

# EDA\_Classfication.R

abhinavmishra

2022-11-11

```
##           A script for exploratory data analysis, and           -
##           classification training four classifiers               -
## to diagnose heart disease based on the Heart Disease Data Set -

#####
##           Author: Abhinav Mishra                               ##
#####

##           Loading/Installing packages required                 -

#install.packages(c("MLeval", caret", "ggvis",
# "skimr", "tidyverse","ggvis", "e1071", "RColorBrewer"))

library(tidyverse)

## -- Attaching packages ----- tidyverse 1.3.2 --
## v ggplot2 3.4.0      v purrr  0.3.5
## v tibble  3.1.8      v dplyr  1.0.10
## v tidyr   1.2.1      v stringr 1.4.1
## v readr   2.1.3      v forcats 0.5.2
## -- Conflicts ----- tidyverse_conflicts() --
## x dplyr::filter() masks stats::filter()
## x dplyr::lag()    masks stats::lag()

library(skimr)
library(ggvis)

##
## Attaching package: 'ggvis'
##
## The following object is masked from 'package:ggplot2':
##
##     resolution

library(caret)

## Loading required package: lattice
##
## Attaching package: 'caret'
##
## The following object is masked from 'package:purrr':
##
##     lift

library(ggvis)
library(caret)
```

```
library(MLeval)
library(RColorBrewer)
```

```
## Loading data from the file -
```

```
processedWithHeader_cleveland <- read_csv("~/Documents/Freie/IFA/WiSe 22-23/Data Science/Week 2/processedWithHeader_cleveland.csv",
                                           show_col_types = FALSE, na = "?")
```

```
heart_data <- na.omit(data.frame(processedWithHeader_cleveland))
data <- data.frame(processedWithHeader_cleveland)
heart_data$ID <- seq.int(nrow(heart_data))
```

```
## Passing as factors -
```

```
#
heart_data$fbs <- as.factor(heart_data$fbs)
heart_data$restecg <- as.factor(heart_data$restecg)
heart_data$exang <- as.factor(heart_data$exang)
heart_data$slope <- as.factor(heart_data$slope)
#heart_data$cp <- as.factor(heart_data$cp)
```

```
## Data wrangling with preparation -
```

```
heart_data[heart_data$sex == 0, ]$sex <- "F"
heart_data[heart_data$sex == 1, ]$sex <- "M"
heart_data$sex <- as.factor(heart_data$sex)

heart_data[heart_data$goal == 0, ]$goal <- "healthy"
heart_data[heart_data$goal == 1, ]$goal <- "unhealthy"
heart_data[heart_data$goal == 2, ]$goal <- "unhealthy"
heart_data[heart_data$goal == 3, ]$goal <- "unhealthy"
heart_data[heart_data$goal == 4, ]$goal <- "unhealthy"
```

```
write.table(heart_data, file = 'heart_data')
table(heart_data$goal)
```

```
##
## healthy unhealthy
##      160      137
```

```
## Descriptive Statistics + NA values omit -
```

```
str(heart_data)
```

```
## 'data.frame': 297 obs. of 15 variables:
## $ age : num 63 67 67 37 41 56 62 57 63 53 ...
## $ sex : Factor w/ 2 levels "F","M": 2 2 2 2 1 2 1 1 2 2 ...
## $ cp : num 1 4 4 3 2 2 4 4 4 4 ...
## $ trestbps: num 145 160 120 130 130 120 140 120 130 140 ...
## $ chol : num 233 286 229 250 204 236 268 354 254 203 ...
## $ fbs : Factor w/ 2 levels "0","1": 2 1 1 1 1 1 1 1 1 2 ...
## $ restecg : Factor w/ 3 levels "0","1","2": 3 3 3 1 3 1 3 1 3 3 ...
## $ thalach : num 150 108 129 187 172 178 160 163 147 155 ...
## $ exang : Factor w/ 2 levels "0","1": 1 2 2 1 1 1 1 2 1 2 ...
```

```
## $ oldpeak : num  2.3 1.5 2.6 3.5 1.4 0.8 3.6 0.6 1.4 3.1 ...
## $ slope   : Factor w/ 3 levels "1","2","3": 3 2 2 3 1 1 3 1 2 3 ...
## $ ca      : num  0 3 2 0 0 0 2 0 1 0 ...
## $ thal    : num  6 3 7 3 3 3 3 7 7 ...
## $ goal    : chr   "healthy" "unhealthy" "unhealthy" "healthy" ...
## $ ID      : int   1 2 3 4 5 6 7 8 9 10 ...
## - attr(*, "na.action")= 'omit' Named int [1:6] 88 167 193 267 288 303
## ..- attr(*, "names")= chr [1:6] "88" "167" "193" "267" ...
```

```
summary(heart_data)
```

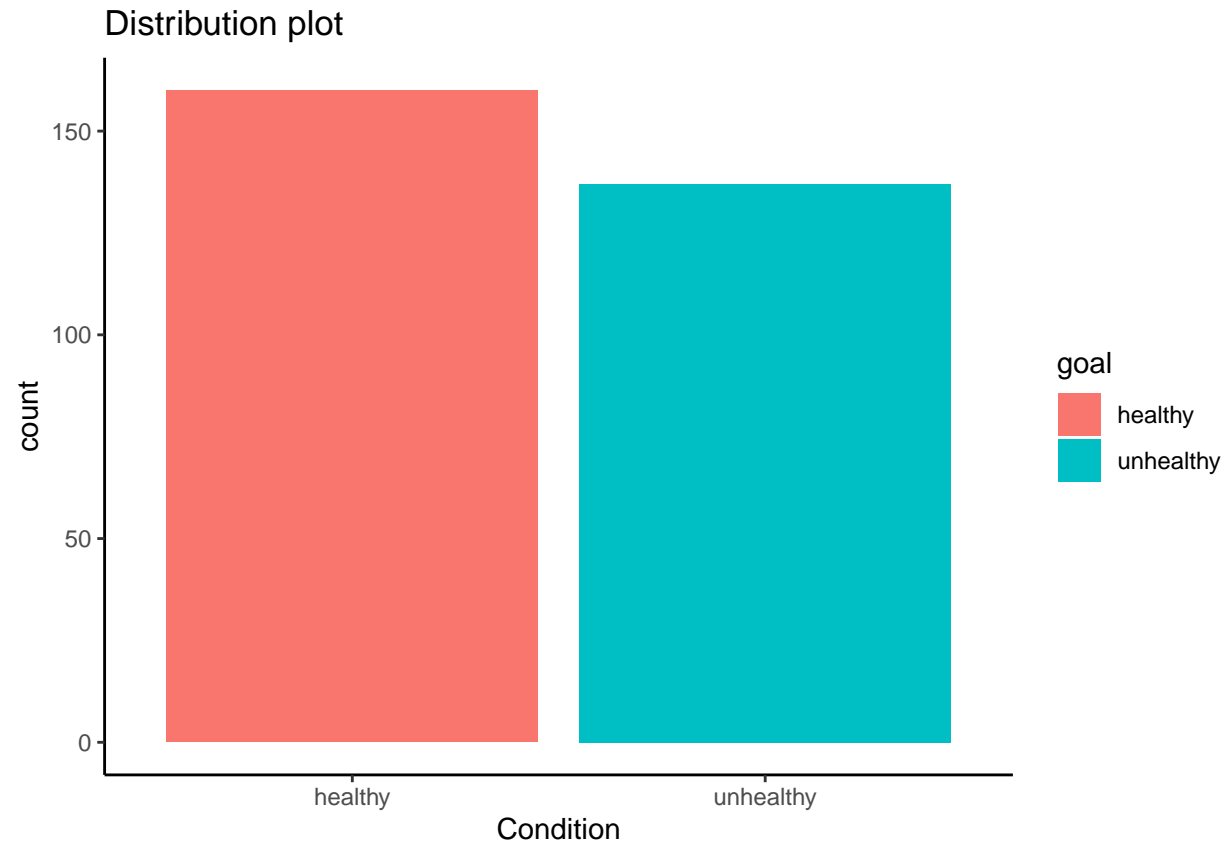
```
##      age      sex      cp      trestbps      chol
## Min.   :29.00  F: 96   Min.   :1.000   Min.   : 94.0   Min.   :126.0
## 1st Qu.:48.00  M:201   1st Qu.:3.000   1st Qu.:120.0   1st Qu.:211.0
## Median :56.00           Median :3.000   Median :130.0   Median :243.0
## Mean   :54.54           Mean   :3.158   Mean   :131.7   Mean   :247.4
## 3rd Qu.:61.00           3rd Qu.:4.000   3rd Qu.:140.0   3rd Qu.:276.0
## Max.   :77.00           Max.   :4.000   Max.   :200.0   Max.   :564.0
## fbs      restecg      thalach      exang      oldpeak      slope
## 0:254    0:147   Min.   : 71.0   0:200   Min.   :0.000   1:139
## 1: 43     1:  4   1st Qu.:133.0   1: 97   1st Qu.:0.000   2:137
##          2:146   Median :153.0           Median :0.800   3: 21
##          Mean   :149.6           Mean   :1.056
##          3rd Qu.:166.0           3rd Qu.:1.600
##          Max.   :202.0           Max.   :6.200
##      ca      thal      goal      ID
## Min.   :0.0000   Min.   :3.000   Length:297   Min.   : 1
## 1st Qu.:0.0000   1st Qu.:3.000   Class :character   1st Qu.: 75
## Median :0.0000   Median :3.000   Mode  :character   Median :149
## Mean   :0.6768   Mean   :4.731           Mean   :149
## 3rd Qu.:1.0000   3rd Qu.:7.000           3rd Qu.:223
## Max.   :3.0000   Max.   :7.000           Max.   :297
```

```
sum(is.na(heart_data))
```

```
## [1] 0
```

```
##      Descriptive plots      -
```

```
ggplot(heart_data, aes(x = goal, fill = goal)) +
  geom_bar() + theme_classic() +
  labs(title='Distribution plot') +
  xlab("Condition")
```



```
table(heart_data$goal)
```

```
##
##   healthy unhealthy
##      160      137
```

```
round(prop.table(table(heart_data$goal)) * 100, digits = 1)
```

```
##
##   healthy unhealthy
##      53.9      46.1
```

```
##           Scatter plots by condition           -
```

```
#heart_data %>% ggvis(~age, ~trestbps, fill = ~goal) %>% layer_points()
```

```
#heart_data %>% ggvis(~age, ~trestbps, fill = ~sex) %>% layer_points()
```

```
#heart_data %>% ggvis(~age, ~trestbps, fill = ~cp) %>% layer_points()
```

```
##           Data partition           -
```

```
test_index <- createDataPartition(y = heart_data$goal, times = 1,
                                   p = 0.2, list= FALSE)
```

```
heart_data$goal <- as.factor(heart_data$goal)
```

```
train_data <- heart_data[-test_index, ]
```

```
test_data <- heart_data[test_index, ]
```

```
##                               Classifiers                               -
##      1. Logistic regression: Fit the logistic regression model,      -
##      that is a GLM using a binomial link using the caret function train() -
```

```
set.seed(112)
log_fit <- train(goal ~.-ID ,
                 data = train_data,
                 method = "glm",
                 family = "binomial")
log_pred <- predict(log_fit, test_data)
confusionMatrix(log_pred, test_data$goal)
```

```
## Confusion Matrix and Statistics
##
##              Reference
## Prediction  healthy unhealthy
## healthy      27          10
## unhealthy     5          18
##
##              Accuracy : 0.75
##              95% CI : (0.6214, 0.8528)
##      No Information Rate : 0.5333
##      P-Value [Acc > NIR] : 0.0004647
##
##              Kappa : 0.4921
##
##  Mcnemar's Test P-Value : 0.3016996
##
##              Sensitivity : 0.8438
##              Specificity : 0.6429
##              Pos Pred Value : 0.7297
##              Neg Pred Value : 0.7826
##              Prevalence : 0.5333
##              Detection Rate : 0.4500
##      Detection Prevalence : 0.6167
##              Balanced Accuracy : 0.7433
##
##      'Positive' Class : healthy
##
```

```
##                               2. Random forest                               -
```

```
set.seed(112)
rf_fit <- train(goal ~.-ID ,
               data = train_data,
               method = "rf")
rf_pred <- predict(rf_fit, test_data)
confusionMatrix(rf_pred, test_data$goal)
```

```
## Confusion Matrix and Statistics
##
##              Reference
## Prediction  healthy unhealthy
## healthy      26           9
```

```
## unhealthy      6      19
##
##           Accuracy : 0.75
##           95% CI : (0.6214, 0.8528)
##       No Information Rate : 0.5333
##       P-Value [Acc > NIR] : 0.0004647
##
##           Kappa : 0.4944
##
## Mcnemar's Test P-Value : 0.6055766
##
##           Sensitivity : 0.8125
##           Specificity : 0.6786
##       Pos Pred Value : 0.7429
##       Neg Pred Value : 0.7600
##           Prevalence : 0.5333
##       Detection Rate : 0.4333
##       Detection Prevalence : 0.5833
##       Balanced Accuracy : 0.7455
##
##       'Positive' Class : healthy
##
```

```
## 3. Boosted logistic regression: using decision stumps (one node decision trees) -
##       as weak learners. It implements a internal version of decision -
##               stump classifier instead of using -
##       calls to rpart. Also, training and testing phases of the -
##       classification process are split into separate functions. -
```

```
set.seed(112)
blog_fit <- train(goal ~.-ID,
                  data = train_data,
                  method = "LogitBoost")
blog_pred <- predict(blog_fit, test_data)
confusionMatrix(blog_pred, test_data$goal)
```

```
## Confusion Matrix and Statistics
##
##           Reference
## Prediction healthy unhealthy
## healthy      25      7
## unhealthy     7     21
##
##           Accuracy : 0.7667
##           95% CI : (0.6396, 0.8662)
##       No Information Rate : 0.5333
##       P-Value [Acc > NIR] : 0.0001655
##
##           Kappa : 0.5313
##
## Mcnemar's Test P-Value : 1.0000000
##
##           Sensitivity : 0.7812
##           Specificity : 0.7500
##       Pos Pred Value : 0.7812
```

```
##          Neg Pred Value : 0.7500
##          Prevalence : 0.5333
##          Detection Rate : 0.4167
##          Detection Prevalence : 0.5333
##          Balanced Accuracy : 0.7656
##
##          'Positive' Class : healthy
##
```

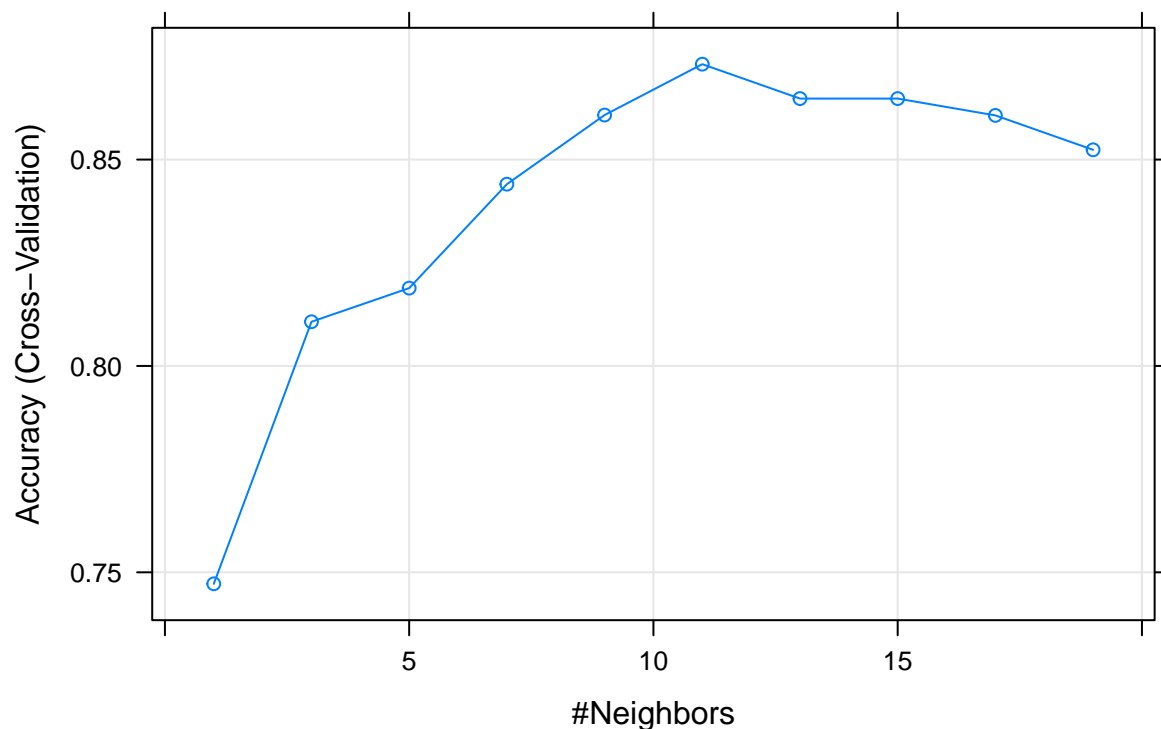
#### ## 4. KNN -

```
ctrl <- trainControl(method = "cv", verboseIter = FALSE, number = 5)

knn_fit <- train(goal ~. -ID , data = train_data,
                 method = "knn", preProcess = c("center","scale"),
                 trControl = ctrl , tuneGrid = expand.grid(k = seq(1, 20, 2)))

plot(knn_fit, main = "K-nearest neighbour")
```

### K-nearest neighbour



```
knn_pred <- predict(knn_fit, test_data)
confusionMatrix(knn_pred, test_data$goal)
```

```
## Confusion Matrix and Statistics
##
##          Reference
## Prediction  healthy unhealthy
## healthy      28      11
## unhealthy     4      17
```

```
##
##           Accuracy : 0.75
##           95% CI : (0.6214, 0.8528)
##      No Information Rate : 0.5333
##      P-Value [Acc > NIR] : 0.0004647
##
##           Kappa : 0.4898
##
## Mcnemar's Test P-Value : 0.1213353
##
##           Sensitivity : 0.8750
##           Specificity : 0.6071
##      Pos Pred Value : 0.7179
##      Neg Pred Value : 0.8095
##           Prevalence : 0.5333
##      Detection Rate : 0.4667
##      Detection Prevalence : 0.6500
##      Balanced Accuracy : 0.7411
##
##      'Positive' Class : healthy
##
```

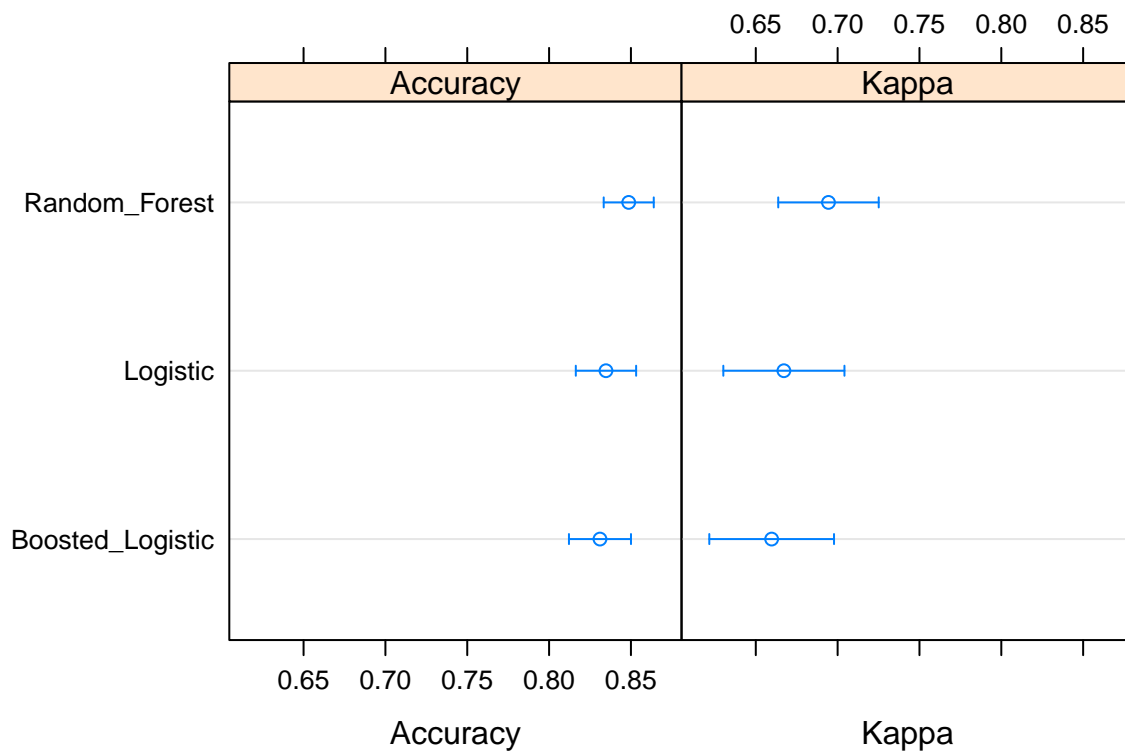
#### ## Performance Comparison -

```
results <- resamples(list(Logistic = log_fit,
                          Random_Forest = rf_fit,
                          Boosted_Logistic = blog_fit))
summary(results)
```

```
##
## Call:
## summary.resamples(object = results)
##
## Models: Logistic, Random_Forest, Boosted_Logistic
## Number of resamples: 25
##
## Accuracy
##           Min.   1st Qu.   Median     Mean   3rd Qu.   Max.
## Logistic      0.7529412 0.8068182 0.8292683 0.8347492 0.8658537 0.9222222
## Random_Forest  0.7888889 0.8160920 0.8461538 0.8486596 0.8764045 0.9176471
## Boosted_Logistic 0.7333333 0.8000000 0.8333333 0.8310928 0.8554217 0.9506173
##           NA's
## Logistic      0
## Random_Forest  0
## Boosted_Logistic 0
##
## Kappa
##           Min.   1st Qu.   Median     Mean   3rd Qu.   Max.
## Logistic      0.4990177 0.5958084 0.6615566 0.6671646 0.7287963 0.8444444
## Random_Forest  0.5761031 0.6270824 0.6923077 0.6944060 0.7463325 0.8321580
## Boosted_Logistic 0.4658754 0.6009852 0.6651660 0.6596931 0.6996016 0.9000000
##           NA's
## Logistic      0
## Random_Forest  0
## Boosted_Logistic 0
```



```
dotplot(results)
```



**Confidence Level: 0.95**

```
cf_rf <- confusionMatrix(rf_pred, test_data$goal)
cf_log <- confusionMatrix(log_pred, test_data$goal)
cf_blog <- confusionMatrix(log_pred, test_data$goal)
cf_knn <- confusionMatrix(knn_pred, test_data$goal)

Accuracy <- 100*rbind(cf_log[["overall"]][["Accuracy"]],
  cf_rf[["overall"]][["Accuracy"]],
  cf_blog[["overall"]][["Accuracy"]],
  cf_knn[["overall"]][["Accuracy"]])

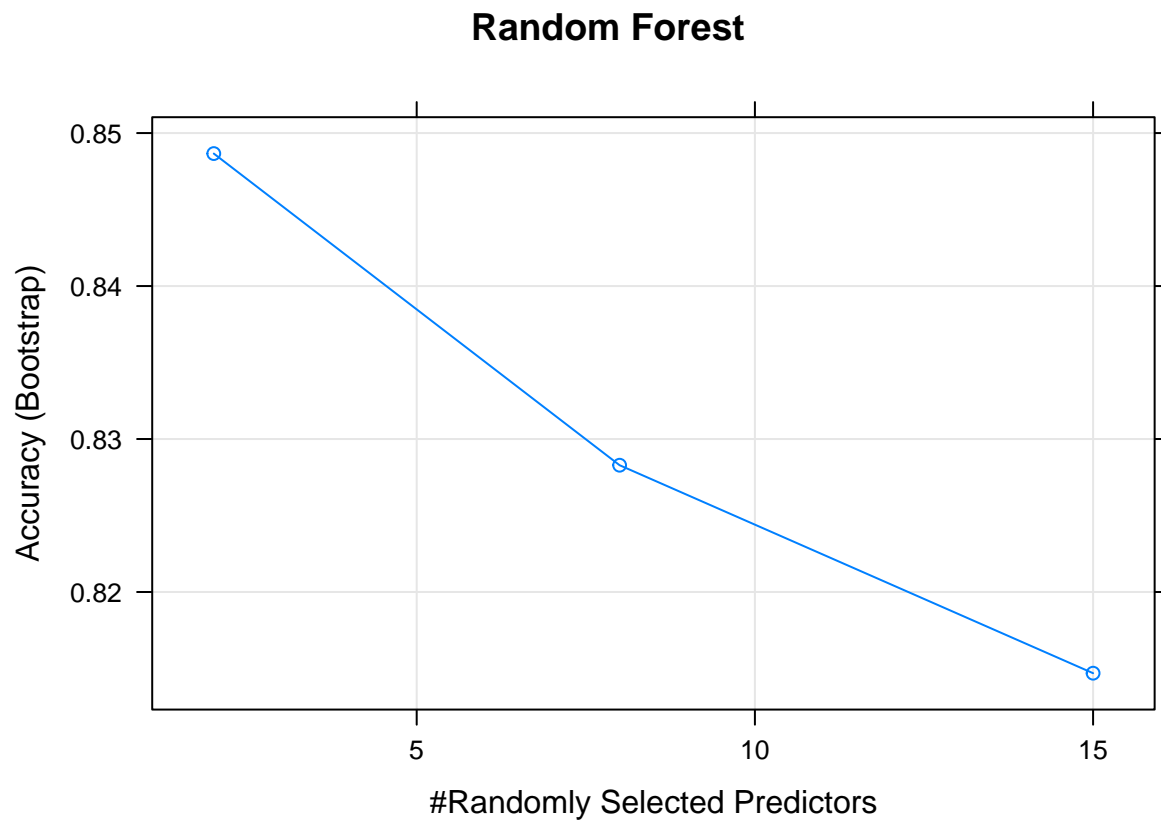
Specificity <- 100*rbind(cf_log[["byClass"]][["Specificity"]],
  cf_rf[["byClass"]][["Specificity"]],
  cf_blog[["byClass"]][["Specificity"]],
  cf_knn[["byClass"]][["Specificity"]])

Sensitivity <- 100*rbind(cf_log[["byClass"]][["Sensitivity"]],
  cf_rf[["byClass"]][["Sensitivity"]],
  cf_blog[["byClass"]][["Sensitivity"]],
  cf_knn[["byClass"]][["Sensitivity"]])

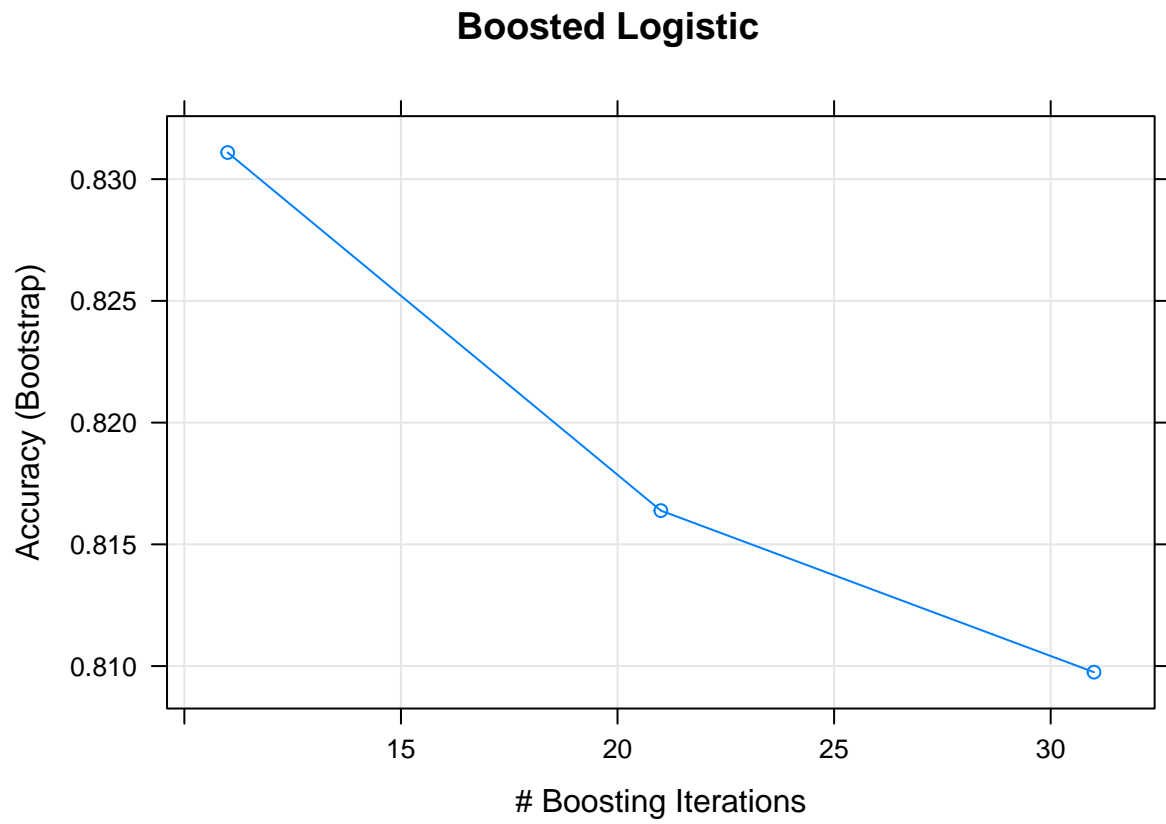
pf_result <- t(data.frame(Accuracy, Specificity, Sensitivity))
colnames(pf_result) <- c("Log", "RF", "LogitB", "KNN")
pf_result <- as.matrix(pf_result)
pf_result
```

```
##           Log      RF      LogitB      KNN
## Accuracy  75.00000 75.00000 75.00000 75.00000
## Specificity 64.28571 67.85714 64.28571 60.71429
## Sensitivity 84.37500 81.25000 84.37500 87.50000
```

```
plot(rf_fit, main = "Random Forest")
```

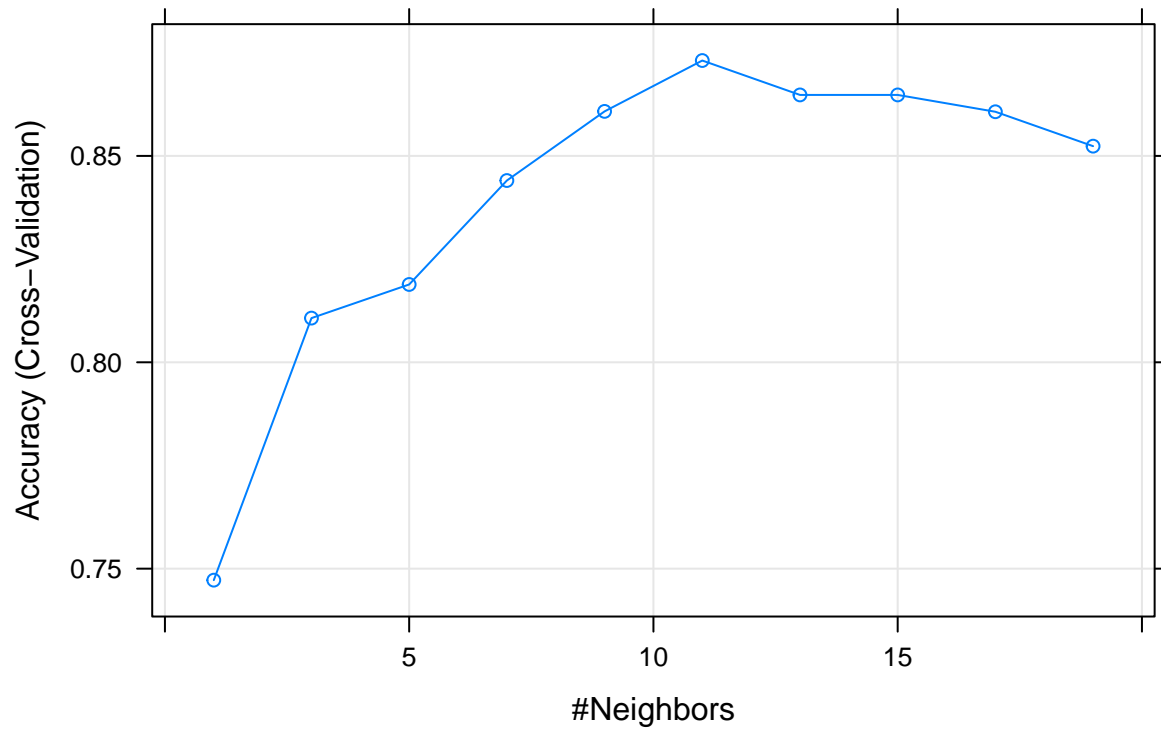


```
plot(blog_fit, main = "Boosted Logistic")
```



```
plot(knn_fit, main = "K-nearest neighbour")
```

## K-nearest neighbour

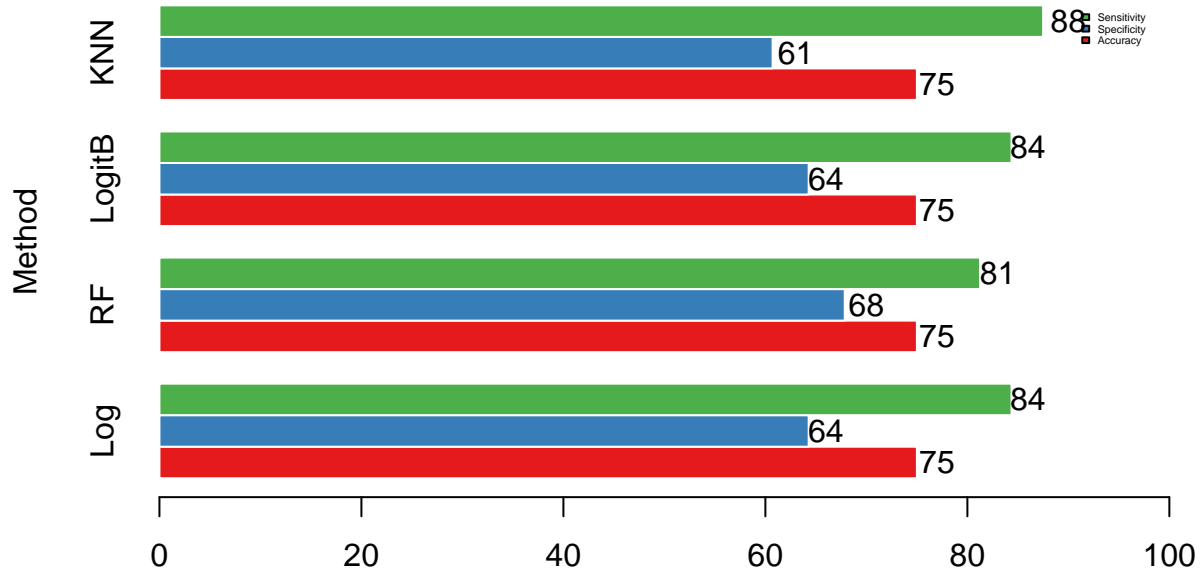


```
y <- barplot(pf_result, beside = TRUE, horiz = TRUE,
             col=brewer.pal(3,"Set1"),border="white",
             legend.text = c("Accuracy", "Specificity", "Sensitivity"),
             args.legend = list(bty = "n", cex = 0.3), xlim=c(0,100),
             main = "Performance Chart", ylab = "Method")

x <- round(pf_result)

text(x+2,y,labels=as.character(x))
```

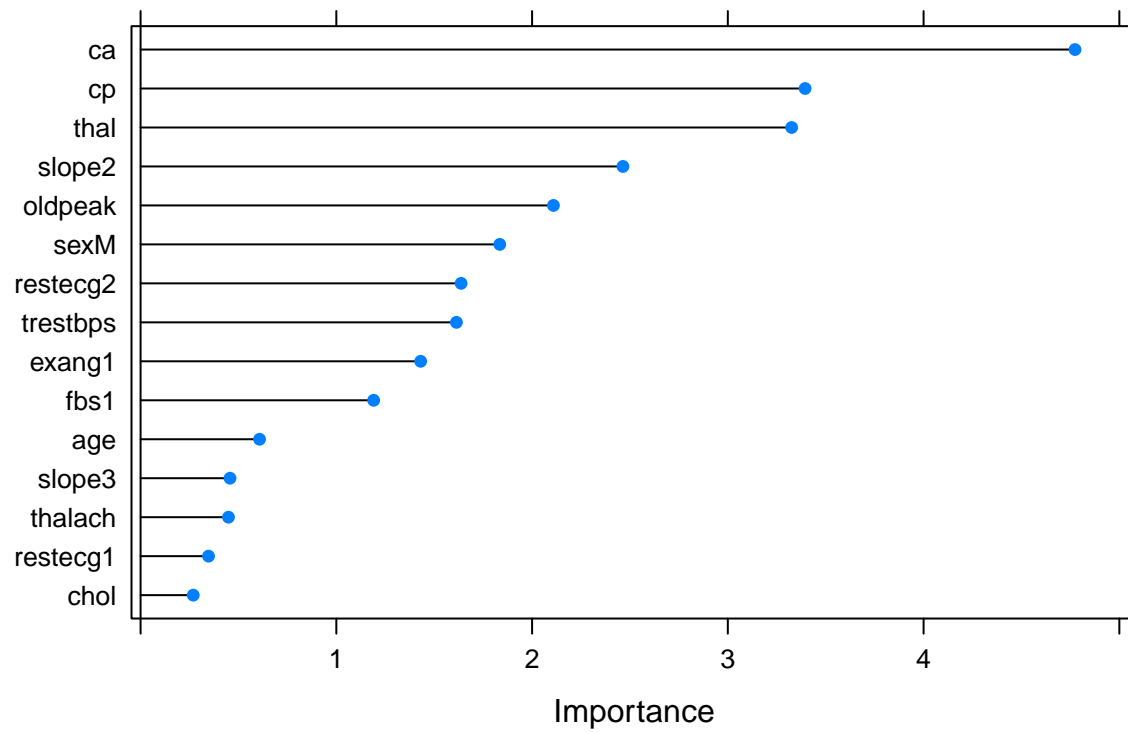
## Performance Chart



```
## Feature extraction -
feat_log <- varImp(log_fit, scale = FALSE)
feat_rf <- varImp(rf_fit, scale = FALSE)
feat_blog <- varImp(blog_fit, scale = FALSE)
feat_knn <- varImp(knn_fit, scale = FALSE)

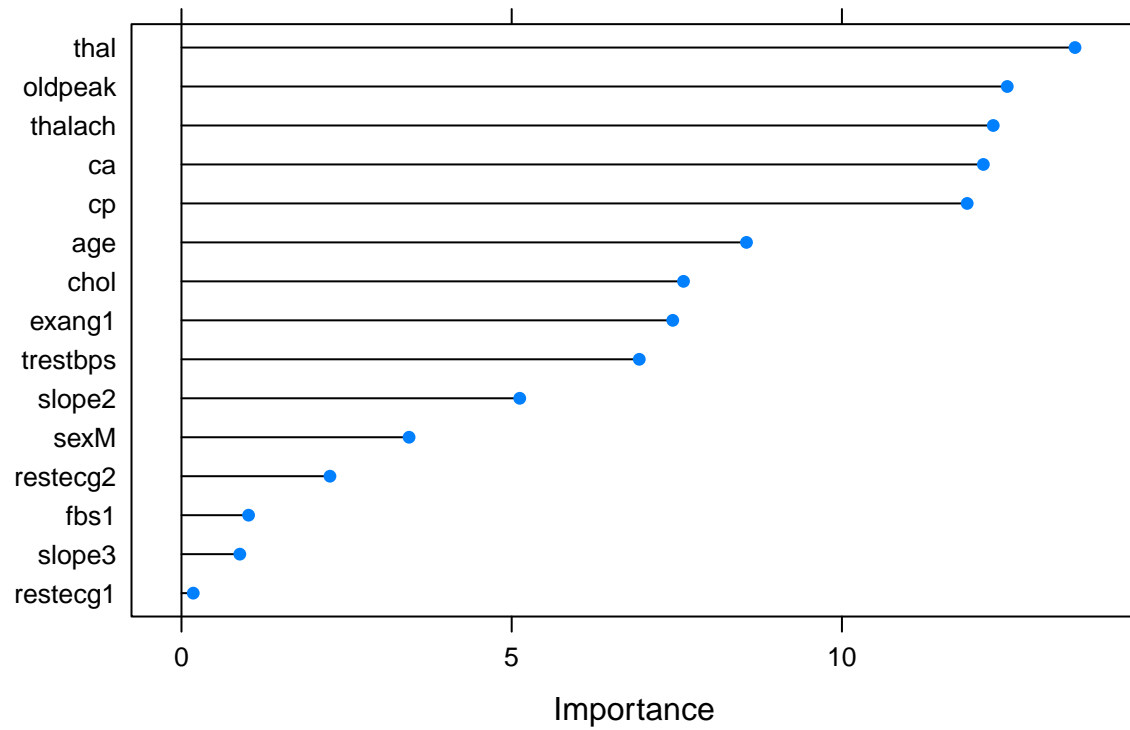
plot(feat_log, main = "Logistic regression: features")
```

## Logistic regression: features



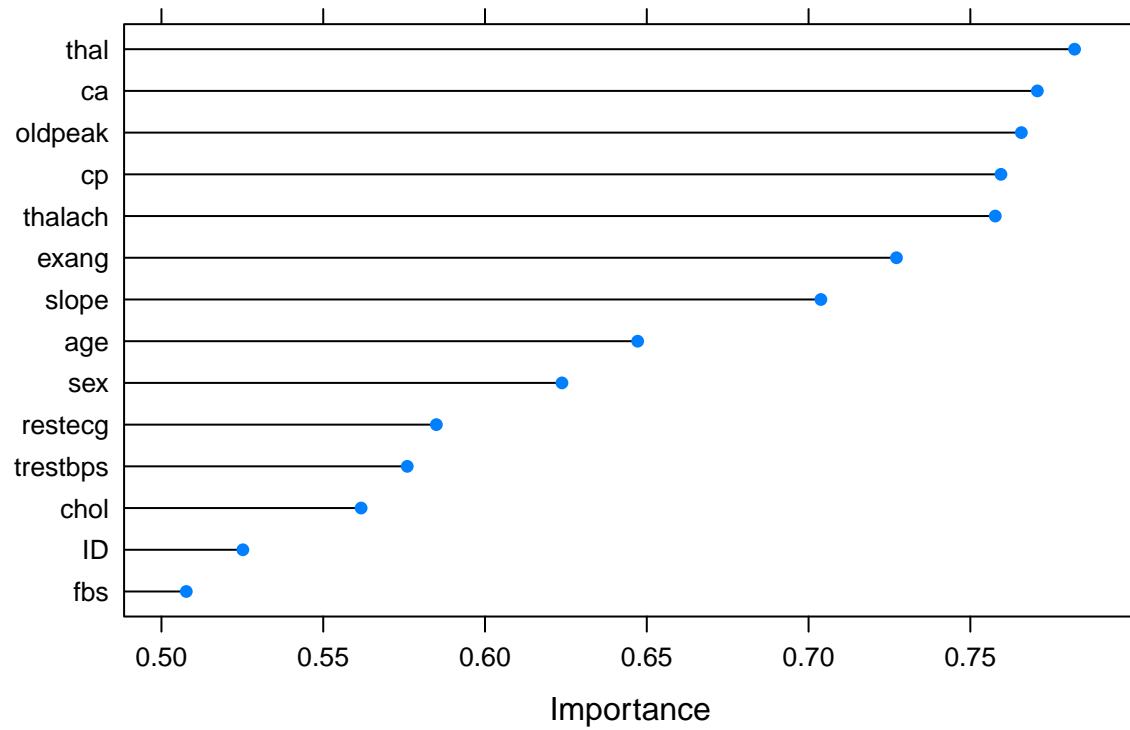
```
plot(feet_rf, main = "Random forest: features")
```

## Random forest: features



```
plot(feet_blog, main = "Boosted Logistic regression: features")
```

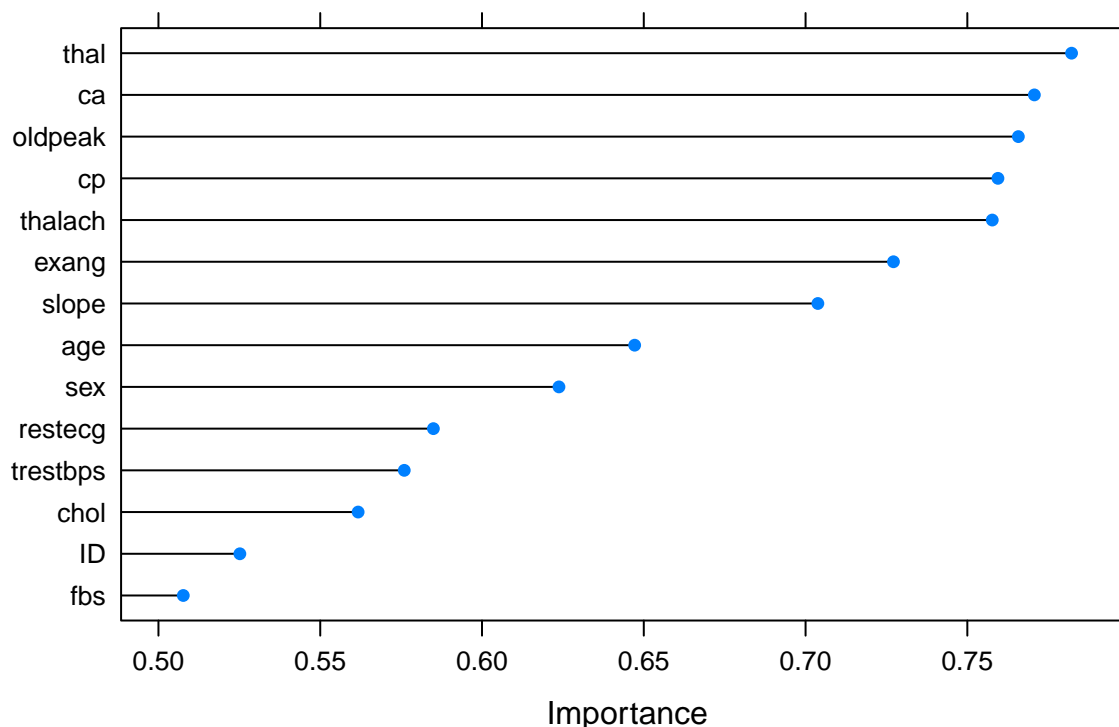
## Boosted Logistic regression: features



```
plot(feet_knn, main = "KNN: features")
```



## KNN: features



```
##      Model Evaluation with ROC, calibration, precision      -
##      recall gain, and Obs vs. Pred probabilities curve     -
```

```
cont <- trainControl(method="cv",
                      summaryFunction=twoClassSummary,
                      classProbs=T,
                      savePredictions = T)
```

```
log <- train(goal ~.-ID ,
             data = train_data,
             method = "glm", preProc=c("center", "scale"),
             family = "binomial", trControl=cont)
```

```
## Warning in train.default(x, y, weights = w, ...): The metric "Accuracy" was not
## in the result set. ROC will be used instead.
```

```
rf <- train(goal ~.-ID ,
            data = train_data, preProc=c("center", "scale"),
            method = "rf", trControl=cont)
```

```
## Warning in train.default(x, y, weights = w, ...): The metric "Accuracy" was not
## in the result set. ROC will be used instead.
```

```
blog <- train(goal ~.-ID,
              data = train_data, preProc=c("center", "scale"),
              method = "LogitBoost", trControl=cont)
```

```
## Warning in train.default(x, y, weights = w, ...): The metric "Accuracy" was not
```

```

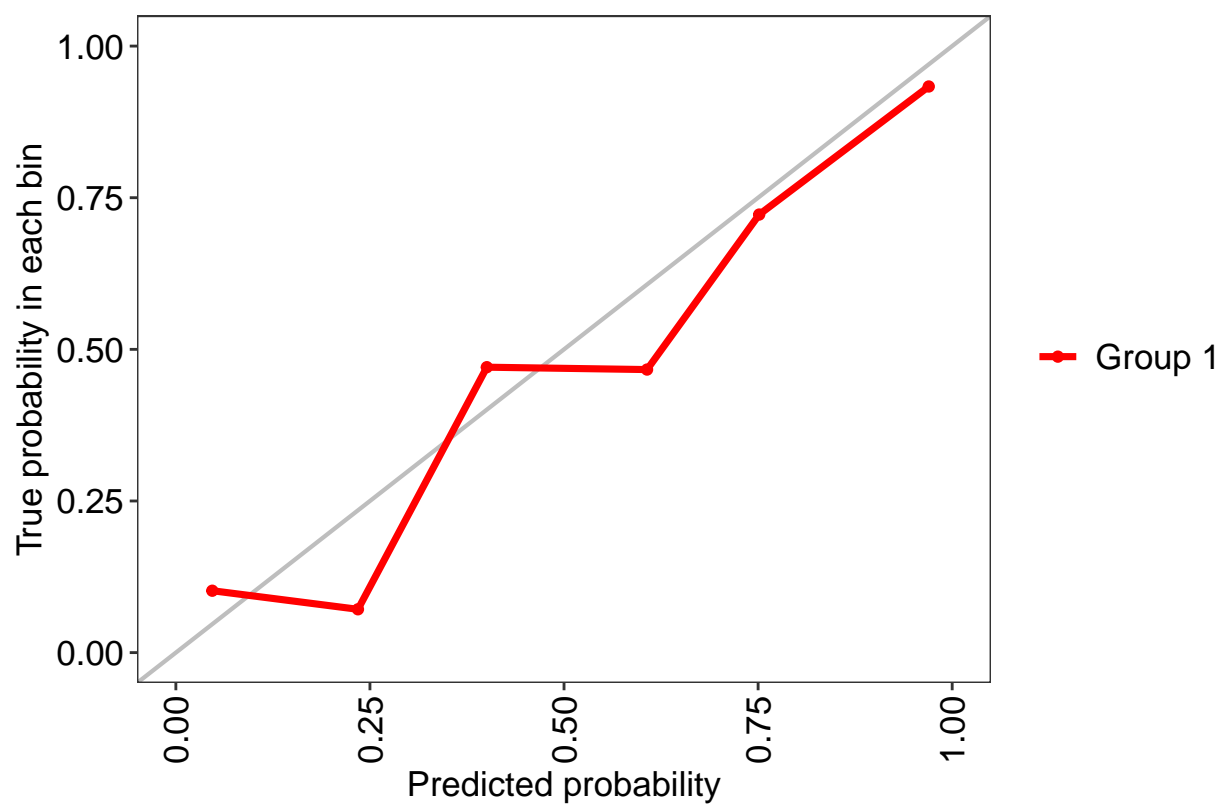
## in the result set. ROC will be used instead.
knn <- train(goal ~. -ID , data = train_data,
             method = "knn", preProcess = c("center","scale"),
             trControl = cont , tuneGrid = expand.grid(k = seq(1, 20, 2)))

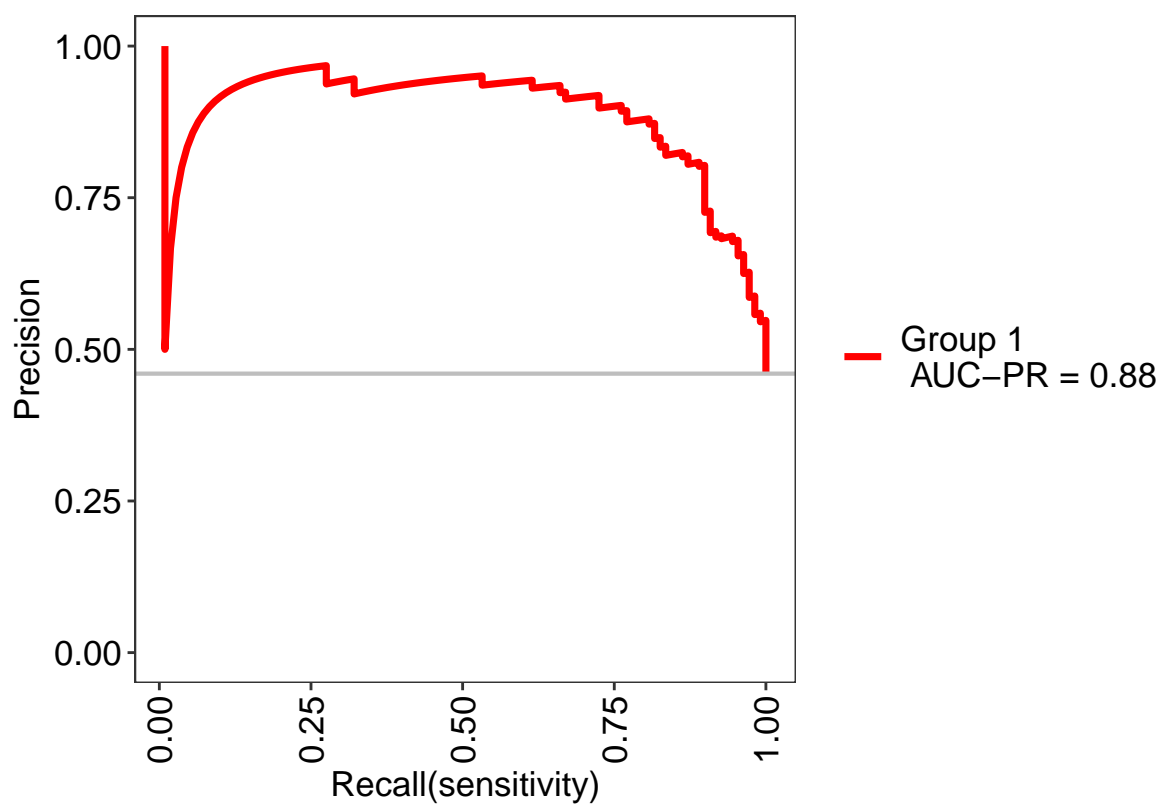
## Warning in train.default(x, y, weights = w, ...): The metric "Accuracy" was not
## in the result set. ROC will be used instead.

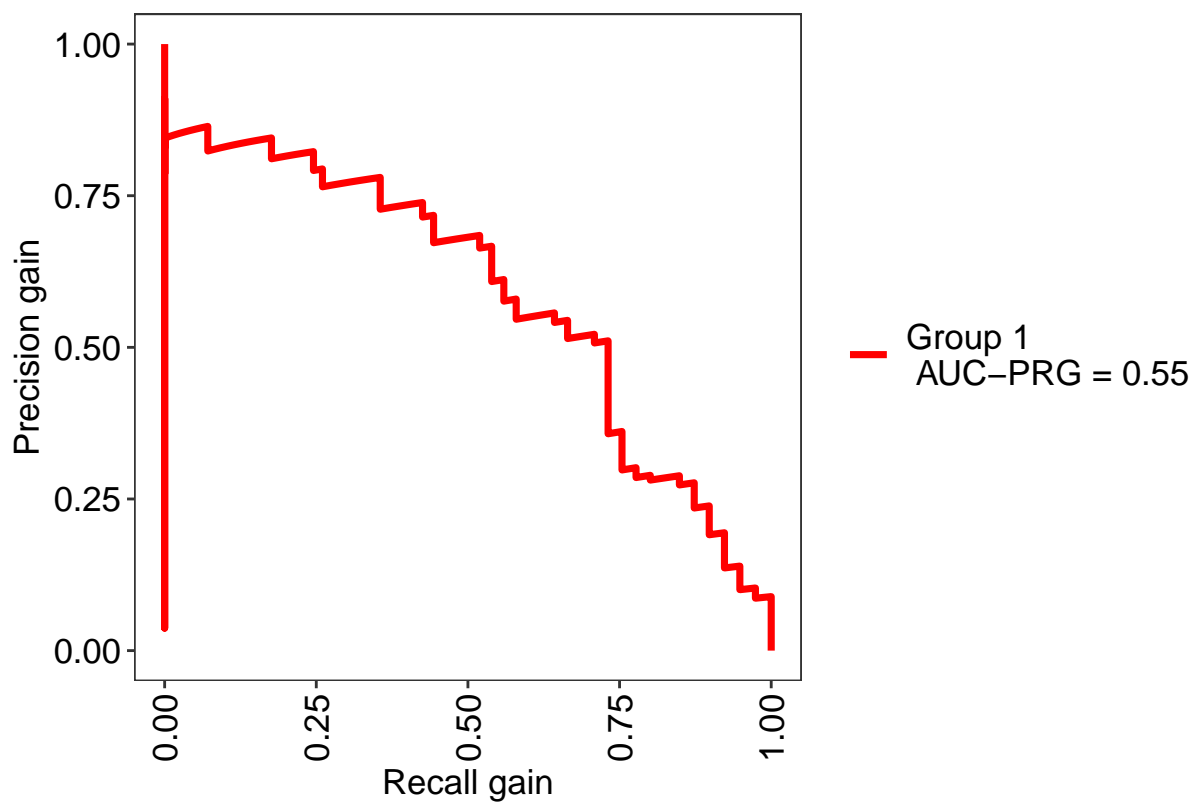
#####
##          Logistic Classifier: Performance Plots          ##
#####
res.log <- evalm(log)

## ***MLeval: Machine Learning Model Evaluation***
## Input: caret train function object
## Not averaging probs.
## Group 1 type: cv
## Observations: 237
## Number of groups: 1
## Observations per group: 237
## Positive: unhealthy
## Negative: healthy
## Group: Group 1
## Positive: 109
## Negative: 128
## ***Performance Metrics***

```

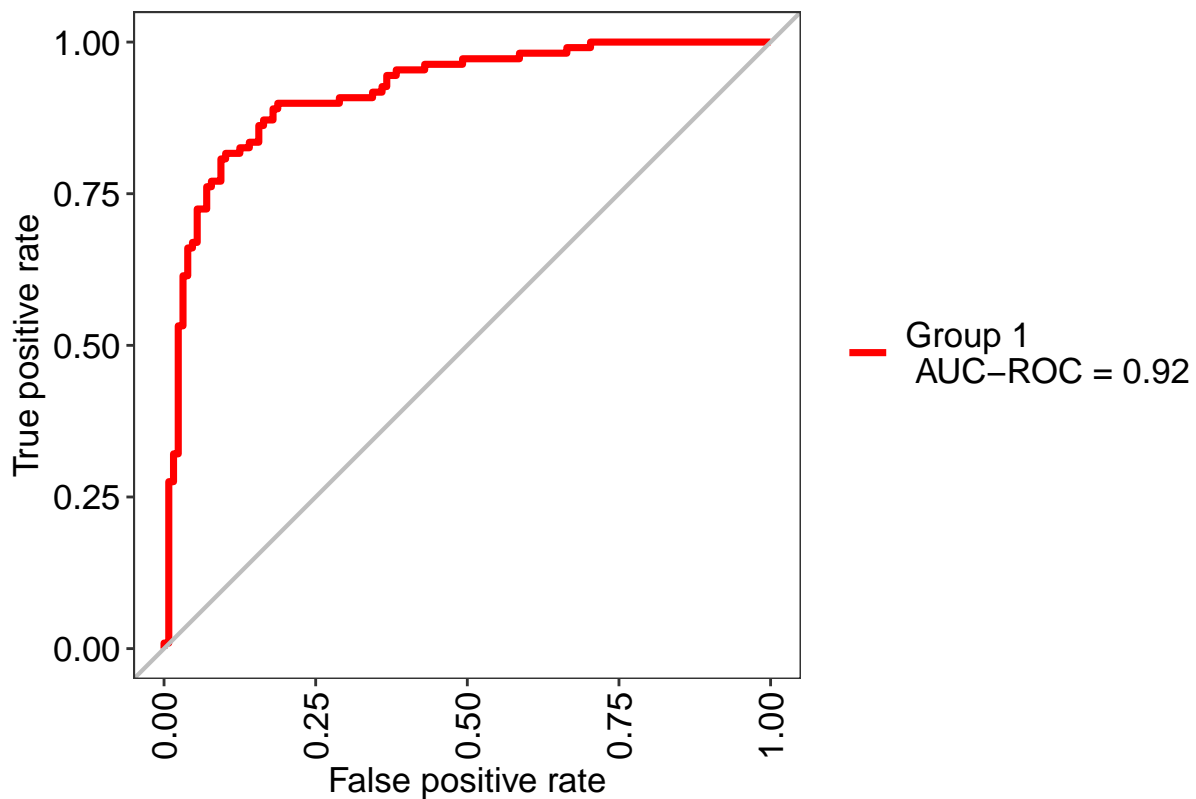






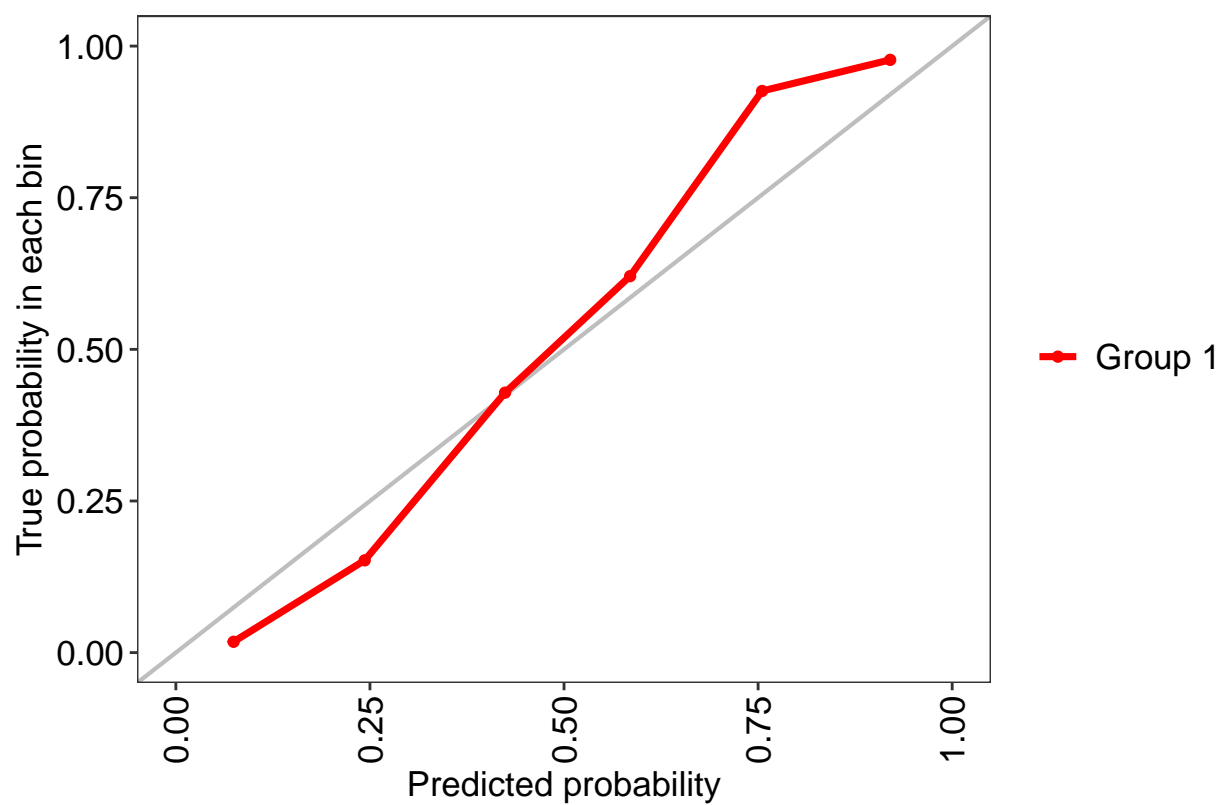
## Group 1 Optimal Informedness = 0.71495126146789

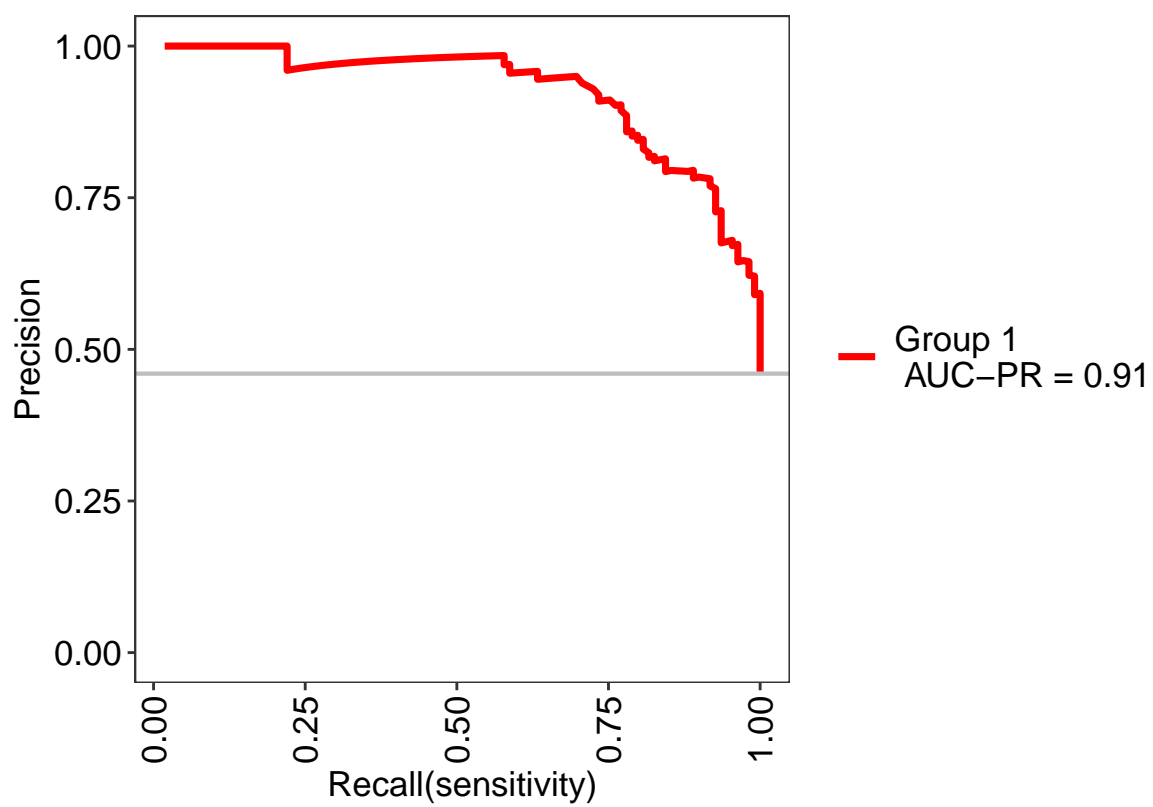
## Group 1 AUC-ROC = 0.92



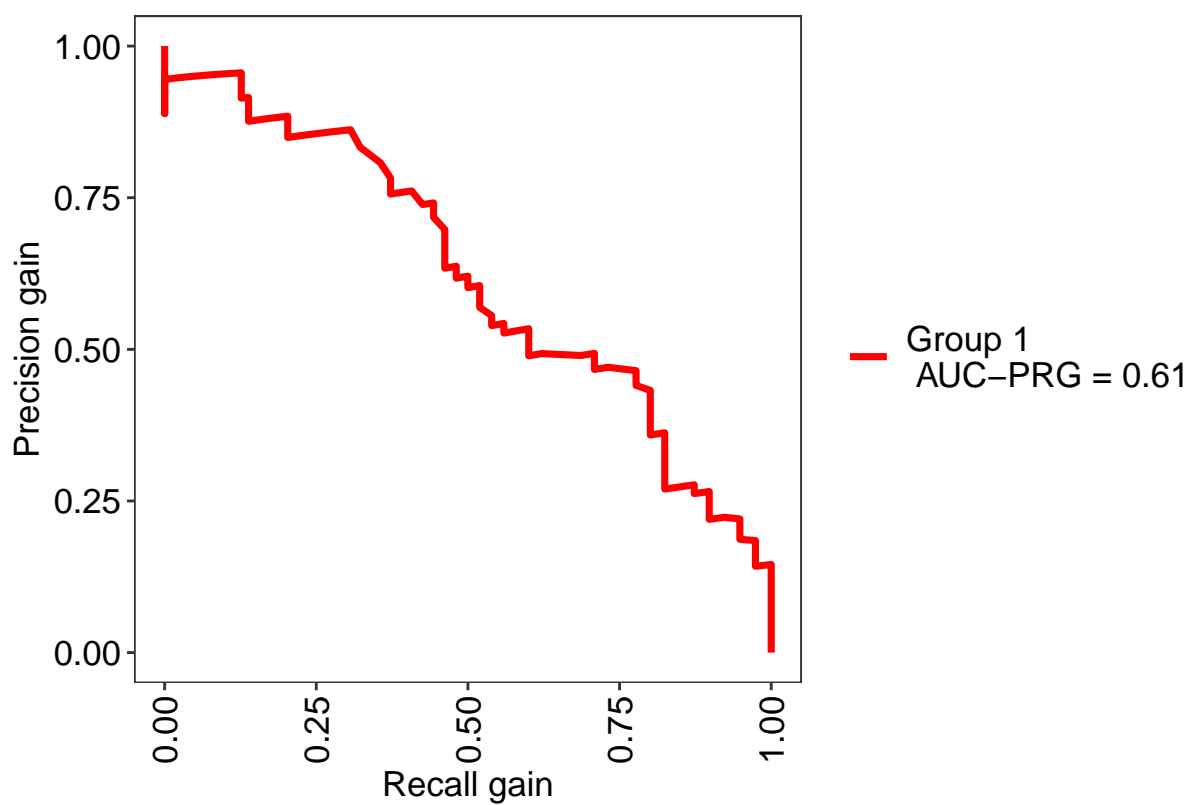
```
#####
##      Random Forest Classifier: Performance Plots      ##
#####
res.rf <- evalm(rf)

## ***MLevel: Machine Learning Model Evaluation***
## Input: caret train function object
## Not averaging probs.
## Group 1 type: cv
## Observations: 237
## Number of groups: 1
## Observations per group: 237
## Positive: unhealthy
## Negative: healthy
## Group: Group 1
## Positive: 109
## Negative: 128
## ***Performance Metrics***
```



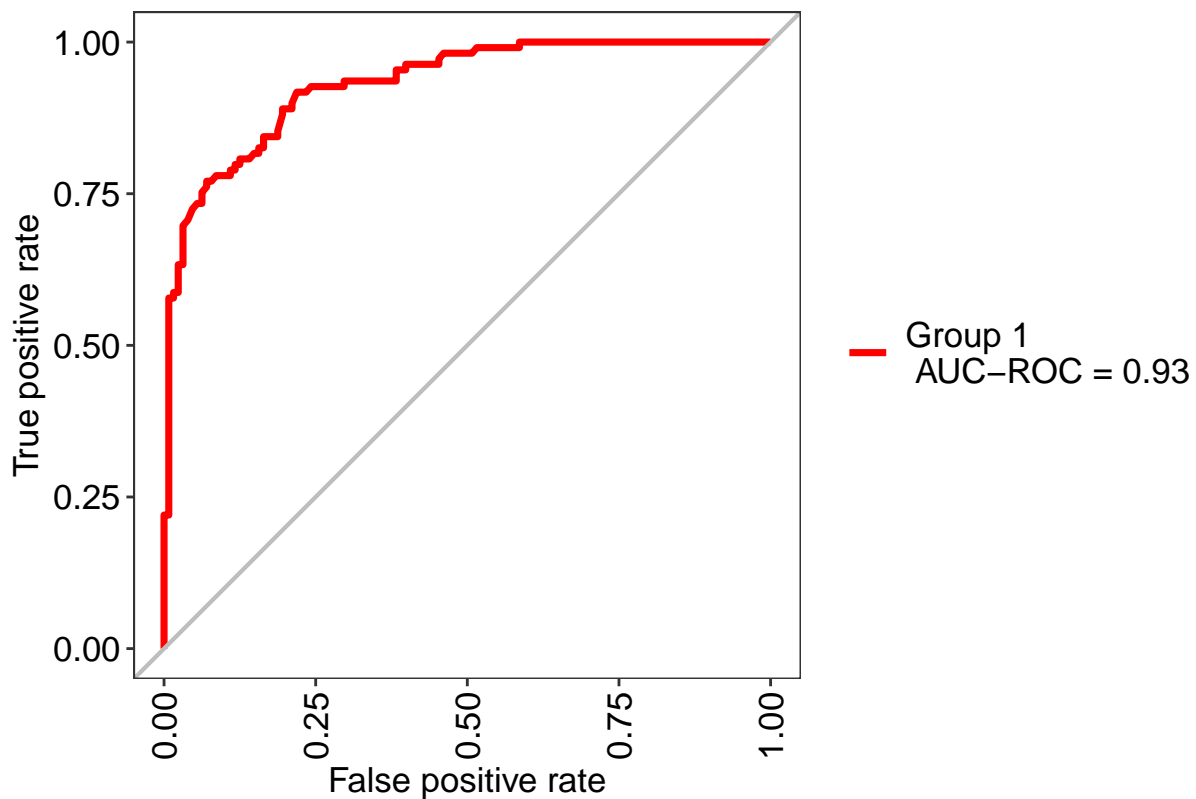






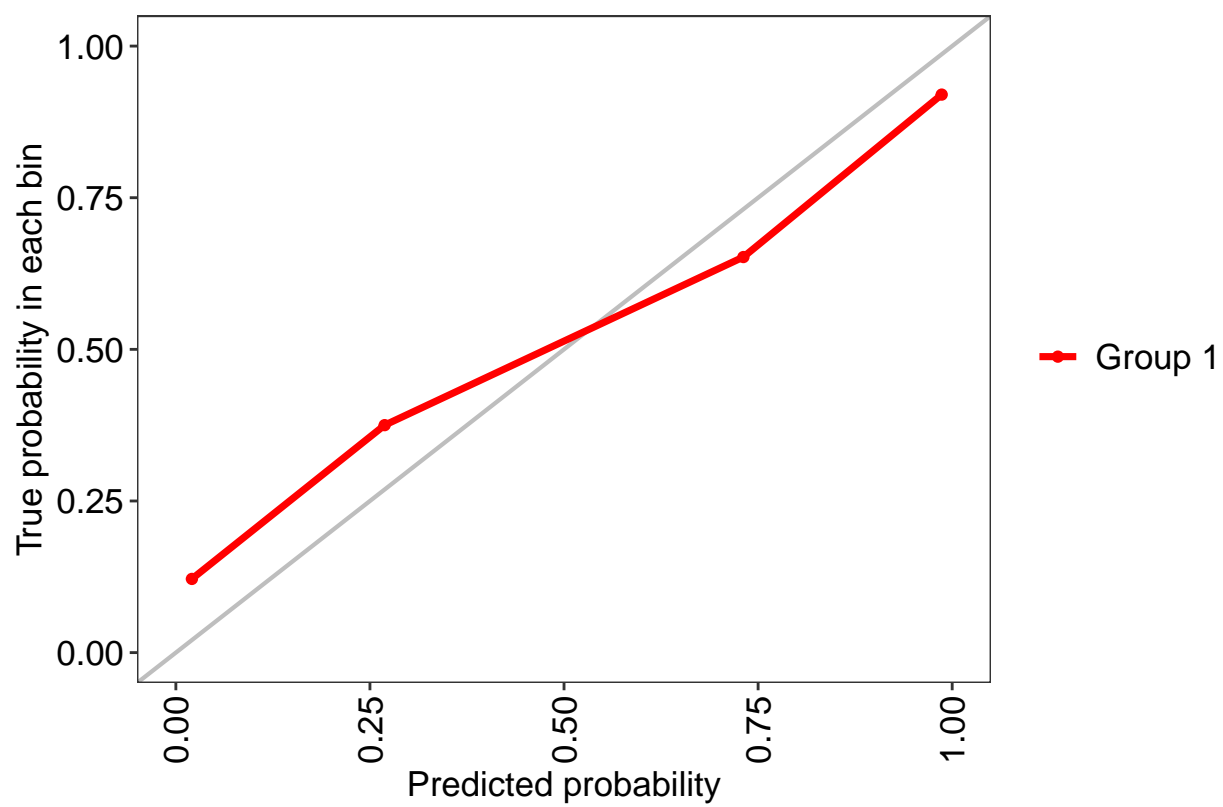
## Group 1 Optimal Informedness = 0.700329701834862

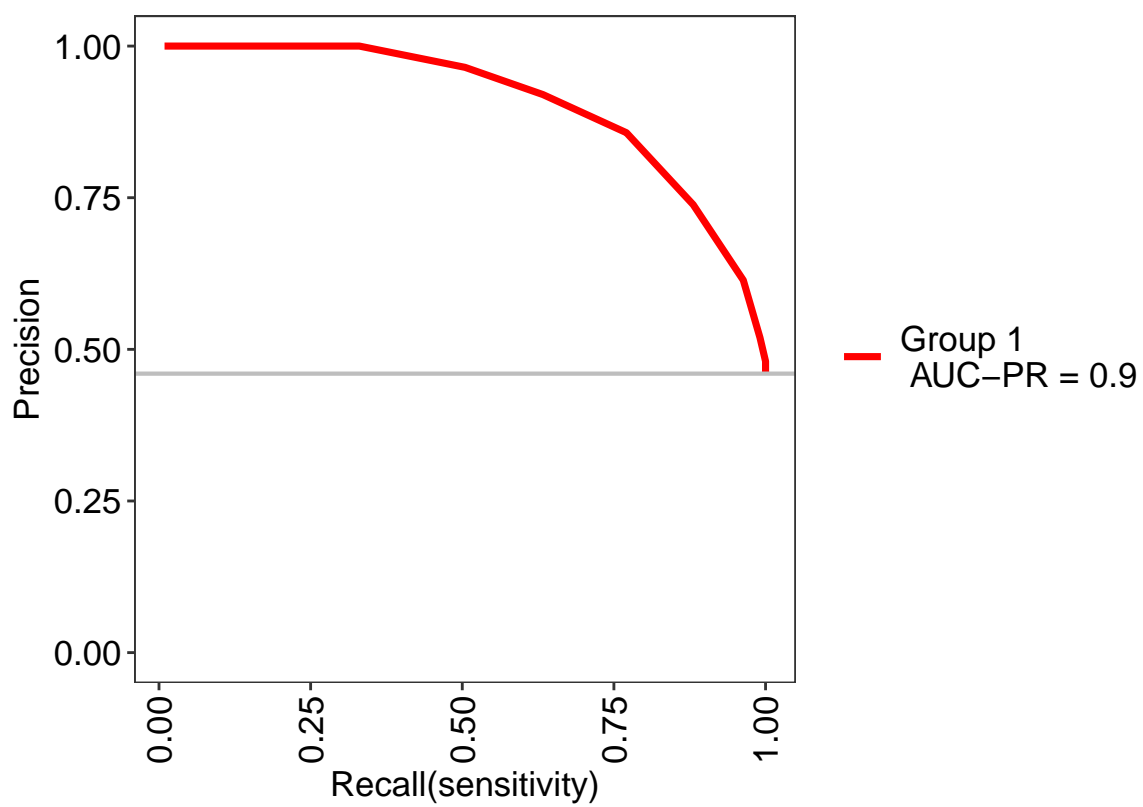
## Group 1 AUC-ROC = 0.93

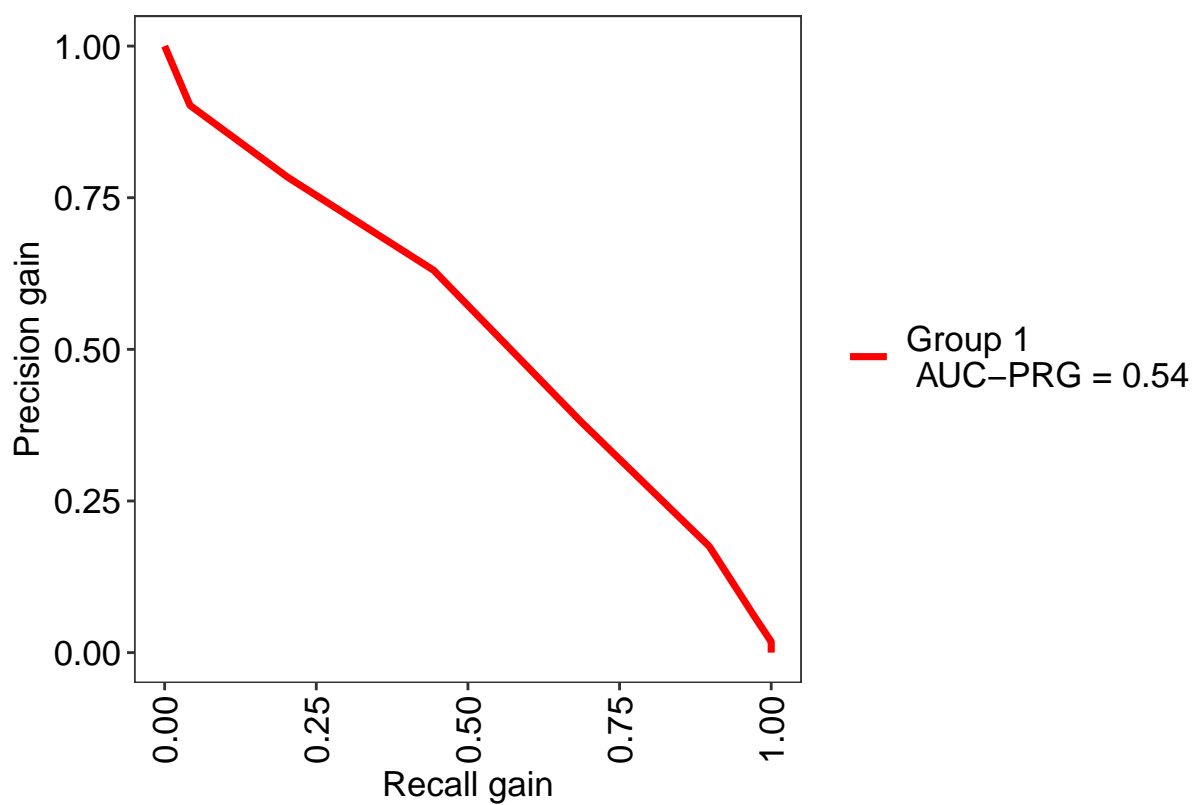


```
#####
##           Boosted logit Classifier: Performance Plots           ##
#####
res.blog <- evalm(blog)

## ***MLevel: Machine Learning Model Evaluation***
## Input: caret train function object
## Not averaging probs.
## Group 1 type: cv
## Observations: 237
## Number of groups: 1
## Observations per group: 237
## Positive: unhealthy
## Negative: healthy
## Group: Group 1
## Positive: 109
## Negative: 128
## ***Performance Metrics***
```

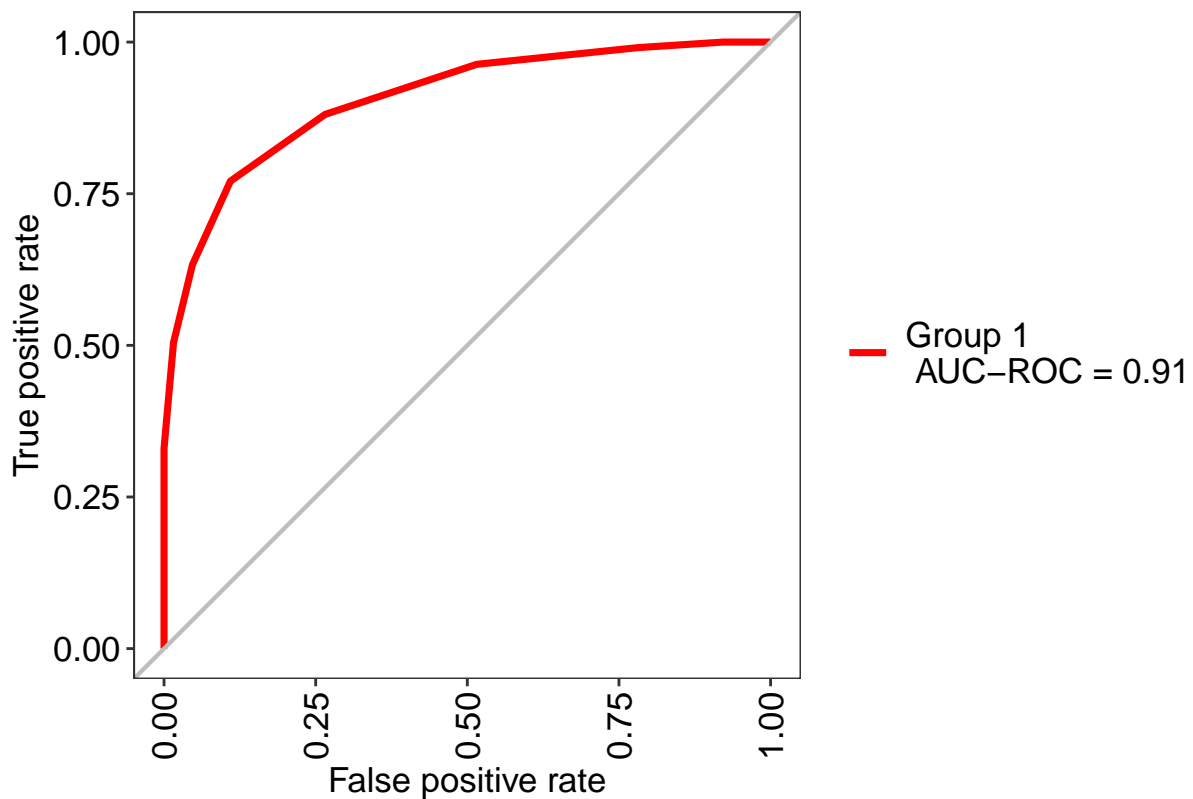






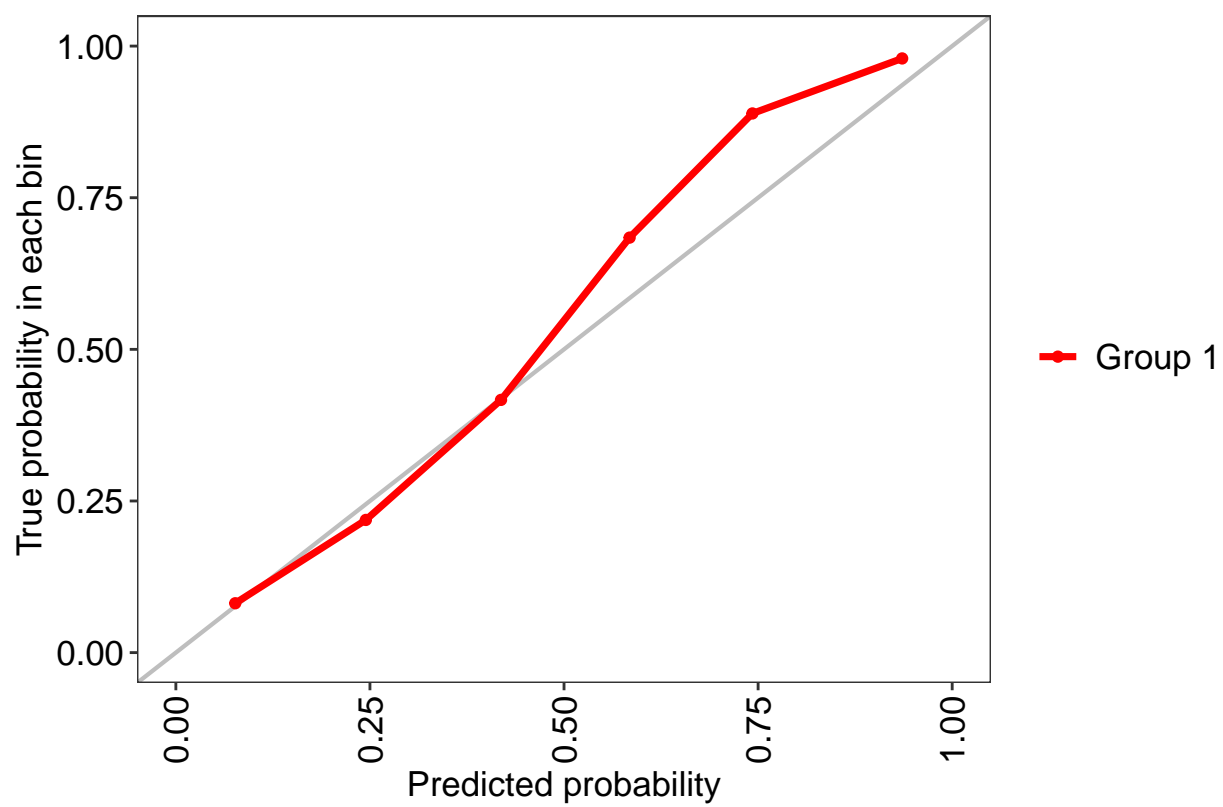
## Group 1 Optimal Informedness = 0.661267201834862

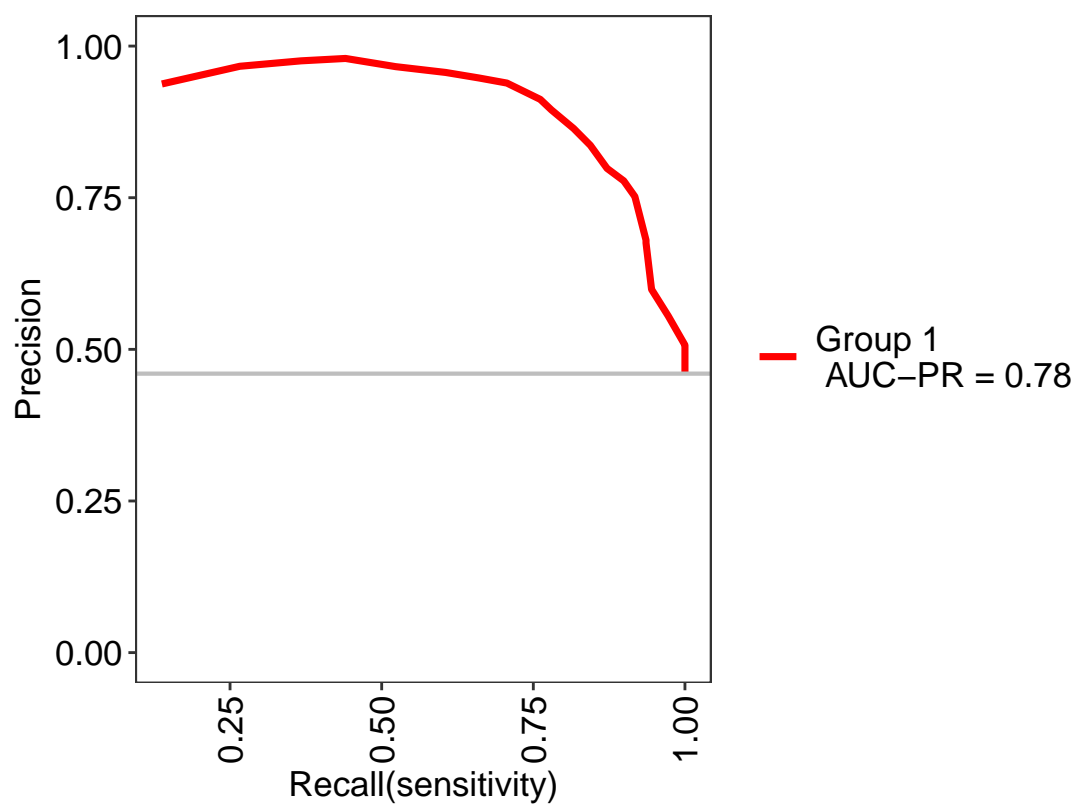
## Group 1 AUC-ROC = 0.91



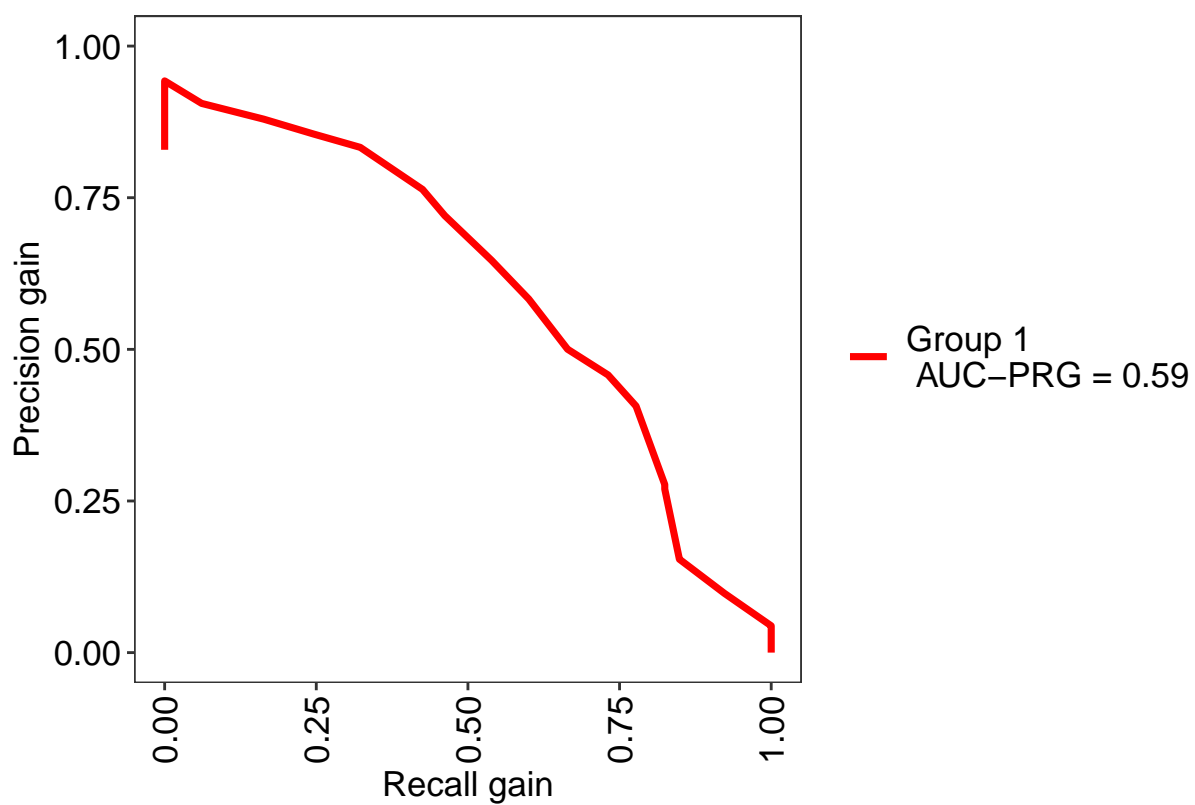
```
#####
##          k-NN Classifier: Performance Plots          ##
#####
res.knn <- evalm(knn)

## ***MLevel: Machine Learning Model Evaluation***
## Input: caret train function object
## Not averaging probs.
## Group 1 type: cv
## Observations: 237
## Number of groups: 1
## Observations per group: 237
## Positive: unhealthy
## Negative: healthy
## Group: Group 1
## Positive: 109
## Negative: 128
## ***Performance Metrics***
```









## Group 1 Optimal Informedness = 0.70713876146789

## Group 1 AUC-ROC = 0.92

