```r
1    ########### Import library necessary
     ######################################################################
     #######
2
3    library(ggplot2)
4    library(tidyverse)
5    library(GGally)
6    library(corrplot)
7    library(nortest)
8    library(pastecs)
9    library(data.table)
10   library(caret)
11   library(dataPreparation)
12   library(caTools)
13   library(microbenchmark)
14   library(ggpubr)
15
16   ########### Import dataset
     ######################################################################
     #################
17   raw <- read.csv('cardio_train.csv',sep = ";")
18
19   ########### Handle Missing Value and Data Formating
     ######################################################################
20
21   # Check for any missing values
22   colSums(is.na(raw))
23
24   # Define categorical list
25   cat.list <- c("id","gender","cholesterol","gluc","smoke","alco","active","cardio")
26
27   # Convert numerical columns to categorical data
28   raw[cat.list ] <- lapply(raw[cat.list ], factor)
29
30   # Define numerical list
31   num.list <- c("age","height","weight","ap_hi","ap_lo")
32
33   # Reduce "age" measurement unit from day to year
34   raw <- mutate(raw, age = age / 365.25)
35
36   ########### Exploratory Data Analysis
     ######################################################################
     ######
37
38   # Check structure of dataframe
39   str(raw)
40
41   # Check descriptive statistic of data
42   summary(raw)
43
44   # Univariate Plots - Plot and export boxplot as jpeg
45   plot.box <- function(x,y){
46     plot <- boxplot(x, xlab=y) + theme_minimal()
47     return(plot)
48   }
49
50   for (i in 1:length(num.list)){
51     jpeg(file=sprintf("plot_box_%s.jpeg",num.list[i]))
52     print(plot.box(raw[num.list][,i],num.list[i]))
53     dev.off()
54     print(i)}
55
56   # Remove outliers of feature "ap_hi" and "ap_lo" and redraw the box plot
57   outliers4 <- boxplot(raw$ap_hi, plot=FALSE)$out
58   raw <- raw[-which(raw$ap_hi %in% outliers4),]
59   boxplot(raw$ap_hi,xlab="ap_hi")
60   outliers5 <- boxplot(raw$ap_lo, plot=FALSE)$out
61   raw <- raw[-which(raw$ap_lo %in% outliers5),]
62   boxplot(raw$ap_lo,xlab="ap_lo")
63
64   # Check descriptive statistic of data again after removed outlier in ap_lo and ap_hi
65   summary(raw[,2:13])
66
```

```r
67    # Check Distribution - Plot and export histogram of numerical variables as jpeg
68    plot.hist <- function(x,y){
69      plot <- ggplot(raw,aes(x)) + geom_histogram(fill='blue',bins = 20,alpha=0.5)+
        labs(x = y) + theme_minimal()
70      return(plot)
71    }
72
73    for (i in 1:length(num.list)){
74      jpeg(file=sprintf("plot_hist_%s.jpeg",num.list[i]))
75      print(plot.hist(raw[num.list][,i],num.list[i]))
76      dev.off()
77      print(i)}
78
79    # Univariate Plots - Plot and export bar chart of categorical variables as jpeg
80    plot.bar <- function(x,y){
81      plot <- ggplot(raw,aes(x)) + geom_bar(aes(fill=factor(x)),alpha=0.5)+ labs(x = y)
        + theme_minimal()+labs(fill = y)
82      return(plot)
83    }
84
85    for (i in 2:length(cat.list)){
86      jpeg(file=sprintf("plot_bar_%s.jpeg",cat.list[i]))
87      print(plot.bar(raw[cat.list][,i],cat.list[i]))
88      dev.off()
89      print(i)}
90
91    # Multivariate Plots - Plot and export scatter plot as jpeg
92    plot.sca <- function(x,y,z,m,n,p){
93      plot <- ggplot(raw,aes(x=x,y=y)) + geom_point(aes(color=factor(z)),alpha=0.5)+
        theme_minimal() + labs(x=m, y=n,col=p)
94      return(plot)
95    }
96
97    for (i in 1:length(num.list)){
98      for (l in 1:length(num.list)){
99        if (num.list[i] != num.list[l]){
100         jpeg(file=sprintf("plot_scatter_%s_%s.jpeg",num.list[i],num.list[l]))
101
          print(plot.sca(raw[num.list][,i],raw[num.list][,l],raw$cardio,num.list[i],num.li
          st[l],"cardio"))
102         dev.off()
103         print(i)
104        }
105      }
106    }
107
108    # Plot and export correlation matrix as jpeg
109    num.cols <- sapply(raw, is.numeric)
110    cor.data <- cor(raw[,num.cols])
111    jpeg(file="corr.jpeg")
112    corrplot(cor.data,method = "number")
113    dev.off()
114
115    ############# Train and Test Set
       ##############################################################################
       ########
116
117    set.seed(123)
118    raw_train <- raw[,2:ncol(raw)]
119    split = sample.split(raw_train $cardio, SplitRatio = 0.75)
120    training_set = subset(raw_train, split == TRUE)
121    test_set = subset(raw_train, split == FALSE)
122
123    summary(training_set)
124
125    X_train = training_set[,1:ncol(raw_train)-1]
126    y_train = training_set[,ncol(raw_train)]
127    X_test = test_set[,1:ncol(raw_train)-1]
128    y_test = test_set[,ncol(raw_train)]
129
130    # Scaling
131    scales <- build_scales(dataSet = training_set, cols = num.list, verbose = TRUE)
132    training_set_s <- fastScale(dataSet = training_set, scales = scales, verbose = TRUE)
```

```r
133    test_set_s <- fastScale(dataSet = test_set, scales = scales, verbose = TRUE)
134    X_train_s <- fastScale(dataSet = X_train, scales = scales, verbose = TRUE)
135    X_test_s <- fastScale(dataSet = X_test, scales = scales, verbose = TRUE)
136
137    summary(training_set_s)
138
139    ############# Test Harness
       ######################################################################################
       #########
140
141    # Run algorithms using 10-fold cross validation
142    control <- trainControl(method="cv", number=10)
143    metric <- "Accuracy"
144
145    ############# Build Models
       ######################################################################################
       #########
146
147    # LDA
148    set.seed(7)
149    fit.lda <- train(cardio~., data=training_set, method="lda", metric=metric,
       trControl=control)
150
151    # CART
152    set.seed(7)
153    fit.cart <- train(cardio~., data=training_set, method="rpart", metric=metric,
       trControl=control)
154
155    # naive bayes
156    set.seed(7)
157    fit.nb <- train(cardio~., data=training_set, method="nb", metric=metric,
       trControl=control)
158
159    # kNN
160    set.seed(7)
161    knn.grid <- expand.grid(k=c(203,253)) # design the parameter tuning grid
162    fit.knn <- train(cardio~., data=training_set_s, method="knn", metric=metric,
       trControl=control, tuneGrid=knn.grid)
163
164    # Random Forest
165    set.seed(7)
166    rf.grid <- expand.grid(mtry=c(2,7,12)) # design the parameter tuning grid
167    fit.rf <- train(cardio~., data=training_set, method="rf", metric=metric,
       trControl=control, tuneGrid=rf.grid)
168
169
170    ############# Measure Training and Testing Time
       ###################################################################################
171
172    # Measure training time
173    knn.grid.final <- expand.grid(k=c(fit.knn$bestTune[[1]]))
174    rf.grid.final <- expand.grid(mtry=c(fit.rf$bestTune[[1]]))
175    mbm_train <- microbenchmark("LDA" = { train1 <- train(cardio~., data=training_set,
       method="lda", metric=metric, trControl=control)},
176                        "CART" = {train2 <- train(cardio~., data=training_set,
                           method="rpart", metric=metric, trControl=control)},
177                        "NB" = {train3 <- train(cardio~., data=training_set,
                           method="nb", metric=metric, trControl=control)},
178                        "KNN" = {train4 <- train(cardio~., data=training_set_s,
                           method="knn", metric=metric, trControl=control,
                           tuneGrid=knn.grid.final)},
179                        "RF" = {train5 <- train(cardio~., data=training_set,
                           method="rf", metric=metric, trControl=control,
                           tuneGrid=rf.grid.final)},
180                        times = 3,unit = "s")
181
182    # Measure testing time
183    mbm_test <- microbenchmark("LDA" = { test1 <- predict(train1, test_set)},
184                        "CART" = {test2 <- predict(train2, test_set)},
185                        "NB" = {test3 <- predict(train3, test_set)},
186                        "KNN" = {test4 <- predict(train4, test_set_s)},
187                        "RF" = {test5 <- predict(train5, test_set)},
188                        times = 3,unit = "s")
```

```r
189
190    ############# Select Best Model
       ##########################################################################
       ####
191
192    # compare accuracy of models
193    results <- resamples(list(lda=train1, cart=train2, nb=train3, kn=train4,rf=train5))
194    summary(results)
195    jpeg(file="model_performance.jpeg")
196    dotplot(results)
197    dev.off()
198
199    # compare training time of models
200    mbm_train
201    jpeg(file="training_time.jpeg")
202    autoplot(mbm_train)
203    dev.off()
204
205    # compare testing time of models
206    mbm_test
207    jpeg(file="testing_time.jpeg")
208    autoplot(mbm_test)
209    dev.off()
210
211    # summarize best model (accuracy)
212    print(train5)
213
214    ############# Make Predictions
       ##########################################################################
       ######
215
216    predictions <- predict(train5, test_set)
217    confusionMatrix(predictions, test_set$cardio)
218
219    ############# Feature Importance
       #########################################################################
220    feature.imp <- varImp(train5)
221    jpeg(file="feature_importance.jpeg")
222    ggplot(data=feature.imp)
223    dev.off()
```