by Hypertension",
245          xlab = "Hypertension",
246          ylab = "Count",
247          col = c("red", "yellow"),
248          legend = c("No Diabetes", "Diabetes"))
249
250
251  # Histogram with distribution of HbA1c Level
252  hist(df$hbA1c_level,
253      main = "Distribution of HbA1c Level",
254      xlab = "HbA1c Level",
255      ylab = "Count",
256      col = "lightblue",
257      breaks = seq(min(df$hbA1c_level),
258                   max(df$hbA1c_level) + 0.5,
259                   by = 1))
282  # Box Plots
283  par(mfrow = c(1, 4))
284  boxplot(df$hbA1c_level, main = "HbA1c Level", col = "lightblue")
285  boxplot(df$bmi, main = "BMI", col = "red")
286  boxplot(df$age, main = "Age", col = "yellow")
287  boxplot(df$blood_glucose_level, main = "Blood Glucose Level", col = "green")
288  par(mfrow = c(1, 1))
```

```r
225  # Per Gender and Age
226  hist(df$age,
227      main = "Age Distribution",
228      xlab = "Age",
229      ylab = "Count",
230      col = "lightblue")
231
232  # Distribution per gender
233  counts <- table(df$is_male)
234  percentages <- round(100 * counts / sum(counts), 1)
235  labels <- paste(c("Female", "Male"), "\n", percentages, "%", sep = "")
236
237  pie(counts,
238      main = "Gender Distribution",
239      labels = labels,
240      col = c("red", "blue"))

261  # BMI distribution
262  plot(density(df$bmi),
263      main = "BMI Distribution",
264      xlab = "BMI",
265      ylab = "Density")
266  polygon(density(df$bmi),
267          col = rgb(1, 0.65, 0, 0.5))
268
269  # Bar plot per type of smoking history
270  smoking_counts <- c(
271    Current = sum(df$smoking_history_current),
272    Never = sum(df$smoking_history_never),
273    Former = sum(df$smoking_history_former)
274  )
275
276  barplot(smoking_counts,
277          main = "Number of People by Smoking Status",
278          xlab = "Smoking Status",
279          ylab = "Count",
280          col = c("red", "green", "yellow"))
```

# Exploratory Data Analysis

- Check Imbalance in dataset, shape, values from attributes, statistics, and missing values.

```
60   ##########################################################################
61   ######################## EXPLORATORY DATA ANALYSIS ######################
62   ##########################################################################
63
64   # Check Imbalance
65   table(data$diabetes) # dataset is highly imbalanced (over 10:1)
66
67   sample_data <- ovun.sample(diabetes~., data=data, p=0.5, seed=42,
68                              method="under")$data
69   table(sample_data$diabetes)
70
71   # Check shape of data
72   nrow(sample_data)
73
74   # Check data
75   str(sample_data)
76
77   # Check unique values from char column that might be good for analysis
78   unique(sample_data$smoking_history)
79
80   # Check Null values
81   sum(is.na(sample_data))
82
83   # Check Summary Statistics
84   summary(sample_data)
```

# Data Cleaning

- Converted gender to numeric, created dummies for smoking history, remove char columns, and fixed naming convention from attributes.

```
90 ▾  ###############################################################################
91 ▾  ############################# DATA CLEANING ############################
92 ▾  ###############################################################################
93
94    # converting gender to numeric
95    sample_data$is_male <- ifelse(sample_data$gender == "Male", 1, 0)
96
97    # Creating Dummies for smoking history column
98    sample_data <- fastDummies::dummy_cols(
99      sample_data,
100     select_columns = "smoking_history",
101     remove_first_dummy = FALSE,
102     remove_selected_columns = TRUE
103   )
104
105   # Removing char columns
106   sample_data <- subset(sample_data, select = -c(location, gender))
107
108   # Replace spaces with underscores
109   colnames(sample_data) <- gsub(" ", "_", colnames(sample_data))
110
111   # Check data after changes
112   str(sample_data)
```

# Creating Train, Test datasets

- We used two sets in our project, one with all attributes named 'full' and one with the selected features named 'fs'. Both with ratio 70:30.

```
115 ▾ ###########################################################################
116 ▾ ################# CREATE TRAIN AND TEST SETS - FULL MODEL #############
117 ▾ ###########################################################################
118
119   sample_data$diabetes <- factor(sample_data$diabetes, levels = c(0,1))
120   df <- sample_data
121
122   # full = full data used
123   sample_split <- sample.split(Y = df$diabetes, SplitRatio = 0.7)
124   full_train_set <- subset(x= df, sample_split == TRUE)
125   full_test_set <- subset(x= df, sample_split == FALSE)
126
192 ▾ ###########################################################################
193 ▾ ############# CREATE TRAIN AND TEST SETS - FEATURE SELECTION #########
194 ▾ ###########################################################################
195
196   df_fs <- sample_data[, final_selected_columns]
197   df_fs$diabetes <- factor(df_fs$diabetes, levels = c(0,1))
198
199
200   # fs = feature selection
201   fs_sample_split <- sample.split(Y = df_fs$diabetes, SplitRatio = 0.7)
202   fs_train_set <- subset(x= df_fs, fs_sample_split == TRUE)
203   fs_test_set <- subset(x= df_fs, fs_sample_split == FALSE)
```

# Feature Selection

- Used to reduce overfitting, and handle multicollinearity

- Used Correlation, LASSO and RIDGE Techniques to decide the better attributes to keep in our model.

```
134 ▾ #### LASSO ####
135   k_train <- model.matrix(diabetes ~ ., full_train_set)[, -1]
136   x_test <- model.matrix(diabetes ~ ., full_test_set)[, -1]
137   y_train <- full_train_set$diabetes
138   y_test <- full_test_set$diabetes
139
140   lasso_cv <- cv.glmnet(x_train, y_train, alpha = 1, nfolds=10,
141                         family = "binomial")
142
143   lasso_coef <- coef(lasso_cv, s = "lambda.1se")
144   selected_lasso <- rownames(lasso_coef)[lasso_coef[,1] != 0]
145
146   selected_lasso
147
148   lasso_features <- c("age", "hypertension", "heart_disease", "bmi",
149                       "hbA1c_level", "blood_glucose_level", "is_male",
150                       "smoking_history_No_Info", "race.Other",
151                       "smoking_history_ever")
```

```
153 ▾ #### RIDGE ####
154
155   ridge_cv <- cv.glmnet(x_train, y_train, alpha = 0, nfolds=10, family = "binomial")
156   ridge_coef <- coef(ridge_cv, s = "lambda.1se")
157
158   ridge_importance <- data.frame(
159     variable = rownames(ridge_coef),
160     coefficient = as.numeric(ridge_coef)
161   )
162
163   ridge_importance <- ridge_importance[order(abs(ridge_importance$coefficient),
164                                         decreasing = TRUE), ]
165
166   ridge_features <- c("hbA1c_level", "hypertension", "heart_disease",
167                       "smoking_history_former", "smoking_history_No_Info",
168                       "smoking_history_ever", "is_male",
169                       "smoking_history_not_current", "smoking_history_current",
170                       "smoking_history_never")
171
```

# Feature Selection

- We did the correlation between attributes selected from LASSO and RIDGE.
- None of the attributes presented correlation higher than 47%, hence all of them could be added.
- but for our study doesn't make sense consider some of the columns even though they could be added.

```
174  # Check Columns chose from LASSO and RIDGE to see if they correlate
175  all_features <- unique(c(lasso_features, ridge_features))
176  all_features
177  feature_data <- full_train_set[, all_features]
178
179  cor_matrix <- cor(feature_data)
180  cor_matrix
181
```

```
186  # Final Selection based in LASSO and RIDGE
187  final_selected_columns <- c("hbA1c_level", "hypertension", "heart_disease",
188                             "smoking_history_former", "is_male", "bmi", "age", "blood_glucose_level",
189                             "smoking_history_current", "smoking_history_never", "diabetes")
190
```

# Models Analysis - Decision Tree

```
> confusionMatrix(full_test_set$diabetes, pred_dt_full)
Confusion Matrix and Statistics

          Reference
Prediction    0    1
         0 2218  334
         1  177 2373

               Accuracy : 0.8998
                 95% CI : (0.8913, 0.9079)
    No Information Rate : 0.5306
    P-Value [Acc > NIR] : < 2.2e-16

                  Kappa : 0.7997

 Mcnemar's Test P-Value : 5.163e-12

            Sensitivity : 0.9261
            Specificity : 0.8766
         Pos Pred Value : 0.8691
         Neg Pred Value : 0.9306
             Prevalence : 0.4694
         Detection Rate : 0.4347
   Detection Prevalence : 0.5002
      Balanced Accuracy : 0.9014

       'Positive' Class : 0
```

```
306 ▾ ###############################################################################
307 ▾ ########################### DECISION TREE ##################################
308 ▾ ###############################################################################
309
310   # Full Model
311   full_model_dt <- rpart(diabetes ~ ., data = full_train_set, method = "class")
312
313   rpart.plot(full_model_dt)
314
315   importances <- varImp(full_model_dt)
316   importances %>% arrange(desc(Overall))
317
318   pred_dt_full <- predict(full_model_dt, newdata = full_test_set, type = "class")
319
320   confusionMatrix(full_test_set$diabetes, pred_dt_full)
321
322   |
323   # Feature Selection Model
324   fs_model_dt <- rpart(diabetes ~ ., data = fs_train_set, method = "class")
325
326   rpart.plot(fs_model_dt)
327
328   importances <- varImp(fs_model_dt)
329   importances %>% arrange(desc(Overall))
330
331   pred_dt_fs <- predict(fs_model_dt, newdata = fs_test_set, type = "class")
332
333   confusionMatrix(fs_test_set$diabetes, pred_dt_fs)
334
```

```
> confusionMatrix(fs_test_set$diabetes, pred_dt_fs)
Confusion Matrix and Statistics

          Reference
Prediction    0    1
         0 2217  335
         1  187 2363

               Accuracy : 0.8977
                 95% CI : (0.889, 0.9059)
    No Information Rate : 0.5288
    P-Value [Acc > NIR] : < 2.2e-16

                  Kappa : 0.7954

 Mcnemar's Test P-Value : 1.243e-10

            Sensitivity : 0.9222
            Specificity : 0.8758
         Pos Pred Value : 0.8687
         Neg Pred Value : 0.9267
             Prevalence : 0.4712
         Detection Rate : 0.4345
   Detection Prevalence : 0.5002
      Balanced Accuracy : 0.8990

       'Positive' Class : 0
```

# Models Analysis - Decision Tree



Decision Tree Full Model

Decision Tree Feature Selection Model

# Models Analysis - Naive Bayes

```
339  #######################################################################
340  ########################### NAIVE BAYES ###########################
341  #######################################################################
342
343  # Full Model
344  full_model_nb <- naiveBayes(diabetes ~ ., data = full_train_set)
345  full_model_nb
346
347  pred_naive_full <- predict(full_model_nb, newdata = full_test_set)
348  table_naive_full <- table(full_test_set$diabetes, pred_naive_full)
349
350  confusionMatrix(table_naive_full)
351
352  # Feature Selection Model
353  fs_model_nb <- naiveBayes(diabetes ~ ., data = fs_train_set)
354
355  pred_naive_fs <- predict(fs_model_nb, newdata = fs_test_set)
356  table_naive_fs <- table(fs_test_set$diabetes, pred_naive_fs)
357
358  confusionMatrix(table_naive_fs)
359
```

```
> confusionMatrix(table_naive_full)
Confusion Matrix and Statistics

    pred_naive_full
        0    1
  0  2171  381
  1   547 2003

               Accuracy : 0.8181
                 95% CI : (0.8072, 0.8286)
    No Information Rate : 0.5327
    P-Value [Acc > NIR] : < 2.2e-16

                  Kappa : 0.6362

 Mcnemar's Test P-Value : 6.081e-08

            Sensitivity : 0.7987
            Specificity : 0.8402
         Pos Pred Value : 0.8507
         Neg Pred Value : 0.7855
             Prevalence : 0.5327
         Detection Rate : 0.4255
   Detection Prevalence : 0.5002
      Balanced Accuracy : 0.8195

       'Positive' Class : 0
```
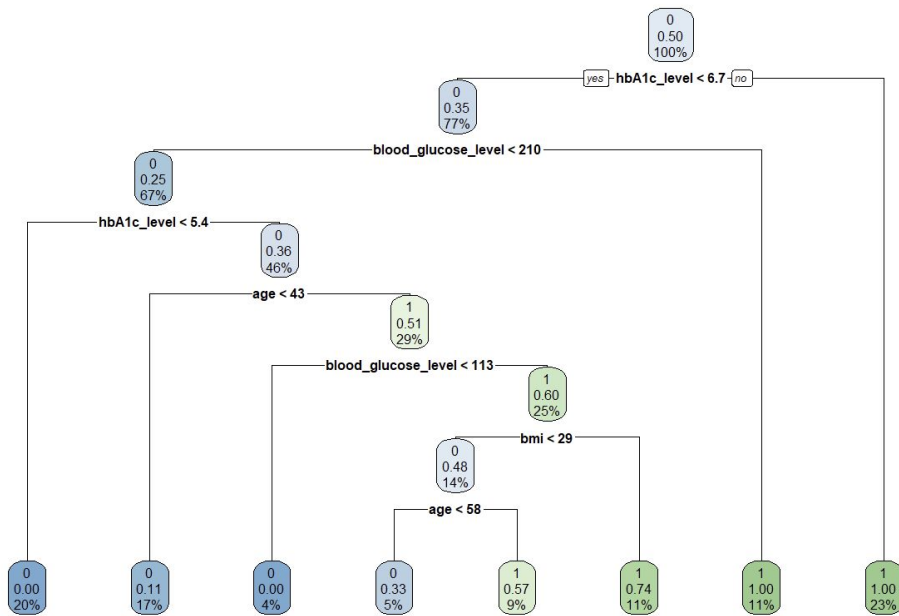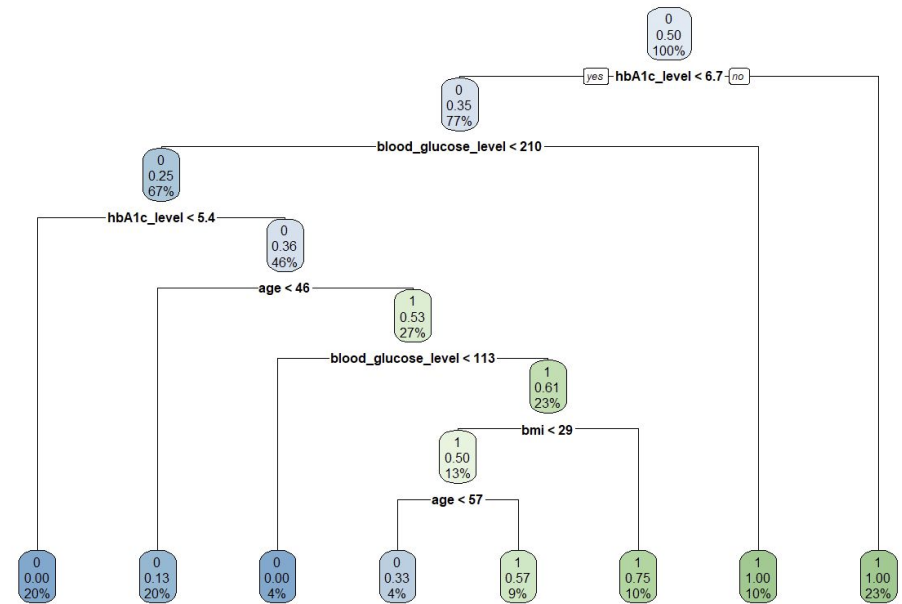
```
> confusionMatrix(table_naive_fs)
Confusion Matrix and Statistics

    pred_naive_fs
        0    1
  0  2245  307
  1   538 2012

               Accuracy : 0.8344
                 95% CI : (0.8239, 0.8445)
    No Information Rate : 0.5455
    P-Value [Acc > NIR] : < 2.2e-16

                  Kappa : 0.6687

 Mcnemar's Test P-Value : 2.528e-15

            Sensitivity : 0.8067
            Specificity : 0.8676
         Pos Pred Value : 0.8797
         Neg Pred Value : 0.7890
             Prevalence : 0.5455
         Detection Rate : 0.4400
   Detection Prevalence : 0.5002
      Balanced Accuracy : 0.8371

       'Positive' Class : 0
```

# Models Analysis - K-Nearest Neighbors

```
366 ########################################################################
367 ########################## K-NEAREST NEIGHBORS #########################
368 ########################################################################
369
370 # Full Model
371 full_train_numeric <- full_train_set %>% select(where(is.numeric))
372 full_test_numeric  <- full_test_set %>% select(where(is.numeric))
373
374 full_classifier_knn <- knn(train = full_train_numeric,
375                            test = full_test_numeric,
376                            cl = full_train_set$diabetes,
377                            k = 3)
378
379 cm_knn_full <- table(full_test_set$diabetes, full_classifier_knn)
380 cm_knn_full
381
382 confusionMatrix(cm_knn_full)
383
384 # Feature Selection Model
385 fs_train_numeric <- fs_train_set %>% select(where(is.numeric))
386 fs_test_numeric  <- fs_test_set %>% select(where(is.numeric))
387
388 fs_classifier_knn <- knn(train = fs_train_numeric,
389                          test = fs_test_numeric,
390                          cl = fs_train_set$diabetes,
391                          k = 3)
392
393 cm_knn_fs <- table(fs_test_set$diabetes, fs_classifier_knn)
394 cm_knn_fs
395
396 confusionMatrix(cm_knn_fs)
```

```
> confusionMatrix(cm_knn_full)
Confusion Matrix and Statistics

       full_classifier_knn
          0    1
     0 2160  392
     1  375 2175

               Accuracy : 0.8497
                 95% CI : (0.8396, 0.8594)
    No Information Rate : 0.5031
    P-Value [Acc > NIR] : <2e-16

                  Kappa : 0.6993

 Mcnemar's Test P-Value : 0.5634

            Sensitivity : 0.8521
            Specificity : 0.8473
         Pos Pred Value : 0.8464
         Neg Pred Value : 0.8529
             Prevalence : 0.4969
         Detection Rate : 0.4234
   Detection Prevalence : 0.5002
      Balanced Accuracy : 0.8497

       'Positive' Class : 0
```

```
> confusionMatrix(cm_knn_fs)
Confusion Matrix and Statistics

       fs_classifier_knn
          0    1
     0 2200  352
     1  402 2148

               Accuracy : 0.8522
                 95% CI : (0.8422, 0.8618)
    No Information Rate : 0.51
    P-Value [Acc > NIR] : < 2e-16

                  Kappa : 0.7044

 Mcnemar's Test P-Value : 0.07435

            Sensitivity : 0.8455
            Specificity : 0.8592
         Pos Pred Value : 0.8621
         Neg Pred Value : 0.8424
             Prevalence : 0.5100
         Detection Rate : 0.4312
   Detection Prevalence : 0.5002
      Balanced Accuracy : 0.8524

       'Positive' Class : 0
```

# Models Analysis - Support Vector Machine

```
405 * ################################################################################
406 * ######################### SUPPORT VECTOR MACHINE ###############################
407 * ################################################################################
408
409   # Full Model
410   full_train_set_svm <- full_train_set
411   full_train_set_svm$diabetes <- factor(full_train_set_svm$diabetes, levels=c(0,1))
412   full_test_set_svm <- full_test_set
413   full_test_set_svm$diabetes  <- factor(full_test_set_svm$diabetes, levels=c(0,1))
414
415   full_model_svm <- svm(diabetes ~ ., data = full_train_set, kernel = "linear")
416
417   full_pred_svm <- predict(full_model_svm, full_test_set)
418
419   confusionMatrix(data=full_pred_svm, reference=full_test_set$diabetes)
420
421   # Feature Selection Model
422   fs_train_set_svm <- fs_train_set
423   fs_train_set_svm$diabetes <- factor(fs_train_set_svm$diabetes, levels=c(0,1))
424   fs_test_set_svm <- fs_test_set
425   fs_test_set_svm$diabetes  <- factor(fs_test_set_svm$diabetes, levels=c(0,1))
426
427   fs_model_svm <- svm(diabetes ~ ., data = fs_train_set, kernel = "linear")
428
429   fs_pred_svm <- predict(fs_model_svm, fs_test_set)
430
431   confusionMatrix(data=fs_pred_svm, reference=fs_test_set$diabetes)
```

```
> confusionMatrix(data=full_pred_svm, reference=full_test_set$diabetes)
Confusion Matrix and Statistics

          Reference
Prediction    0    1
         0 2282  313
         1  270 2237

               Accuracy : 0.8857
                 95% CI : (0.8767, 0.8943)
    No Information Rate : 0.5002
    P-Value [Acc > NIR] : < 2e-16

                  Kappa : 0.7715

 Mcnemar's Test P-Value : 0.08195

            Sensitivity : 0.8942
            Specificity : 0.8773
         Pos Pred Value : 0.8794
         Neg Pred Value : 0.8923
             Prevalence : 0.5002
         Detection Rate : 0.4473
   Detection Prevalence : 0.5086
      Balanced Accuracy : 0.8857

       'Positive' Class : 0

> confusionMatrix(data=fs_pred_svm, reference=fs_test_set$diabetes)
Confusion Matrix and Statistics

          Reference
Prediction    0    1
         0 2288  285
         1  264 2265

               Accuracy : 0.8924
                 95% CI : (0.8836, 0.9008)
    No Information Rate : 0.5002
    P-Value [Acc > NIR] : <2e-16

                  Kappa : 0.7848

 Mcnemar's Test P-Value : 0.3933

            Sensitivity : 0.8966
            Specificity : 0.8882
         Pos Pred Value : 0.8892
         Neg Pred Value : 0.8956
             Prevalence : 0.5002
         Detection Rate : 0.4485
   Detection Prevalence : 0.5043
      Balanced Accuracy : 0.8924

       'Positive' Class : 0
```

# Models Analysis - OLS Regression

```
452 ▾ ##################################################################################
453 ▾ ########################### OLS REGRESSION #####################################
454 ▾ ##################################################################################
455
456   # Full Model
457   full_train_set_lm <- full_train_set
458   full_train_set_lm$diabetes <- as.numeric(as.character(full_train_set_lm$diabetes))
459   full_test_set_lm <- full_test_set
460   full_test_set_lm$diabetes <- as.numeric(as.character(full_test_set$diabetes))
461
462   full_ols_fit <- lm(diabetes ~., data=full_train_set_lm)
463   full_pred_ols <- predict(full_ols_fit, newdata = full_test_set_lm)
464
465   # Convert numeric predictions to 0/1 (threshold = 0.5)
466   full_ols_class <- ifelse(full_pred_ols > 0.5, 1, 0)
467   full_ols_class <- factor(full_ols_class, levels = c(0,1))
468
469   # True labels as factor
470   full_y_test_factor <- factor(full_test_set_lm$diabetes, levels = c(0,1))
471
472   # Confusion matrix
473   confusionMatrix(data = full_ols_class, reference = full_y_test_factor)
474
475   # Feature Selection Model
476   fs_train_set_lm <- fs_train_set
477   fs_train_set_lm$diabetes <- as.numeric(as.character(fs_train_set_lm$diabetes))
478   fs_test_set_lm <- fs_test_set
479   fs_test_set_lm$diabetes <- as.numeric(as.character(fs_test_set$diabetes))
480
481   fs_ols_fit <- lm(diabetes ~., data=fs_train_set_lm)
482   fs_pred_ols <- predict(fs_ols_fit, newdata = fs_test_set_lm)
483
484   # Convert numeric predictions to 0/1 (threshold = 0.5)
485   fs_ols_class <- ifelse(fs_pred_ols > 0.5, 1, 0)
486   fs_ols_class <- factor(fs_ols_class, levels = c(0,1))
487
488   # True labels as factor
489   fs_y_test_factor <- factor(fs_test_set_lm$diabetes, levels = c(0,1))
490
491   # Confusion matrix
492   confusionMatrix(data = fs_ols_class, reference = fs_y_test_factor)
```

```
> confusionMatrix(data = full_ols_class, reference = full_y_test_factor)
Confusion Matrix and Statistics

          Reference
Prediction    0    1
         0 2308  333
         1  244 2217

               Accuracy : 0.8869
                 95% CI : (0.8779, 0.8955)
    No Information Rate : 0.5002
    P-Value [Acc > NIR] : < 2.2e-16

                  Kappa : 0.7738

 Mcnemar's Test P-Value : 0.0002488

            Sensitivity : 0.9044
            Specificity : 0.8694
         Pos Pred Value : 0.8739
         Neg Pred Value : 0.9009
             Prevalence : 0.5002
         Detection Rate : 0.4524
   Detection Prevalence : 0.5176
      Balanced Accuracy : 0.8869

       'Positive' Class : 0

> confusionMatrix(data = fs_ols_class, reference = fs_y_test_factor)
Confusion Matrix and Statistics

          Reference
Prediction    0    1
         0 2298  312
         1  254 2238

               Accuracy : 0.8891
                 95% CI : (0.8801, 0.8976)
    No Information Rate : 0.5002
    P-Value [Acc > NIR] : < 2e-16

                  Kappa : 0.7781

 Mcnemar's Test P-Value : 0.01658

            Sensitivity : 0.9005
            Specificity : 0.8776
         Pos Pred Value : 0.8805
         Neg Pred Value : 0.8981
             Prevalence : 0.5002
         Detection Rate : 0.4504
   Detection Prevalence : 0.5116
      Balanced Accuracy : 0.8891

       'Positive' Class : 0
```

# Models Analysis - GLM Regression

```
508 ▾ ######################################################################
509 ▾ ############################# GLM REGRESSION ##########################
510 ▾ ######################################################################
511
512   # Full Model
513   full_model_glm <- glm(diabetes~., data=full_train_set_lm, family=binomial())
514
515   full_pred_glm <- predict(full_model_glm, newdata=full_test_set_lm, type="response")
516
517   full_pred_class_glm <- ifelse(full_pred_glm > 0.5, 1, 0)
518   full_pred_class_glm <- factor(full_pred_class_glm, levels=c(0,1))
519
520   # True labels
521   full_y_test_factor <- factor(full_test_set_lm$diabetes, levels=c(0,1))
522
523   # Confusion matrix
524   confusionMatrix(data=full_pred_class_glm, reference=full_y_test_factor)
525
526
527   # Feature Selection Model
528
529   fs_model_glm <- glm(diabetes~., data=fs_train_set_lm, family=binomial())
530
531   fs_pred_glm <- predict(fs_model_glm, newdata=fs_test_set_lm, type="response")
532
533   fs_pred_class_glm <- ifelse(fs_pred_glm > 0.5, 1, 0)
534   fs_pred_class_glm <- factor(fs_pred_class_glm, levels=c(0,1))
535
536   # True labels
537   fs_y_test_factor <- factor(fs_test_set_lm$diabetes, levels=c(0,1))
538
539   # Confusion matrix
540   confusionMatrix(data=fs_pred_class_glm, reference=fs_y_test_factor)
```

```
> confusionMatrix(data=full_pred_class_glm, reference=full_y_test_factor)
Confusion Matrix and Statistics

          Reference
Prediction    0    1
         0 2274  315
         1  278 2235

               Accuracy : 0.8838
                 95% CI : (0.8747, 0.8924)
    No Information Rate : 0.5002
    P-Value [Acc > NIR] : <2e-16

                  Kappa : 0.7675

 Mcnemar's Test P-Value : 0.1393

            Sensitivity : 0.8911
            Specificity : 0.8765
         Pos Pred Value : 0.8783
         Neg Pred Value : 0.8894
             Prevalence : 0.5002
         Detection Rate : 0.4457
   Detection Prevalence : 0.5074
      Balanced Accuracy : 0.8838

       'Positive' Class : 0


> confusionMatrix(data=fs_pred_class_glm, reference=fs_y_test_factor)
Confusion Matrix and Statistics

          Reference
Prediction    0    1
         0 2275  292
         1  277 2258

               Accuracy : 0.8885
                 95% CI : (0.8795, 0.897)
    No Information Rate : 0.5002
    P-Value [Acc > NIR] : <2e-16

                  Kappa : 0.7769

 Mcnemar's Test P-Value : 0.5573

            Sensitivity : 0.8915
            Specificity : 0.8855
         Pos Pred Value : 0.8862
         Neg Pred Value : 0.8907
             Prevalence : 0.5002
         Detection Rate : 0.4459
   Detection Prevalence : 0.5031
      Balanced Accuracy : 0.8885

       'Positive' Class : 0
```

# Models Analysis - Bagging

```
586 ▾ ##############################################################################
587 ▾ ############################## BAGGING ######################################
588 ▾ ##############################################################################
589
590    # Full Model
591    full_model_bag <- bagging(formula = diabetes ~ .,
592                              data=full_train_set, nbagg = 50, coob=TRUE,
593                              control=rpart.control(minsplit=2, cp=0, min_depth=2))
594
595    full_pred_bag <- predict(full_model_bag, newdata = full_test_set)
596
597    confusionMatrix(full_test_set$diabetes, full_pred_bag)
598
599    # Feature Selection Model
600
601    fs_model_bag <- bagging(formula = diabetes ~ .,
602                            data=fs_train_set, nbagg = 50, coob=TRUE,
603                            control=rpart.control(minsplit=2, cp=0, min_depth=2))
604
605    fs_pred_bag <- predict(fs_model_bag, newdata = fs_test_set)
606
607    confusionMatrix(full_test_set$diabetes, fs_pred_bag)
608
```

```
> confusionMatrix(full_test_set$diabetes, full_pred_bag)
Confusion Matrix and Statistics

          Reference
Prediction    0    1
         0 2278  274
         1  269 2281

               Accuracy : 0.8936
                 95% CI : (0.8848, 0.9019)
    No Information Rate : 0.5008
    P-Value [Acc > NIR] : <2e-16

                  Kappa : 0.7871

 Mcnemar's Test P-Value : 0.8637

            Sensitivity : 0.8944
            Specificity : 0.8928
         Pos Pred Value : 0.8926
         Neg Pred Value : 0.8945
             Prevalence : 0.4992
         Detection Rate : 0.4465
   Detection Prevalence : 0.5002
      Balanced Accuracy : 0.8936

       'Positive' Class : 0

> confusionMatrix(fs_test_set$diabetes, fs_pred_bag)
Confusion Matrix and Statistics

          Reference
Prediction    0    1
         0 2287  265
         1  236 2314

               Accuracy : 0.9018
                 95% CI : (0.8933, 0.9098)
    No Information Rate : 0.5055
    P-Value [Acc > NIR] : <2e-16

                  Kappa : 0.8036

 Mcnemar's Test P-Value : 0.211

            Sensitivity : 0.9065
            Specificity : 0.8972
         Pos Pred Value : 0.8962
         Neg Pred Value : 0.9075
             Prevalence : 0.4945
         Detection Rate : 0.4483
   Detection Prevalence : 0.5002
      Balanced Accuracy : 0.9019

       'Positive' Class : 0
```

# Models Analysis - XGBoost

```
634 ▾ ####################################################################################
635 ▾ ############################### XGBOOST ###########################################
636 ▾ ####################################################################################
637   full_x_train <- as.matrix(full_train_set[, setdiff(names(full_train_set), "diabetes")])
638   full_y_train <- as.numeric(as.character(full_train_set$diabetes))
639   full_x_test <- as.matrix(full_test_set[, setdiff(names(full_test_set), "diabetes")])
640   full_y_test  <- as.numeric(as.character(full_test_set$diabetes))
641
642   fs_x_train <- as.matrix(fs_train_set[, setdiff(names(fs_train_set), "diabetes")])
643   fs_y_train   <- as.numeric(as.character(fs_train_set$diabetes))
644   fs_x_test <- as.matrix(fs_test_set[, setdiff(names(fs_test_set), "diabetes")])
645   fs_y_test    <- as.numeric(as.character(fs_test_set$diabetes))
646
647   unique(full_y_train)
648   unique(full_y_test)
649
650   # Full Model
651   full_model_boost <- xgboost::xgboost(data=full_x_train, label=full_y_train,
652                      max.depth=2, eta=0.3, nthread=2, nrounds=2,
653                      objective="binary:logistic")
654
655   full_pred_boost <- predict(full_model_boost, newdata = full_x_test)
656
657   full_pred_class <- ifelse(full_pred_boost > 0.5, 1, 0)
658   full_pred_class <- factor(full_pred_class, levels = c(0,1))
659   full_y_test_factor <- factor(full_y_test, levels = c(0,1))
660
661   confusionMatrix(data = full_pred_class, reference = full_y_test_factor)
662
663   # Feature Selection Model
664   fs_model_boost <- xgboost::xgboost(data=fs_x_train, label=fs_y_train,
665                      max.depth=2, eta=0.3, nthread=2, nrounds=2,
666                      objective="binary:logistic")
667
668   fs_pred_boost <- predict(fs_model_boost, newdata = fs_x_test)
669
670   fs_pred_class <- ifelse(fs_pred_boost > 0.5, 1, 0)
671   fs_pred_class <- factor(fs_pred_class, levels = c(0,1))
672   fs_y_test_factor <- factor(fs_y_test, levels = c(0,1))
673
674   confusionMatrix(data = fs_pred_class, reference = fs_y_test_factor)
675
```

```
> confusionMatrix(data = full_pred_class, reference = full_y_test_factor)
Confusion Matrix and Statistics

          Reference
Prediction    0    1
         0 2552  846
         1    0 1704

               Accuracy : 0.8342
                 95% CI : (0.8237, 0.8443)
    No Information Rate : 0.5002
    P-Value [Acc > NIR] : < 2.2e-16

                  Kappa : 0.6683

 Mcnemar's Test P-Value : < 2.2e-16

            Sensitivity : 1.0000
            Specificity : 0.6682
         Pos Pred Value : 0.7510
         Neg Pred Value : 1.0000
             Prevalence : 0.5002
         Detection Rate : 0.5002
   Detection Prevalence : 0.6660
      Balanced Accuracy : 0.8341

       'Positive' Class : 0

> confusionMatrix(data = fs_pred_class, reference = fs_y_test_factor)
Confusion Matrix and Statistics

          Reference
Prediction    0    1
         0 2552  836
         1    0 1714

               Accuracy : 0.8361
                 95% CI : (0.8257, 0.8462)
    No Information Rate : 0.5002
    P-Value [Acc > NIR] : < 2.2e-16

                  Kappa : 0.6722

 Mcnemar's Test P-Value : < 2.2e-16

            Sensitivity : 1.0000
            Specificity : 0.6722
         Pos Pred Value : 0.7532
         Neg Pred Value : 1.0000
             Prevalence : 0.5002
         Detection Rate : 0.5002
   Detection Prevalence : 0.6641
      Balanced Accuracy : 0.8361

       'Positive' Class : 0
```

# Models Analysis - Random Forest

```
546 ▾ ###################################################################################
547 ▾ ############################## RANDOM FOREST #####################################
548 ▾ ###################################################################################
549
550   # Full Model
551   full_model_RF <- randomForest(diabetes~., data=full_train_set,
552                                  ntree=500, ntry=6, importance=TRUE,
553                                  na.action = randomForest::na.roughfix, replace = FALSE)
554   varImpPlot(full_model_RF, col=3)
555
556   full_pred_RF <- predict(full_model_RF, newdata = full_test_set)
557
558   confusionMatrix(full_test_set$diabetes, full_pred_RF)
559
560   # Feature Selection Model
561   fs_model_RF <- randomForest(diabetes~., data=fs_train_set,
562                                ntree=500, ntry=6, importance=TRUE,
563                                na.action = randomForest::na.roughfix, replace = FALSE)
564   varImpPlot(fs_model_RF, col=3)
565
566   fs_pred_RF <- predict(fs_model_RF, newdata = fs_test_set)
567
568   confusionMatrix(fs_test_set$diabetes, fs_pred_RF)
```

```
> confusionMatrix(full_test_set$diabetes, full_pred_RF)
Confusion Matrix and Statistics

          Reference
Prediction    0    1
         0 2334  218
         1  254 2296

               Accuracy : 0.9075
                 95% CI : (0.8992, 0.9153)
    No Information Rate : 0.5073
    P-Value [Acc > NIR] : <2e-16

                  Kappa : 0.815

 Mcnemar's Test P-Value : 0.1072

            Sensitivity : 0.9019
            Specificity : 0.9133
         Pos Pred Value : 0.9146
         Neg Pred Value : 0.9004
             Prevalence : 0.5073
         Detection Rate : 0.4575
   Detection Prevalence : 0.5002
      Balanced Accuracy : 0.9076

       'Positive' Class : 0

> confusionMatrix(fs_test_set$diabetes, fs_pred_RF)
Confusion Matrix and Statistics

          Reference
Prediction    0    1
         0 2314  238
         1  229 2321

               Accuracy : 0.9085
                 95% CI : (0.9002, 0.9162)
    No Information Rate : 0.5016
    P-Value [Acc > NIR] : <2e-16

                  Kappa : 0.8169

 Mcnemar's Test P-Value : 0.7112

            Sensitivity : 0.9099
            Specificity : 0.9070
         Pos Pred Value : 0.9067
         Neg Pred Value : 0.9102
             Prevalence : 0.4984
         Detection Rate : 0.4535
   Detection Prevalence : 0.5002
      Balanced Accuracy : 0.9085

       'Positive' Class : 0
```
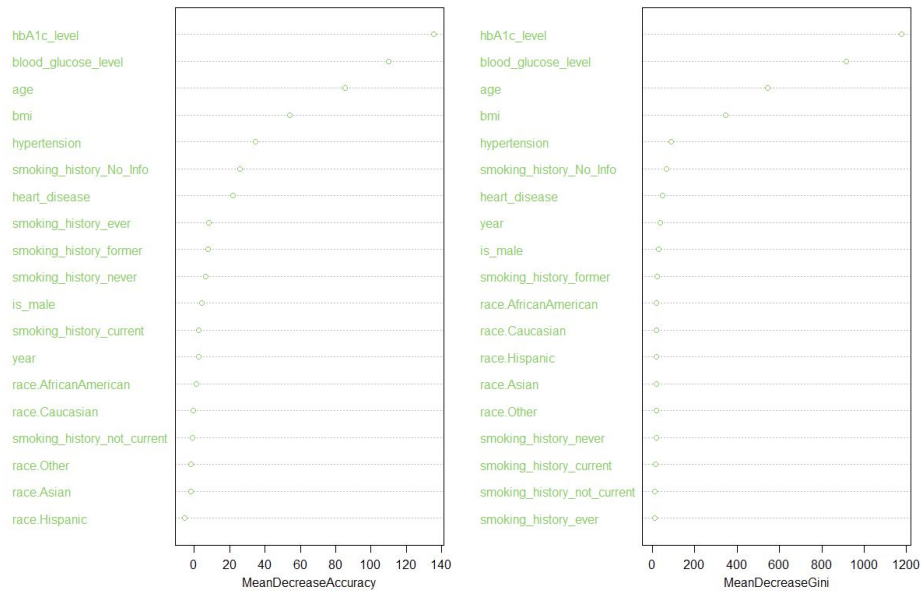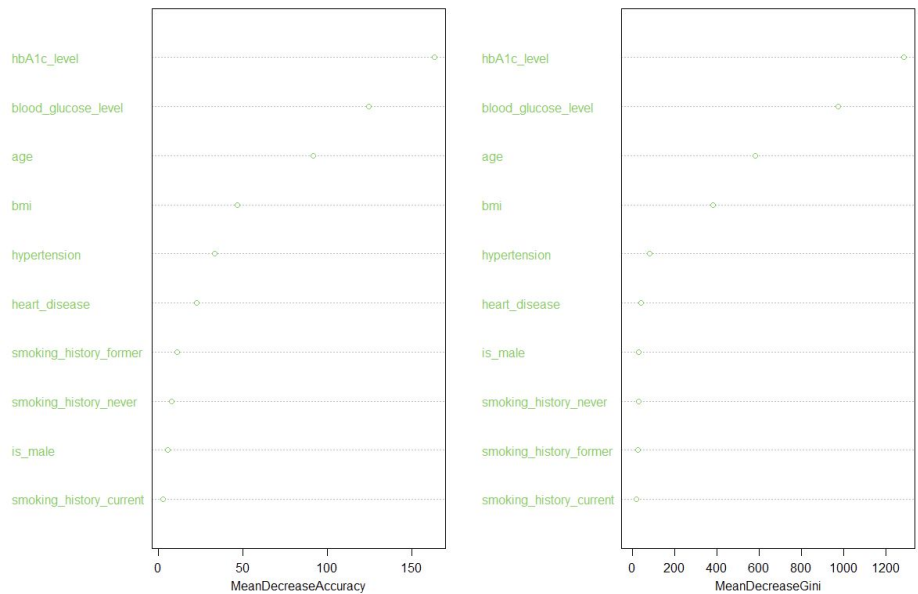
# Models Analysis - Random Forest

# Model Comparison

| Model | Accuracy | Kappa |
|---|---|---|
| Decision Tree - Full | 89.98% | 0.7997 |
| Decision Tree - FS | 89.77% | 0.7954 |
| Naive Bayes - Full | 81.81% | 0.6362 |
| Naive Bayes - FS | 83.44% | 0.6687 |
| KNN - Full | 84.97% | 0.6993 |
| KNN - FS | 85.22% | 0.7044 |
| SVM - Full | 88.57% | 0.7715 |
| SVM - FS | 89.24% | 0.7848 |
| Linear Regression - Full | 88.69% | 0.7738 |
| Linear Regression - FS | 88.91% | 0.7781 |

| Model | Accuracy | Kappa |
|---|---|---|
| Logistic Regression - Full | 88.38% | 0.7675 |
| GLM - FS | 88.85% | 0.7769 |
| RF - Full | 90.75% | 0.815 |
| RF - FS | 90.85% | 0.8169 |
| Bagging - Full | 89.36% | 0.7871 |
| Bagging - FS | 90.18% | 0.8036 |
| XGBoost - Full | 83.42% | 0.6683 |
| XGBoost - FS | 83.61% | 0.6722 |

# Research Findings

- Random Forest was the model with best accuracy (~91%)

- The model describes correctly 82% of the data according to the kappa metric

- Accordingly to the feature importance plot, hemoglobin A1c and Blood glucose level heavy indicates if the patience has diabetes or not, followed by age and bmi

- It provides a good balance for between catching disease cases while minimizing unnecessary follow-up tests

# Research Findings

- The model performance is statistically significant as it has extremely low p-value

- The model has a good recall value which is important for this study case where we want to avoid false negatives as much as possible

- Random Forest was the best one due to its capacity to handling complex non-linear relationships, less hyperparameter sensitivity, and robust to overfitting

# Bibliography

- Classes Notes

- https://www.geeksforgeeks.org/machine-learning/what-are-the-advantages-and-disadvantages-of-random-forest/

- https://www.geeksforgeeks.org/machine-learning/confusion-matrix-machine-learning/

Group Members:

'Aafrin Zahid Memon
Akintunde Akinro
Bruno do Nascimento Beserra