

# Model Selection - death\_2year

Eduardo Yuki Yada

## Global parameters

```
k <- params$k # Number of folds for cross validation
grid_size <- params$grid_size # Number of parameter combination to tune on each model
repeats <- params$repeats
```

Minutes to run: 0

## Imports

```
library(tidyverse)
library(yaml)
library(tidymodels)
library(usemodels)
library(vip)
library(bonsai)
library(lightgbm)
library(caret)
library(pROC)

source("aux_functions.R")
```

Minutes to run: 0

## Loading data

```
load('dataset/processed_data.RData')
load('dataset/processed_dictionary.RData')

columns_list <- yaml.load_file("../auxiliar/columns_list.yaml")

outcome_column <- params$outcome_column
features_list <- params$features_list

df <- mutate(df, across(where(is.character), as.factor))
```

Minutes to run: 0.005

```
dir.create(file.path("../auxiliar/model_selection/hyperparameters/"),
           showWarnings = FALSE,
           recursive = TRUE)

dir.create(file.path("../auxiliar/model_selection/performance/"),
           showWarnings = FALSE,
           recursive = TRUE)
```

Minutes to run: 0

## Eligible features

```
cat_features_list = readRDS(sprintf(
  "./auxiliar/significant_columns/categorical_%s.rds",
  outcome_column
))

num_features_list = readRDS(sprintf(
  "./auxiliar/significant_columns/numerical_%s.rds",
  outcome_column
))

features_list = c(cat_features_list, num_features_list)
```

Minutes to run: 0

```
eligible_columns = df_names %>%
  filter(momento.aquisicao == 'Admissão t0') %>%
  .$variable.name

exception_columns = c('death_intraop', 'death_intraop_1', 'disch_outcomes_t0')

correlated_columns = c('year_procedure_1', # com year_adm_t0
  'age_surgery_1', # com age
  'admission_t0', # com admission_pre_t0_count
  'atb', # com meds_antimicrobianos
  'classe_meds_cardio_qtde', # com classe_meds_qtde
  'suporte_hemod', # com proced_invasivos_qtde,
  'radiografia', # com exames_imagem_qtde
  'ecg' # com metodos_graficos_qtde
)

eligible_features = eligible_columns %>%
  base::intersect(c(columns_list$categories_columns, columns_list$numerical_columns)) %>%
  setdiff(c(exception_columns, correlated_columns))

if (is.null(features_list)) {
  features = eligible_features
} else {
  features = base::intersect(eligible_features, features_list)
}

gluedown::md_order(features, seq = TRUE, pad = TRUE)

## 01. sex
## 02. age
## 03. education_level
## 04. underlying_heart_disease
## 05. heart_disease
## 06. nyha_basal
## 07. hypertension
## 08. prior_mi
## 09. heart_failure
## 10. af
## 11. cardiac_arrest
## 12. valvopathy
## 13. diabetes
## 14. renal_failure
## 15. hemodialysis
## 16. stroke
## 17. copd
## 18. comorbidities_count
```

## 19. procedure\_type\_1  
## 20. reop\_type\_1  
## 21. procedure\_type\_new  
## 22. cied\_final\_1  
## 23. cied\_final\_group\_1  
## 24. admission\_pre\_t0\_count  
## 25. admission\_pre\_t0\_180d  
## 26. year\_adm\_t0  
## 27. icu\_t0  
## 28. dialysis\_t0  
## 29. admission\_t0\_emergency  
## 30. aco  
## 31. antiarritmico  
## 32. ieca\_bra  
## 33. dva  
## 34. digoxina  
## 35. estatina  
## 36. diuretico  
## 37. vasodilatador  
## 38. insuf\_cardiaca  
## 39. espirolactona  
## 40. antiplaquetario\_ev  
## 41. insulina  
## 42. psicofarmacos  
## 43. antifungico  
## 44. antiviral  
## 45. classe\_meds\_qtde  
## 46. meds\_cardiovasc\_qtde  
## 47. meds\_antimicrobianos  
## 48. vni  
## 49. ventilacao\_mecanica  
## 50. transplante\_cardiaco  
## 51. outros\_proced\_cirurgicos  
## 52. icp  
## 53. angioplastia  
## 54. cateterismo  
## 55. cateter\_venoso\_central  
## 56. proced\_invasivos\_qtde  
## 57. transfusao  
## 58. interconsulta  
## 59. equipe\_multiprof  
## 60. holter  
## 61. teste\_esforco  
## 62. tilt\_teste  
## 63. metodos\_graficos\_qtde  
## 64. laboratorio  
## 65. cultura  
## 66. analises\_clinicas\_qtde  
## 67. citologia  
## 68. histopatologia\_qtde  
## 69. angio\_tc  
## 70. cintilografia  
## 71. ecocardiograma  
## 72. endoscopia  
## 73. flebografia  
## 74. pet\_ct  
## 75. ultrassom  
## 76. tomografia  
## 77. ressonancia  
## 78. exames\_imagem\_qtde  
## 79. bic

```
## 80. hospital_stay
```

Minutes to run: 0

## Train test split (70%/30%)

```
set.seed(42)

if (outcome_column == 'readmission_30d') {
  df_split <- readRDS("./dataset/split_object.rds")
} else {
  df_split <- initial_split(df, prop = .7, strata = all_of(outcome_column))
}

df_train <- training(df_split) %>% dplyr::select(all_of(c(features, outcome_column)))
df_test <- testing(df_split) %>% dplyr::select(all_of(c(features, outcome_column)))

df_folds <- vfold_cv(df_train, v = k,
                     strata = all_of(outcome_column))
```

Minutes to run: 0.001

## Boosted Tree (XGBoost)

```
xgboost_recipe <-
  recipe(formula = sprintf("%s ~ .", outcome_column) %>% as.formula, data = df_train) %>%
  step_novel(all_nominal_predictors()) %>%
  step_unknown(all_nominal_predictors()) %>%
  step_other(all_nominal_predictors(), threshold = 0.05, other = ".merged") %>%
  step_dummy(all_nominal_predictors())

xgboost_spec <- boost_tree(
  trees = tune(),
  min_n = tune(),
  tree_depth = tune(),
  learn_rate = tune(),
  loss_reduction = tune(),
  sample_size = tune()
) %>%
  set_engine("xgboost",
             nthread = 8) %>%
  set_mode("classification")

xgboost_grid <- grid_latin_hypercube(
  trees(range = c(50L, 300L)),
  min_n(),
  tree_depth(),
  learn_rate(range = c(0.01, 0.3), trans = NULL),
  loss_reduction(),
  sample_prop(range = c(1/10, 1), trans = NULL),
  size = grid_size
)

xgboost_workflow <-
  workflow() %>%
  add_recipe(xgboost_recipe) %>%
  add_model(xgboost_spec)

xgboost_tune <-
  xgboost_workflow %>%
```

```

tune_grid(resamples = df_folds,
          grid = xgboost_grid)

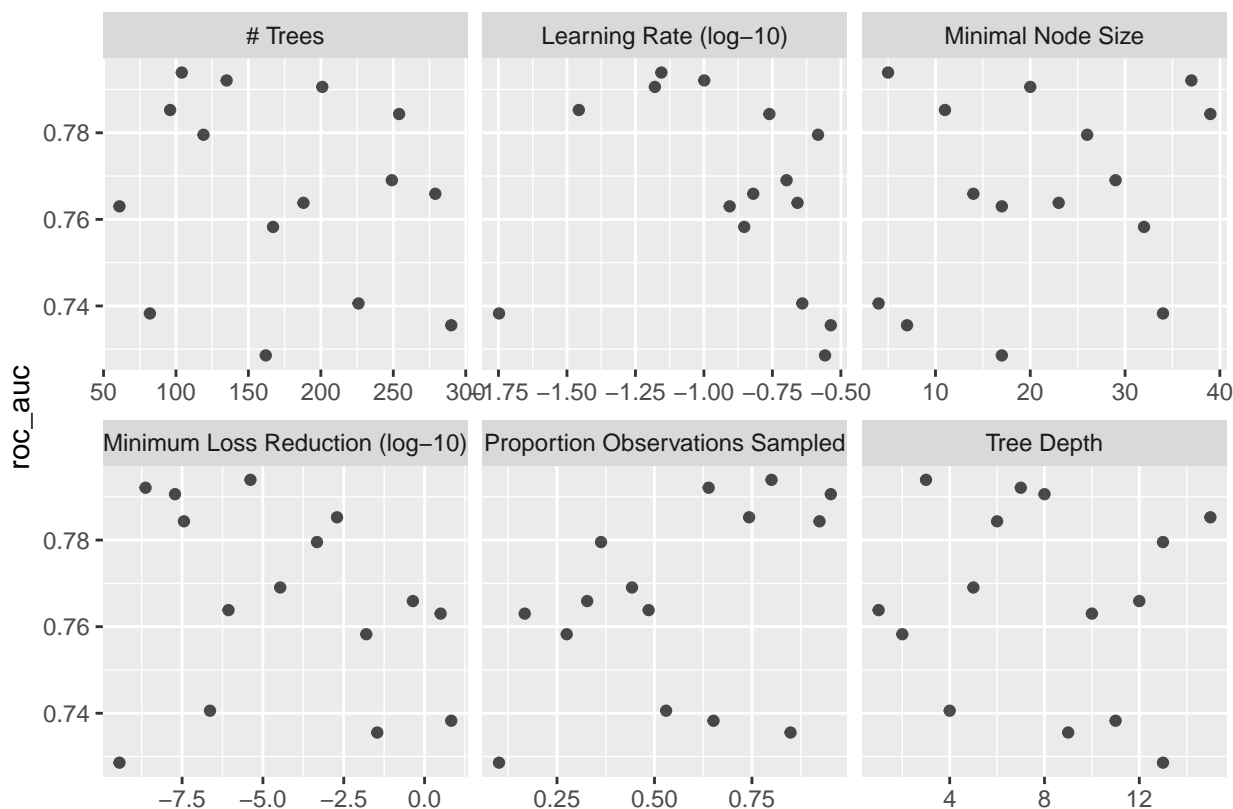
xgboost_tune %>%
  show_best("roc_auc")

## # A tibble: 5 x 12
##   trees min_n tree_depth learn_rate loss_reduction sample~1 .metric .esti~2 mean      n std_err .config
##   <int> <int>    <int>    <dbl>      <dbl>      <dbl> <chr>  <chr>  <dbl> <int>  <dbl> <chr>
## 1   104     5        3    0.0699  0.00000412    0.800 roc_auc binary  0.794     4 0.0105 Prepro~
## 2   135    37        7    0.100  0.00000000234  0.639 roc_auc binary  0.792     4 0.0104 Prepro~
## 3   201    20        8    0.0662  0.0000000191  0.952 roc_auc binary  0.791     4 0.0111 Prepro~
## 4    96    11       15    0.0349  0.00196      0.743 roc_auc binary  0.785     4 0.00966 Prepro~
## 5   254    39        6    0.173  0.0000000361  0.923 roc_auc binary  0.784     4 0.0106 Prepro~
## # ... with abbreviated variable names 1: sample_size, 2: .estimator

best_xgboost <- xgboost_tune %>%
  select_best("roc_auc")

autoplot(xgboost_tune, metric = "roc_auc")

```



```

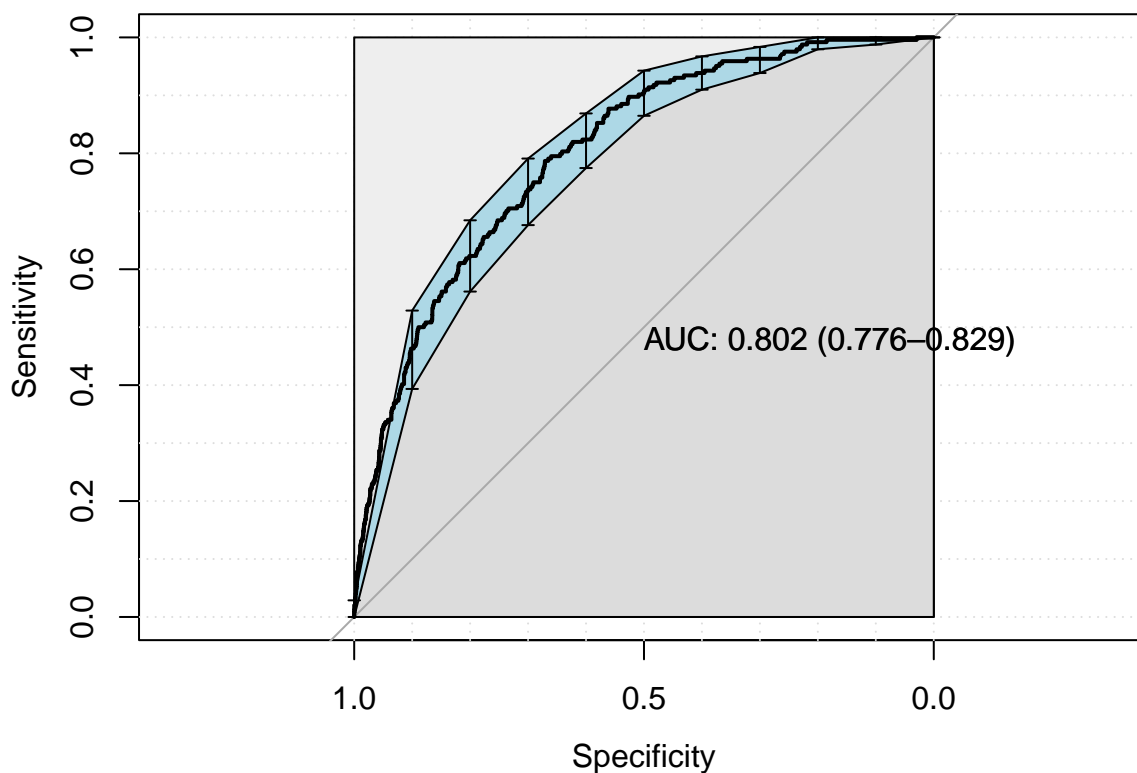
final_xgboost_workflow <-
  xgboost_workflow %>%
  finalize_workflow(best_xgboost)

last_xgboost_fit <-
  final_xgboost_workflow %>%
  last_fit(df_split)

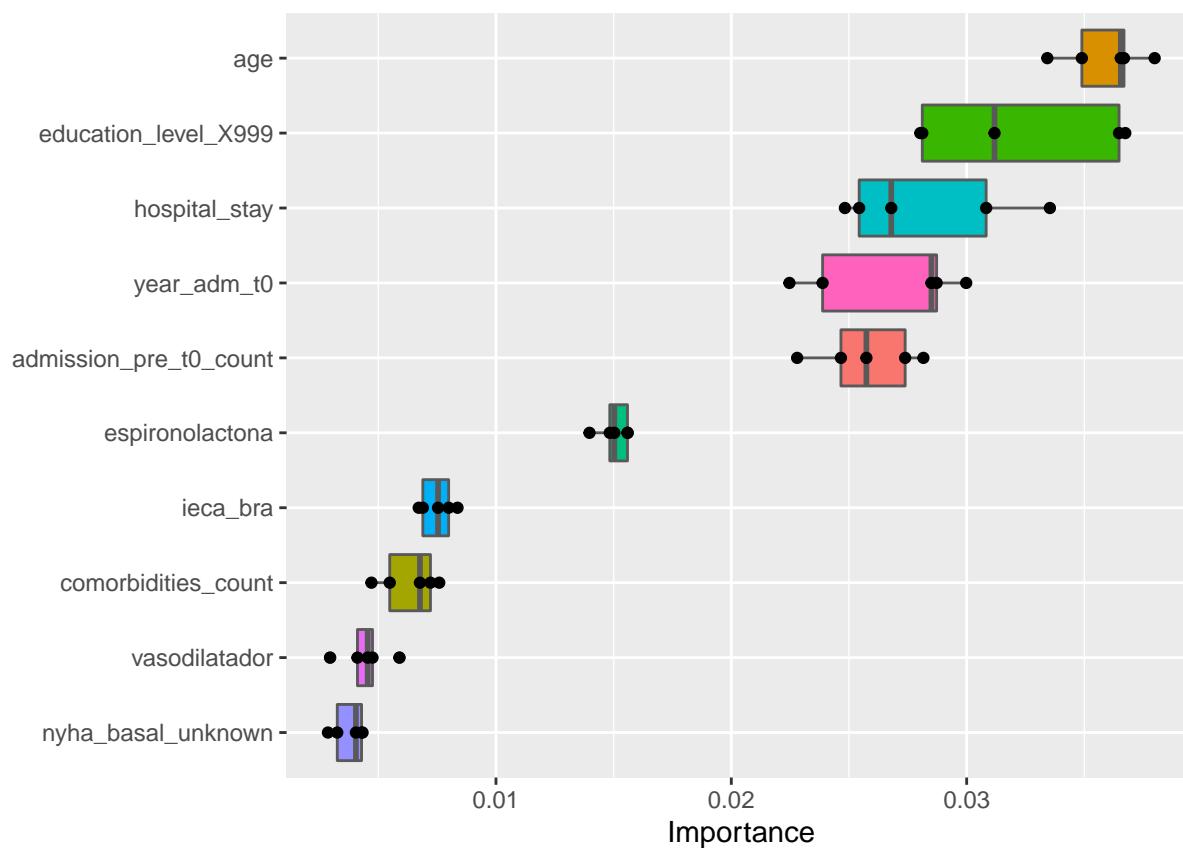
final_xgboost_fit <- extract_workflow(last_xgboost_fit)

xgboost_auc <- validation(final_xgboost_fit, df_test)

```



```
## [1] "Optimal Threshold: 0.04"
## Confusion Matrix and Statistics
##
##      reference
## data    0    1
##    0 3010   52
##    1 1476  192
##
##              Accuracy : 0.677
##              95% CI   : (0.6634, 0.6903)
##    No Information Rate : 0.9484
##    P-Value [Acc > NIR] : 1
##
##              Kappa   : 0.1218
##
## Mcnemar's Test P-Value : <2e-16
##
##              Sensitivity : 0.6710
##              Specificity : 0.7869
##              Pos Pred Value : 0.9830
##              Neg Pred Value : 0.1151
##              Prevalence : 0.9484
##              Detection Rate : 0.6364
##              Detection Prevalence : 0.6474
##              Balanced Accuracy : 0.7289
##
##              'Positive' Class : 0
##
extract_vip(final_xgboost_fit, pred_wrapper = predict,
            reference_class = "0")
```

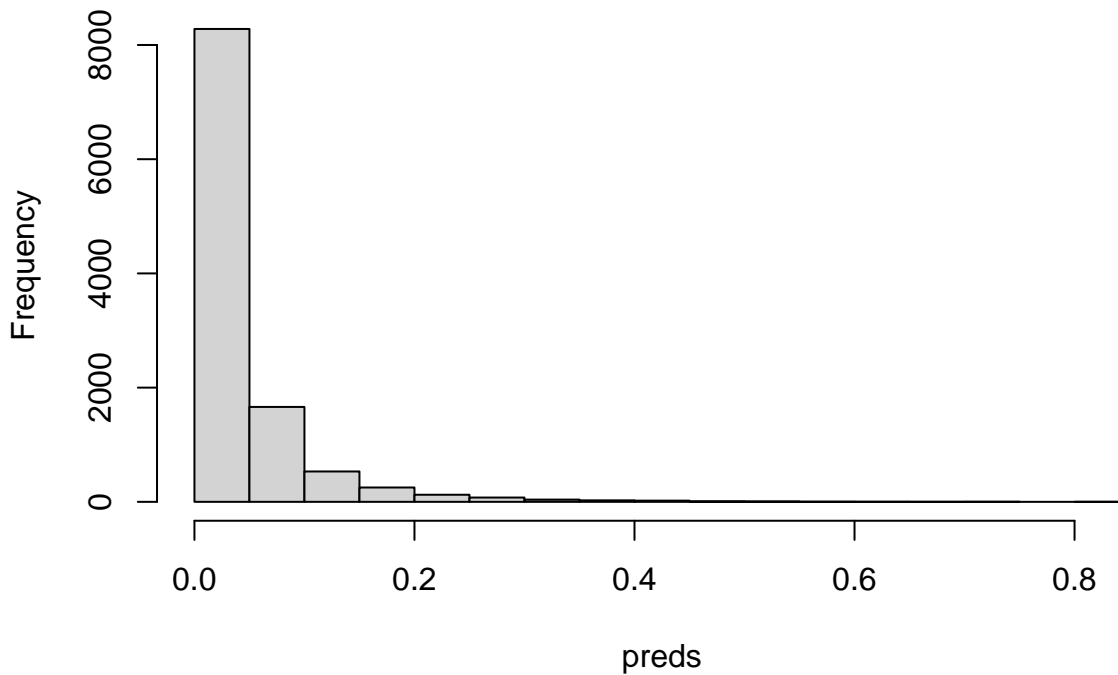


```
xgboost_parameters <- xgboost_tune %>%
  show_best("roc_auc", n = 1) %>%
  select(trees, min_n, tree_depth, learn_rate, loss_reduction) %>%
  as.list

saveRDS(
  xgboost_parameters,
  file = sprintf(
    "./auxiliar/model_selection/hyperparameters/xgboost_%s.rds",
    outcome_column
  )
)
```

Minutes to run: 2.15

## Histogram of preds



Minutes to run:

0.006

## Boosted Tree (LightGBM)

```
lightgbm_recipe <-  
  recipe(formula = sprintf("%s ~ .", outcome_column) %>% as.formula, data = df_train) %>%  
  step_novel(all_nominal_predictors()) %>%  
  step_unknown(all_nominal_predictors()) %>%  
  step_other(all_nominal_predictors(), threshold = 0.05, other = ".merged") %>%  
  step_dummy(all_nominal_predictors())  
  
lightgbm_spec <- boost_tree(  
  trees = tune(),  
  min_n = tune(),  
  tree_depth = tune(),  
  learn_rate = tune(),  
  loss_reduction = tune(),  
  sample_size = 1  
) %>%  
  set_engine("lightgbm",  
    nthread = 8) %>%  
  set_mode("classification")  
  
lightgbm_grid <- grid_latin_hypercube(  
  trees(range = c(50L, 300L)),  
  min_n(),  
  tree_depth(),  
  learn_rate(range = c(0.01, 0.3), trans = NULL),  
  loss_reduction(),  
  size = grid_size  
)  
  
lightgbm_workflow <-
```



```
workflow() %>%
  add_recipe(lightgbm_recipe) %>%
  add_model(lightgbm_spec)
```

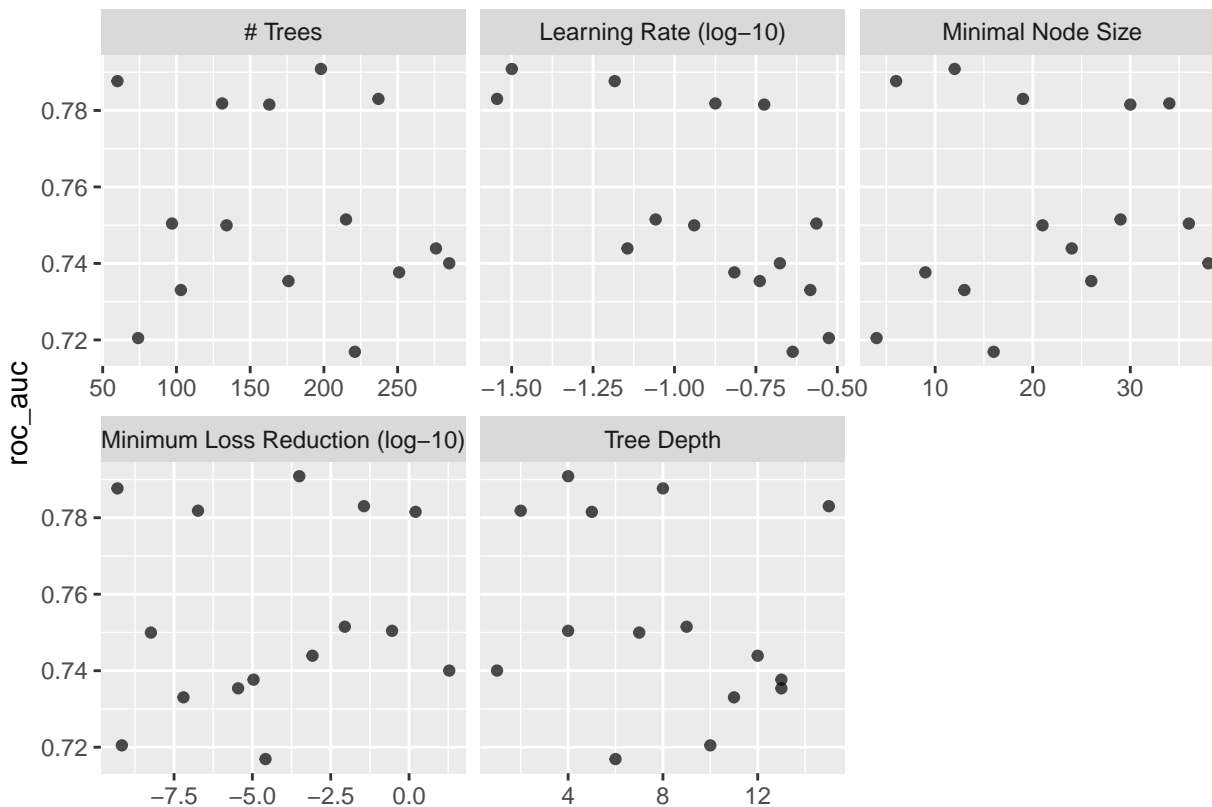
```
lightgbm_tune <-
  lightgbm_workflow %>%
  tune_grid(resamples = df_folds,
            grid = lightgbm_grid)
```

```
lightgbm_tune %>%
  show_best("roc_auc")
```

```
## # A tibble: 5 x 11
##   trees min_n tree_depth learn_rate loss_reduction .metric .estimator mean n std_err .config
##   <int> <int>   <int>   <dbl>         <dbl> <chr>   <chr>   <dbl> <int> <dbl> <chr>
## 1   198    12      4    0.0317      3.14e- 4 roc_auc binary   0.791     4  0.0119 Preprocessor~
## 2    60     6      8    0.0655      4.95e-10 roc_auc binary   0.788     4  0.0112 Preprocessor~
## 3   237    19     15    0.0285      3.64e- 2 roc_auc binary   0.783     4  0.0111 Preprocessor~
## 4   131    34      2    0.134       1.84e- 7 roc_auc binary   0.782     4  0.0143 Preprocessor~
## 5   163    30      5    0.189       1.61e+ 0 roc_auc binary   0.782     4  0.0134 Preprocessor~
```

```
best_lightgbm <- lightgbm_tune %>%
  select_best("roc_auc")
```

```
autoplot(lightgbm_tune, metric = "roc_auc")
```

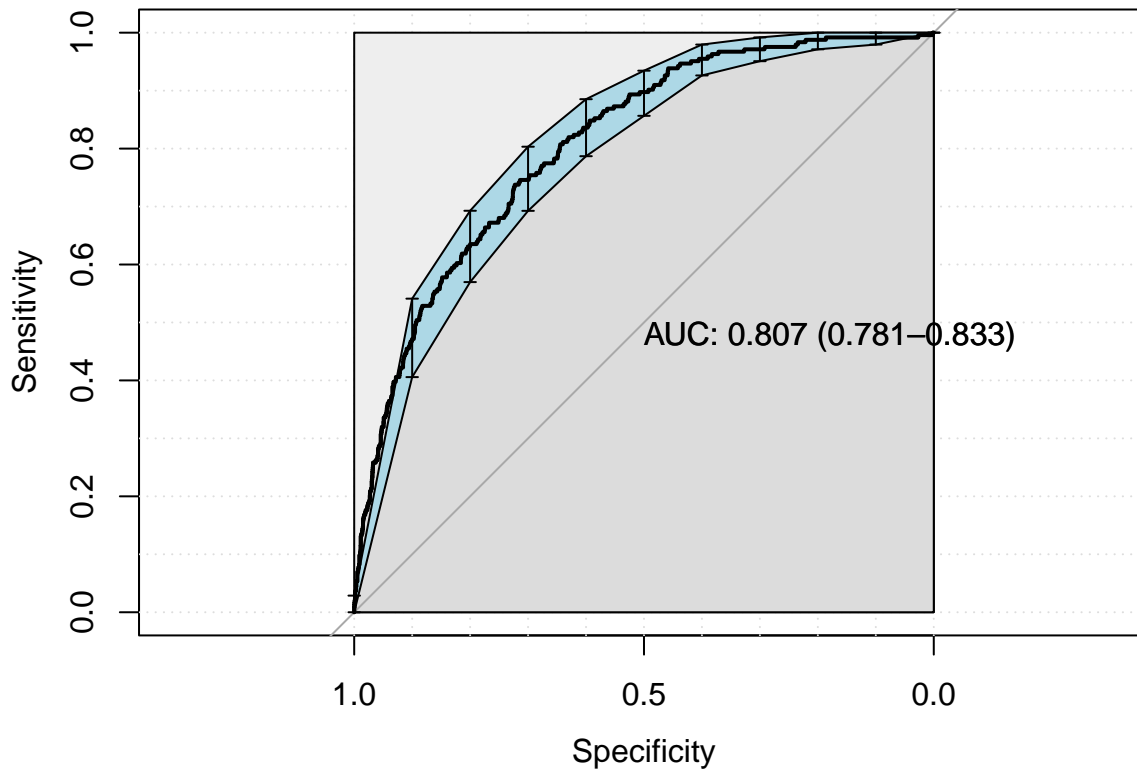


```
final_lightgbm_workflow <-
  lightgbm_workflow %>%
  finalize_workflow(best_lightgbm)
```

```
last_lightgbm_fit <-
  final_lightgbm_workflow %>%
  last_fit(df_split)
```

```
final_lightgbm_fit <- extract_workflow(last_lightgbm_fit)

lightgbm_auc <- validation(final_lightgbm_fit, df_test)
```



```
## [1] "Optimal Threshold: 0.04"
## Confusion Matrix and Statistics
##
##      reference
## data    0    1
## 0 3243   64
## 1 1243  180
##
##              Accuracy : 0.7237
##              95% CI   : (0.7107, 0.7364)
##    No Information Rate : 0.9484
##    P-Value [Acc > NIR] : 1
##
##              Kappa   : 0.1402
##
## Mcnemar's Test P-Value : <2e-16
##
##              Sensitivity : 0.7229
##              Specificity : 0.7377
##              Pos Pred Value : 0.9806
##              Neg Pred Value : 0.1265
##              Prevalence : 0.9484
##              Detection Rate : 0.6856
##              Detection Prevalence : 0.6992
##              Balanced Accuracy : 0.7303
##
##              'Positive' Class : 0
##
```

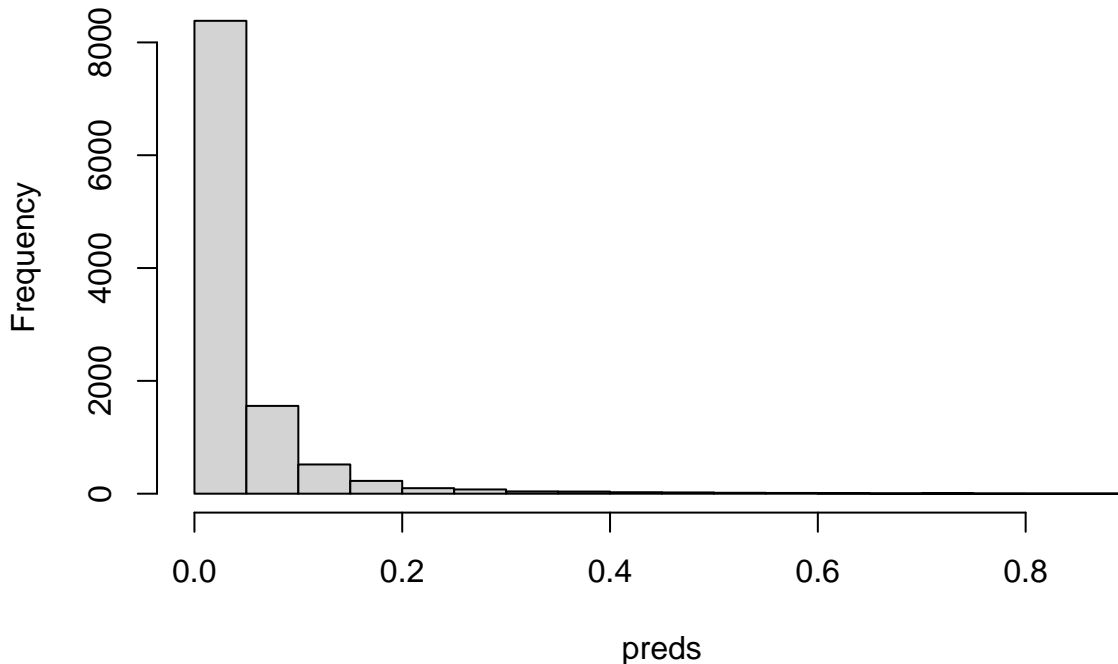
```
# pfun_lightgbm <- function(object, newdata) predict(object, data = newdata)
# extract_vip(final_lightgbm_fit, pred_wrapper = pfun_lightgbm,
#             reference_class = "1")

lightgbm_parameters <- lightgbm_tune %>%
  show_best("roc_auc", n = 1) %>%
  select(trees, min_n, tree_depth, learn_rate, loss_reduction) %>%
  as.list

saveRDS(
  lightgbm_parameters,
  file = sprintf(
    "./auxiliar/model_selection/hyperparameters/lightgbm_%s.rds",
    outcome_column
  )
)
```

Minutes to run: 1.968

## Histogram of preds



Minutes to run:

0.007

## GLM

```
glmnet_recipe <-
  recipe(formula = sprintf("%s ~ .", outcome_column) %>% as.formula, data = df_train) %>%
  step_novel(all_nominal_predictors()) %>%
  step_unknown(all_nominal_predictors()) %>%
  step_other(all_nominal_predictors(), threshold = 0.05, other = ".merged") %>%
  step_dummy(all_nominal_predictors()) %>%
  step_zv(all_predictors()) %>%
  step_normalize(all_numeric_predictors())

glmnet_spec <-
```

```

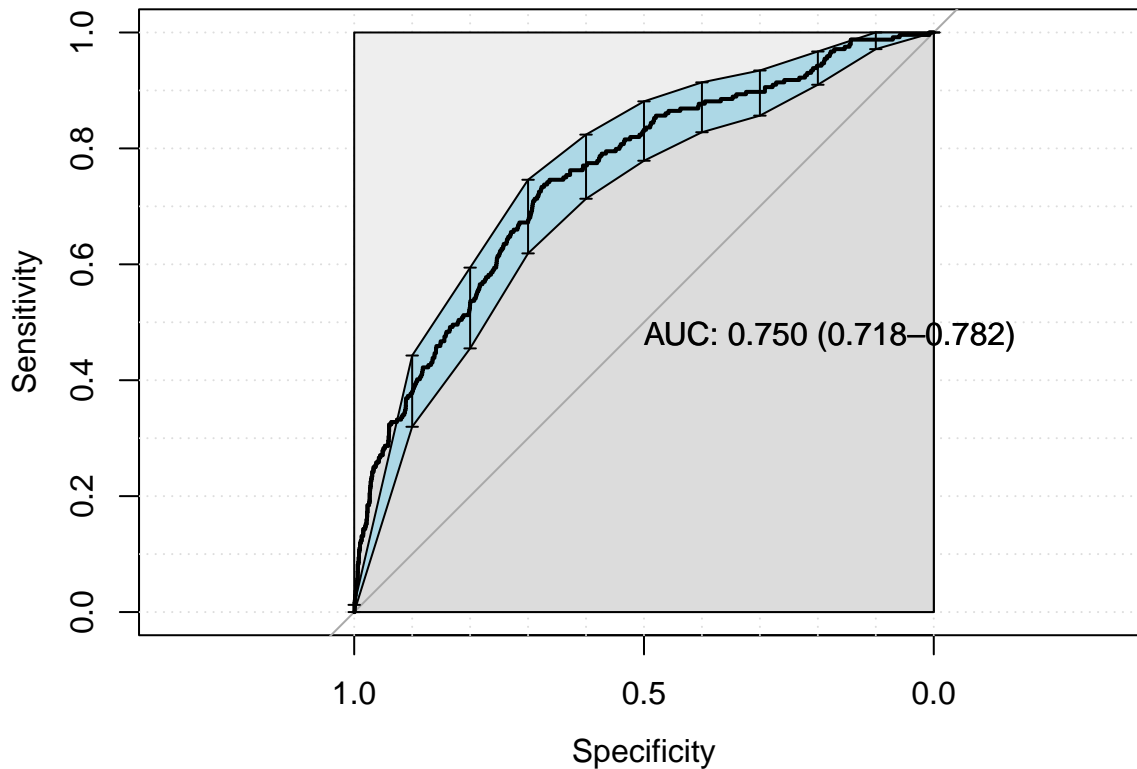
logistic_reg(penalty = 0) %>%
set_mode("classification") %>%
set_engine("glmnet")

glmnet_workflow <-
  workflow() %>%
  add_recipe(glmnet_recipe) %>%
  add_model(glmnet_spec)

glm_fit <- glmnet_workflow %>%
  fit(df_train)

glmnet_auc <- validation(glm_fit, df_test)

```



```

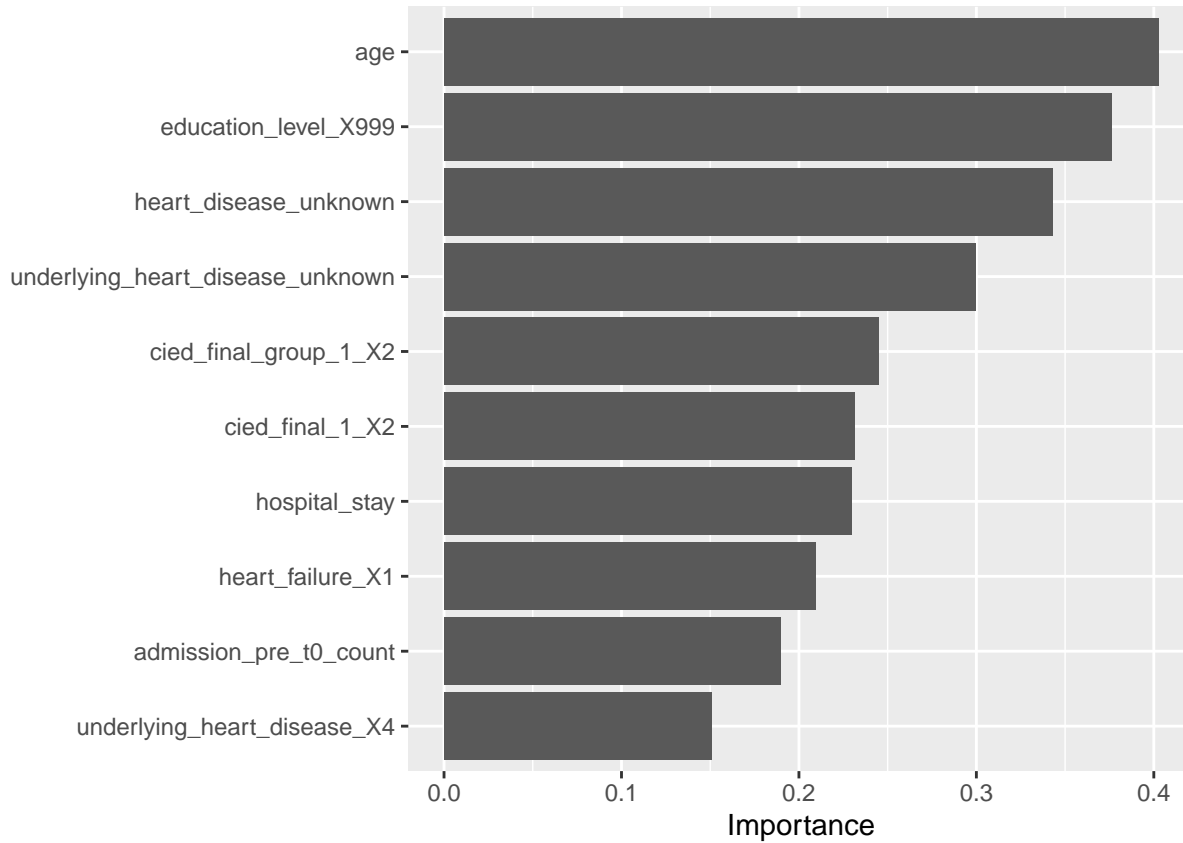
## [1] "Optimal Threshold: 0.04"
## Confusion Matrix and Statistics
##
##      reference
## data    0    1
## 0 3036   65
## 1 1450  179
##
##              Accuracy : 0.6797
##              95% CI   : (0.6662, 0.693)
##    No Information Rate : 0.9484
##    P-Value [Acc > NIR] : 1
##
##              Kappa   : 0.1114
##
##  Mcnemar's Test P-Value : <2e-16
##
##              Sensitivity : 0.6768
##              Specificity : 0.7336

```

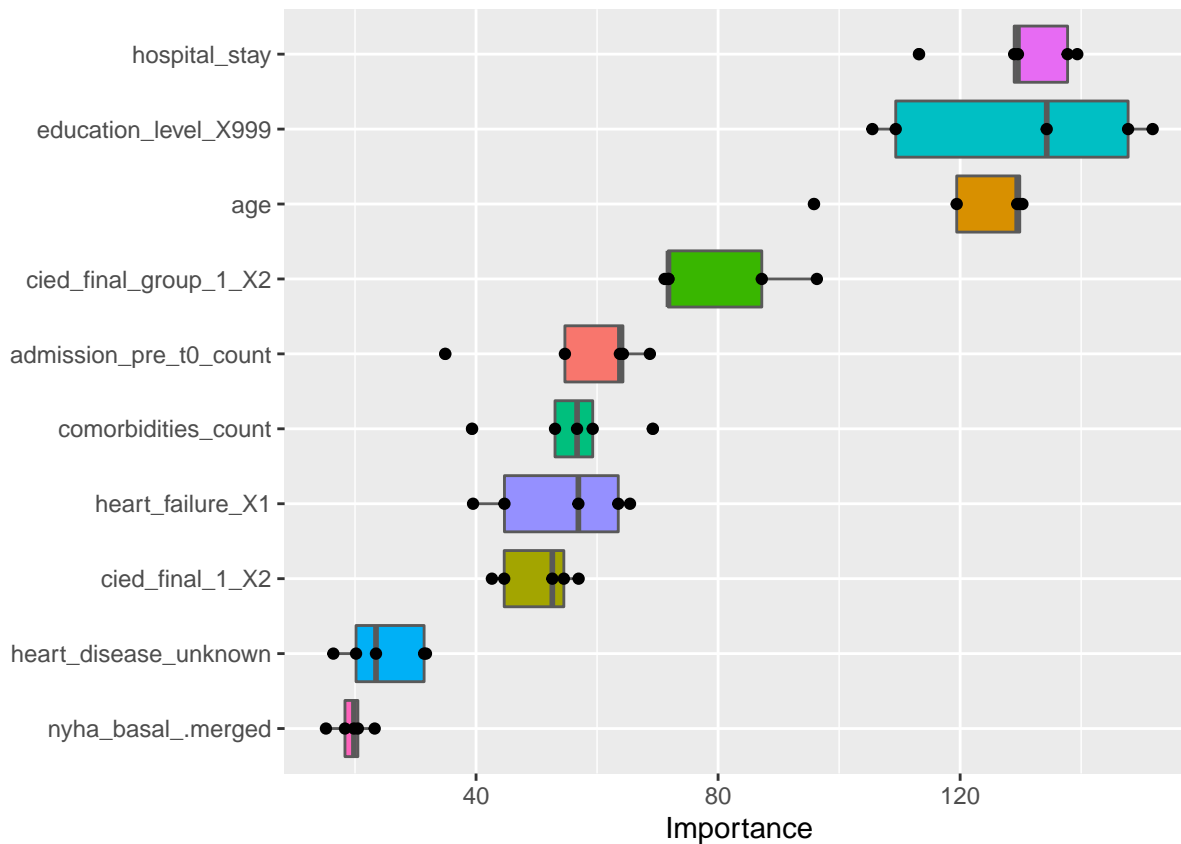
```
##          Pos Pred Value : 0.9790
##          Neg Pred Value : 0.1099
##          Prevalence : 0.9484
##          Detection Rate : 0.6419
##          Detection Prevalence : 0.6556
##          Balanced Accuracy : 0.7052
##
##          'Positive' Class : 0
##
```

```
pfun_glmnet <- function(object, newdata) predict(object, newx = newdata)
```

```
extract_vip(glm_fit, pred_wrapper = pfun_glmnet,
             reference_class = "1", method = 'model')
```



```
extract_vip(glm_fit, pred_wrapper = pfun_glmnet,
             reference_class = "1", method = 'permute')
```



Minutes to run:

1.982

## Decision Tree

```
tree_recipe <-
  recipe(formula = sprintf("%s ~ .", outcome_column) %>% as.formula, data = df_train) %>%
  step_novel(all_nominal_predictors()) %>%
  step_unknown(all_nominal_predictors()) %>%
  step_other(all_nominal_predictors(), threshold = 0.05, other = ".merged") %>%
  step_dummy(all_nominal_predictors()) %>%
  step_zv(all_predictors())

tree_spec <-
  decision_tree(cost_complexity = tune(),
                tree_depth = tune(),
                min_n = tune()) %>%
  set_mode("classification") %>%
  set_engine("rpart")

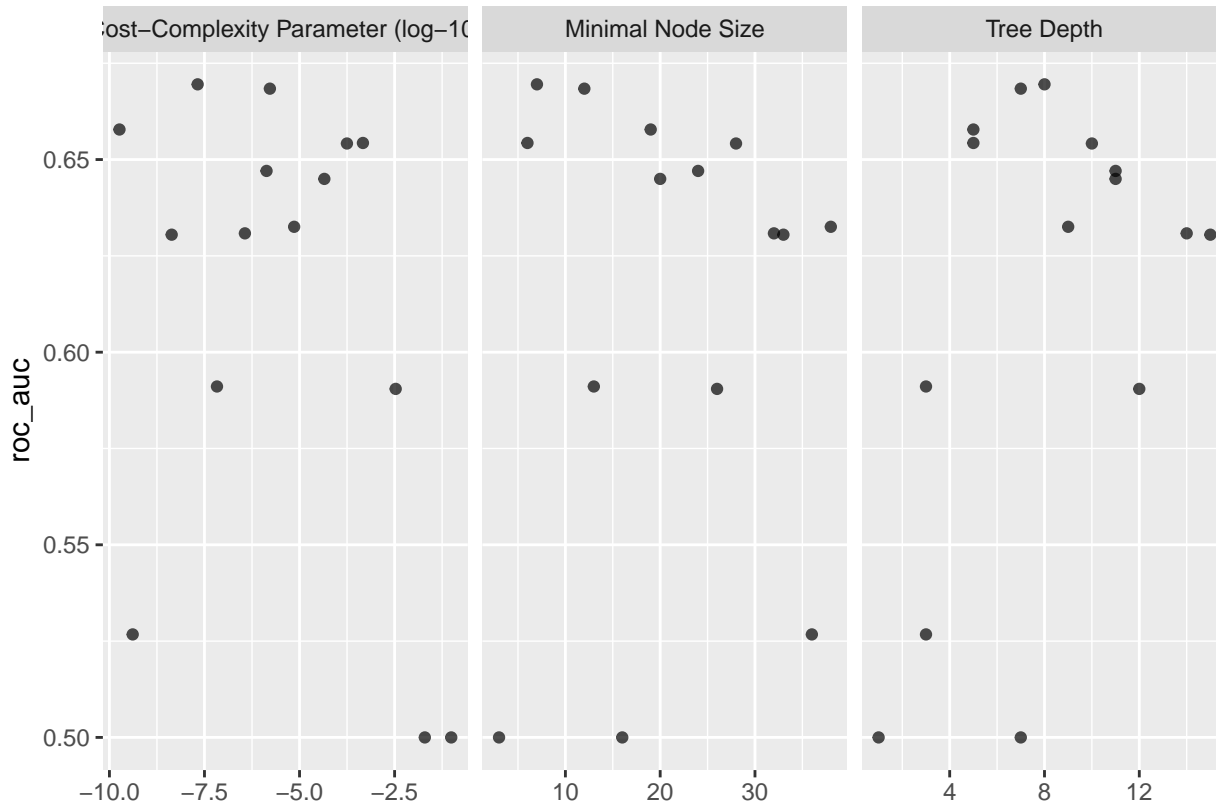
tree_grid <- grid_latin_hypercube(cost_complexity(),
                                  tree_depth(),
                                  min_n(),
                                  size = grid_size)

tree_workflow <-
  workflow() %>%
  add_recipe(tree_recipe) %>%
  add_model(tree_spec)

tree_tune <-
  tree_workflow %>%
  tune_grid(resamples = df_folds,
            grid = tree_grid)
```

```
tree_tune %>%
  collect_metrics()

autoplot(tree_tune, metric = "roc_auc")
```



```
tree_tune %>%
  show_best("roc_auc")

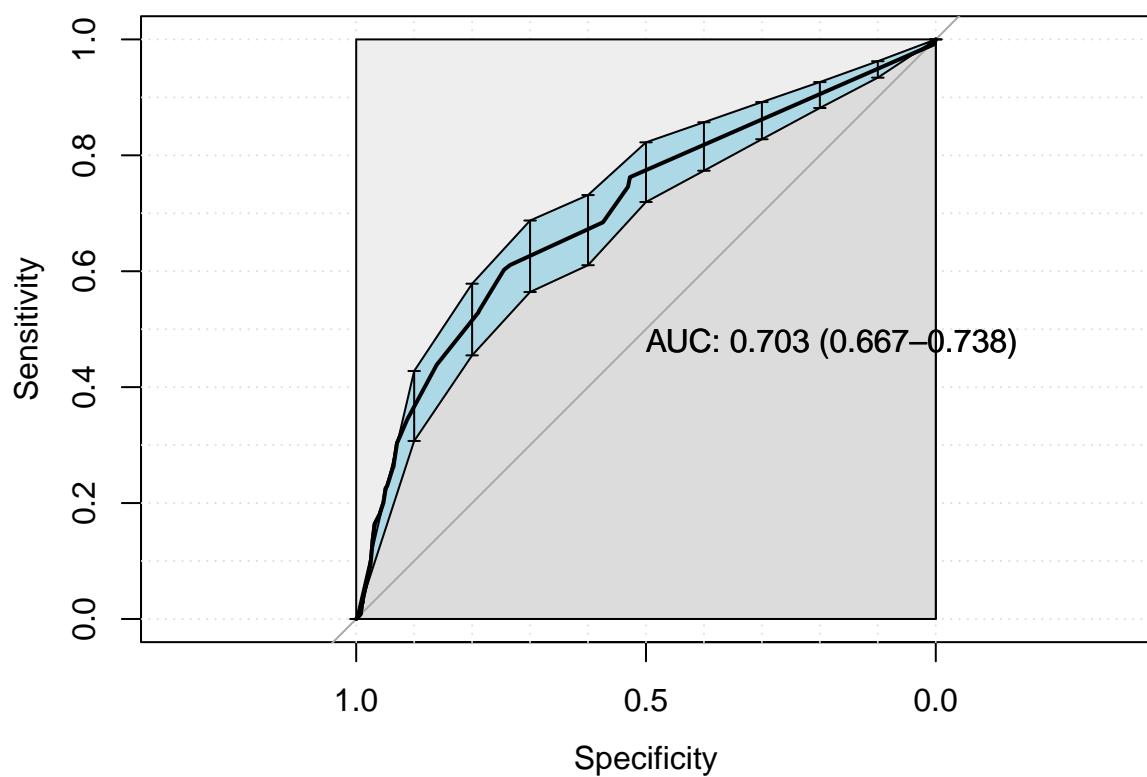
best_tree <- tree_tune %>%
  select_best("roc_auc")

final_tree_workflow <-
  tree_workflow %>%
  finalize_workflow(best_tree)

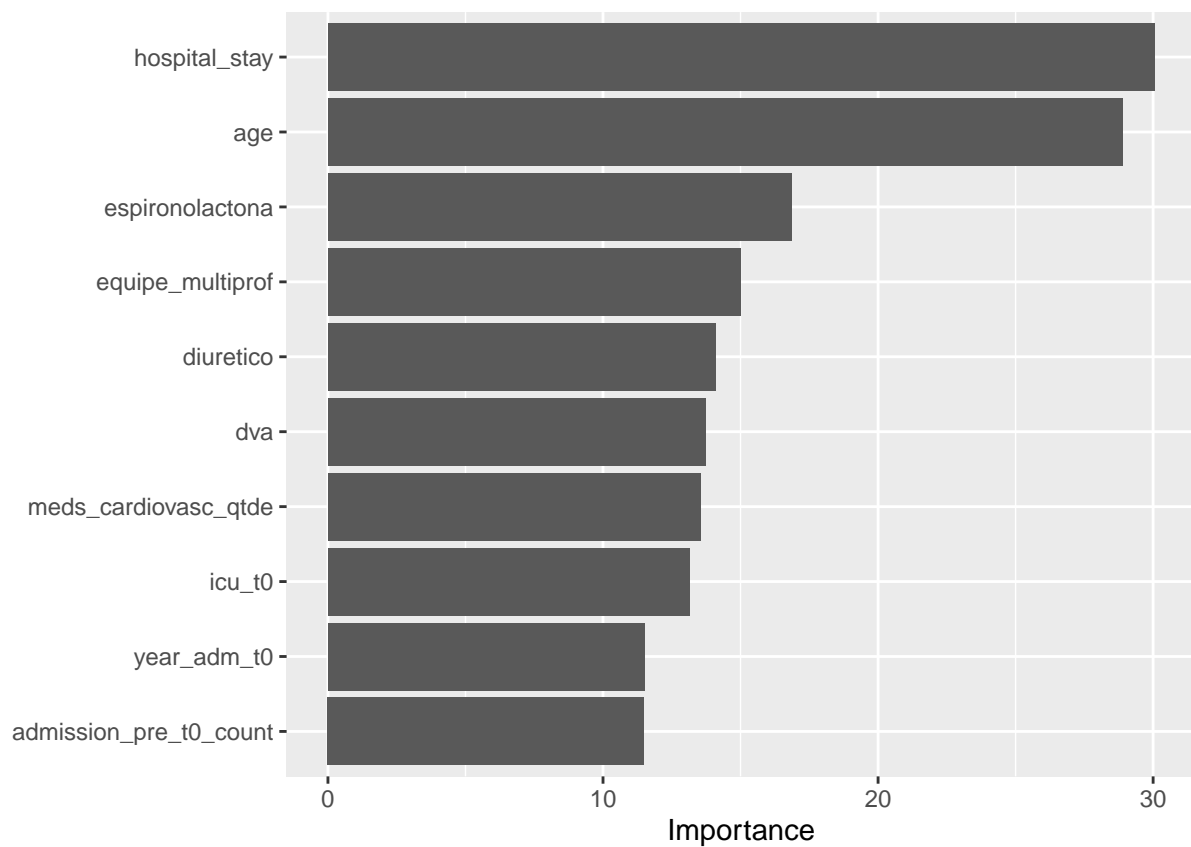
last_tree_fit <-
  final_tree_workflow %>%
  last_fit(df_split)

final_tree_fit <- extract_workflow(last_tree_fit)

tree_auc <- validation(final_tree_fit, df_test)
```



```
extract_vip(final_tree_fit, pred_wrapper = predict,
            reference_class = "0", use_matrix = FALSE,
            method = 'model')
```



```
# extract_vip(final_tree_fit, pred_wrapper = predict,
#             reference_class = "1", use_matrix = FALSE,
```

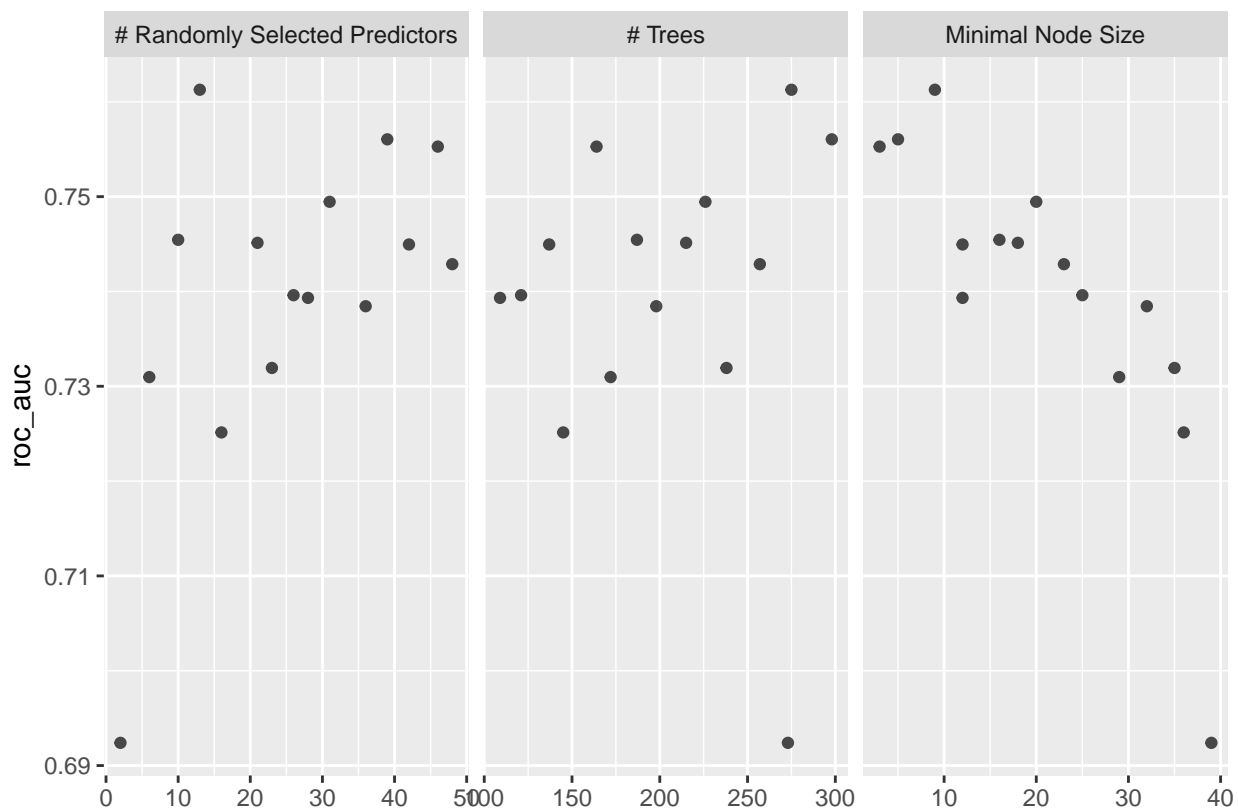


```
# method = 'permute')
```

Minutes to run: 1.381

## Random Forest

```
rf_recipe <-  
  recipe(formula = sprintf("%s ~ .", outcome_column) %>% as.formula,  
    data = df_train) %>%  
  step_nominal(all_nominal_predictors()) %>%  
  step_unknown(all_nominal_predictors()) %>%  
  step_other(all_nominal_predictors(), threshold = 0.05, other = ".merged") %>%  
  step_dummy(all_nominal_predictors()) %>%  
  step_zv(all_predictors()) %>%  
  step_impute_mean(all_numeric_predictors())  
  
rf_spec <-  
  rand_forest(mtry = tune(),  
    trees = tune(),  
    min_n = tune()) %>%  
  set_mode("classification") %>%  
  set_engine("randomForest",  
    probability = TRUE,  
    nthread = 8)  
  
rf_grid <- grid_latin_hypercube(mtry(range = c(1L, 50L)),  
  trees(range = c(100L, 300L)),  
  min_n(),  
  size = grid_size)  
  
rf_workflow <-  
  workflow() %>%  
  add_recipe(rf_recipe) %>%  
  add_model(rf_spec)  
  
rf_tune <-  
  rf_workflow %>%  
  tune_grid(resamples = df_folds,  
    grid = rf_grid)  
  
rf_tune %>%  
  collect_metrics()  
  
autoplot(rf_tune, metric = "roc_auc")
```



```
rf_tune %>%
  show_best("roc_auc")

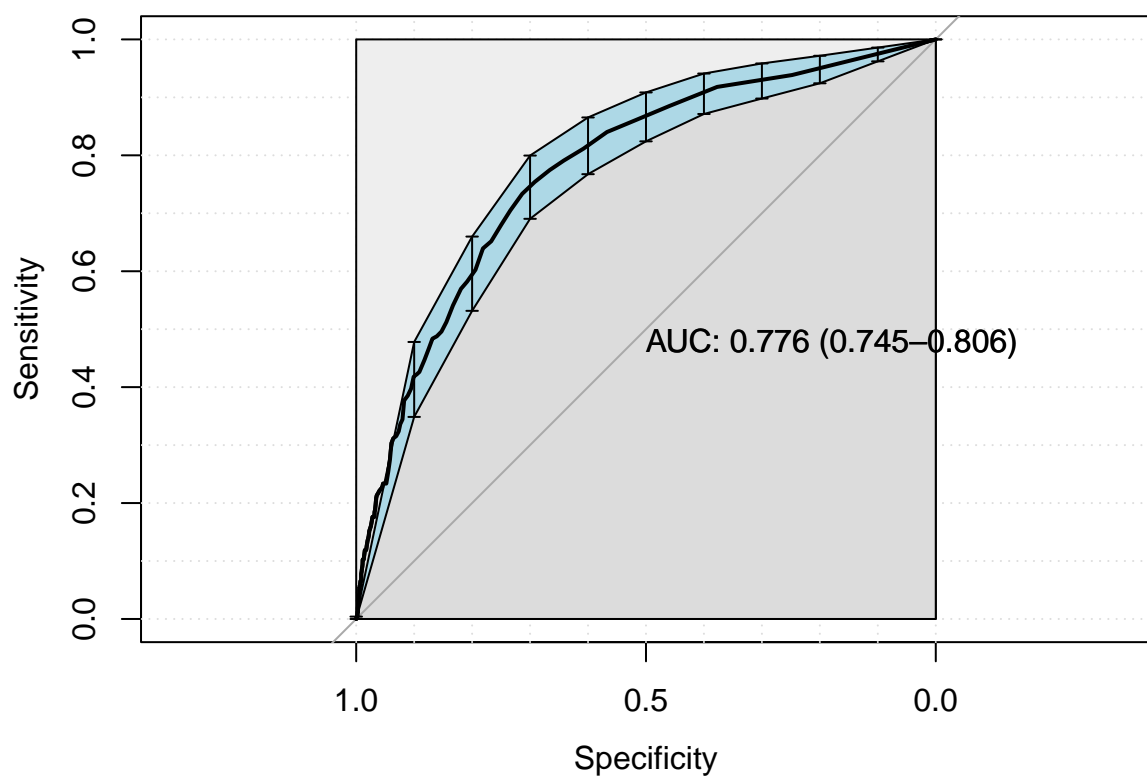
best_rf <- rf_tune %>%
  select_best("roc_auc")

final_rf_workflow <-
  rf_workflow %>%
  finalize_workflow(best_rf)

last_rf_fit <-
  final_rf_workflow %>%
  last_fit(df_split)

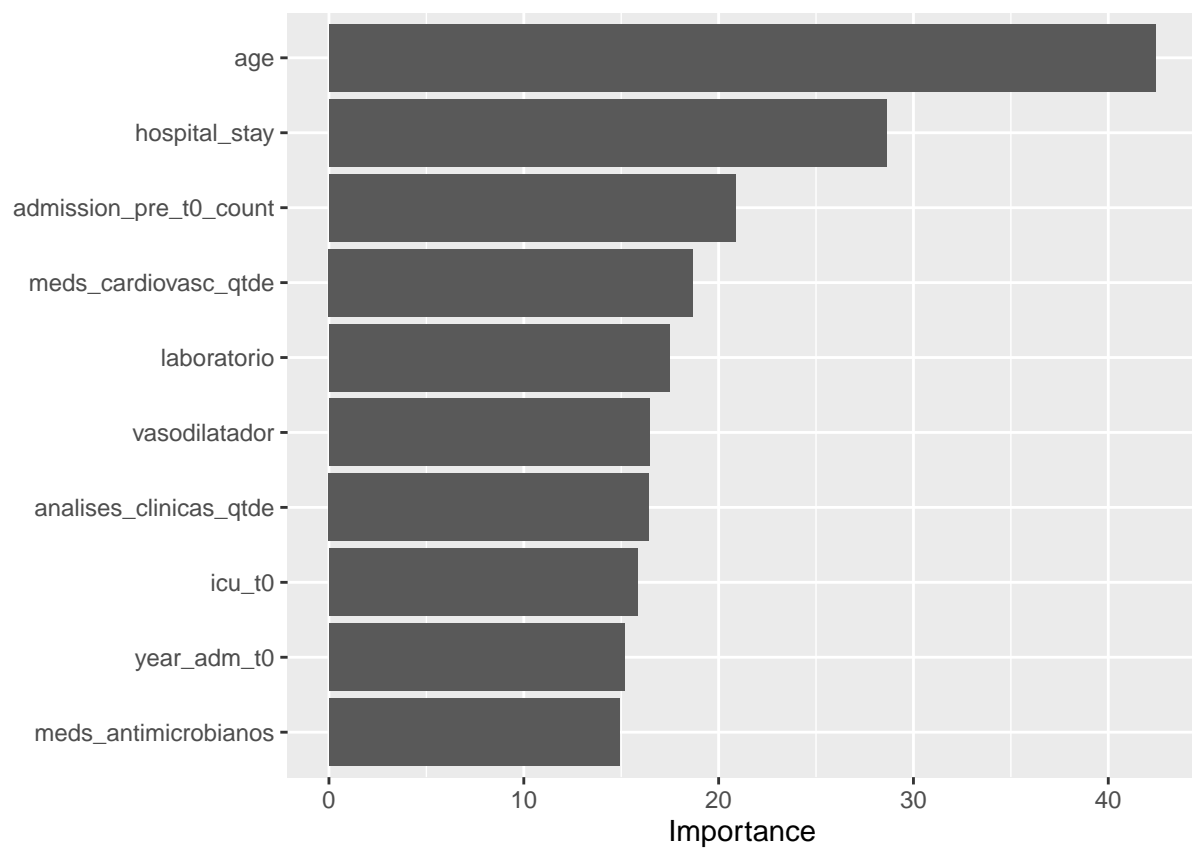
final_rf_fit <- extract_workflow(last_rf_fit)

rf_auc <- validation(final_rf_fit, df_test)
```



```
pfun_rf <- function(object, newdata) predict(object, data = newdata)

extract_vip(final_rf_fit, pred_wrapper = predict,
            reference_class = "1", use_matrix = FALSE,
            method = 'model')
```



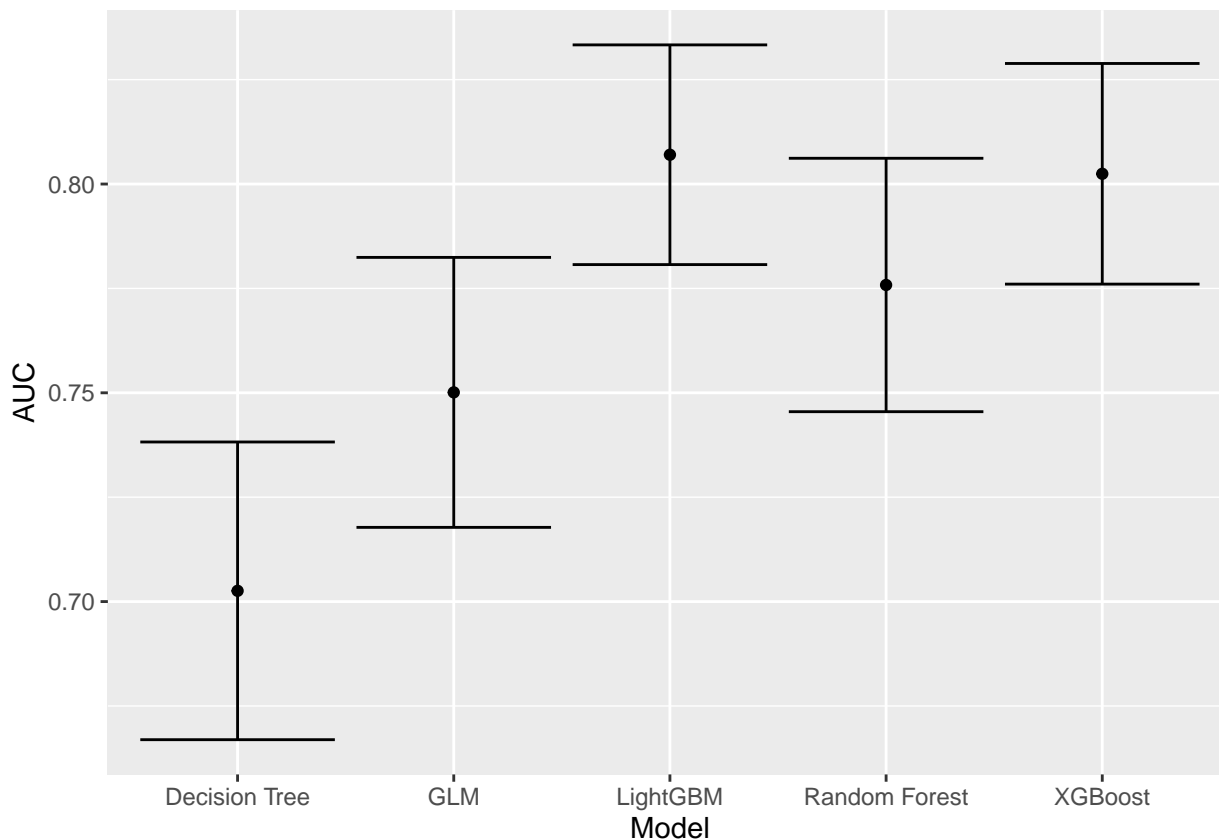
```
# extract_vip(final_rf_fit, pred_wrapper = predict,
#             reference_class = "1", use_matrix = FALSE,
#             method = 'permute')
```

Minutes to run: 16.267

## Models Comparison

```
df_auc <- tibble::tribble(
  ~Model, ~`AUC`, ~`Lower Limit`, ~`Upper Limit`,
  'XGBoost', as.numeric(xgboost_auc$auc), xgboost_auc$ci[1], xgboost_auc$ci[3],
  'LightGBM', as.numeric(lightgbm_auc$auc), lightgbm_auc$ci[1], lightgbm_auc$ci[3],
  'GLM', as.numeric(glmnet_auc$auc), glmnet_auc$ci[1], glmnet_auc$ci[3],
  'Decision Tree', as.numeric(tree_auc$auc), tree_auc$ci[1], tree_auc$ci[3],
  'Random Forest', as.numeric(rf_auc$auc), rf_auc$ci[1], rf_auc$ci[3]
) %>%
  mutate(Target = outcome_column)

df_auc %>%
  ggplot(aes(x = Model, y = AUC, ymin = `Lower Limit`, ymax = `Upper Limit`)) +
    geom_point() +
    geom_errorbar()
```



```
saveRDS(df_auc, sprintf("./auxiliar/model_selection/performance/%s.RData", outcome_column))
```

Minutes to run: 0.003