

# Model Selection - readmission\_1year

Eduardo Yuki Yada

## Global parameters

```
k = 5 # Number of folds for cross validation
grid_size = 15 # Number of parameter combination to tune on each model
```

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: 0.006

```
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

```
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))

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. race
## 04. education_level
## 05. patient_state
## 06. underlying_heart_disease
## 07. heart_disease
## 08. nyha_basal
## 09. prior_mi
## 10. heart_failure
## 11. af
## 12. cardiac_arrest
## 13. transplant
## 14. valvopathy
## 15. endocardites
## 16. diabetes
## 17. renal_failure
## 18. hemodialysis
## 19. copd
## 20. comorbidities_count
## 21. procedure_type_1
## 22. reop_type_1
## 23. procedure_type_new
## 24. cied_final_1
## 25. cied_final_group_1
## 26. admission_pre_t0_count
## 27. admission_pre_t0_180d
## 28. year_adm_t0
## 29. icu_t0
## 30. dialysis_t0
## 31. admission_t0_emergency
```

## 32. aco  
## 33. antiarritmico  
## 34. betabloqueador  
## 35. ieca\_bra  
## 36. dva  
## 37. digoxina  
## 38. estatina  
## 39. diuretico  
## 40. vasodilatador  
## 41. insuf\_cardiaca  
## 42. espironolactona  
## 43. bloq\_calcio  
## 44. antiplaquetario\_ev  
## 45. insulina  
## 46. anticonvulsivante  
## 47. psicofarmacos  
## 48. antifungico  
## 49. antiviral  
## 50. antiretroviral  
## 51. classe\_meds\_qtde  
## 52. meds\_cardiovasc\_qtde  
## 53. meds\_antimicrobianos  
## 54. ventilacao\_mecanica  
## 55. cec  
## 56. transplante\_cardiaco  
## 57. cir\_toracica  
## 58. outros\_proced\_cirurgicos  
## 59. icp  
## 60. intervencao\_cv  
## 61. angioplastia  
## 62. cateterismo  
## 63. eletrofisiologia  
## 64. cateter\_venoso\_central  
## 65. proced\_invasivos\_qtde  
## 66. cve\_desf  
## 67. transfusao  
## 68. interconsulta  
## 69. equipe\_multiprof  
## 70. holter  
## 71. teste\_esforco  
## 72. espiro\_ergoespiro  
## 73. tilt\_teste  
## 74. metodos\_graficos\_qtde  
## 75. laboratorio  
## 76. cultura  
## 77. analises\_clinicas\_qtde  
## 78. citologia  
## 79. biopsia  
## 80. histopatologia\_qtde  
## 81. angio\_rm  
## 82. angio\_tc  
## 83. aortografia  
## 84. arteriografia  
## 85. cintilografia  
## 86. ecocardiograma  
## 87. endoscopia  
## 88. flebografia  
## 89. pet\_ct  
## 90. ultrassom  
## 91. tomografia  
## 92. ressonancia

```
## 93. exames_imagem_qtde
## 94. dieta_parenteral
## 95. bic
## 96. mpp
## 97. 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(
  mtry = tune(),
  trees = tune(),
  min_n = tune(),
  tree_depth = tune(),
  learn_rate = tune(),
  loss_reduction = tune()
) %>%
  set_engine("xgboost",
             nthread = 8) %>%
  set_mode("classification")

xgboost_grid <- grid_latin_hypercube(
  finalize(mtry(), df_train),
  trees(range = c(100L, 300L)),
  min_n(),
  tree_depth(),
  learn_rate(),
  loss_reduction(),
  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)

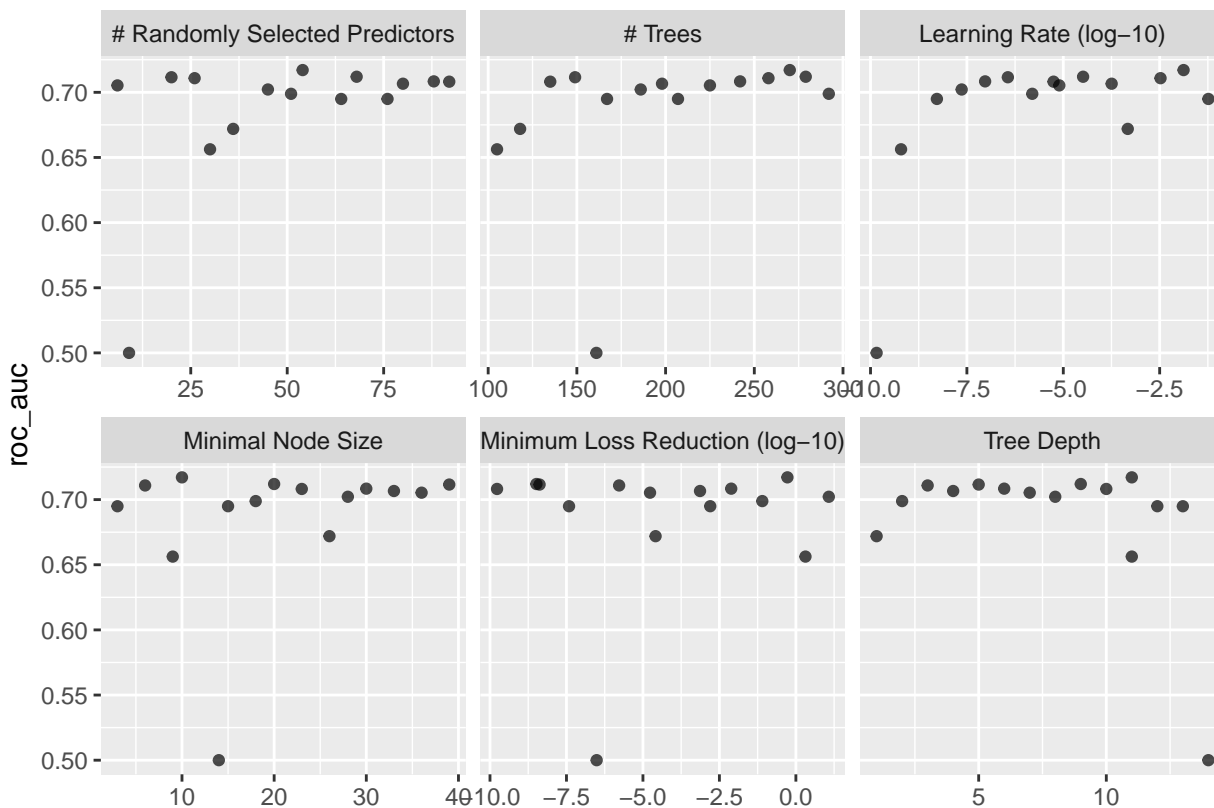
xgboost_tune %>%
  show_best("roc_auc")

## # A tibble: 5 x 12
##   mtry trees min_n tree_depth learn_rate loss_reduction .metric .estimator mean n std_err .config
##   <int> <int> <int>    <int>    <dbl>      <dbl> <chr>    <chr>    <dbl> <int>  <dbl> <chr>
## 1    54   270    10      11 0.0132      0.529  roc_auc  binary  0.717     5 0.00614 Prepro~
## 2    68   279    20       9 0.0000327  0.00000000329 roc_auc  binary  0.712     5 0.00656 Prepro~
## 3    20   149    39       5 0.000000365 0.00000000418 roc_auc  binary  0.712     5 0.00761 Prepro~
## 4    26   258     6       3 0.00333    0.00000167  roc_auc  binary  0.711     5 0.00609 Prepro~
## 5    88   242    30       6 0.0000000939 0.00770    roc_auc  binary  0.708     5 0.00729 Prepro~
## # ... with abbreviated variable name 1: .estimator

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

autoplot(xgboost_tune, metric = "roc_auc")

```



```

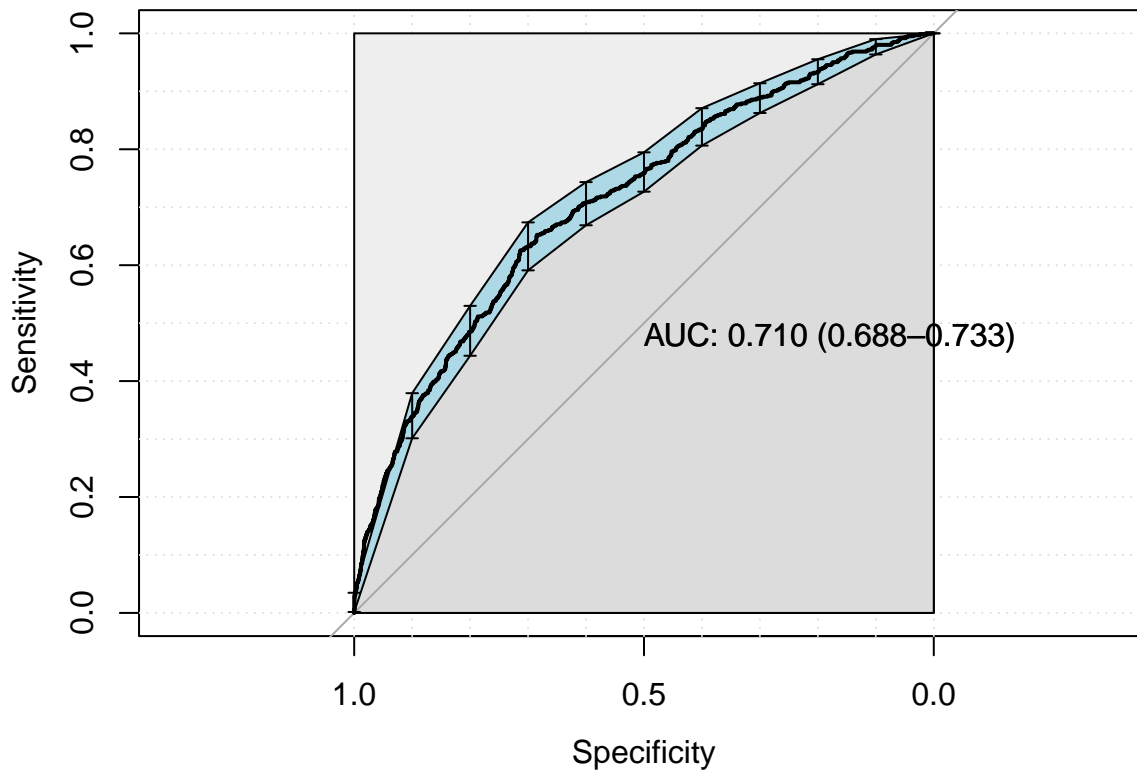
final_xgboost_workflow <-
  xgboost_workflow %>%
  finalize_workflow(best_xgboost)

last_xgboost_fit <-
  final_xgboost_workflow %>%
  last_fit(df_split)

final_xgboost_fit <- extract_workflow(last_xgboost_fit)

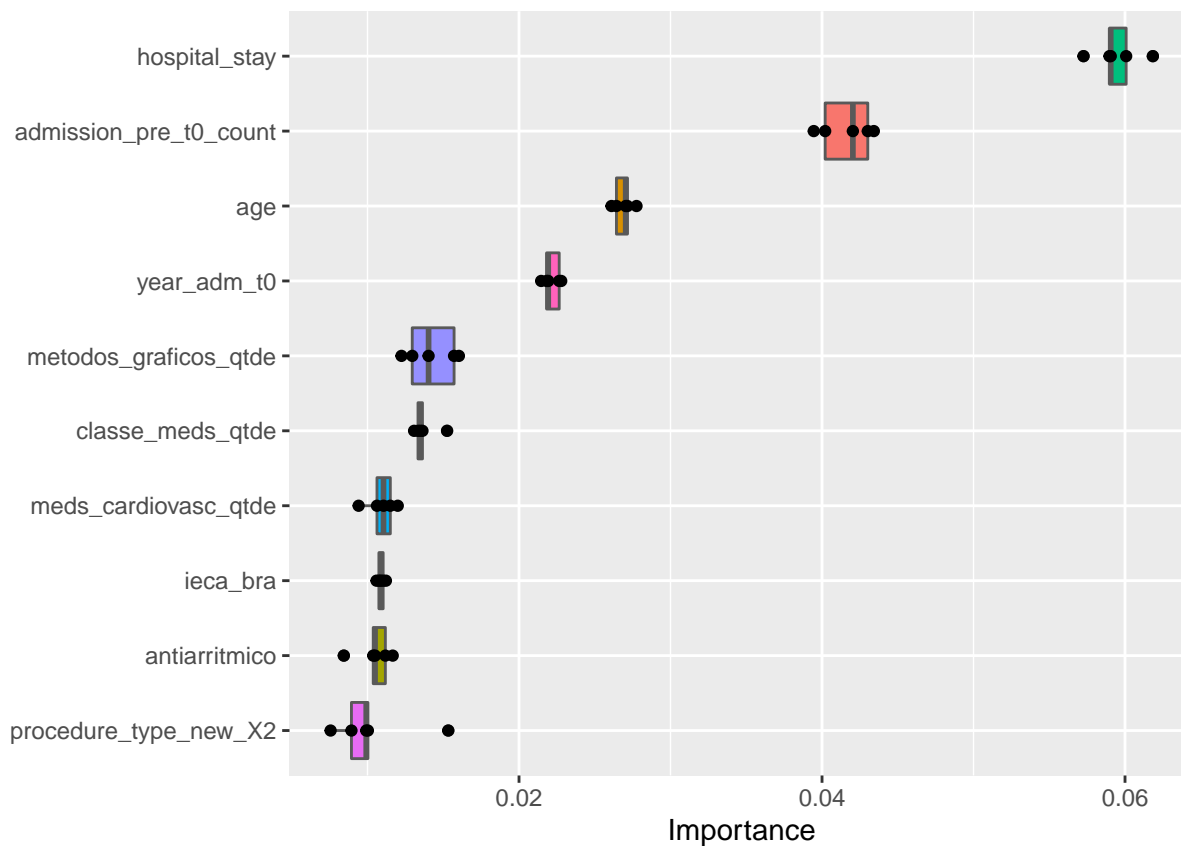
```

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



```
## [1] "Optimal Threshold: 0.13"
## Confusion Matrix and Statistics
##
##      reference
## data    0    1
## 0 2943  226
## 1 1184  378
##
##              Accuracy : 0.702
##              95% CI   : (0.6887, 0.715)
##      No Information Rate : 0.8723
##      P-Value [Acc > NIR] : 1
##
##              Kappa   : 0.2021
##
##      McNemar's Test P-Value : <2e-16
##
##              Sensitivity : 0.7131
##              Specificity : 0.6258
##              Pos Pred Value : 0.9287
##              Neg Pred Value : 0.2420
##              Prevalence : 0.8723
##              Detection Rate : 0.6221
##              Detection Prevalence : 0.6698
##              Balanced Accuracy : 0.6695
##
##              '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, mtry, 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: 3.003

## 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(
  mtry = tune(),
  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(
  finalize(mtry(), df_train),
  trees(range = c(100L, 300L)),
  min_n(),
  tree_depth(),
  learn_rate(),
  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 12
##   mtry trees min_n tree_depth   learn_rate loss_reduction .metric .estim~1 mean     n std_err .config
##   <int> <int> <int>     <int>         <dbl>         <dbl> <chr>   <chr>   <dbl> <int>   <dbl> <chr>
## 1    47   113    32         11 0.000521         0.0279   roc_auc binary  0.716     5 0.00719 Prepro~
## 2    46   133    36         14 0.000000167       0.00000492   roc_auc binary  0.714     5 0.00762 Prepro~
## 3    30   105    19          7 0.000000945       0.00322      roc_auc binary  0.713     5 0.00691 Prepro~
## 4    61   141    10         13 0.0000000137      0.000000278   roc_auc binary  0.712     5 0.00752 Prepro~
## 5    57   288     7          8 0.00000000173     0.00000353   roc_auc binary  0.712     5 0.00703 Prepro~
## # ... with abbreviated variable name 1: .estimator

```

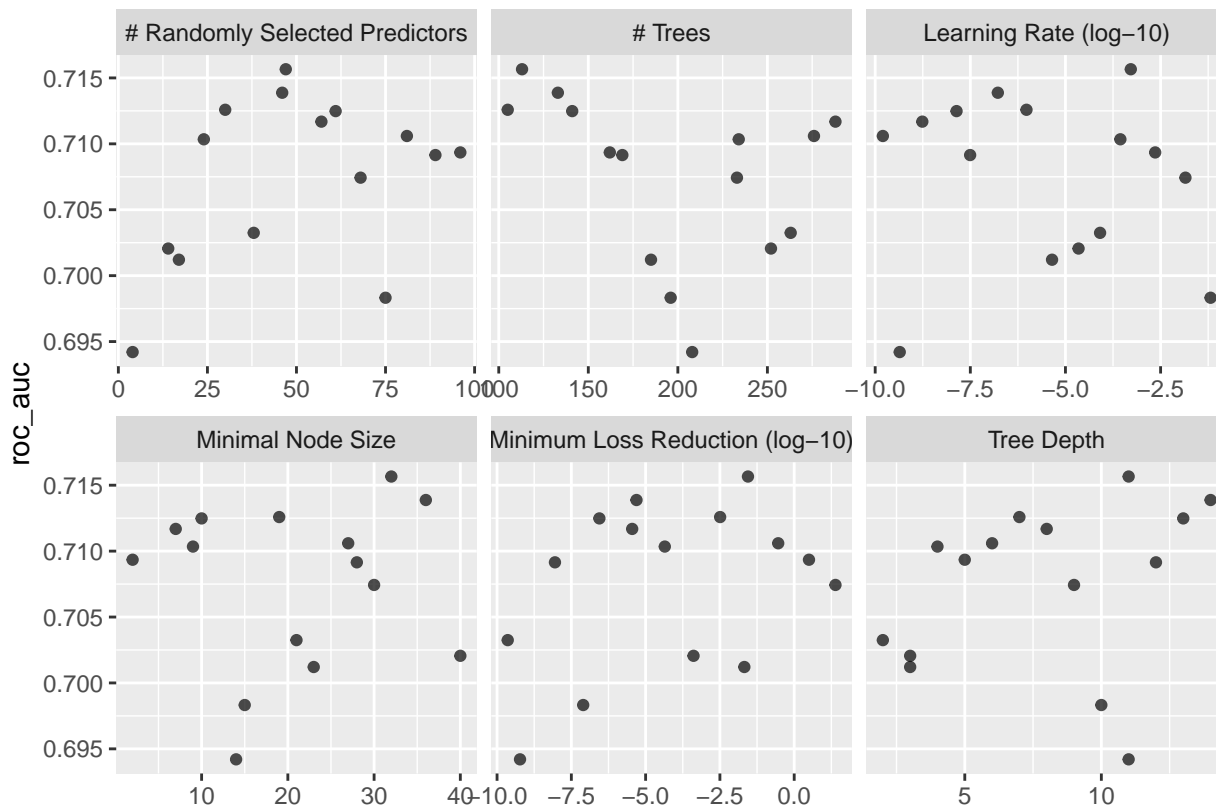
```

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

autoplot(lightgbm_tune, metric = "roc_auc")

```



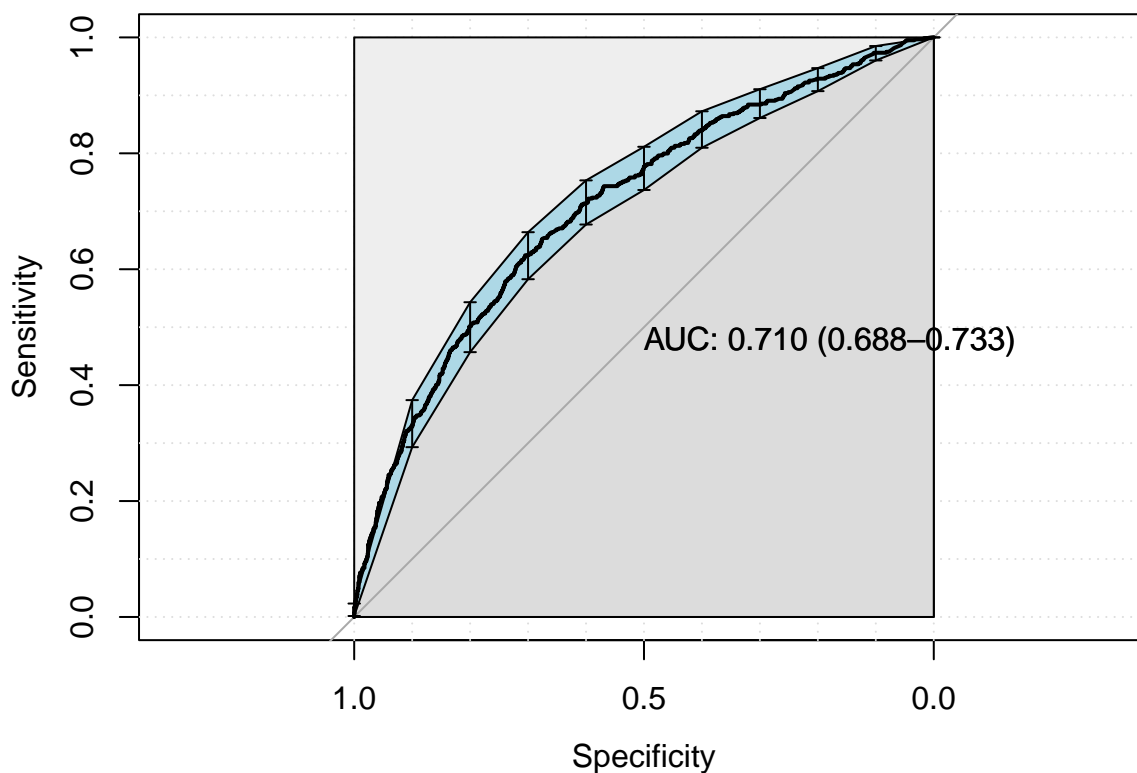


```
final_lightgbm_workflow <-
  lightgbm_workflow %>%
  finalize_workflow(best_lightgbm)

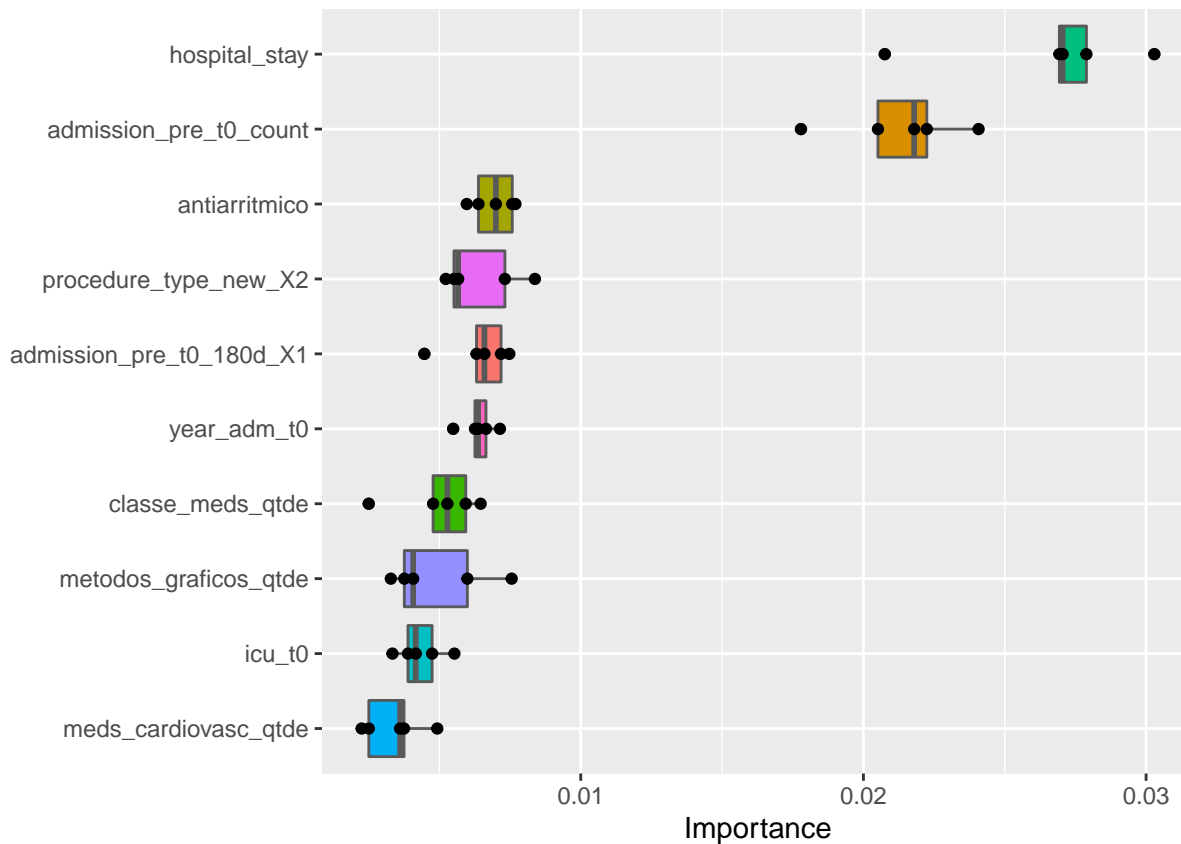
last_lightgbm_fit <-
  final_lightgbm_workflow %>%
  last_fit(df_split)

final_lightgbm_fit <- extract_workflow(last_lightgbm_fit)

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



```
## [1] "Optimal Threshold: 0.13"
## Confusion Matrix and Statistics
##
##      reference
## data    0    1
##    0 2909  227
##    1 1218  377
##
##              Accuracy : 0.6946
##              95% CI   : (0.6812, 0.7077)
##    No Information Rate : 0.8723
##    P-Value [Acc > NIR] : 1
##
##              Kappa   : 0.1935
##
## Mcnemar's Test P-Value : <2e-16
##
##              Sensitivity : 0.7049
##              Specificity : 0.6242
##              Pos Pred Value : 0.9276
##              Neg Pred Value : 0.2364
##              Prevalence : 0.8723
##              Detection Rate : 0.6149
##              Detection Prevalence : 0.6629
##              Balanced Accuracy : 0.6645
##
##              '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, mtry, 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: 2.615

## 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 = tune(), mixture = tune()) %>%
  set_mode("classification") %>%
  set_engine("glmnet")

glmnet_grid <- grid_latin_hypercube(penalty(),
  mixture(),
```

```
size = grid_size)
```

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

```
glmnet_tune <-  
  glmnet_workflow %>%  
  tune_grid(resamples = df_folds,  
            grid = glmnet_grid)
```

```
glmnet_tune %>%  
  collect_metrics()
```

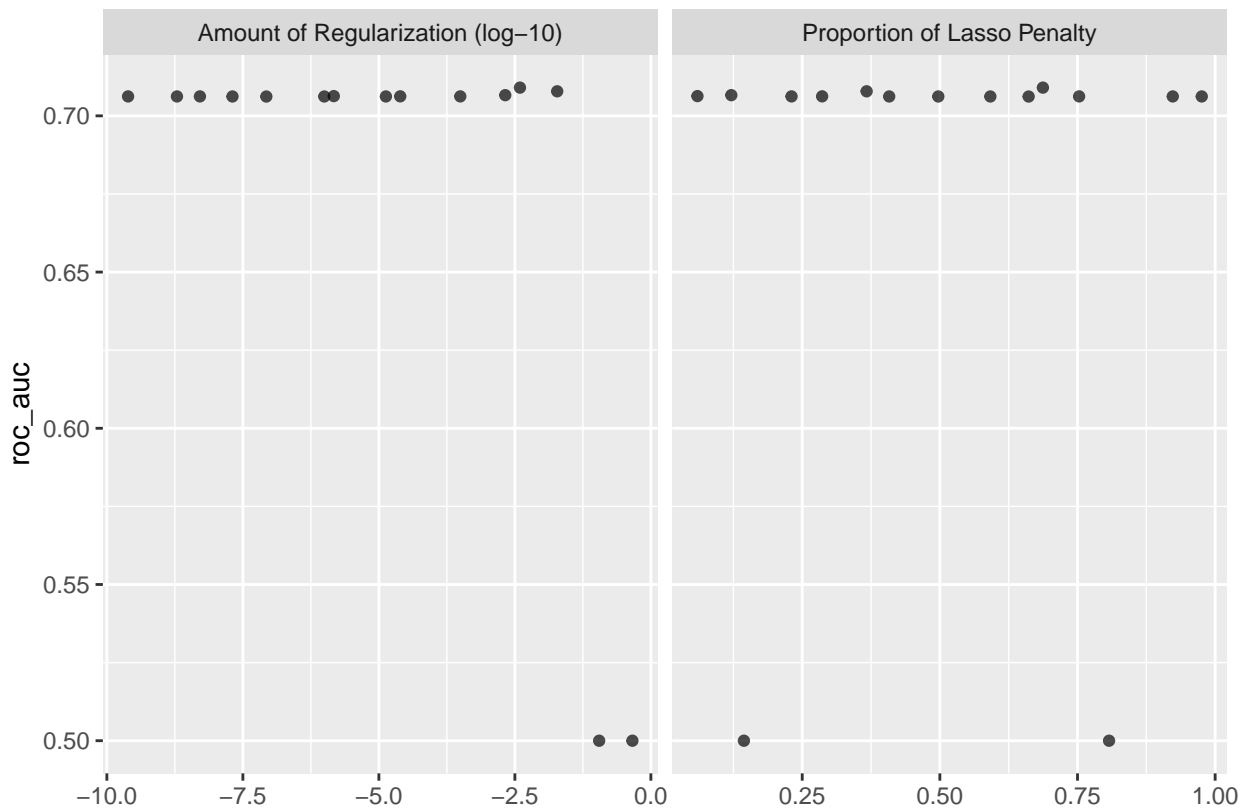
```
## # A tibble: 30 x 8
```

##	penalty	mixture	.metric	.estimator	mean	n	std_err	.config
##	<dbl>	<dbl>	<chr>	<chr>	<dbl>	<int>	<dbl>	<chr>
##	1	0.00000148	0.0597	accuracy	binary	0.874	5 0.00140	Preprocessor1_Model01
##	2	0.00000148	0.0597	roc_auc	binary	0.706	5 0.0100	Preprocessor1_Model01
##	3	0.00210	0.121	accuracy	binary	0.874	5 0.00144	Preprocessor1_Model02
##	4	0.00210	0.121	roc_auc	binary	0.707	5 0.00992	Preprocessor1_Model02
##	5	0.457	0.144	accuracy	binary	0.872	5 0.0000880	Preprocessor1_Model03
##	6	0.457	0.144	roc_auc	binary	0.5	5 0	Preprocessor1_Model03
##	7	0.000314	0.231	accuracy	binary	0.874	5 0.00148	Preprocessor1_Model04
##	8	0.000314	0.231	roc_auc	binary	0.706	5 0.0100	Preprocessor1_Model04
##	9	0.0000247	0.286	accuracy	binary	0.874	5 0.00140	Preprocessor1_Model05
##	10	0.0000247	0.286	roc_auc	binary	0.706	5 0.00999	Preprocessor1_Model05

```
## # ... with 20 more rows
```

```
## # i Use 'print(n = ...)' to see more rows
```

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



```

glmnet_tune %>%
  show_best("roc_auc")

## # A tibble: 5 x 8
##   penalty mixture .metric .estimator mean      n std_err .config
##   <dbl>   <dbl> <chr>   <chr>   <dbl> <int>   <dbl> <chr>
## 1 0.00392     0.687 roc_auc binary    0.709     5 0.00874 Preprocessor1_Model11
## 2 0.0190      0.367 roc_auc binary    0.708     5 0.00800 Preprocessor1_Model06
## 3 0.00210     0.121 roc_auc binary    0.707     5 0.00992 Preprocessor1_Model02
## 4 0.00000148 0.0597 roc_auc binary    0.706     5 0.0100  Preprocessor1_Model01
## 5 0.0000247  0.286 roc_auc binary    0.706     5 0.00999 Preprocessor1_Model05

best_glmnet <- glmnet_tune %>%
  select_best("roc_auc")

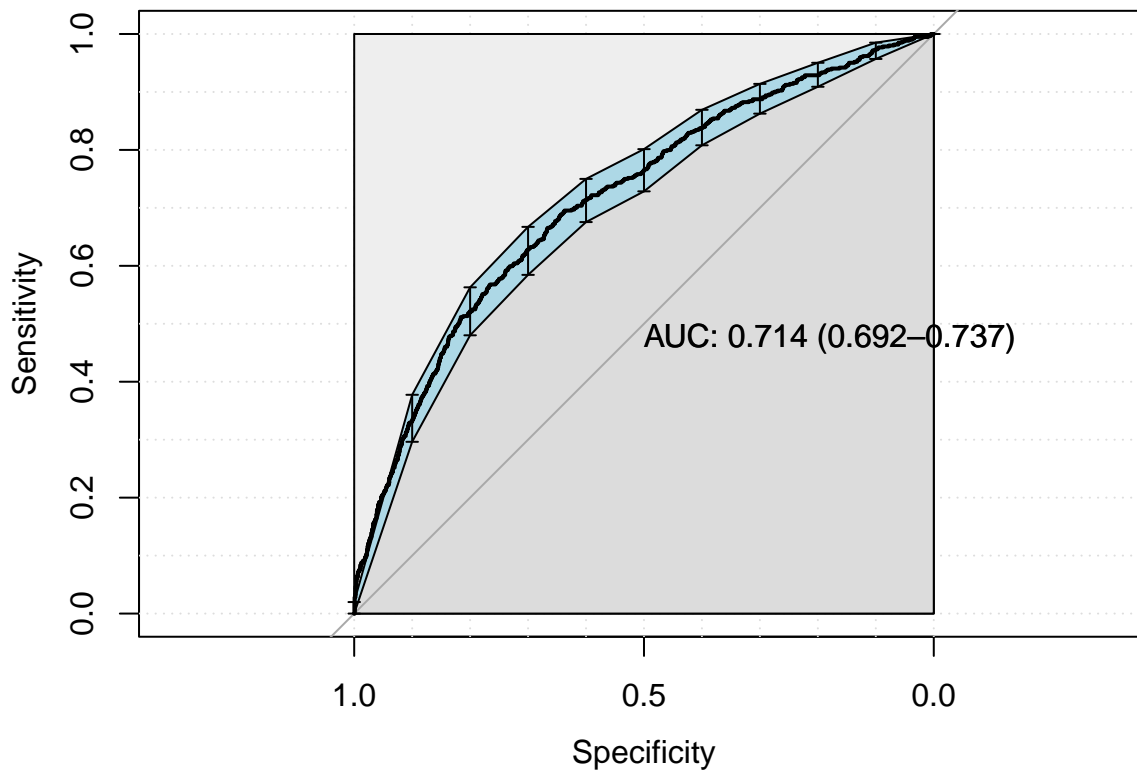
final_glmnet_workflow <-
  glmnet_workflow %>%
  finalize_workflow(best_glmnet)

last_glmnet_fit <-
  final_glmnet_workflow %>%
  last_fit(df_split)

final_glmnet_fit <- extract_workflow(last_glmnet_fit)

glmnet_auc <- validation(final_glmnet_fit, df_test)

```



```

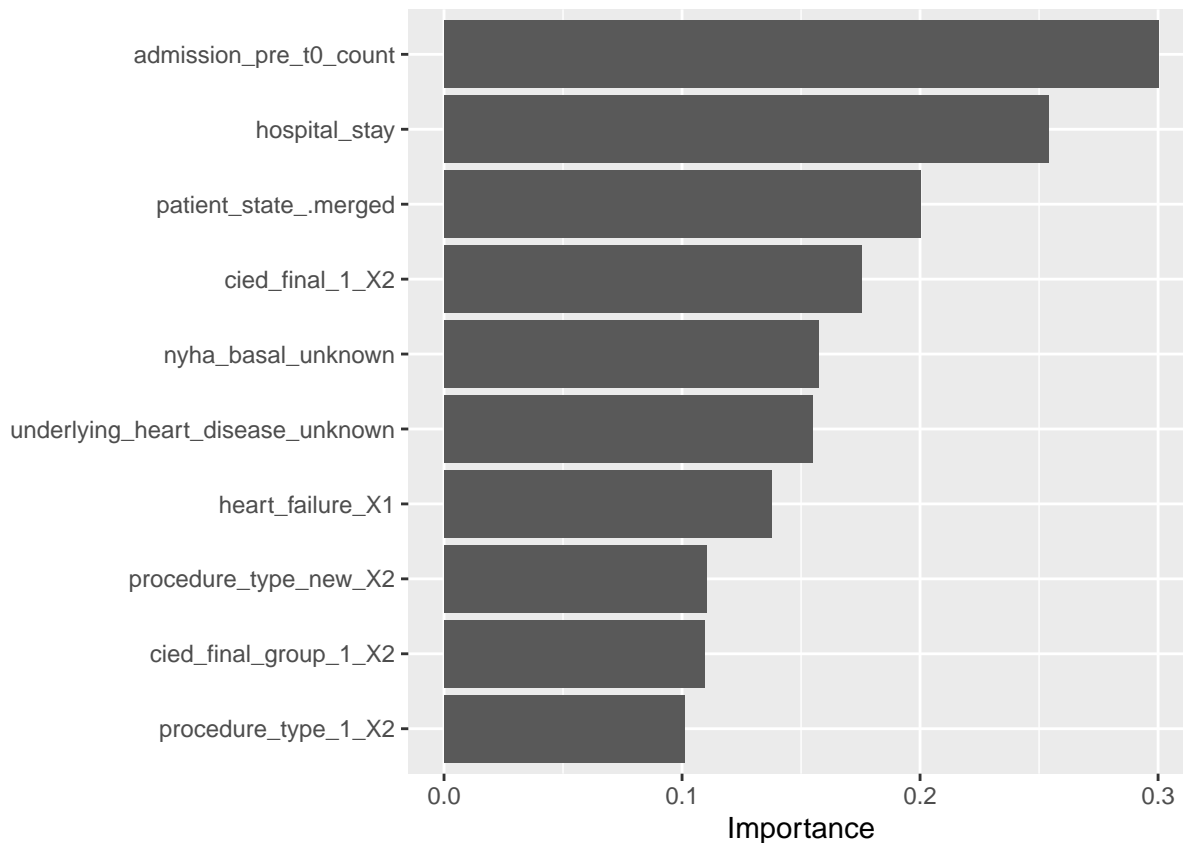
## [1] "Optimal Threshold: 0.14"
## Confusion Matrix and Statistics
##
##   reference
## data    0    1
##   0 3159 261
##   1  968 343

```

```
##
##           Accuracy : 0.7402
##           95% CI : (0.7275, 0.7527)
##           No Information Rate : 0.8723
##           P-Value [Acc > NIR] : 1
##
##           Kappa : 0.2223
##
## Mcnemar's Test P-Value : <2e-16
##
##           Sensitivity : 0.7654
##           Specificity : 0.5679
##           Pos Pred Value : 0.9237
##           Neg Pred Value : 0.2616
##           Prevalence : 0.8723
##           Detection Rate : 0.6677
##           Detection Prevalence : 0.7229
##           Balanced Accuracy : 0.6667
##
##           'Positive' Class : 0
##
```

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

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



1.077

Minutes to run:

## 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)

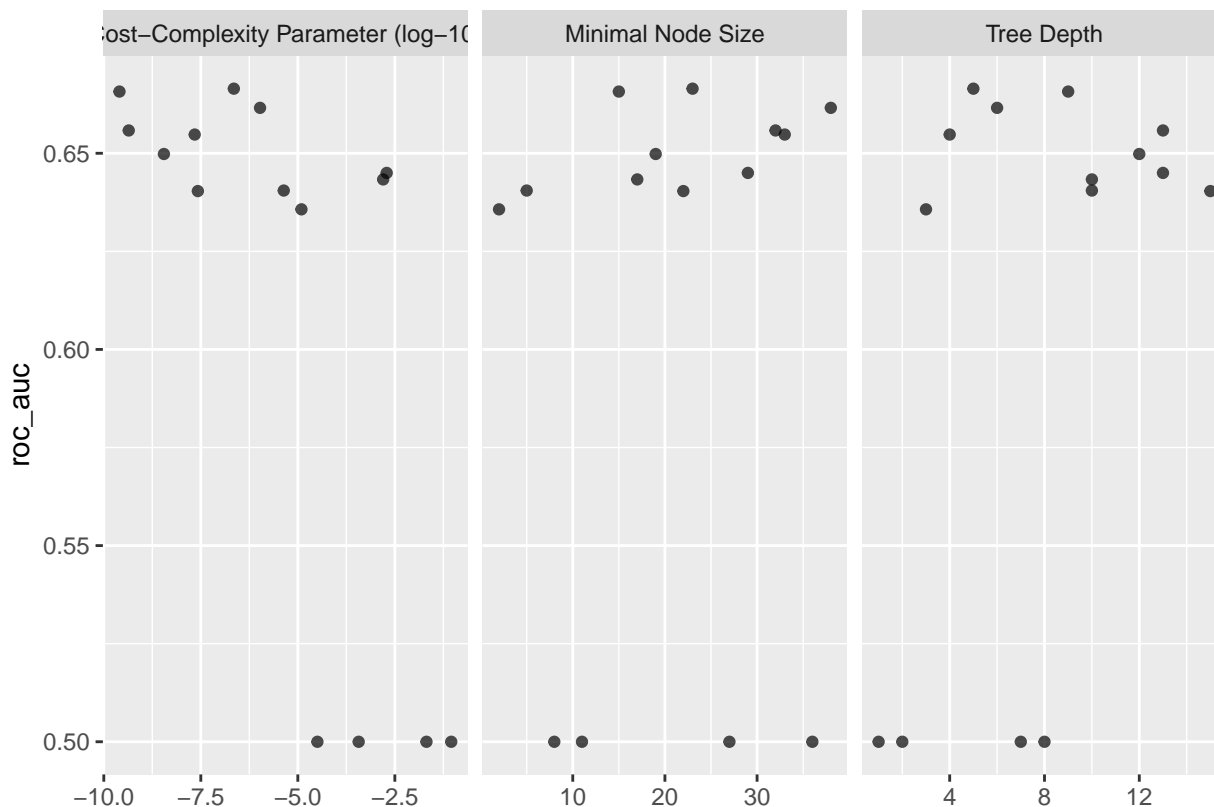
tree_workflow <-
  workflow() %>%
  add_recipe(tree_recipe) %>%
  add_model(tree_spec)

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")

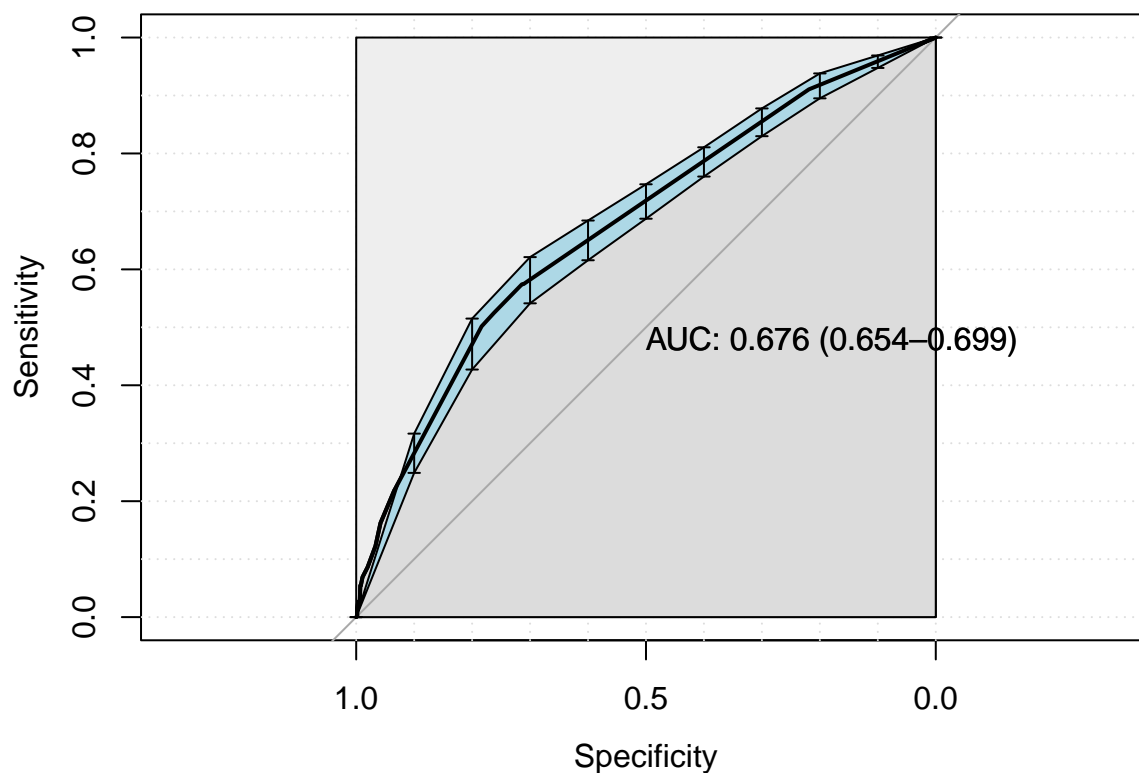
final_tree_workflow <-
  tree_workflow %>%
  finalize_workflow(best_tree)

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)

```

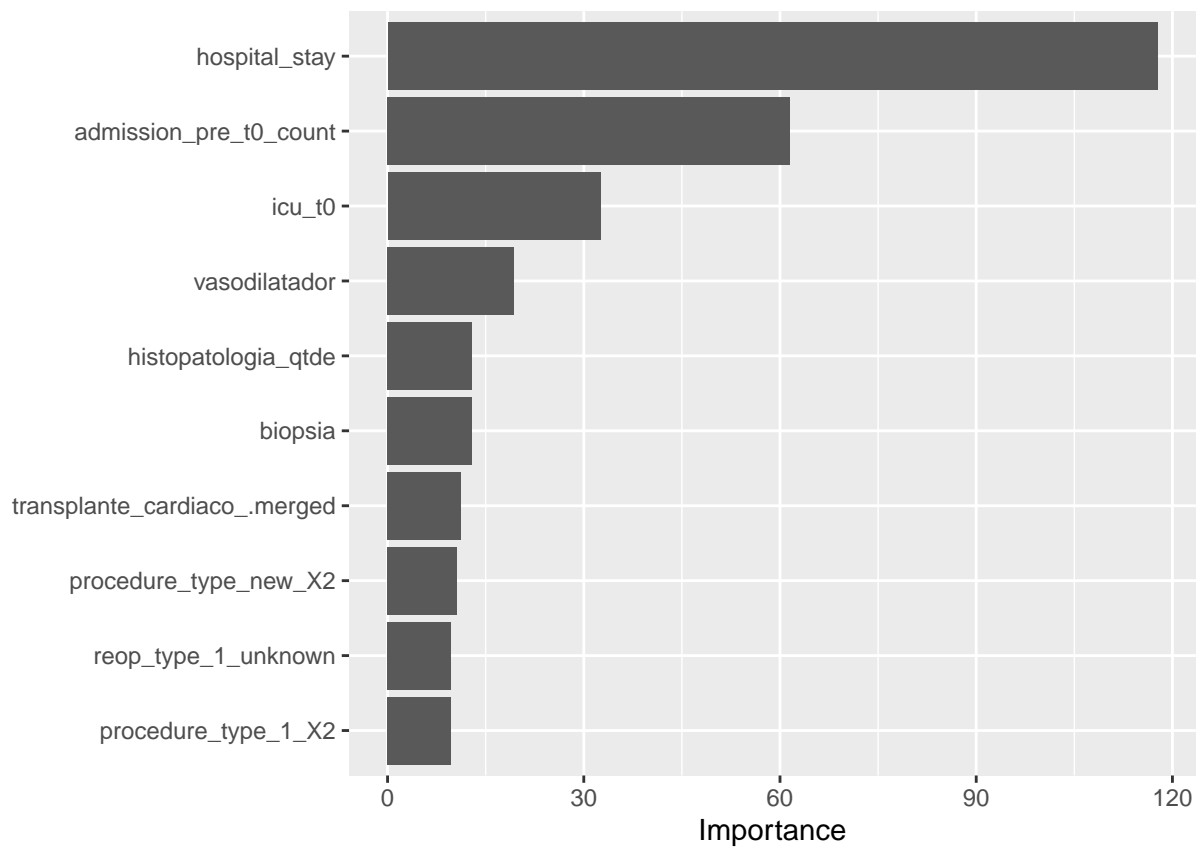


```

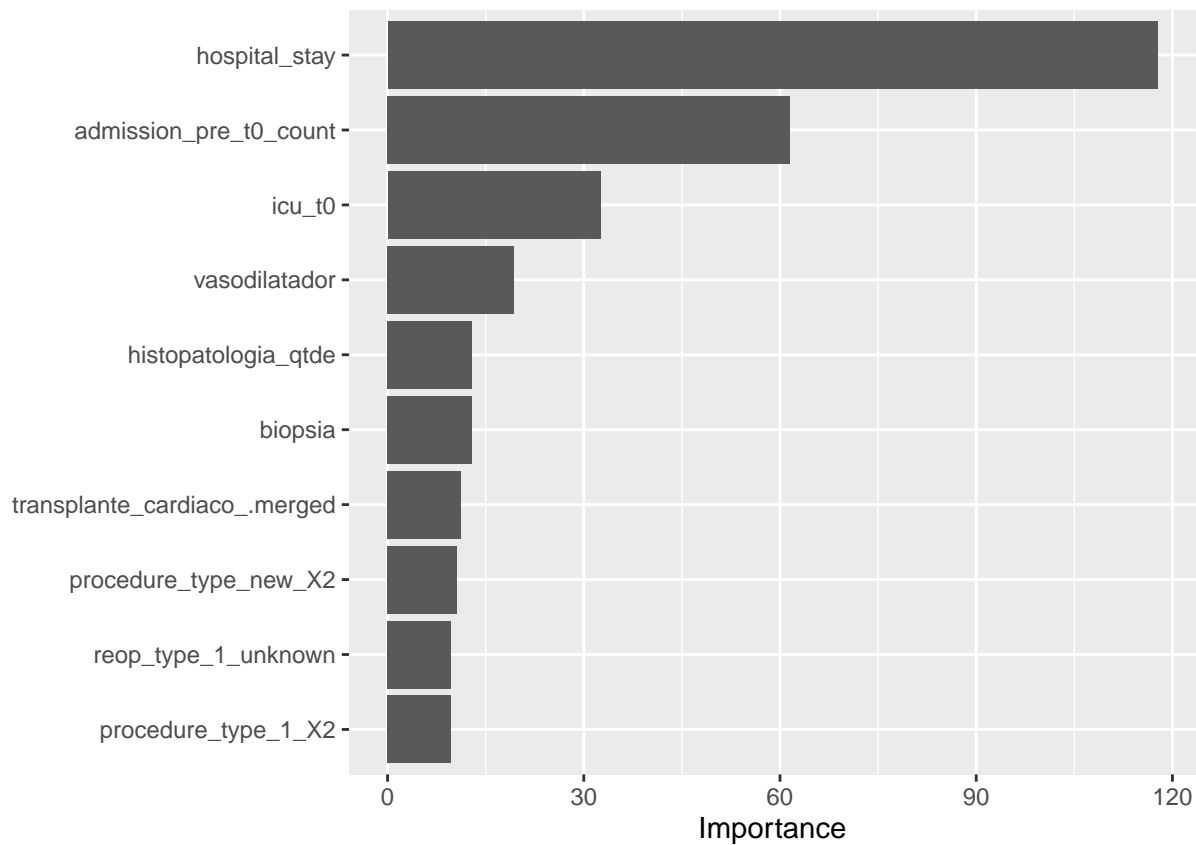
if (tree_auc$auc > 0.55) {
  final_tree_fit %>%
    extract_fit_parsnip() %>%
    vip()
}

```





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

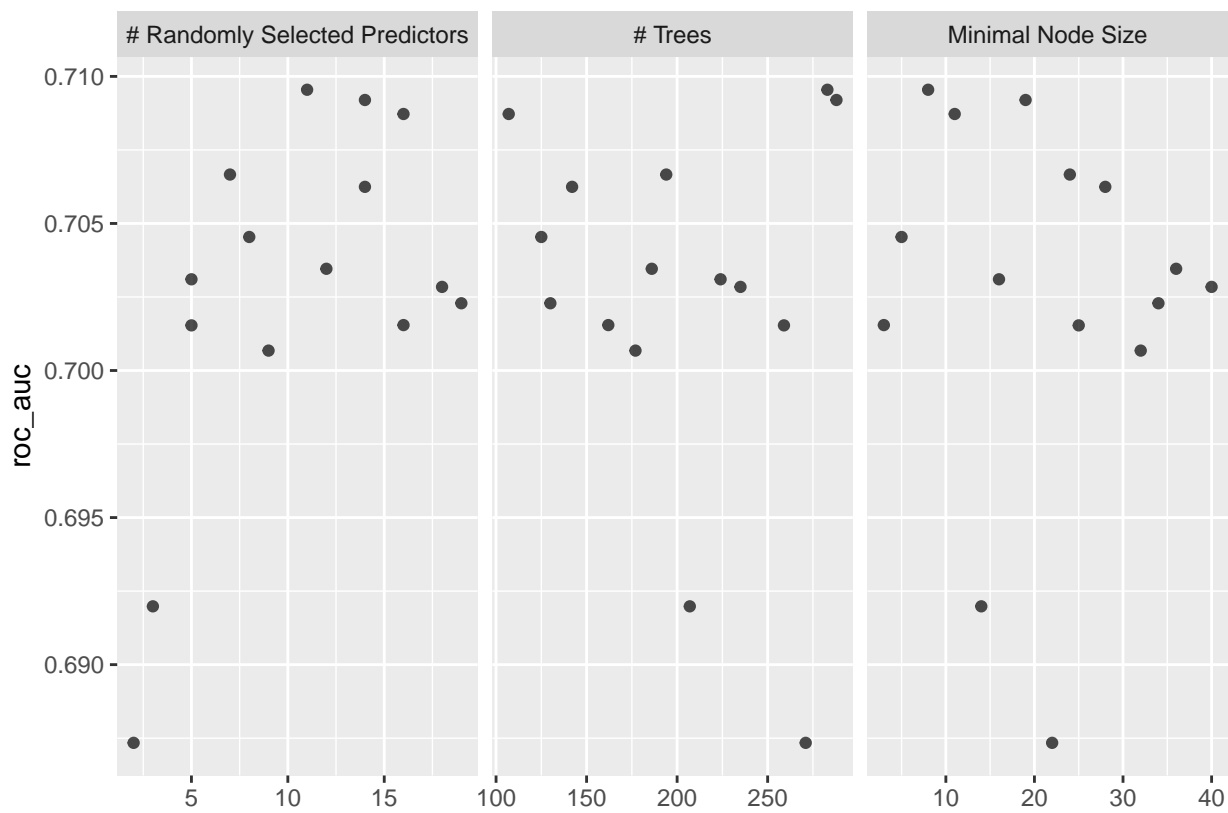


1.946

Minutes to run:

# 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, 20L)),  
  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")

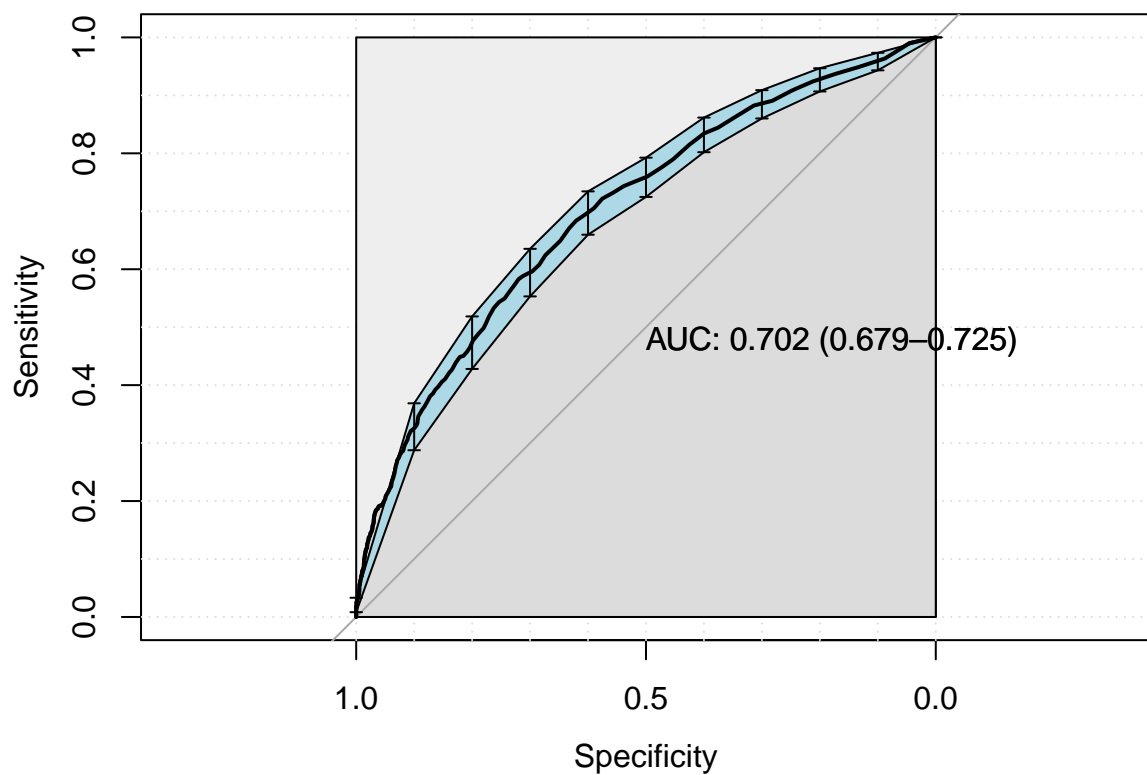
best_rf <- rf_tune %>%
  select_best("roc_auc")

final_rf_workflow <-
  rf_workflow %>%
  finalize_workflow(best_rf)

last_rf_fit <-
  final_rf_workflow %>%
  last_fit(df_split)

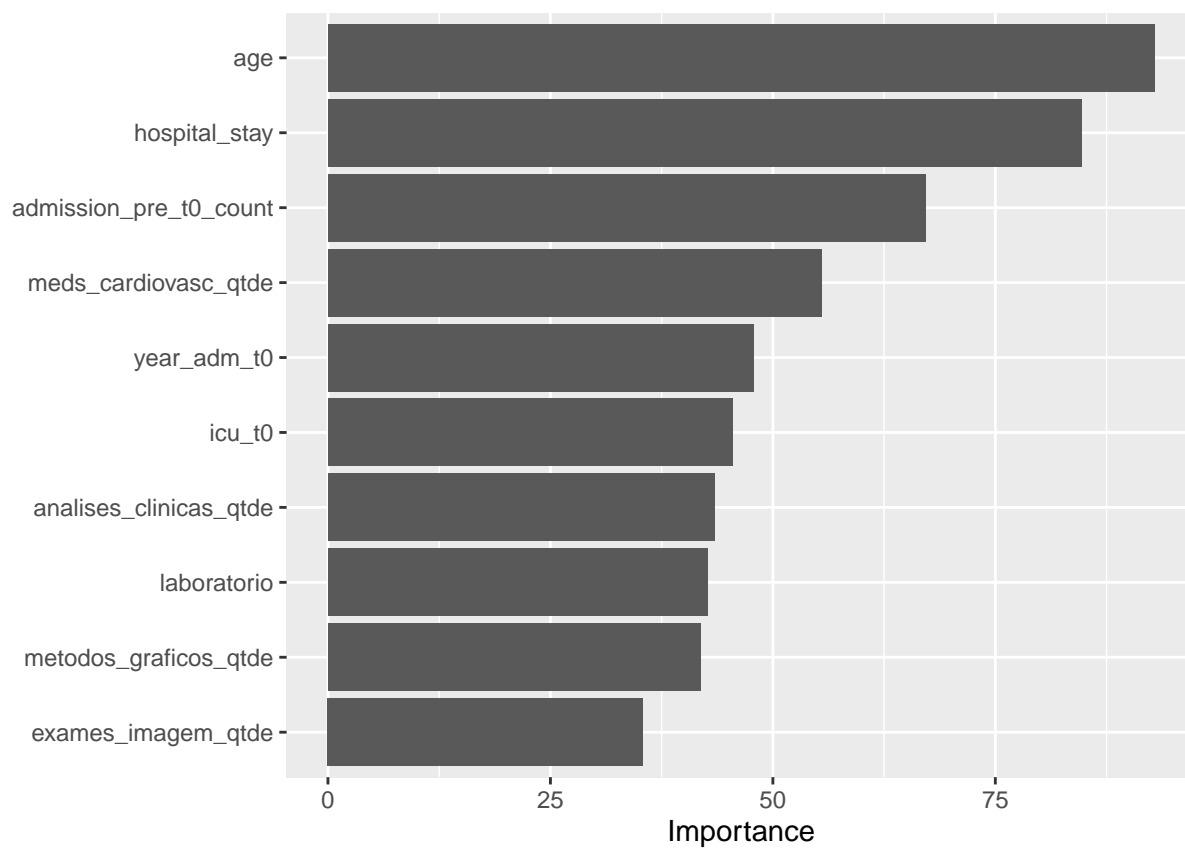
final_rf_fit <- extract_workflow(last_rf_fit)

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')
```



## KNN

```
# knn_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())  
#  
# knn_spec <-  
#   nearest_neighbor(neighbors = tune(),  
#                     weight_func = tune(),  
#                     dist_power = tune()) %>%  
#   set_mode("classification") %>%  
#   set_engine("kknn")  
#  
# knn_grid <- grid_latin_hypercube(neighbors(),  
#                                   weight_func(),  
#                                   dist_power(),  
#                                   size = grid_size)  
#  
# knn_workflow <-  
#   workflow() %>%  
#   add_recipe(knn_recipe) %>%  
#   add_model(knn_spec)  
#  
# knn_tune <-  
#   knn_workflow %>%  
#   tune_grid(resamples = df_folds,  
#             grid = knn_grid)  
#  
# knn_tune %>%  
#   collect_metrics()  
#  
# autoplot(knn_tune, metric = "roc_auc")  
#  
# knn_tune %>%  
#   show_best("roc_auc")  
#  
# best_knn <- knn_tune %>%  
#   select_best("roc_auc")  
#  
# final_knn_workflow <-  
#   knn_workflow %>%  
#   finalize_workflow(best_knn)  
#  
# last_knn_fit <-  
#   final_knn_workflow %>%  
#   last_fit(df_split)  
#  
# final_knn_fit <- extract_workflow(last_knn_fit)  
#  
# knn_auc = validation(final_knn_fit, df_test)
```

Minutes to run: 0

# SVM

```
# sum_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())
#
# sum_spec <-
#   sum_rbf(cost = tune(), rbf_sigma = tune()) %>%
#   set_mode("classification") %>%
#   set_engine("kernlab")
#
# sum_grid <- grid_latin_hypercube(cost(),
#                                   rbf_sigma(),
#                                   size = grid_size)
#
# sum_workflow <-
#   workflow() %>%
#   add_recipe(sum_recipe) %>%
#   add_model(sum_spec)
#
# sum_tune <-
#   sum_workflow %>%
#   tune_grid(resamples = df_folds,
#             grid = grid_size)
#
# sum_tune %>%
#   collect_metrics()
#
# autoplot(sum_tune, metric = "roc_auc")
#
# sum_tune %>%
#   show_best("roc_auc")
#
# best_sum <- sum_tune %>%
#   select_best("roc_auc")
#
# final_sum_workflow <-
#   sum_workflow %>%
#   finalize_workflow(best_sum)
#
# last_sum_fit <-
#   final_sum_workflow %>%
#   last_fit(df_split)
#
# final_sum_fit <- extract_workflow(last_sum_fit)
#
# sum_auc = validation(final_sum_fit, df_test)
```

Minutes to run: 0

## Models Comparison

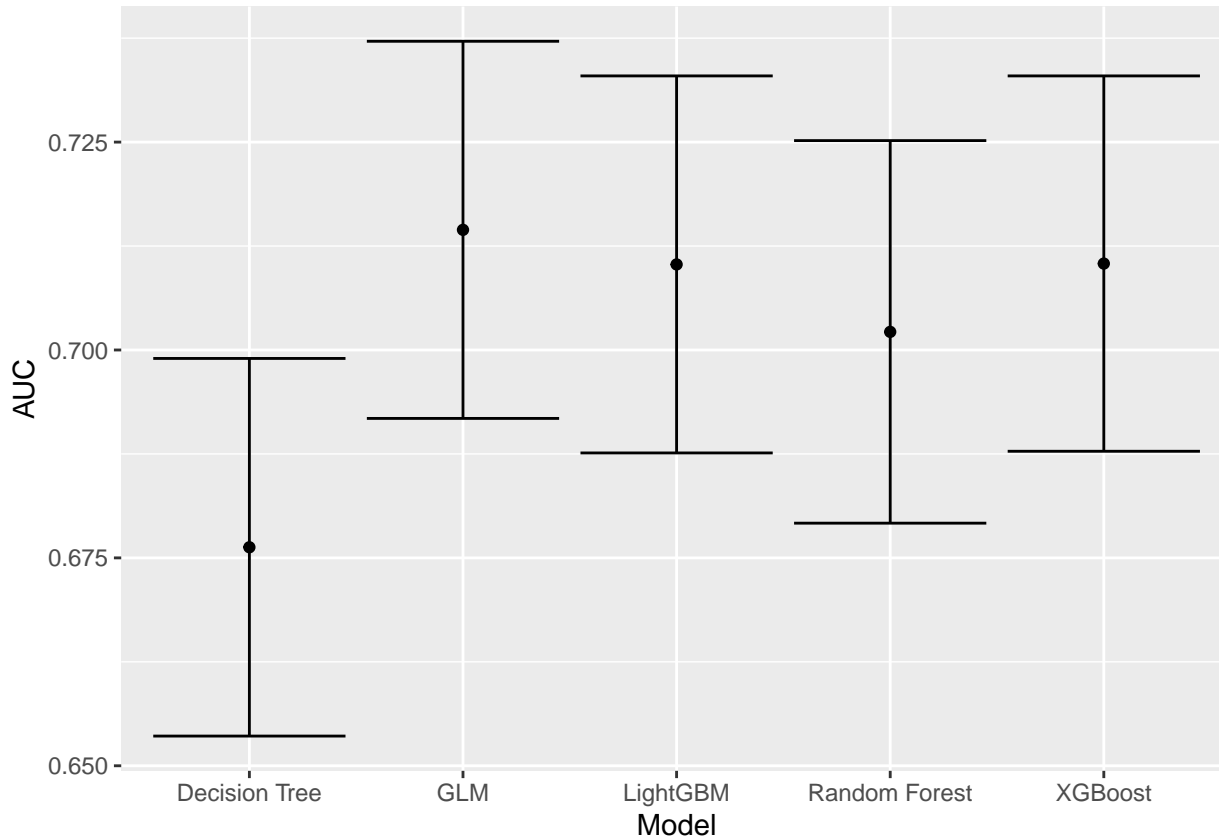
```
df_auc <- tibble::tribble(
  ~Model, ~`AUC`, ~`Lower Limit`, ~`Upper Limit`,
  'XGBoost', as.numeric(xgboost_auc$auc), xgboost_auc$ci[1], xgboost_auc$ci[3],
```

```

'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: 0.002