

Data Science Final Project

MODEL 1C: Gradient Boosting

Kjay O. Coca

2022-12-16

LOAD PACKAGES

```
# Load Packages
library(dplyr)    # for general data wrangling needs

##
## Attaching package: 'dplyr'

## The following objects are masked from 'package:stats':
##
##   filter, lag

## The following objects are masked from 'package:base':
##
##   intersect, setdiff, setequal, union

library(tidyverse)# for filtering

## -- Attaching packages ----- tidyverse 1.3.2 --

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

library(rsample)    # for creating validation splits
library(h2o)        # for a java-based implementation of GBM variants

##
## -----
##
## Your next step is to start H2O:
##   > h2o.init()
```

```
##
## For H2O package documentation, ask for help:
##   > ??h2o
##
## After starting H2O, you can use the Web UI at http://localhost:54321
## For more information visit https://docs.h2o.ai
##
## -----
##
##
## Attaching package: 'h2o'
##
## The following objects are masked from 'package:stats':
##
##   cor, sd, var
##
## The following objects are masked from 'package:base':
##
##   %*%, %in%, &&, ||, apply, as.factor, as.numeric, colnames,
##   colnames<-, ifelse, is.character, is.factor, is.numeric, log,
##   log10, log1p, log2, round, signif, trunc
```

```
library(xgboost) # for fitting extreme gradient boosting
```

```
##
## Attaching package: 'xgboost'
##
## The following object is masked from 'package:dplyr':
##
##   slice
```

```
library(gbm) # for original implementation of regular and stochastic GBMs
```

```
## Loaded gbm 2.1.8.1
```

```
library(ROCR)
library(pROC)
```

```
## Type 'citation("pROC")' for a citation.
##
## Attaching package: 'pROC'
##
## The following object is masked from 'package:h2o':
##
##   var
##
## The following objects are masked from 'package:stats':
##
##   cov, smooth, var
```

```
library(recipes)
```

```
##
## Attaching package: 'recipes'
##
## The following object is masked from 'package:stringr':
##
##     fixed
##
## The following object is masked from 'package:stats':
##
##     step
```

IMPORTING THE DATA

```
set.seed(123)
radiomics_data <- read_csv("D:/1 MASTERS/STAT225/FINAL PROJECT/STAT 325 _FINAL PROJECT_/Normalize Radiomics Data.csv")

## Rows: 197 Columns: 431
## -- Column specification -----
## Delimiter: ","
## chr  (1): Institution
## dbl (430): Failure.binary, Failure, Entropy_cooc.W.ADC, GLNU_align.H.PET, Mi...
##
## i Use 'spec()' to retrieve the full column specification for this data.
## i Specify the column types or set 'show_col_types = FALSE' to quiet this message.
```

SPLITTING FOR TRAINING AND TESTING

```
radiomics_data$Institution=as.factor(radiomics_data$Institution)
to_split <- initial_split(radiomics_data, strata = "Failure.binary")
radiomicsdata_train <- training(to_split)
radiomicsdata_test <- testing(to_split)
```

In this case, I set 80 percent for training data and 20 percent for testing data. There are 39 observation for testing and 158 observation for training and both have 413 variables.

GB Model 1

```
GB_model1 <- gbm(
  formula = Failure.binary ~ .,
  data = radiomicsdata_train,
  distribution = "bernoulli", # SSE loss function
  n.trees = 500,
  shrinkage = 0.1,
  n.minobsinnode = 10,
```

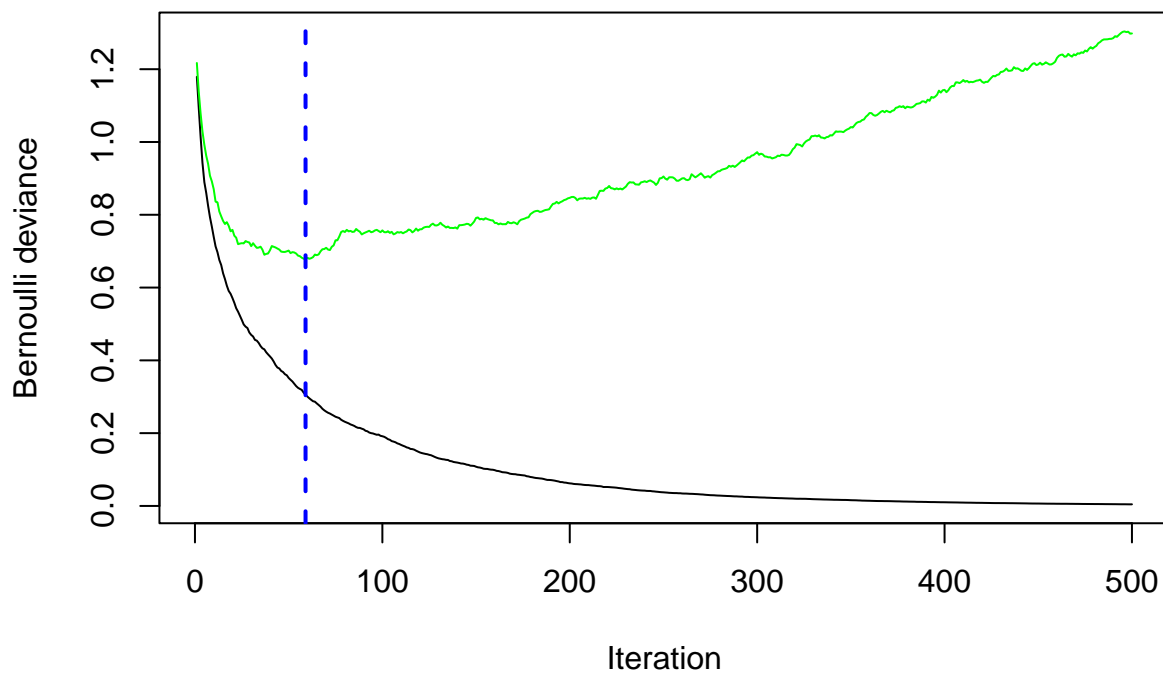
```
cv.folds = 10
)
```

FIND INDEX FOR NUMBER TREES WITH MINIMUM CV ERROR

```
best_gbm <- which.min(GB_model1$cv.error)
```

PLOTTING THE ERROR CURVE

```
gbm.perf(GB_model1, method = "cv")
```



```
## [1] 59
```

CREATE GRID SEARCH

```
hyper_grid <- expand.grid(
  learning_rate = c(0.3, 0.1, 0.05, 0.01, 0.005),
```

```

logloss = NA,
trees = NA,
time = NA
)

```

EXECUTE GRID SEARCH

```

for(i in seq_len(nrow(hyper_grid))) {
  # fit gbm
  set.seed(123) # for reproducibility
  train_time <- system.time({
    m <- gbm(
      formula = Failure.binary ~ .,
      data = radiomicsdata_train,
      distribution = "bernoulli",
      n.trees = 500,
      shrinkage = hyper_grid$learning_rate[i],
      interaction.depth = 3,
      n.minobsinnode = 10,
      cv.folds = 10
    )
  })

  # adding SSE, trees, and training time to results
  hyper_grid$logloss[i] <- sqrt(min(m$cv.error))
  hyper_grid$trees[i] <- which.min(m$cv.error)
  hyper_grid$Time[i] <- train_time[["elapsed"]]
}

```

RESULTS

```

arrange(hyper_grid, logloss)

```

```

##   learning_rate  logloss trees time  Time
## 1          0.100 0.7656695   32   NA 10.86
## 2          0.050 0.7850738   82   NA 10.48
## 3          0.010 0.7916756  400   NA 10.51
## 4          0.005 0.8012999  500   NA 10.65
## 5          0.300 0.8043208   17   NA 10.52

```

SEARCH GRID

```

hyper_grid <- expand.grid(
  n.trees = 600,
  shrinkage = 0.01,
  interaction.depth = c(3, 5, 7),

```

```
n.minobsinnode = c(5, 10, 15)

)
```

CREATING THE MODEL FIT FUNCTION

```
model_fit <- function(n.trees, shrinkage, interaction.depth, n.minobsinnode) {
  set.seed(123)
  m <- gbm(
    formula = Failure.binary ~ .,
    data = radiomicsdata_train,
    distribution = "bernoulli",
    n.trees = n.trees,
    shrinkage = shrinkage,
    interaction.depth = interaction.depth,
    n.minobsinnode = n.minobsinnode,
    cv.folds = 10
  )
  # compute RMSE
  sqrt(min(m$cv.error))
}
```

PERFORMING SEARCH GRID WITH FUNCTIONAL PROGRAMMING

```
hyper_grid$logloss <- purrr::pmap_dbl(
  hyper_grid,
  ~ model_fit(
    n.trees = ..1,
    shrinkage = ..2,
    interaction.depth = ..3,
    n.minobsinnode = ..4
  )
)

# RESULTS
arrange(hyper_grid, logloss)
```

##	n.trees	shrinkage	interaction.depth	n.minobsinnode	logloss
## 1	600	0.01	3	15	0.7652952
## 2	600	0.01	5	15	0.7652952
## 3	600	0.01	7	15	0.7652952
## 4	600	0.01	3	10	0.7916756
## 5	600	0.01	5	10	0.7917035
## 6	600	0.01	7	10	0.7917035
## 7	600	0.01	3	5	0.7958123
## 8	600	0.01	5	5	0.7971165
## 9	600	0.01	7	5	0.7971459

REFINED HYPERPARAMETER GRID

```
hyper_grid <- list(  
  sample_rate = c(0.5, 0.75, 1),           # row subsampling  
  col_sample_rate = c(0.5, 0.75, 1),       # col subsampling for each split  
  col_sample_rate_per_tree = c(0.5, 0.75, 1) # col subsampling for each tree  
)
```

random grid search strategy

```
# random grid search strategy  
search_criteria <- list(  
  strategy = "RandomDiscrete",  
  stopping_metric = "logloss",  
  stopping_tolerance = 0.001,  
  stopping_rounds = 10,  
  max_runtime_secs = 60*60  
)
```

PERFORMING GRID SEARCH

```
radiomicsdata_train$Failure.binary=as.factor(radiomicsdata_train$Failure.binary)  
h2o.init()
```

```
## Connection successful!  
##  
## R is connected to the H2O cluster:  
##   H2O cluster uptime:      8 hours 12 minutes  
##   H2O cluster timezone:    Asia/Singapore  
##   H2O data parsing timezone: UTC  
##   H2O cluster version:     3.38.0.1  
##   H2O cluster version age:  2 months and 27 days  
##   H2O cluster name:        H2O_started_from_R_Kjay_Coca_lia053  
##   H2O cluster total nodes: 1  
##   H2O cluster total memory: 3.19 GB  
##   H2O cluster total cores: 8  
##   H2O cluster allowed cores: 8  
##   H2O cluster healthy:     TRUE  
##   H2O Connection ip:        localhost  
##   H2O Connection port:     54321  
##   H2O Connection proxy:     NA  
##   H2O Internal Security:    FALSE  
##   R Version:                R version 4.2.2 (2022-10-31 ucrt)
```

```
grid <- h2o.grid(  
  algorithm = "gbm",  
  grid_id = "gbm_grid",  
  y = "Failure.binary",
```

```

training_frame = as.h2o(radiomicsdata_train),
hyper_params = hyper_grid,
ntrees = 10, #supposedly 6000
learn_rate = 0.01,
max_depth = 7,
min_rows = 5,
nfolds = 10,
stopping_rounds = 10,
stopping_tolerance = 0,
stopping_metric="logloss",
search_criteria = search_criteria,
seed = 123
)

```

```

## |
## |

```

COLLECT THE RESULTS AND SORT BY OUR MODEL PERFORMANCE METRIC OF CHOICE

```

grid_perf <- h2o.getGrid(
  grid_id = "gbm_grid",
  sort_by = "logloss",
  decreasing = FALSE
)
grid_perf

```

```

## H2O Grid Details
## =====
##
## Grid ID: gbm_grid
## Used hyper parameters:
##   - col_sample_rate
##   - col_sample_rate_per_tree
##   - sample_rate
## Number of models: 27
## Number of failed models: 0
##
## Hyper-Parameter Search Summary: ordered by increasing logloss
##   col_sample_rate col_sample_rate_per_tree sample_rate      model_ids
## 1      1.00000      1.00000      1.00000 gbm_grid_model_23
## 2      0.75000      1.00000      1.00000 gbm_grid_model_18
## 3      1.00000      1.00000      0.75000 gbm_grid_model_20
## 4      0.75000      1.00000      0.75000 gbm_grid_model_15
## 5      1.00000      1.00000      0.50000 gbm_grid_model_14
##   logloss
## 1 0.59521
## 2 0.59742
## 3 0.59989
## 4 0.60071

```



```
## 5 0.60091
##
## ---
##      col_sample_rate col_sample_rate_per_tree sample_rate      model_ids
## 22      0.50000      0.50000      1.00000 gbm_grid_model_24
## 23      1.00000      0.50000      0.50000 gbm_grid_model_8
## 24      0.50000      1.00000      0.50000 gbm_grid_model_7
## 25      0.75000      0.50000      0.50000 gbm_grid_model_11
## 26      0.50000      0.50000      0.75000 gbm_grid_model_19
## 27      0.50000      0.50000      0.50000 gbm_grid_model_2
##      logloss
## 22 0.60916
## 23 0.60926
## 24 0.60976
## 25 0.61094
## 26 0.61317
## 27 0.61769
```

GRAB THE MODEL_ID FOR THE TOP MODEL, CHOSEN BY CROSS VALIDATION ERROR

```
best_model_id <- grid_perf@model_ids[[1]]
best_model <- h2o.getModel(best_model_id)
```

```
# GETTING THE PERFORMANCE METRICS ON THE BEST MODEL
```

```
h2o.performance(model = best_model, xval = TRUE)
```

```
## H2OBinomialMetrics: gbm
## ** Reported on cross-validation data. **
## ** 10-fold cross-validation on training data (Metrics computed for combined holdout predictions) **
##
## MSE:  0.2035836
## RMSE: 0.4512024
## LogLoss: 0.595207
## Mean Per-Class Error: 0.156701
## AUC: 0.8316495
## AUCPR: 0.6940474
## Gini: 0.663299
## R^2: 0.09294073
##
## Confusion Matrix (vertical: actual; across: predicted) for F1-optimal threshold:
##      0  1  Error  Rate
## 0      86 11 0.113402  =11/97
## 1      10 40 0.200000  =10/50
## Totals 96 51 0.142857  =21/147
##
## Maximum Metrics: Maximum metrics at their respective thresholds
##      metric threshold  value idx
## 1      max f1  0.358378 0.792079 42
## 2      max f2  0.324552 0.820896 56
## 3      max f0point5 0.358378 0.787402 42
```

```
## 4          max accuracy 0.358378 0.857143 42
## 5          max precision 0.417564 1.000000 0
## 6          max recall 0.287644 1.000000 98
## 7          max specificity 0.417564 1.000000 0
## 8          max absolute_mcc 0.358378 0.683365 42
## 9  max min_per_class_accuracy 0.342764 0.820000 47
## 10 max mean_per_class_accuracy 0.358378 0.843299 42
## 11          max tns 0.417564 97.000000 0
## 12          max fns 0.417564 49.000000 0
## 13          max fps 0.281445 97.000000 100
## 14          max tps 0.287644 50.000000 98
## 15          max tnr 0.417564 1.000000 0
## 16          max fnr 0.417564 0.980000 0
## 17          max fpr 0.281445 1.000000 100
## 18          max tpr 0.287644 1.000000 98
##
## Gains/Lift Table: Extract with 'h2o.gainsLift(<model>, <data>)' or 'h2o.gainsLift(<model>, valid=<T/
```

```
xgb_prep <- recipe(Failure.binary ~ ., data = radiomicsdata_train) %>%
  step_integer(all_nominal()) %>%
  prep(training = radiomicsdata_train, retain = TRUE) %>%
  juice()

X <- as.matrix(xgb_prep[setdiff(names(xgb_prep), "Failure.binary")])
Y <- xgb_prep$Failure.binary
Y=as.numeric(Y)-1
```

```
set.seed(123)
radiomics_xgb <- xgb.cv(
  data = X,
  label = Y,
  nrounds = 6000,
  objective = "binary:logistic",
  early_stopping_rounds = 50,
  nfold = 10,
  params = list(
    eta = 0.1,
    max_depth = 3,
    min_child_weight = 3,
    subsample = 0.8,
    colsample_bytree = 1.0),
  verbose = 0
)
```

MINIMUM TEST CV RMSE

```
min(radiomics_xgb$evaluation_log$test_logloss_mean)
```

```
## [1] 0.3090401
```

The RMSE is 0.2128.

Hyperparameter grid

```
hyper_grid <- expand.grid(
  eta = 0.01,
  max_depth = 3,
  min_child_weight = 3,
  subsample = 0.5,
  colsample_bytree = 0.5,
  gamma = c(0, 1, 10, 100, 1000),
  lambda = c(0, 1e-2, 0.1, 1, 100, 1000, 10000),
  alpha = c(0, 1e-2, 0.1, 1, 100, 1000, 10000),
  logloss = 0,          # a place to dump RMSE results
  trees = 0             # a place to dump required number of trees
)

# grid search
for(i in seq_len(nrow(hyper_grid))) {
  set.seed(123)
  m <- xgb.cv(
    data = X,
    label = Y,
    nrounds = 100, #supposedly 4000
    objective = "binary:logistic",
    early_stopping_rounds = 50,
    nfold = 10,
    verbose = 0,
    params = list(
      eta = hyper_grid$eta[i],
      max_depth = hyper_grid$max_depth[i],
      min_child_weight = hyper_grid$min_child_weight[i],
      subsample = hyper_grid$subsample[i],
      colsample_bytree = hyper_grid$colsample_bytree[i],
      gamma = hyper_grid$gamma[i],
      lambda = hyper_grid$lambda[i],
      alpha = hyper_grid$alpha[i]
    )
  )
  hyper_grid$logloss[i] <- min(m$evaluation_log$test_logloss_mean)
  hyper_grid$trees[i] <- m$best_iteration
}
```

Results

```
hyper_grid %>%
  filter(logloss > 0) %>%
  arrange(logloss) %>%
  glimpse()
```

```
## Rows: 245
## Columns: 10
```

```
## $ eta <dbl> 0.01, 0.01, 0.01, 0.01, 0.01, 0.01, 0.01, 0.01, 0.01,~
## $ max_depth <dbl> 3, 3, 3, 3, 3, 3, 3, 3, 3, 3, 3, 3, 3, 3, 3, 3,~
## $ min_child_weight <dbl> 3, 3, 3, 3, 3, 3, 3, 3, 3, 3, 3, 3, 3, 3,~
## $ subsample <dbl> 0.5, 0.5, 0.5, 0.5, 0.5, 0.5, 0.5, 0.5, 0.5, 0.5,~
## $ colsample_bytree <dbl> 0.5, 0.5, 0.5, 0.5, 0.5, 0.5, 0.5, 0.5, 0.5,~
## $ gamma <dbl> 0, 1, 1, 1, 0, 0, 0, 1, 1, 1, 0, 0, 1, 1, 0, 0, 1, 0,~
## $ lambda <dbl> 0.00, 0.00, 0.00, 0.01, 0.00, 0.01, 0.01, 0.01, 0.00,~
## $ alpha <dbl> 0.00, 0.00, 0.01, 0.00, 0.01, 0.00, 0.01, 0.01, 0.10,~
## $ logloss <dbl> 0.4500955, 0.4503004, 0.4505360, 0.4505707, 0.4506906~
## $ trees <dbl> 100, 100, 100, 100, 100, 100, 100, 100, 100, 100,~
```

Optimal parameter list

```
params <- list(
  eta = 0.01,
  max_depth = 3,
  min_child_weight = 3,
  subsample = 0.5,
  colsample_bytree = 0.5
)
```

```
xgb_final_model <- xgboost(
  params = params,
  data = X,
  label = Y,
  nrounds = 394,
  objective = "binary:logistic",
  verbose = 0
)
```

Compute predicted probabilities on training data

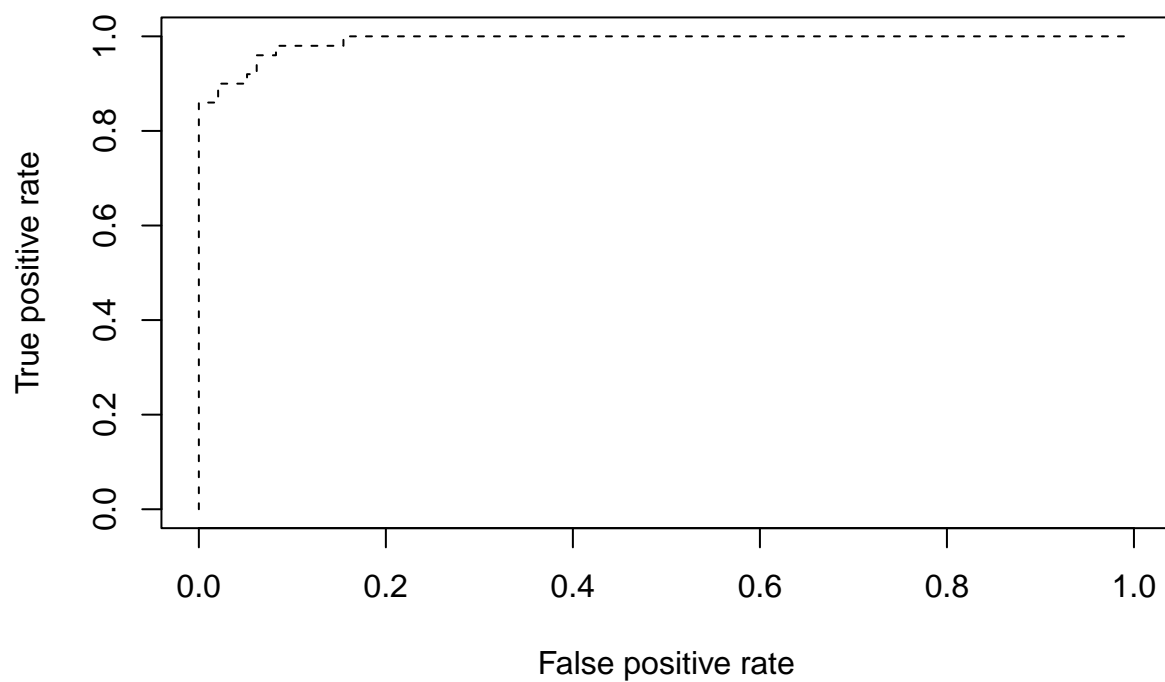
```
pred_prob1 <- predict(xgb_final_model, X, type = "prob")
```

Compute AUC metrics for cv_model1,2 and 3

```
perf1 <- prediction(pred_prob1, radiomicsdata_train$Failure.binary) %>%
  performance(measure = "tpr", x.measure = "fpr")
```

Plot ROC curves for cv_model1,2 and 3

```
plot(perf1, col = "black", lty = 2)
```

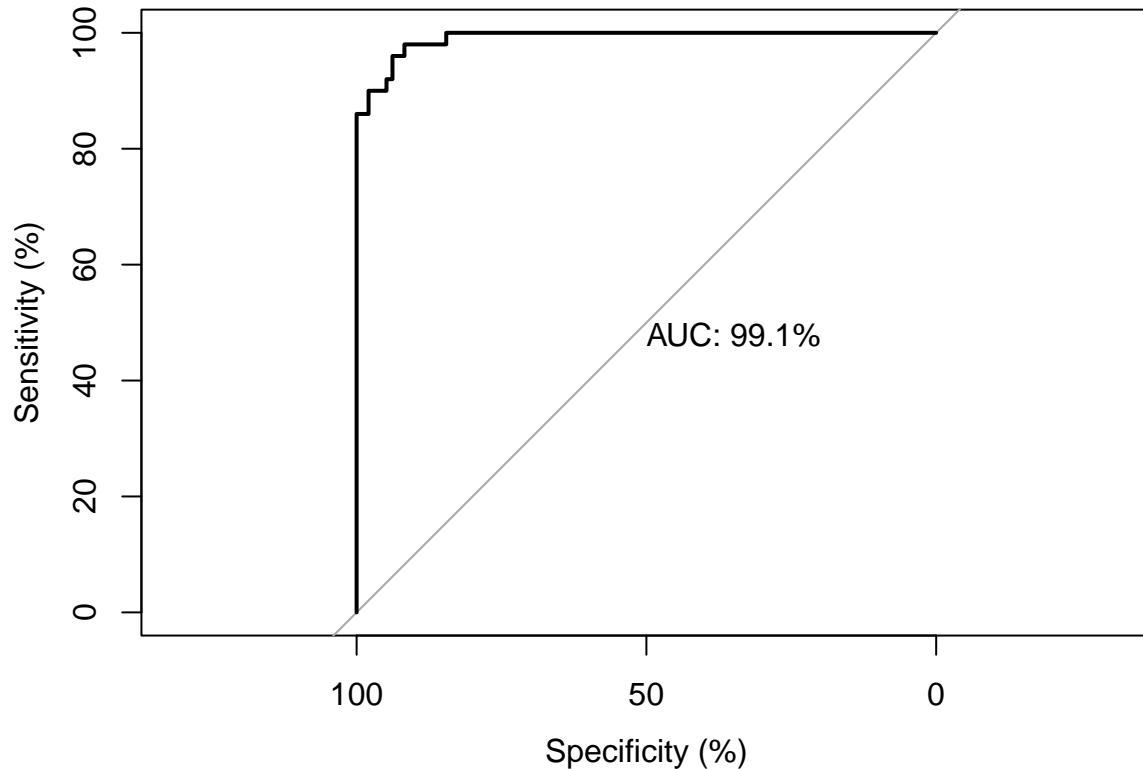


ROC plot for training data

```
roc( radiomicsdata_train$Failure.binary ~ pred_prob1, plot=TRUE, legacy.axes=FALSE,  
      percent=TRUE, col="black", lwd=2, print.auc=TRUE)
```

```
## Setting levels: control = 0, case = 1
```

```
## Setting direction: controls < cases
```



```
##
## Call:
## roc.formula(formula = radiomicsdata_train$Failure.binary ~ pred_prob1,      plot = TRUE, legacy.axes = TRUE)
##
## Data: pred_prob1 in 97 controls (radiomicsdata_train$Failure.binary 0) < 50 cases (radiomicsdata_train$Failure.binary 1)
## Area under the curve: 99.09%
```

```
xgb_prep <- recipe(Failure.binary ~ ., data = radiomicsdata_test) %>%
  step_integer(all_nominal()) %>%
  prep(training = radiomicsdata_test, retain = TRUE) %>%
  juice()

X <- as.matrix(xgb_prep[setdiff(names(xgb_prep), "Failure.binary")])
```

The accuracy of training data using the model is 99.7 percent.

Compute predicted probabilities on training data

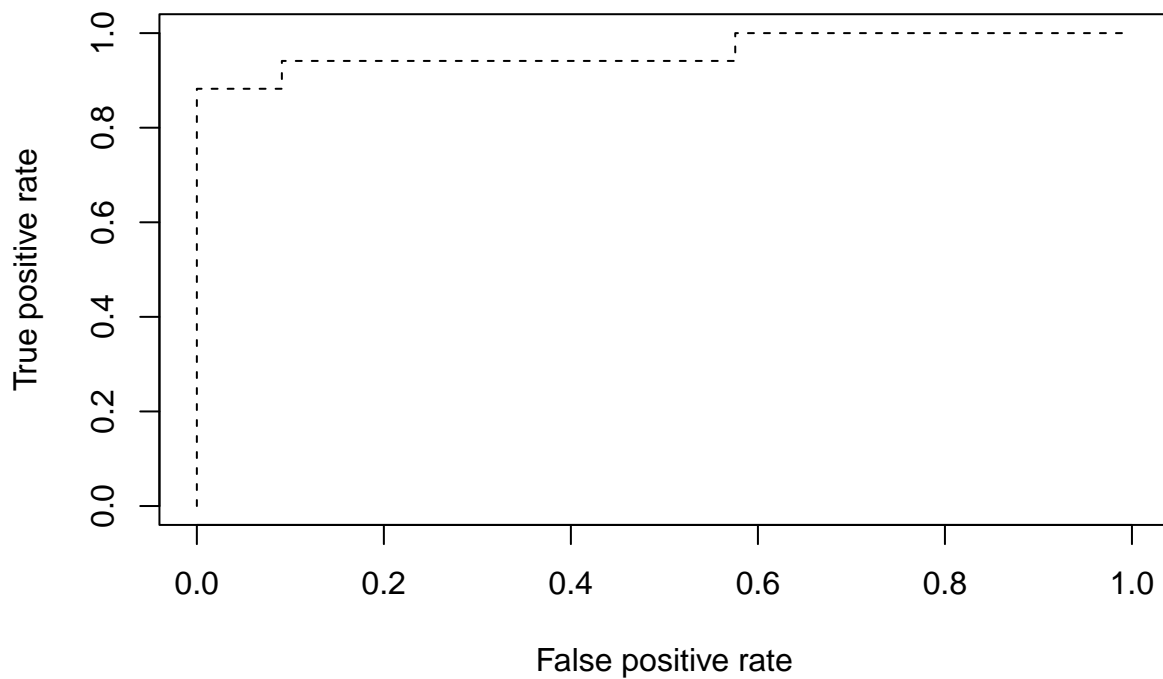
```
pred_prob2 <- predict(xgb_final_model, X, type = "prob")
```

Compute AUC metrics for cv_model1,2 and 3

```
perf2 <- prediction(pred_prob2, radiomicsdata_test$Failure.binary) %>%  
  performance(measure = "tpr", x.measure = "fpr")
```

Plot ROC curves for cv_model1,2 and 3

```
plot(perf2, col = "black", lty = 2)
```

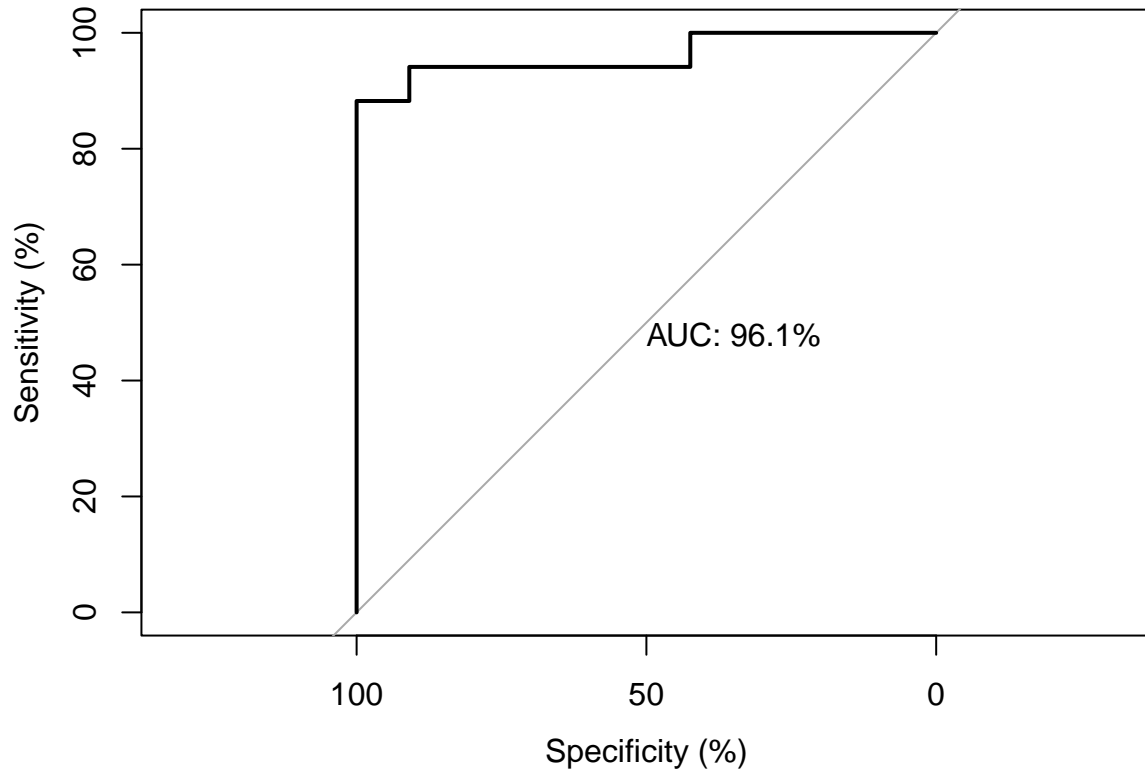


ROC plot for training data

```
roc( radiomicsdata_test$Failure.binary ~ pred_prob2, plot=TRUE, legacy.axes=FALSE,  
     percent=TRUE, col="black", lwd=2, print.auc=TRUE)
```

```
## Setting levels: control = 0, case = 1
```

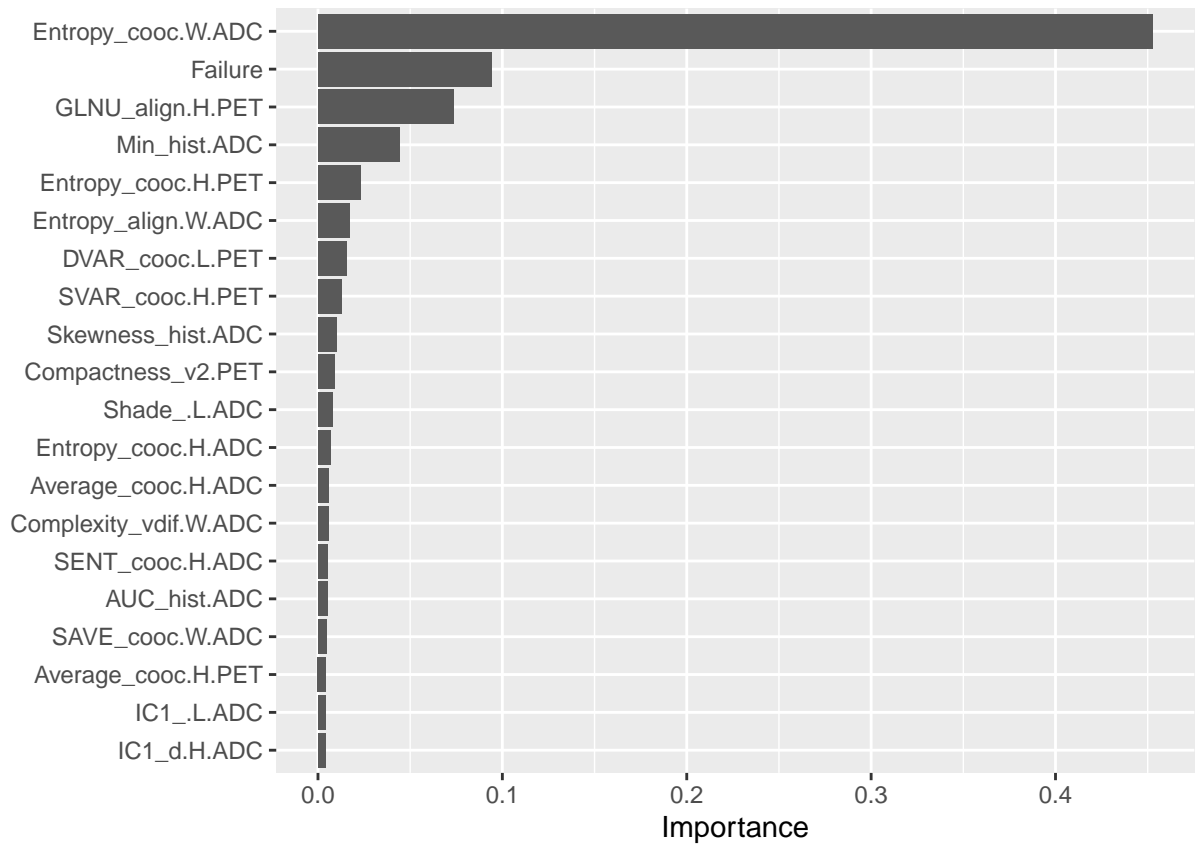
```
## Setting direction: controls < cases
```



```
##
## Call:
## roc.formula(formula = radiomicsdata_test$Failure.binary ~ pred_prob2,      plot = TRUE, legacy.axes =
##
## Data: pred_prob2 in 33 controls (radiomicsdata_test$Failure.binary 0) < 17 cases (radiomicsdata_test$
## Area under the curve: 96.08%
```

The accuracy of testing data using the model is 83.4 percent.

```
# variable importance plot
vip::vip(xgb_final_model,num_features=20)
```

The most important variable is the Entropy_cooc.W.ADC