

Model Selection - death_1year

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$repeats
RUN_ALL_MODELS <- params$RUN_ALL_MODELS
Hmisc::list.tree(params)

## params = list 5 (952 bytes)
## . outcome_column = character 1= death_1year
## . k = double 1= 10
## . grid_size = double 1= 30
## . repeats = double 1= 2
## . RUN_ALL_MODELS = logical 1= TRUE
```

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: 0.008

```
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$categorical_columns, columns_list$numerical_columns)) %>%
  setdiff(c(exception_columns, correlated_columns))

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. cancer
19. comorbidities_count
20. procedure_type_1
21. reop_type_1
22. procedure_type_new
23. cied_final_1
24. cied_final_group_1
25. admission_pre_t0_count
26. admission_pre_t0_180d
27. year_adm_t0
28. icu_t0
29. dialysis_t0
30. admission_t0_emergency
31. aco
32. antiarritmico
33. ieca_bra
34. dva
35. digoxina
36. estatina
37. diuretico
38. vasodilatador
39. insuf_cardiaca
40. espirolactona
41. antiplaquetario_ev
42. insulina
43. psicofarmacos
44. antifungico
45. antiviral
46. classe_meds_qtde
47. meds_cardiovasc_qtde
48. meds_antimicrobianos
49. vni
50. ventilacao_mecanica
51. transplante_cardiaco
52. cir_toracica
53. outros_proced_cirurgicos
54. icp
55. cateterismo
56. cateter_venoso_central
57. proced_invasivos_qtde
58. transfusao
59. interconsulta
60. equipe_multiprof
61. holter
62. teste_esforco
63. tilt_teste
64. metodos_graficos_qtde
65. laboratorio
66. cultura

```
## 67. analises_clinicas_qtde
## 68. citologia
## 69. histopatologia_qtde
## 70. angio_tc
## 71. angiografia
## 72. aortografia
## 73. cintilografia
## 74. ecocardiograma
## 75. endoscopia
## 76. flebografia
## 77. pet_ct
## 78. ultrassom
## 79. tomografia
## 80. ressonancia
## 81. exames_imagem_qtde
## 82. bic
## 83. 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 <- rsample::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, 200L)),
  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_size .metric .estimator mean      n std_err .config
##   <int> <int>    <int>      <dbl>         <dbl>      <dbl> <chr>   <chr>    <dbl> <int>  <dbl> <chr>
## 1   154    14        11    0.0568         6.26e- 5    0.865 roc_auc binary  0.820    10 0.0108 Preproc
## 2   101     9         7    0.0928         9.25e-10    0.699 roc_auc binary  0.816    10 0.00946 Preproc
## 3   116    17         2    0.152          6.97e- 6    0.478 roc_auc binary  0.812    10 0.0130 Preproc
## 4   134    33         8    0.0684         6.27e+ 0    0.895 roc_auc binary  0.812    10 0.0148 Preproc
## 5    55    38         9    0.245          8.68e- 7    0.996 roc_auc binary  0.810    10 0.0129 Preproc

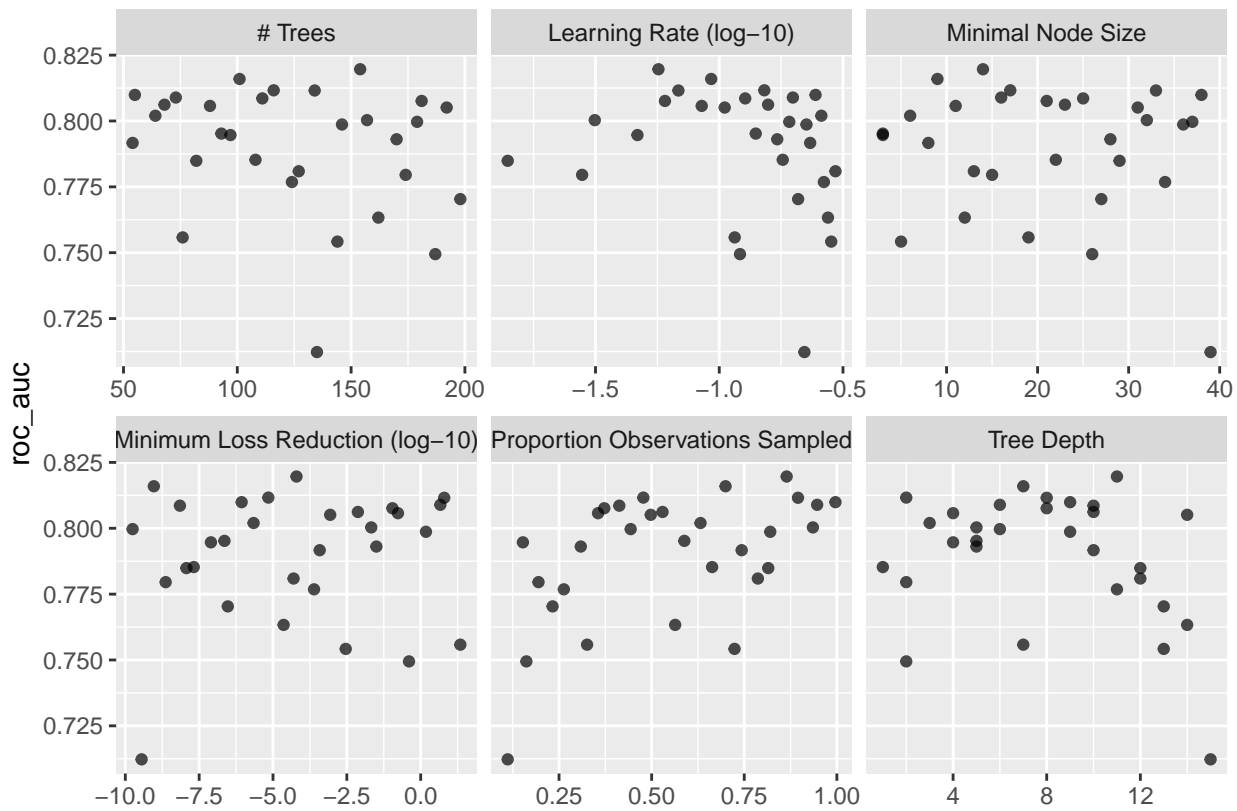
```

```

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

autoplot(xgboost_tune, metric = "roc_auc")

```

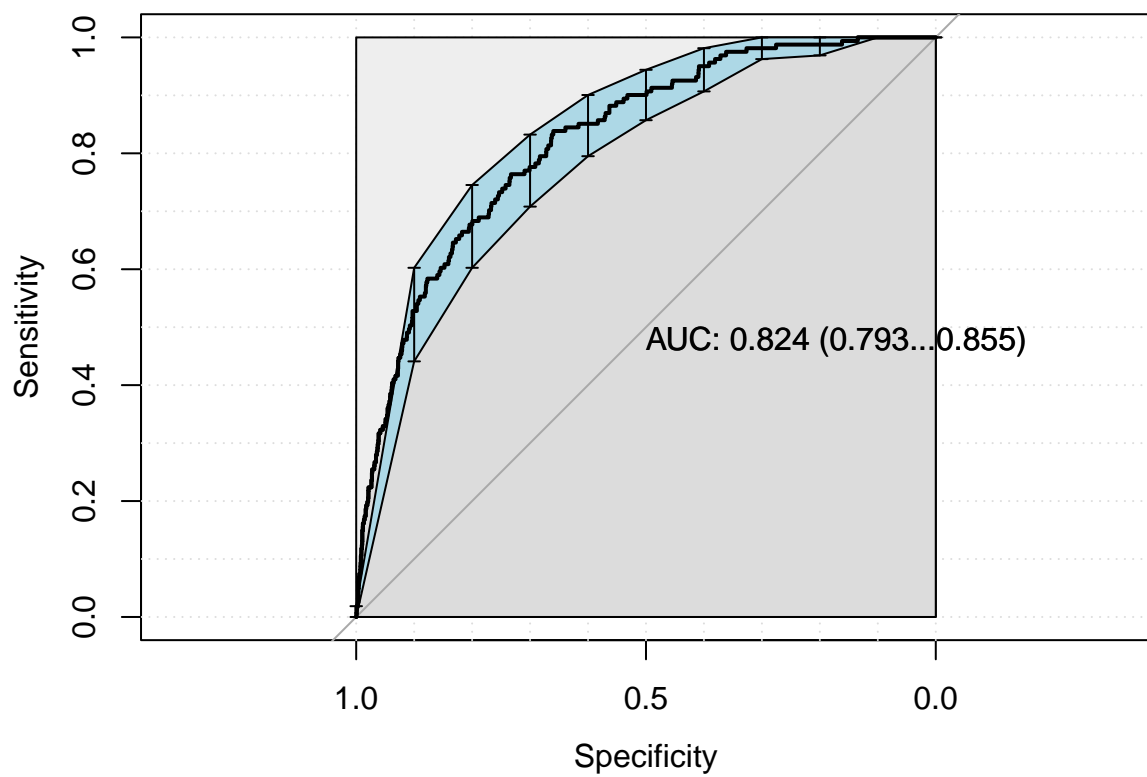


```
final_xgboost_workflow <-
  xgboost_workflow %>%
  finalize_workflow(best_xgboost)

last_xgboost_fit <-
  final_xgboost_workflow %>%
  last_fit(df_split)

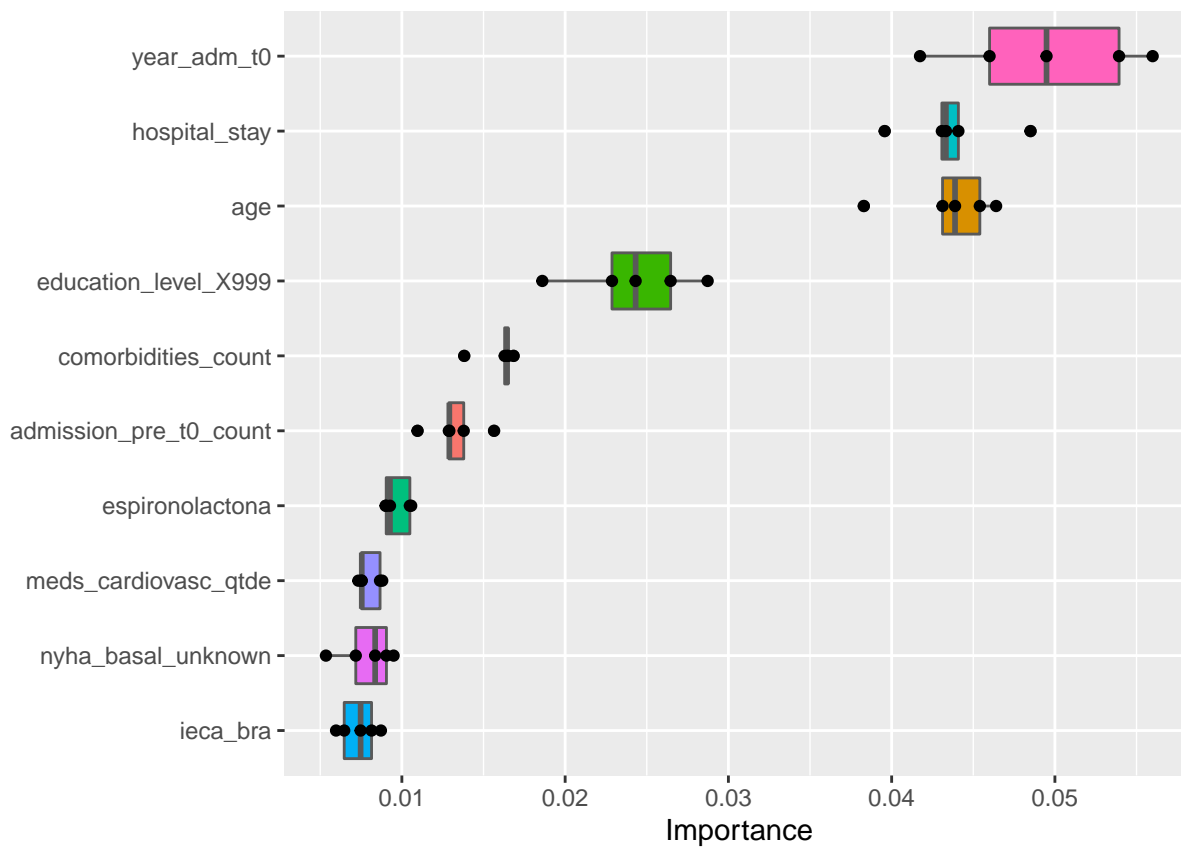
final_xgboost_fit <- extract_workflow(last_xgboost_fit)

xgboost_auc <- validation(final_xgboost_fit, df_test)
```



|

```
extract_vip(final_xgboost_fit, pred_wrapper = predict,  
            reference_class = "0")
```



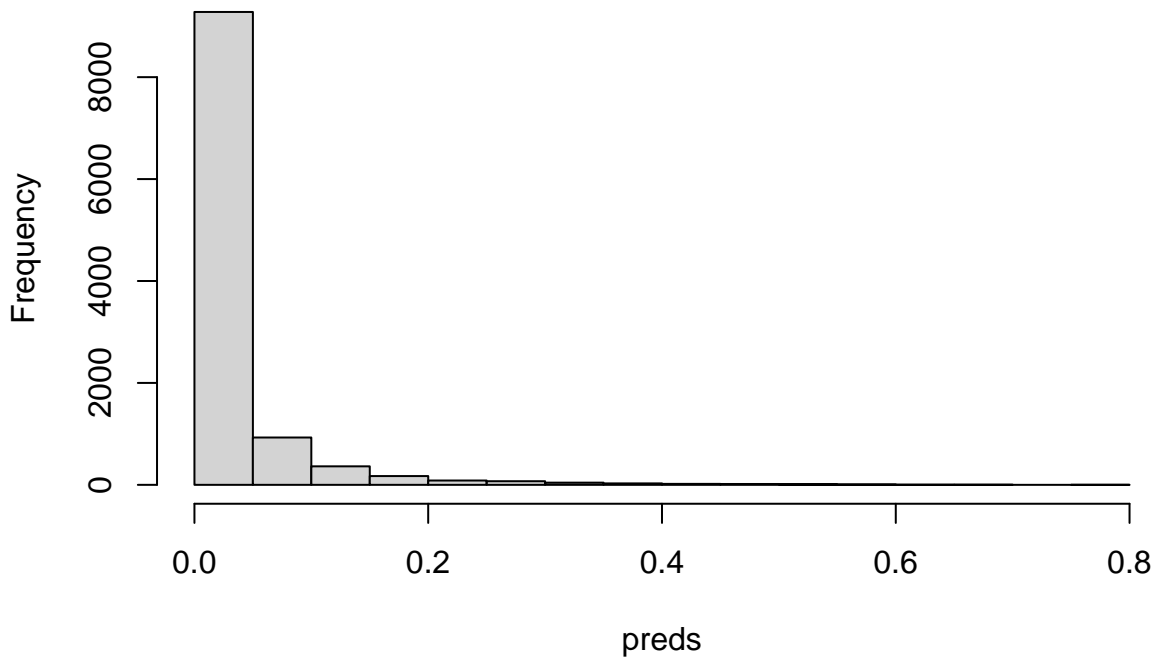
```
xgboost_parameters <- xgboost_tune %>%
  show_best("roc_auc", n = 1) %>%
  select(trees, min_n, tree_depth, learn_rate, loss_reduction) %>%
  as.list

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

preds <- predict(final_xgboost_fit, new_data = df_train, type = "prob") %>%
  rename_at(vars(starts_with(".pred_")), ~ str_remove(., ".pred_")) %>%
  .$`1`

hist(preds)
```


Histogram of preds



Minutes to run:

11.219

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(),  
  sample_size = 1  
) %>%  
  set_engine("lightgbm",  
    nthread = 8) %>%  
  set_mode("classification")  
  
lightgbm_grid <- grid_latin_hypercube(  
  trees(range = c(25L, 150L)),  
  min_n(range = c(2L, 100L)),  
  tree_depth(range = c(5L, 15L)),  
  learn_rate(range = c(-3, -1), trans = log10_trans()),  
  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 10
##   trees min_n tree_depth learn_rate .metric .estimator mean      n std_err .config
##   <int> <int>    <int>    <dbl> <chr>   <chr>    <dbl> <int>  <dbl> <chr>
## 1    56    37      13    0.0495 roc_auc binary  0.815    10  0.0125 Preprocessor1_Model11
## 2   123     6      14    0.0328 roc_auc binary  0.813    10  0.0107 Preprocessor1_Model102
## 3    99    81       5    0.0597 roc_auc binary  0.811    10  0.0126 Preprocessor1_Model125
## 4   141    87      12    0.0201 roc_auc binary  0.809    10  0.0119 Preprocessor1_Model126
## 5   135    65       8    0.0349 roc_auc binary  0.807    10  0.0107 Preprocessor1_Model120

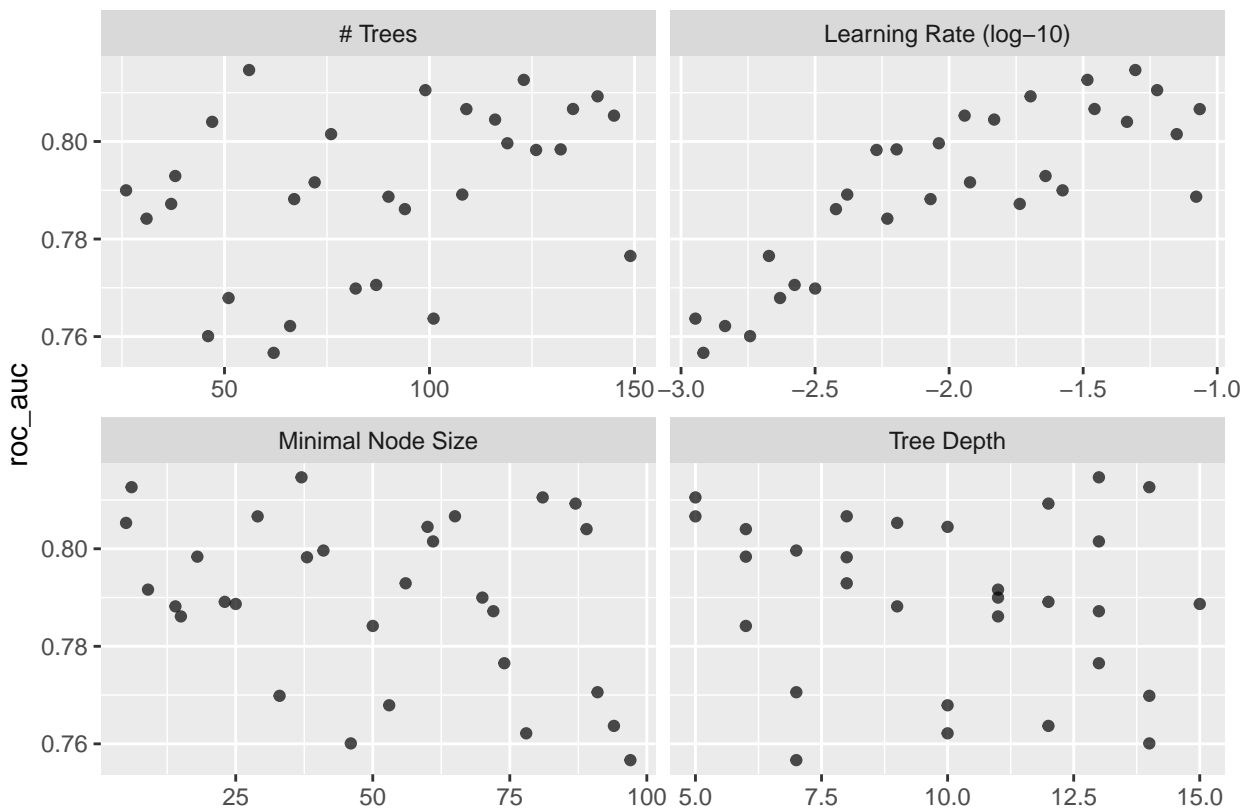
```

```

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

autoplot(lightgbm_tune, metric = "roc_auc")

```



```

final_lightgbm_workflow <-
  lightgbm_workflow %>%
  finalize_workflow(best_lightgbm)

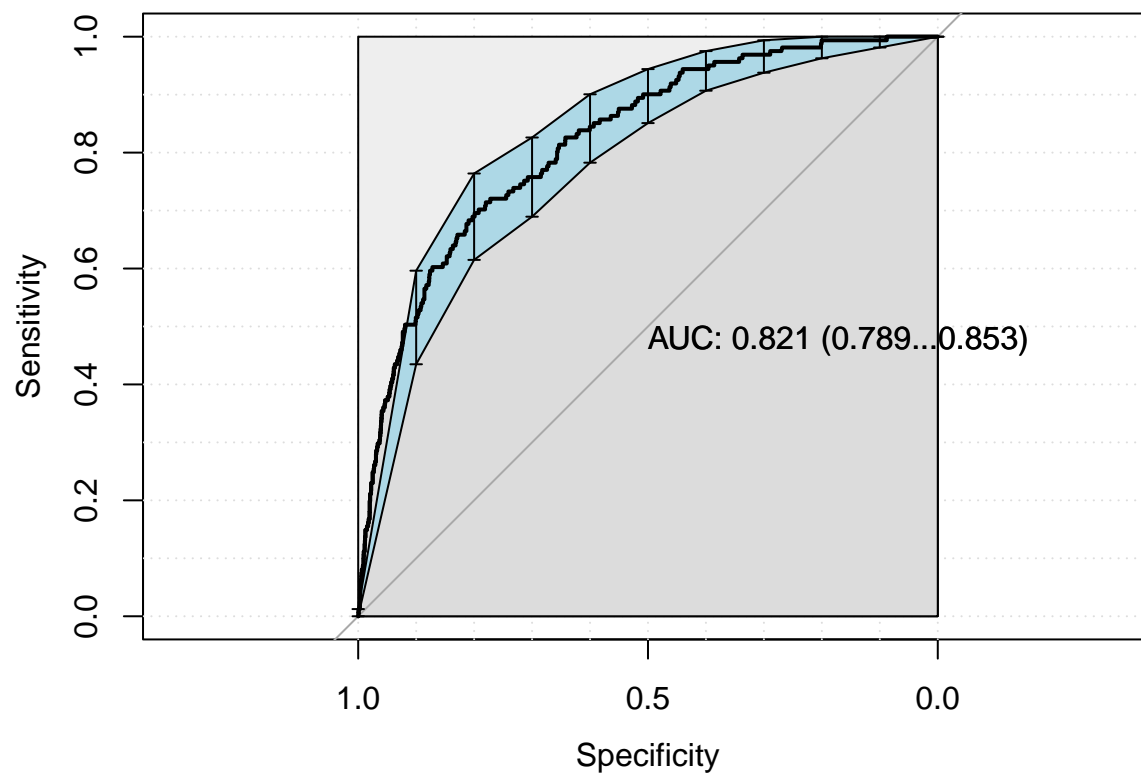
last_lightgbm_fit <-
  final_lightgbm_workflow %>%

```

```
last_fit(df_split)

final_lightgbm_fit <- extract_workflow(last_lightgbm_fit)

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



|

```
##
##           'Positive' Class : 0
##

lightgbm_parameters <- lightgbm_tune %>%
  show_best("roc_auc", n = 1) %>%
  select(-.metric, -.estimator, -.config, -mean, -n, -std_err) %>%
  as.list

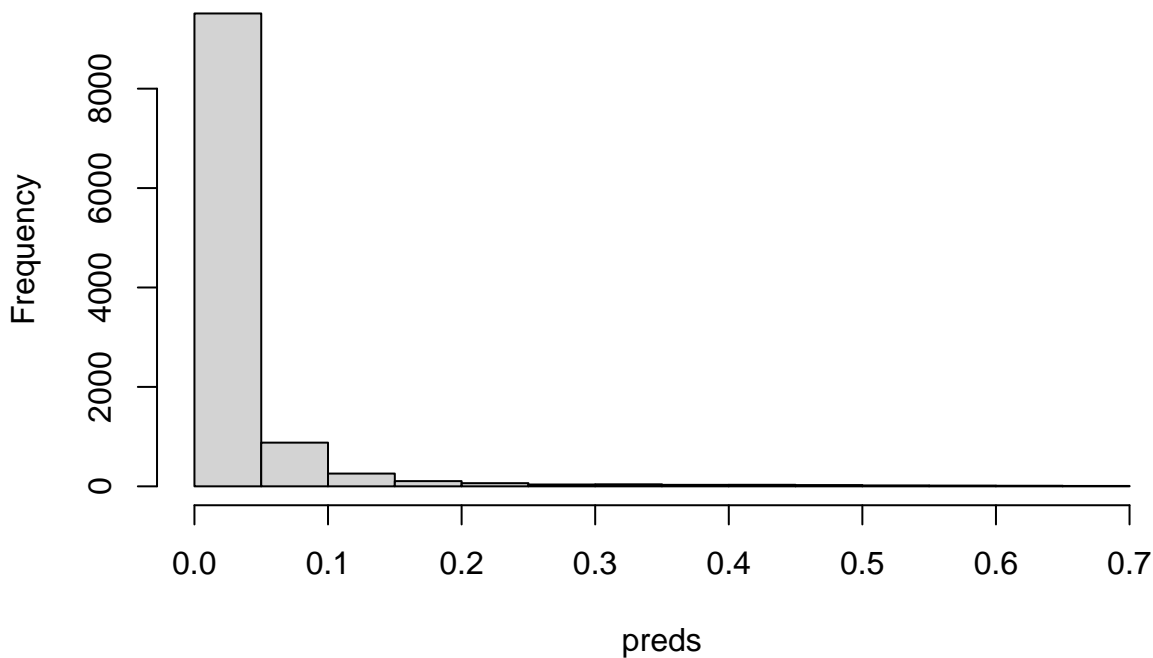
Hmisc::list.tree(lightgbm_parameters)
```

```
## lightgbm_parameters = list 4 (736 bytes)
## . trees = integer 1= 56
## . min_n = integer 1= 37
## . tree_depth = integer 1= 13
## . learn_rate = double 1= 0.049504
```

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

Minutes to run: 3.585

Histogram of preds



0.005

Minutes to run:

GLM

```

glmnet_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_normalize(all_numeric_predictors())

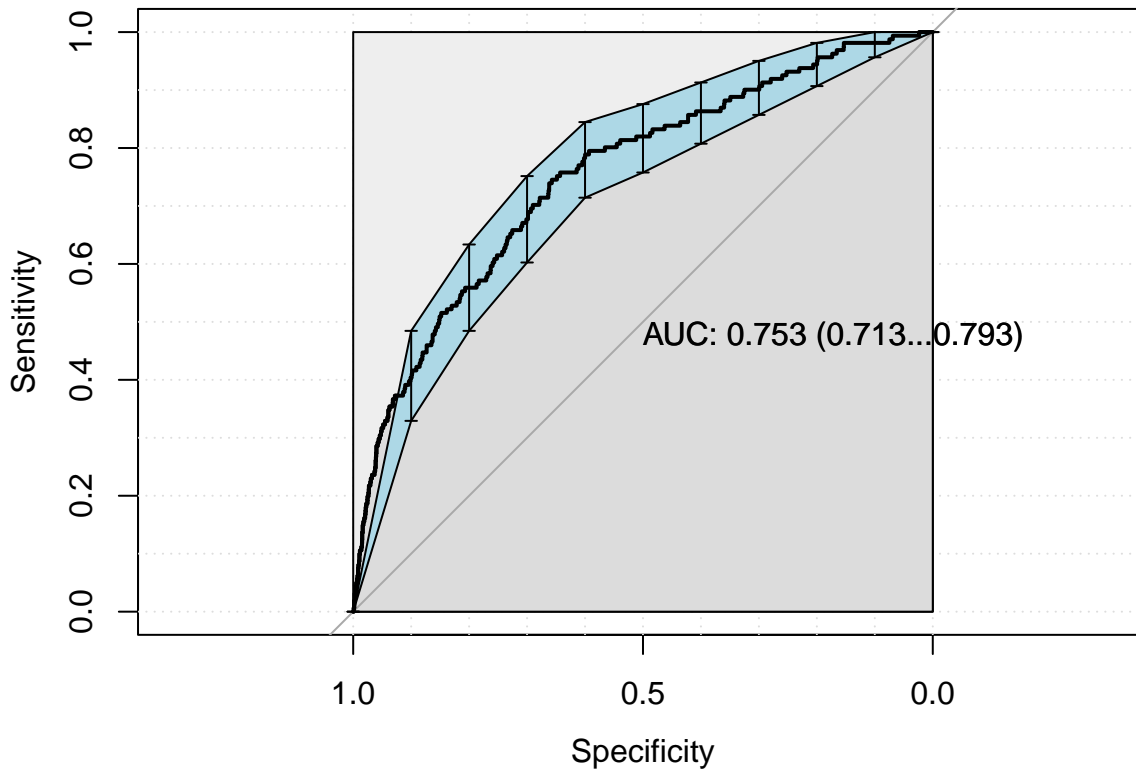
glmnet_spec <-
  logistic_reg(penalty = 0) %>%
  set_mode("classification") %>%
  set_engine("glmnet")

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)

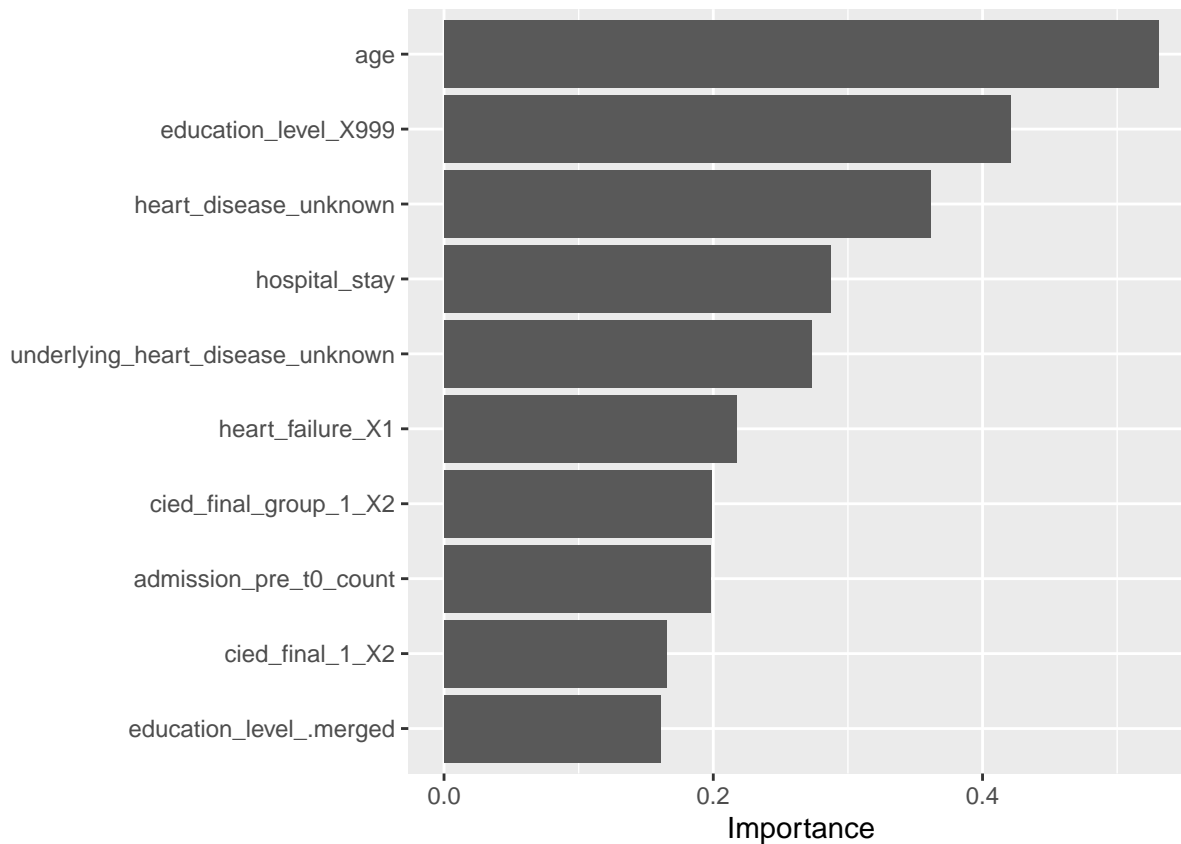
```



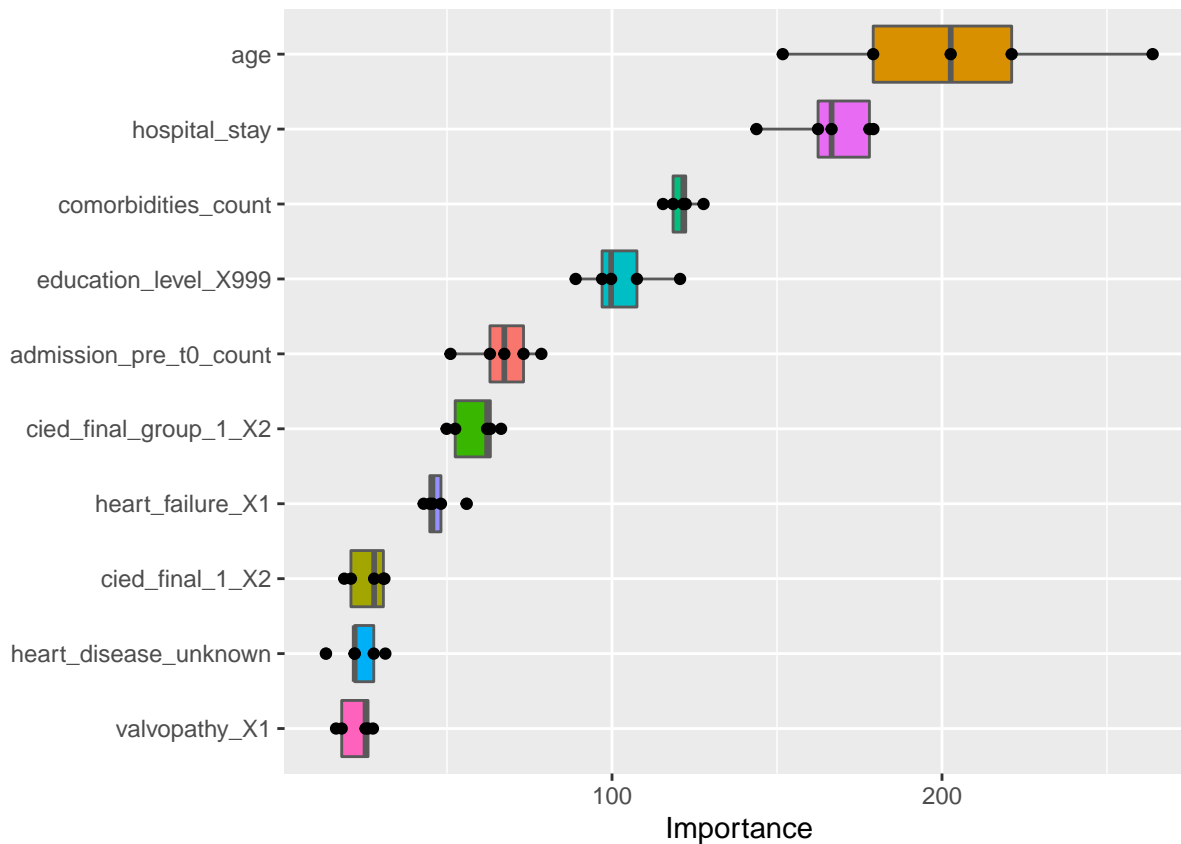
```
##
##          Accuracy : 0.6603
##          95% CI   : (0.6466, 0.6738)
## No Information Rate : 0.966
## P-Value [Acc > NIR] : 1
##
##          Kappa   : 0.0723
##
## Mcnemar's Test P-Value : <2e-16
##
##          Sensitivity : 0.65726
##          Specificity : 0.74534
##          Pos Pred Value : 0.98653
##          Neg Pred Value : 0.07117
##          Prevalence : 0.96596
##          Detection Rate : 0.63488
##          Detection Prevalence : 0.64355
##          Balanced Accuracy : 0.70130
##
##          '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:

3.189

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)

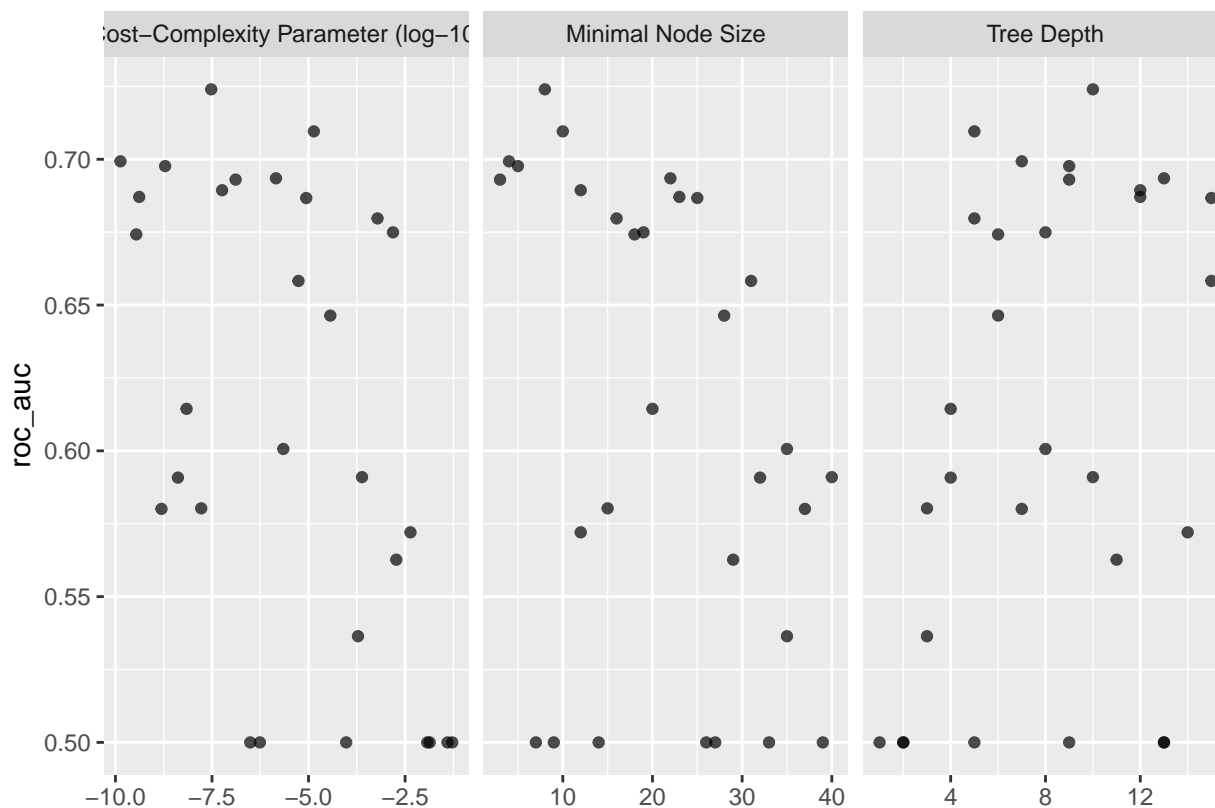
tree_tune <-
  tree_workflow %>%
```

```
tune_grid(resamples = df_folds,
          grid = tree_grid)
```

```
tree_tune %>%
  collect_metrics()
```

```
## # A tibble: 60 x 9
##   cost_complexity tree_depth min_n .metric .estimator mean      n std_err .config
##   <dbl>          <int> <int> <chr>   <chr>      <dbl> <int>   <dbl> <chr>
## 1      3.47e-10           6    18 accuracy binary    0.965    10 0.00179 Preprocessor1_Model101
## 2      3.47e-10           6    18 roc_auc  binary    0.674    10 0.0176  Preprocessor1_Model101
## 3      5.49e- 6          15    31 accuracy binary    0.966    10 0.00152 Preprocessor1_Model102
## 4      5.49e- 6          15    31 roc_auc  binary    0.658    10 0.0227  Preprocessor1_Model102
## 5      1.56e- 9           7    37 accuracy binary    0.966    10 0.00161 Preprocessor1_Model103
## 6      1.56e- 9           7    37 roc_auc  binary    0.580    10 0.0281  Preprocessor1_Model103
## 7      1.43e- 6          13    22 accuracy binary    0.962    10 0.00146 Preprocessor1_Model104
## 8      1.43e- 6          13    22 roc_auc  binary    0.693    10 0.0109  Preprocessor1_Model104
## 9      4.15e- 9           4    32 accuracy binary    0.967    10 0.00157 Preprocessor1_Model105
## 10     4.15e- 9           4    32 roc_auc  binary    0.591    10 0.0226  Preprocessor1_Model105
## # ... with 50 more rows
```

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



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

```
## # A tibble: 5 x 9
##   cost_complexity tree_depth min_n .metric .estimator mean      n std_err .config
##   <dbl>          <int> <int> <chr>   <chr>      <dbl> <int>   <dbl> <chr>
## 1      3.02e- 8          10     8 roc_auc binary    0.724    10 0.0141  Preprocessor1_Model107
## 2      1.37e- 5           5    10 roc_auc binary    0.710    10 0.0136  Preprocessor1_Model110
```


## 3	1.35e-10	7	4	roc_auc	binary	0.699	10	0.0158	Preprocessor1_Model130
## 4	1.93e- 9	9	5	roc_auc	binary	0.698	10	0.0157	Preprocessor1_Model117
## 5	1.43e- 6	13	22	roc_auc	binary	0.693	10	0.0109	Preprocessor1_Model104

```

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

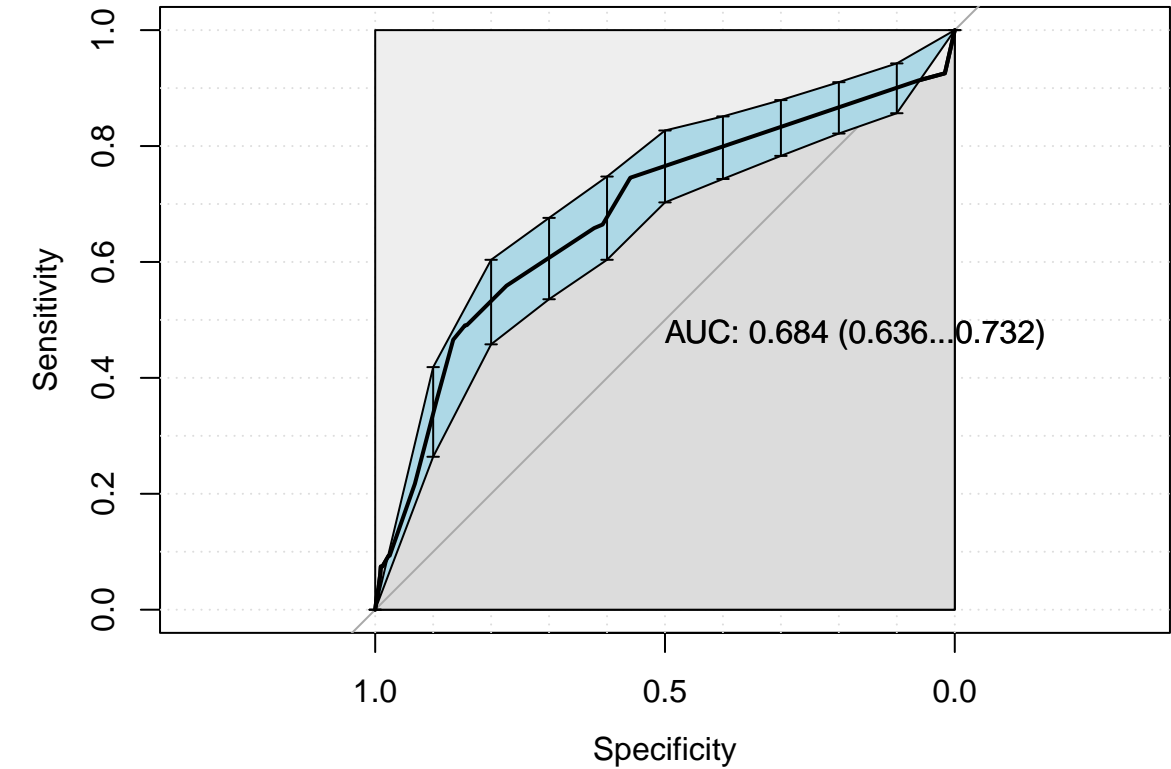
final_tree_workflow <-
  tree_workflow %>%
  finalize_workflow(best_tree)

last_tree_fit <-
  final_tree_workflow %>%
  last_fit(df_split)

final_tree_fit <- extract_workflow(last_tree_fit)

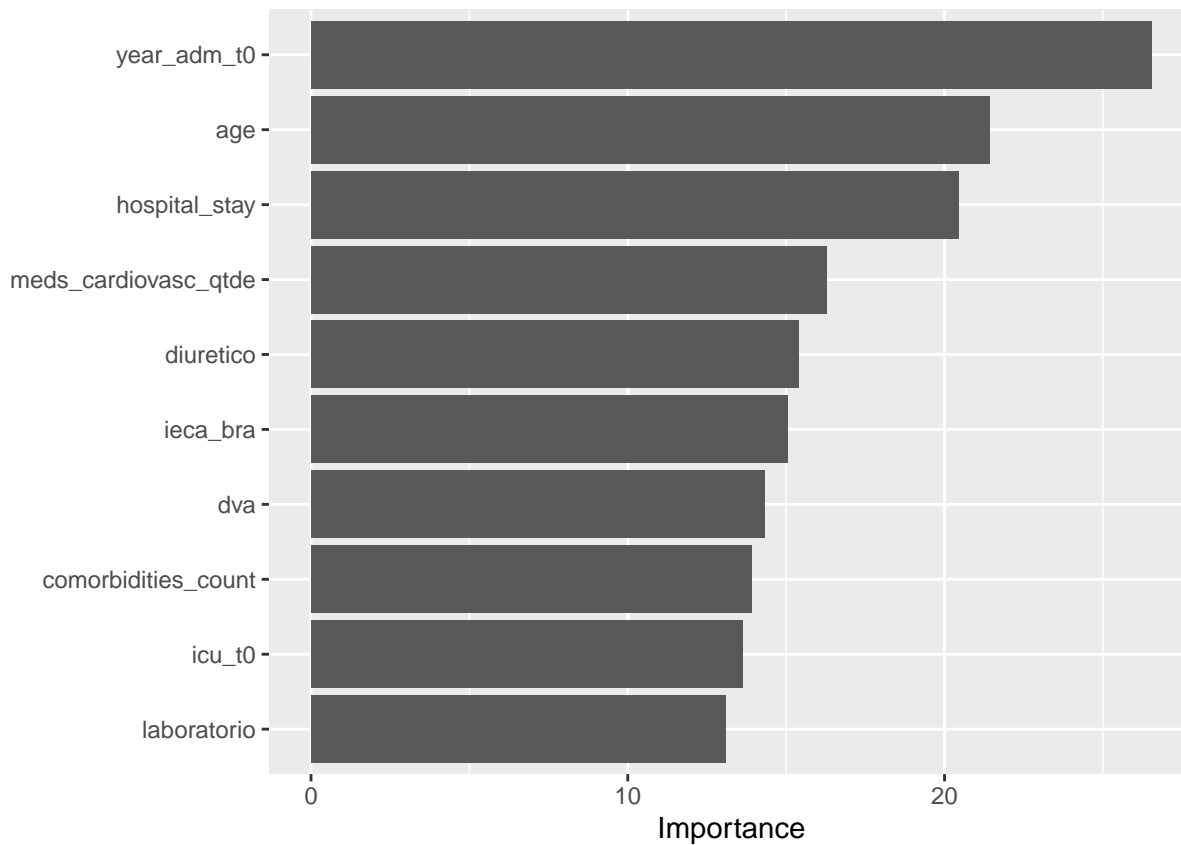
tree_auc <- validation(final_tree_fit, df_test)

```



```
##          Kappa : 0.1168
##
##  McNemar's Test P-Value : <2e-16
##
##          Sensitivity : 0.8450
##          Specificity : 0.4907
##          Pos Pred Value : 0.9792
##          Neg Pred Value : 0.1004
##          Prevalence : 0.9660
##          Detection Rate : 0.8163
##          Detection Prevalence : 0.8336
##          Balanced Accuracy : 0.6679
##
##          'Positive' Class : 0
##
```

```
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: 8.67

Random Forest

```

rf_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_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()

```

```

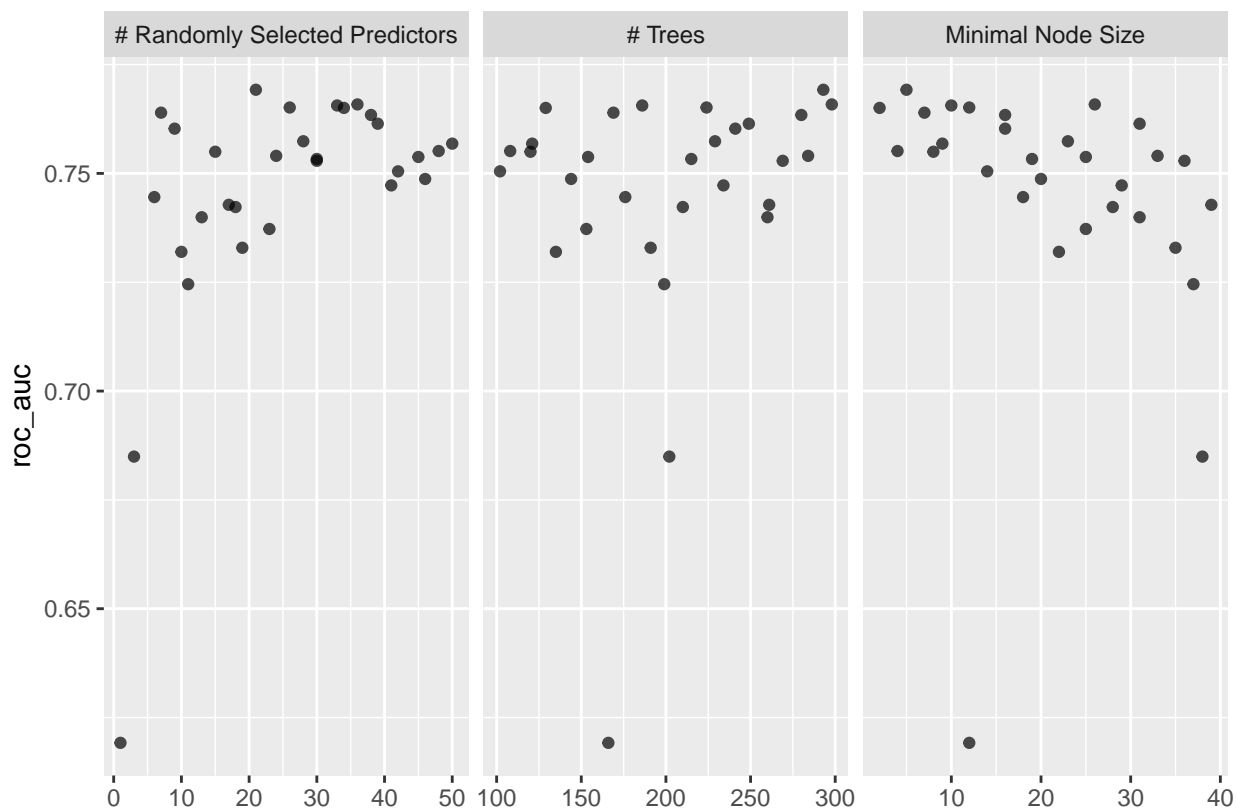
## # A tibble: 60 x 9
##   mtry trees min_n .metric .estimator mean      n std_err .config
##   <int> <int> <int> <chr>   <chr>   <dbl> <int>   <dbl> <chr>
## 1     48   108     4 accuracy binary    0.967    10 0.00135 Preprocessor1_Model01
## 2     48   108     4 roc_auc  binary    0.755    10 0.0136  Preprocessor1_Model01
## 3      6   176    18 accuracy binary    0.968    10 0.00143 Preprocessor1_Model02
## 4      6   176    18 roc_auc  binary    0.745    10 0.0133  Preprocessor1_Model02
## 5     26   224    12 accuracy binary    0.968    10 0.00140 Preprocessor1_Model03
## 6     26   224    12 roc_auc  binary    0.765    10 0.0116  Preprocessor1_Model03
## 7     15   120     8 accuracy binary    0.968    10 0.00143 Preprocessor1_Model04
## 8     15   120     8 roc_auc  binary    0.755    10 0.0157  Preprocessor1_Model04
## 9     30   215    19 accuracy binary    0.968    10 0.00140 Preprocessor1_Model05
## 10    30   215    19 roc_auc  binary    0.753    10 0.0169  Preprocessor1_Model05
## # ... with 50 more rows

```

```

autoplot(rf_tune, metric = "roc_auc")

```



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

```
## # A tibble: 5 x 9
##   mtry trees min_n .metric .estimator mean      n std_err .config
##   <int> <int> <int> <chr>   <chr>   <dbl> <int>   <dbl> <chr>
## 1    21   293     5 roc_auc binary  0.769    10  0.0143 Preprocessor1_Model26
## 2    36   298    26 roc_auc binary  0.766    10  0.0147 Preprocessor1_Model18
## 3    33   186    10 roc_auc binary  0.766    10  0.0133 Preprocessor1_Model20
## 4    26   224    12 roc_auc binary  0.765    10  0.0116 Preprocessor1_Model03
## 5    34   129     2 roc_auc binary  0.765    10  0.0124 Preprocessor1_Model08
```

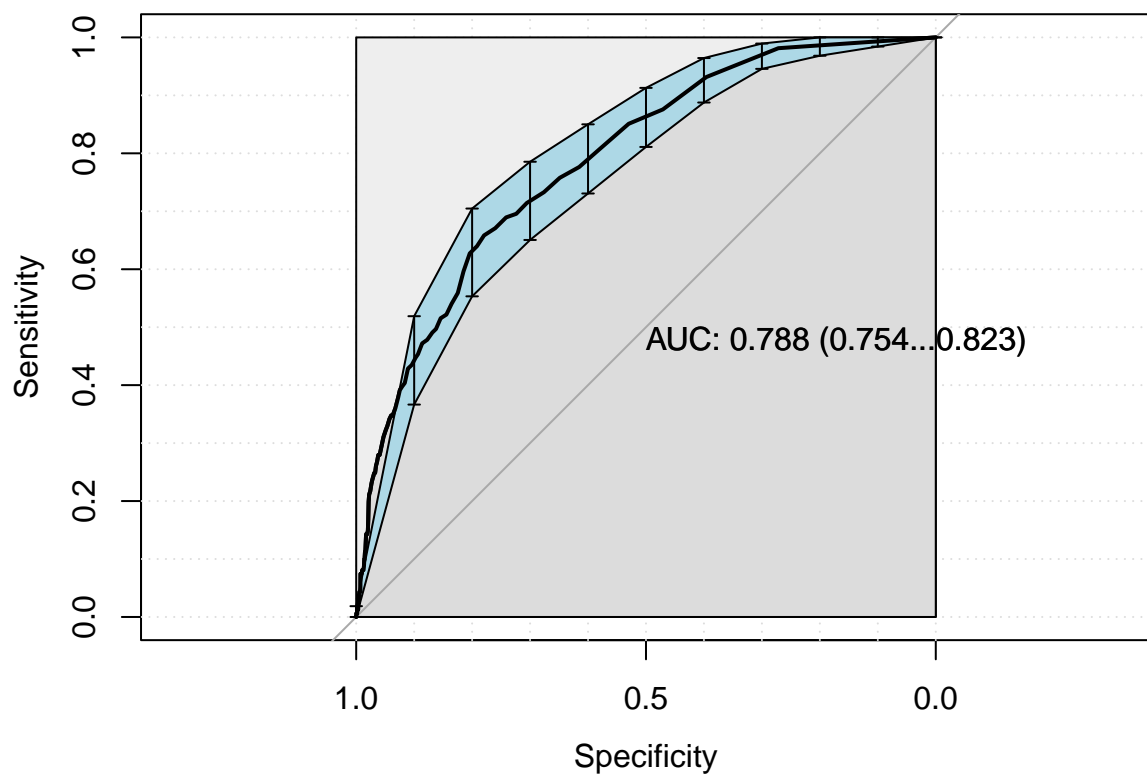
```
best_rf <- rf_tune %>%
  select_best("roc_auc")

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)

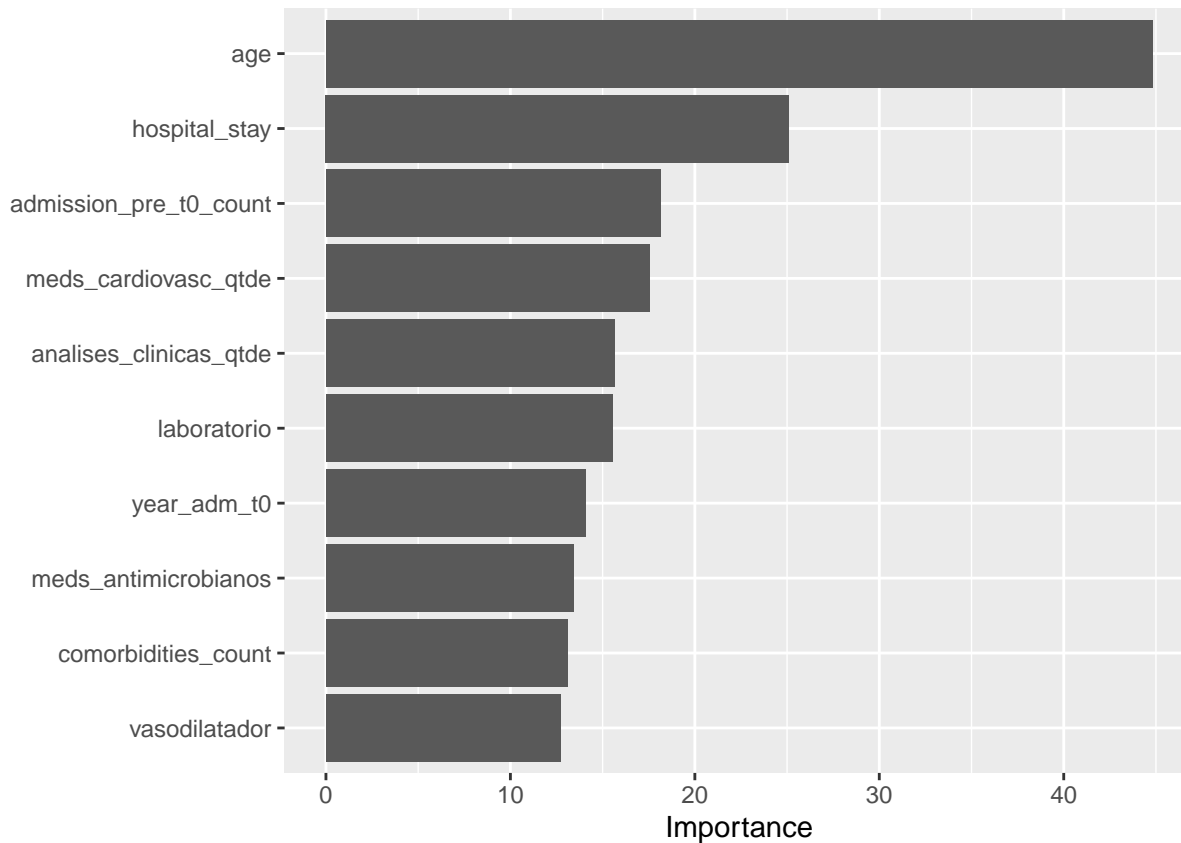
rf_auc <- validation(final_rf_fit, df_test)
```



|

```
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')
```



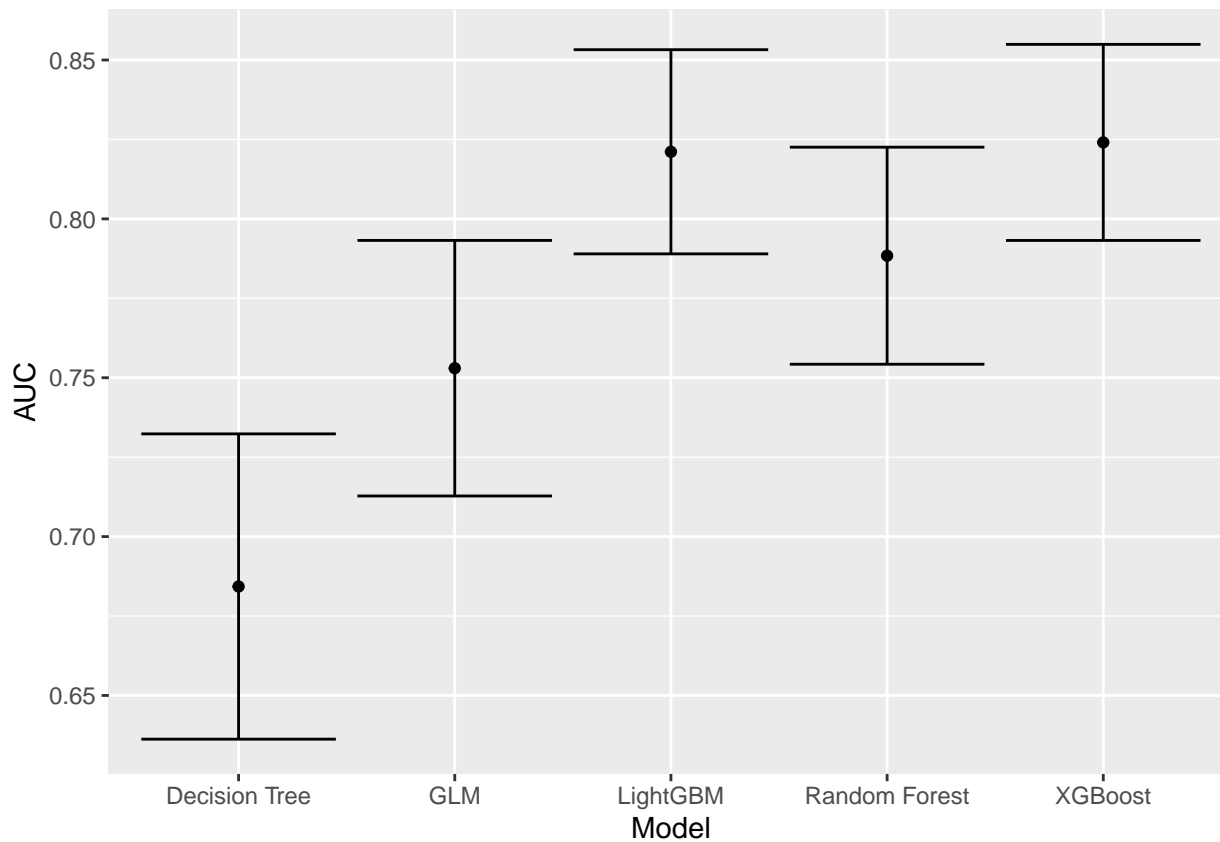
```
# extract_vip(final_rf_fit, pred_wrapper = predict,
#             reference_class = "1", use_matrix = FALSE,
#             method = 'permute')
```

Minutes to run: 141.815

Models Comparison

```
if (RUN_ALL_MODELS) {
  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$ci[3]
  ) %>%
    mutate(Target = outcome_column)
} else {
  df_auc <- tibble::tribble(
    ~Model, ~`AUC`, ~`Lower Limit`, ~`Upper Limit`,
    'LightGBM', as.numeric(lightgbm_auc$auc), lightgbm_auc$ci[1], lightgbm_auc$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: 0.006