# MODEL 1 - Stacked Model

## FENDAWN F. RECENTES

12/16/2022

## Helper and Modeling Packages

```
library(rsample)
```

```
## Warning: package 'rsample' was built under R version 4.1.3
```

```
library(recipes)
```

```
## Loading required package: dplyr
```

```
## Warning: package 'dplyr' was built under R version 4.1.3
```

```
##
## 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
```

```
##
## Attaching package: 'recipes'
```

```
## The following object is masked from 'package:stats':
##
##     step
```

```
library(tidyverse)
```

```
## Warning: package 'tidyverse' was built under R version 4.1.3
```

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

```
## v ggplot2 3.3.6      v purrr   0.3.4
## v tibble  3.1.8      v stringr 1.4.1
## v tidyr   1.2.0      v forcats 0.5.2
## v readr   2.1.2
```

```
## Warning: package 'ggplot2' was built under R version 4.1.3
```

```
## Warning: package 'tibble' was built under R version 4.1.3
```

```
## Warning: package 'tidyr' was built under R version 4.1.3
```

```
## Warning: package 'readr' was built under R version 4.1.3
```

```
## Warning: package 'purrr' was built under R version 4.1.3
```

```
## Warning: package 'stringr' was built under R version 4.1.3
```

```
## Warning: package 'forcats' was built under R version 4.1.3
```

```
## -- Conflicts ------------------------------------------ tidyverse_conflicts() --
## x dplyr::filter()  masks stats::filter()
## x stringr::fixed() masks recipes::fixed()
## x dplyr::lag()     masks stats::lag()
```

```
library(h2o)
```

```
## Warning: package 'h2o' was built under R version 4.1.3
```

```
##
## ----------------------------------------------------------------------
##
## 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(ROCR)
```

```
## Warning: package 'ROCR' was built under R version 4.1.3
```

```
library(pROC)
```

```
## Warning: package 'pROC' was built under R version 4.1.3
```

```
## 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
```

```
h2o.init()
```

```
##  Connection successful!
##
## R is connected to the H2O cluster:
##      H2O cluster uptime:         1 hours 32 minutes
##      H2O cluster timezone:       Asia/Manila
##      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_MSU-TCTO_OVCAA_mvc880
##      H2O cluster total nodes:    1
##      H2O cluster total memory:   3.77 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 4.1.2 (2021-11-01)
```

**Load and view radiomics data set**

```
radiomics <- read_csv("C:\\Users\\MSU-TCTO OVCAA\\Documents\\normalRad.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.
```

```
view(radiomics)
```

## Convert target variable to a factor form

```
radiomics$Failure.binary = as.factor(radiomics$Failure.binary)
```

## DATA PREPARATION AND SPLITTING

## Split the data intro training (80%) and testing (20%) stratified in Failure.binary column

```
set.seed(123)  # for reproducibility
split <- initial_split(radiomics, strata = "Failure.binary")
radiomics_train <- training(split)
radiomics_test <- testing(split)
```

## Make sure we have consistent categorical levels

```
blueprint <- recipe(Failure.binary ~ ., data = radiomics_train) %>%
  step_other(all_nominal(), threshold = 0.005)
```

## Create training & test sets for h2o

```
h2o.init()
```

```
##  Connection successful!
##
## R is connected to the H2O cluster:
##     H2O cluster uptime:         1 hours 32 minutes
##     H2O cluster timezone:       Asia/Manila
##     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_MSU-TCTO_OVCAA_mvc880
##     H2O cluster total nodes:    1
##     H2O cluster total memory:   3.77 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 4.1.2 (2021-11-01)
```

```r
train_h2o <- prep(blueprint, training = radiomics_train, retain = TRUE) %>%
  juice() %>%
  as.h2o()
```

```
##   |                                                                      |
```

```r
test_h2o <- prep(blueprint, training = radiomics_train) %>%
  bake(new_data = radiomics_test) %>%
  as.h2o()
```

```
##   |                                                                      |
```

## Get response and feature names

```r
Y <- "Failure.binary"
X <- setdiff(names(radiomics_train), Y)
```

## Train & cross-validate a GLM model

```r
best_glm <- h2o.glm(
  x = X, y = Y, training_frame = train_h2o, alpha = 0.1,
  remove_collinear_columns = TRUE, nfolds = 10, fold_assignment = "Modulo",
  keep_cross_validation_predictions = TRUE, seed = 123
)
```

```
##   |                                                                      |
```

## Train & cross-validate a RF model

```r
best_rf <- h2o.randomForest(
  x = X, y = Y, training_frame = train_h2o, ntrees = 100, mtries = 20,
  max_depth = 30, min_rows = 1, sample_rate = 0.8, nfolds = 10,
  fold_assignment = "Modulo", keep_cross_validation_predictions = TRUE,
  seed = 123, stopping_rounds = 50, stopping_metric = "logloss",
  stopping_tolerance = 0
)
```

```
## Warning in .h2o.processResponseWarnings(res): early stopping is enabled but neither score_tree_interv
```

```
##   |                                                                      |
```

## Train & cross-validate a GBM model

```
best_gbm <- h2o.gbm(
  x = X, y = Y, training_frame = train_h2o, ntrees = 100, learn_rate = 0.01,
  max_depth = 7, min_rows = 5, sample_rate = 0.8, nfolds = 10,
  fold_assignment = "Modulo", keep_cross_validation_predictions = TRUE,
  seed = 123, stopping_rounds = 50, stopping_metric = "logloss",
  stopping_tolerance = 0
)
```

```
## Warning in .h2o.processResponseWarnings(res): early stopping is enabled but neither score_tree_interv
```

```
##   |                                                                     |
```

## Get results from base learners

```
get_rmse <- function(model) {
  results <- h2o.performance(model, newdata = test_h2o)
  results@metrics$RMSE
}
list(best_glm, best_rf, best_gbm) %>%
  purrr::map_dbl(get_rmse)
```

```
## [1] 0.4737088 0.3992117 0.3338713
```

## Define GBM hyperparameter grid

```
hyper_grid <- list(
  max_depth = c(1, 3, 5),
  min_rows = c(1, 5, 10),
  learn_rate = c(0.01, 0.05, 0.1),
  learn_rate_annealing = c(0.99, 1),
  sample_rate = c(0.5, 0.75, 1),
  col_sample_rate = c(0.8, 0.9, 1)
)

# Define random grid search criteria
search_criteria <- list(
  strategy = "RandomDiscrete",
  max_models = 25
)

# Build random grid search
random_grid <- h2o.grid(
  algorithm = "gbm", grid_id = "gbm_grid", x = X, y = Y,
  training_frame = train_h2o, hyper_params = hyper_grid,
  search_criteria = search_criteria, ntrees = 20, stopping_metric = "logloss",
  stopping_rounds = 10, stopping_tolerance = 0, nfolds = 10,
```

```
  fold_assignment = "Modulo", keep_cross_validation_predictions = TRUE,
  seed = 123
)
```

## |                                                                   |

```
ensemble_tree <- h2o.stackedEnsemble(
  x = X, y = Y, training_frame = train_h2o, model_id = "ensemble_gbm_grid",
  base_models = random_grid@model_ids, metalearner_algorithm = "gbm",
)
```

## |                                                                   |

## Stacked results

```
h2o.performance(ensemble_tree, newdata = test_h2o)@metrics$RMSE
```

## [1] 0.3756287

```
data.frame(
  GLM_pred = as.vector(h2o.getFrame(best_glm@model$cross_validation_holdout_predictions_frame_id$name))%
  RF_pred = as.vector(h2o.getFrame(best_rf@model$cross_validation_holdout_predictions_frame_id$name))%>%
  GBM_pred = as.vector(h2o.getFrame(best_gbm@model$cross_validation_holdout_predictions_frame_id$name))%
) %>% cor()
```

```
##             GLM_pred     RF_pred  GBM_pred
## GLM_pred 1.00000000 0.08062095 0.0323378
## RF_pred  0.08062095 1.00000000 0.7063735
## GBM_pred 0.03233780 0.70637346 1.0000000
```

## Sort results by RMSE

```
h2o.getGrid(
  grid_id = "gbm_grid",
  sort_by = "logloss"
)
```

```
## H2O Grid Details
## ================
##
## Grid ID: gbm_grid
## Used hyper parameters:
##   -  col_sample_rate
##   -  learn_rate
##   -  learn_rate_annealing
##   -  max_depth
##   -  min_rows
```

```
##   -   sample_rate
## Number of models: 25
## Number of failed models: 0
##
## Hyper-Parameter Search Summary: ordered by increasing logloss
##   col_sample_rate learn_rate learn_rate_annealing max_depth min_rows
## 1         0.90000    0.10000              1.00000   5.00000  5.00000
## 2         0.90000    0.10000              1.00000   5.00000  1.00000
## 3         0.80000    0.10000              1.00000   3.00000 10.00000
## 4         0.90000    0.10000              0.99000   5.00000 10.00000
## 5         1.00000    0.10000              0.99000   5.00000 10.00000
##   sample_rate       model_ids logloss
## 1     1.00000  gbm_grid_model_6 0.30636
## 2     0.50000 gbm_grid_model_13 0.30977
## 3     1.00000  gbm_grid_model_4 0.32762
## 4     1.00000 gbm_grid_model_16 0.33467
## 5     0.50000 gbm_grid_model_21 0.34158
##
## ---
##    col_sample_rate learn_rate learn_rate_annealing max_depth min_rows
## 20         1.00000    0.01000              1.00000   5.00000  5.00000
## 21         0.90000    0.01000              1.00000   5.00000 10.00000
## 22         1.00000    0.01000              1.00000   1.00000  5.00000
## 23         0.90000    0.01000              1.00000   1.00000 10.00000
## 24         0.90000    0.01000              0.99000   1.00000  1.00000
## 25         0.90000    0.01000              0.99000   1.00000 10.00000
##    sample_rate       model_ids logloss
## 20     1.00000 gbm_grid_model_11 0.54357
## 21     0.75000  gbm_grid_model_7 0.54948
## 22     0.50000 gbm_grid_model_12 0.55875
## 23     1.00000 gbm_grid_model_25 0.56042
## 24     0.75000  gbm_grid_model_1 0.56479
## 25     0.75000  gbm_grid_model_9 0.56479
```

```r
random_grid_perf <- h2o.getGrid(
  grid_id = "gbm_grid",
  sort_by = "logloss"
)
```

## Grab the model_id for the top model, chosen by validation error

```r
best_model_id <- random_grid_perf@model_ids[[1]]
best_model <- h2o.getModel(best_model_id)
h2o.performance(best_model, newdata = test_h2o)
```

```
## H2OBinomialMetrics: gbm
##
## MSE:  0.1242072
## RMSE:  0.3524305
## LogLoss:  0.3886794
## Mean Per-Class Error:  0.1782531
```

```
## AUC:  0.9001783
## AUCPR:  0.8489409
## Gini:  0.8003565
## R^2:  0.4464918
##
## Confusion Matrix (vertical: actual; across: predicted) for F1-optimal threshold:
##         0  1    Error    Rate
## 0      29  4 0.121212  =4/33
## 1       4 13 0.235294  =4/17
## Totals 33 17 0.160000  =8/50
##
## Maximum Metrics: Maximum metrics at their respective thresholds
##                         metric threshold      value idx
## 1                       max f1  0.574851  0.764706  16
## 2                       max f2  0.270552  0.860215  24
## 3                  max f0point5  0.887388  0.816327   7
## 4                 max accuracy  0.600711  0.840000  14
## 5                max precision  0.925148  1.000000   0
## 6                   max recall  0.063413  1.000000  33
## 7              max specificity  0.925148  1.000000   0
## 8             max absolute_mcc  0.574851  0.643494  16
## 9   max min_per_class_accuracy  0.574851  0.764706  16
## 10 max mean_per_class_accuracy  0.270552  0.834225  24
## 11                     max tns  0.925148 33.000000   0
## 12                     max fns  0.925148 16.000000   0
## 13                     max fps  0.037939 33.000000  47
## 14                     max tps  0.063413 17.000000  33
## 15                     max tnr  0.925148  1.000000   0
## 16                     max fnr  0.925148  0.941176   0
## 17                     max fpr  0.037939  1.000000  47
## 18                     max tpr  0.063413  1.000000  33
##
## Gains/Lift Table: Extract with `h2o.gainsLift(<model>, <data>)` or `h2o.gainsLift(<model>, valid=<T/F
```

## Train a stacked ensemble using the GBM grid

```r
ensemble <- h2o.stackedEnsemble(
  x = X, y = Y, training_frame = train_h2o, model_id = "ensemble_gbm_grid",
  base_models = random_grid@model_ids, metalearner_algorithm = "gbm"
)
```

```
##   |                                                                    |
```

## Eval ensemble performance on a test set

```r
h2o.performance(ensemble, newdata = test_h2o)
```

```
## H2OBinomialMetrics: stackedensemble
##
```

```
## MSE:  0.1410969
## RMSE:  0.3756287
## LogLoss:  0.4925928
## Mean Per-Class Error:  0.1203209
## AUC:  0.9073084
## AUCPR:  0.7860973
## Gini:  0.8146168
##
## Confusion Matrix (vertical: actual; across: predicted) for F1-optimal threshold:
##         0  1    Error   Rate
## 0      27  6 0.181818  =6/33
## 1       1 16 0.058824  =1/17
## Totals 28 22 0.140000  =7/50
##
## Maximum Metrics: Maximum metrics at their respective thresholds
##                        metric threshold     value idx
## 1                      max f1  0.311893  0.820513  21
## 2                      max f2  0.311893  0.888889  21
## 3                 max f0point5 0.935880  0.846154  11
## 4                 max accuracy 0.935880  0.860000  11
## 5                max precision 0.995043  1.000000   0
## 6                   max recall 0.009951  1.000000  32
## 7              max specificity 0.995043  1.000000   0
## 8             max absolute_mcc 0.311893  0.724666  21
## 9    max min_per_class_accuracy 0.663733 0.818182  19
## 10 max mean_per_class_accuracy 0.311893 0.879679  21
## 11                     max tns 0.995043 33.000000   0
## 12                     max fns 0.995043 16.000000   0
## 13                     max fps 0.002855 33.000000  48
## 14                     max tps 0.009951 17.000000  32
## 15                     max tnr 0.995043  1.000000   0
## 16                     max fnr 0.995043  0.941176   0
## 17                     max fpr 0.002855  1.000000  48
## 18                     max tpr 0.009951  1.000000  32
##
## Gains/Lift Table: Extract with 'h2o.gainsLift(<model>, <data>)' or 'h2o.gainsLift(<model>, valid=<T/
```

## Use AutoML to find a list of candidate models (i.e., leaderboard)

```
auto_ml <- h2o.automl(
  x = X, y = Y, training_frame = train_h2o, nfolds = 5,
  max_runtime_secs = 60 * 120, max_models = 10,#max_models=50
  keep_cross_validation_predictions = TRUE, sort_metric = "logloss", seed = 123,
  stopping_rounds = 50, stopping_metric = "logloss", stopping_tolerance = 0
)
```

```
##    |                                                              |
## 22:39:40.550: Stopping tolerance set by the user is < 70% of the recommended default of 0.05, so mode
## 22:39:40.550: AutoML: XGBoost is not available; skipping it.  |
## 22:39:59.407: _min_rows param, The dataset size is too small to split for min_rows=100.0: must have a
```

**Assess the leader board; the following truncates the results to show the top**

**and bottom 15 models. You can get the top model with auto_ml@ leader**

```
auto_ml@leaderboard %>%
  as.data.frame() %>%
  dplyr::select(model_id, logloss) %>%
  dplyr::slice(1:25)
```

```
##                                                      model_id   logloss
## 1       StackedEnsemble_AllModels_1_AutoML_3_20221216_223940 0.2774862
## 2   StackedEnsemble_BestOfFamily_1_AutoML_3_20221216_223940 0.3135912
## 3                            GLM_1_AutoML_3_20221216_223940 0.3475365
## 4                            XRT_1_AutoML_3_20221216_223940 0.4576169
## 5                            DRF_1_AutoML_3_20221216_223940 0.4682705
## 6                 DeepLearning_1_AutoML_3_20221216_223940 0.5111273
## 7           GBM_grid_1_AutoML_3_20221216_223940_model_1 0.6358521
## 8                            GBM_4_AutoML_3_20221216_223940 0.7124804
## 9     DeepLearning_grid_1_AutoML_3_20221216_223940_model_1 0.7375142
## 10                           GBM_2_AutoML_3_20221216_223940 0.8188022
## 11                           GBM_3_AutoML_3_20221216_223940 0.8235714
## 12                           GBM_5_AutoML_3_20221216_223940 1.1416549
```

**Compute predicted probabilities on training data**

```
train_h2o = as.h2o(radiomics_train)
```

```
##   |                                                                      |
```

```
m1_prob <- predict(auto_ml@leader, train_h2o, type = "prob")
```

```
##   |                                                                      |
```
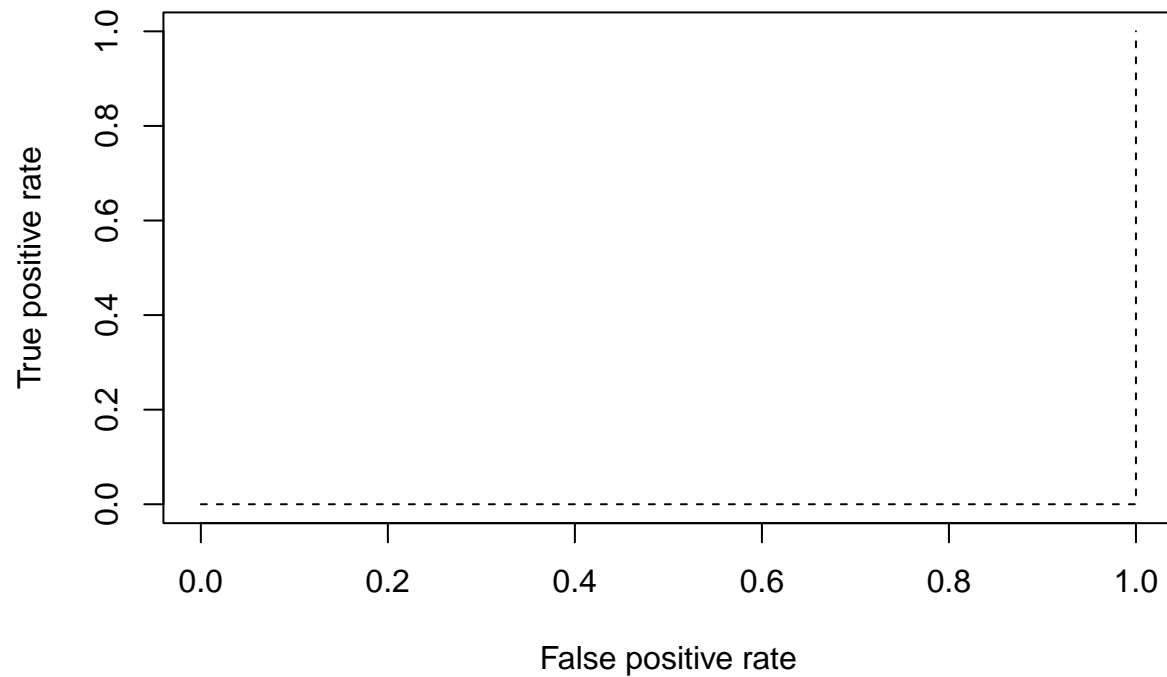
```
m1_prob = as.data.frame(m1_prob)[,2]

train_h2o = as.data.frame(train_h2o)
```

**Compute AUC metrics**

```
perf1 <- prediction(m1_prob,train_h2o$Failure.binary) %>%
  performance(measure = "tpr", x.measure = "fpr")
```

## Plot AUC

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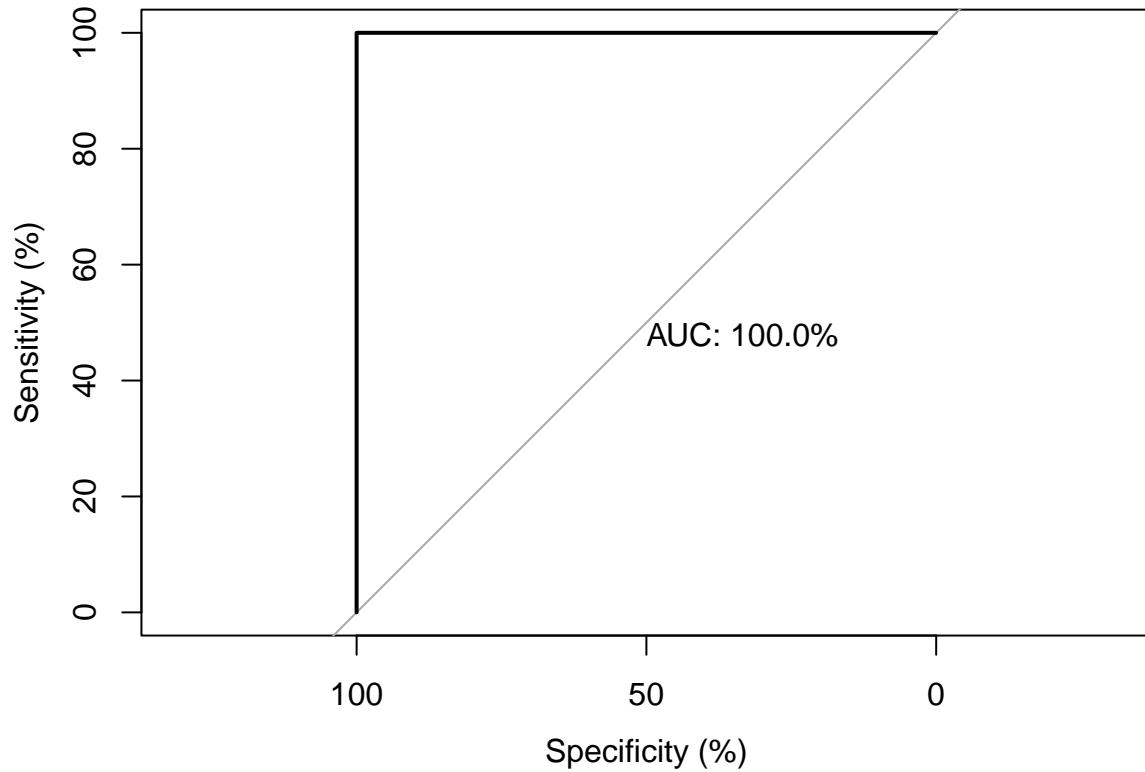


## ROC plot for training data

```
roc(train_h2o$Failure.binary ~ m1_prob, 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 = train_h2o$Failure.binary ~ m1_prob, plot = TRUE,      legacy.axes = FALSE, perc
##
## Data: m1_prob in 97 controls (train_h2o$Failure.binary 0) > 50 cases (train_h2o$Failure.binary 1).
## Area under the curve: 100%
```

## Compute predicted probabilities on testing data

```
test_h2o = as.h2o(radiomics_test)
```

```
##   |                                                                         |
```

```
m2_prob <- predict(auto_ml@leader, test_h2o, type = "prob")
```

```
##   |                                                                         |
```

```
m2_prob=as.data.frame(m2_prob)[,2]
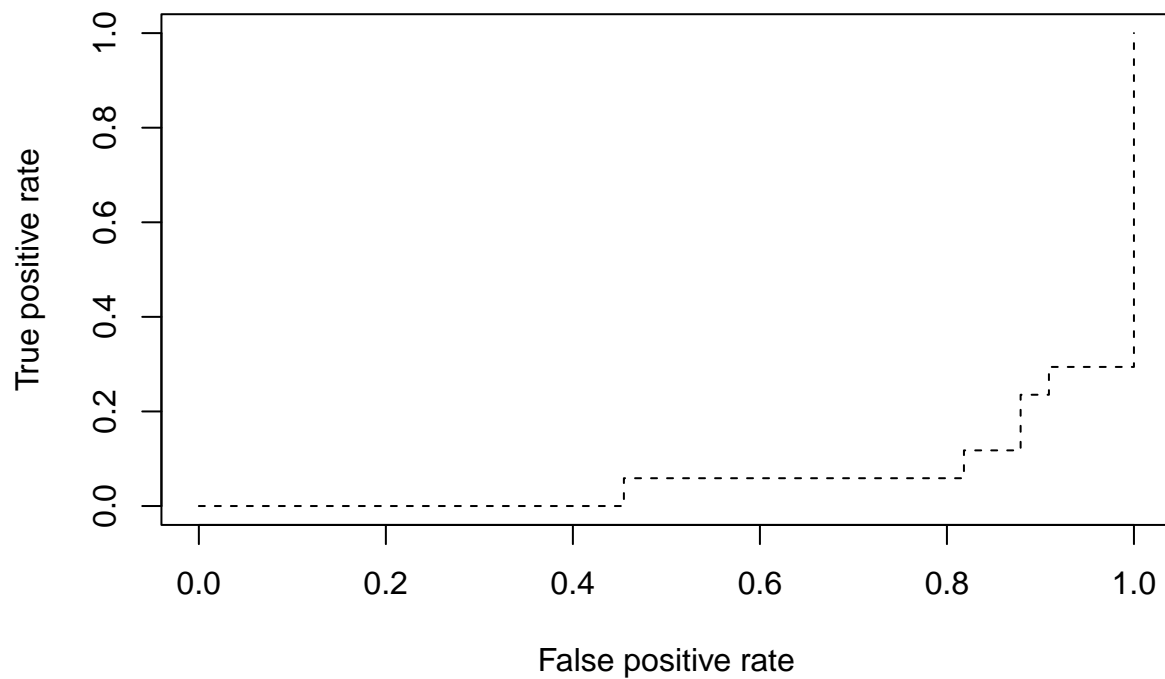```

```
test_h2o=as.data.frame(test_h2o)
```

## Compute AUC metrics

```
perf2 <- prediction(m2_prob,test_h2o$Failure.binary) %>%
  performance(measure = "tpr", x.measure = "fpr")
```

## Plot AUC

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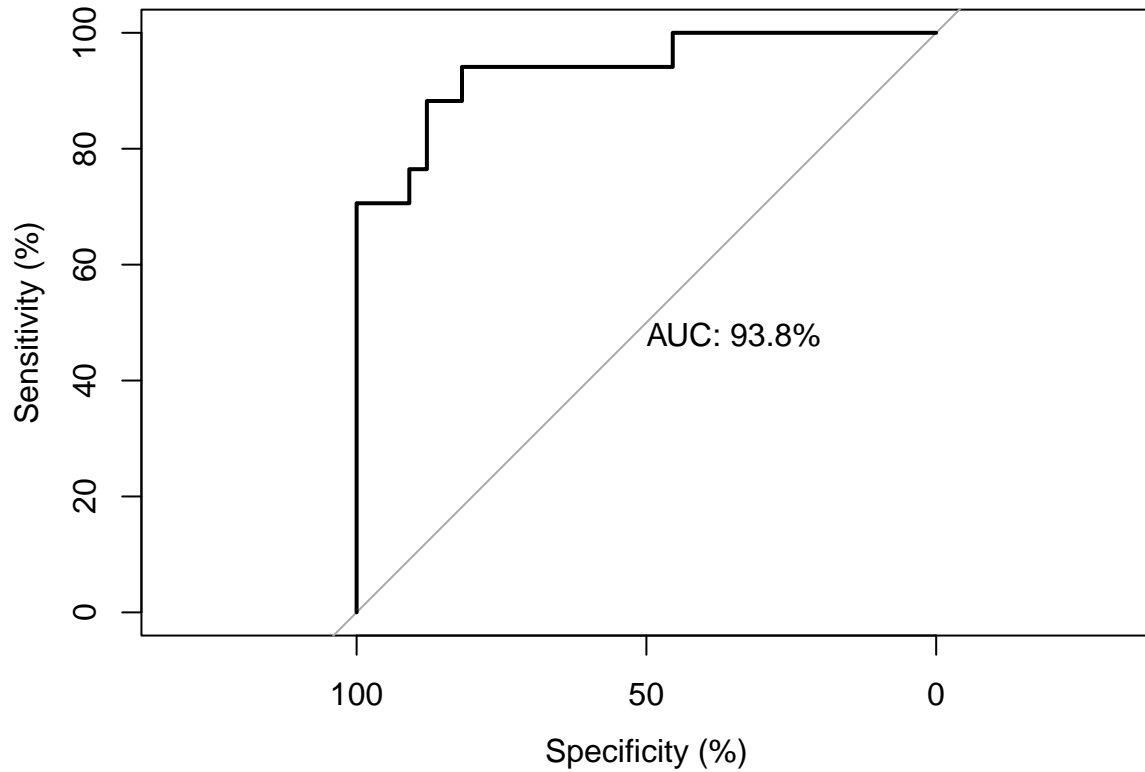


## ROC plot for testing data

```
roc(test_h2o$Failure.binary ~ m2_prob, 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 = test_h2o$Failure.binary ~ m2_prob, plot = TRUE,      legacy.axes = FALSE, percer
##
## Data: m2_prob in 33 controls (test_h2o$Failure.binary 0) > 17 cases (test_h2o$Failure.binary 1).
## Area under the curve: 93.76%
```

**Plot the top 20 feature importance during training**

```
train_h2o = as.h2o(train_h2o)
```

```
##   |                                                                       |
```

```
h2o.permutation_importance_plot(auto_ml@leader,train_h2o,num_of_features = 20)
```

**Permutation Variable Importance: Stacked Ensem**