

Final Model

Eduardo Yuki Yada

Imports

```
library(tidyverse)
library(yaml)
library(tidymodels)
library(usemodels)
library(vip)

library(SHAPforxgboost)
library(xgboost)
library(Matrix)
library(mltools)
library(bonsai)
library(lightgbm)
```

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
```

Minutes to run: 0.001

Filtering eligible pacients

```
df = df %>%
  filter(disch_outcomes_t0 == 0)

df %>% dim

## [1] 15766   239
```

Minutes to run: 0.007

Eligible features

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

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

```

correlated_columns = c('year_procedure_1', # com year_adm_t0
                     'age_surgery_1', # com age
                     'admission_pre_t0_count', # com admission_t0
                     'atb', # com meds_antimicrobianos
                     'classe_meds_cardio_qtde', # com classe_meds_qtde
                     'suporte_hemod' # com proced_invasivos_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. age
## 02. education_level
## 03. underlying_heart_disease
## 04. heart_disease
## 05. nyha_basal
## 06. prior_mi
## 07. heart_failure
## 08. af
## 09. cardiac_arrest
## 10. transplant
## 11. valvopathy
## 12. diabetes
## 13. hemodialysis
## 14. comorbidities_count
## 15. procedure_type_1
## 16. reop_type_1
## 17. procedure_type_new
## 18. cied_final_1
## 19. cied_final_group_1
## 20. admission_t0
## 21. admission_pre_t0_180d
## 22. icu_t0
## 23. dialysis_t0
## 24. disch_outcomes_t0
## 25. admission_t0_emergency
## 26. aco
## 27. antiarritmico
## 28. betabloqueador
## 29. ieca_bra
## 30. dva
## 31. digoxina
## 32. estatina
## 33. diuretico
## 34. vasodilatador
## 35. insuf_cardiaca
## 36. espironolactona
## 37. bloq_calcio
## 38. antiplaquetario_ev
## 39. insulina
## 40. anticonvulsivante
## 41. psicofarmacos

```

```

## 42. antifungico
## 43. antiviral
## 44. classe_meds_qtde
## 45. meds_cardiovasc_qtde
## 46. meds_antimicrobianos
## 47. cec
## 48. transplante_cardiaco
## 49. cir_toracica
## 50. outros_proced_cirurgicos
## 51. icp
## 52. angioplastia
## 53. cateterismo
## 54. eletrofisiologia
## 55. cateter_venoso_central
## 56. proced_invasivos_qtde
## 57. cve_desf
## 58. transfusao
## 59. equipe_multiprof
## 60. ecg
## 61. holter
## 62. teste_esforco
## 63. espiro_ergoespiro
## 64. tilt_teste
## 65. metodos_graficos_qtde
## 66. laboratorio
## 67. cultura
## 68. analises_clinicas_qtde
## 69. citologia
## 70. biopsia
## 71. histopatologia_qtde
## 72. angio_rm
## 73. angio_tc
## 74. arteriografia
## 75. cintilografia
## 76. ecocardiograma
## 77. endoscopia
## 78. pet_ct
## 79. ultrassom
## 80. tomografia
## 81. radiografia
## 82. ressonancia
## 83. exames_imagem_qtde
## 84. bic

```

Minutes to run: 0

Train test split (70%/30%)

```

set.seed(42)

df[columns_list$outcome_columns] <- lapply(df[columns_list$outcome_columns], factor)
df <- mutate(df, across(where(is.character), as.factor))

df_split <- initial_split(df %>% dplyr::select(all_of(c(features, outcome_column))), 
                           prop = .7, strata = all_of(outcome_column))
df_train <- training(df_split)
df_test <- testing(df_split)

dim(df_train)[1] / dim(df)[1]

## [1] 0.6999873

```

```
dim(df_test)[1] / dim(df)[1]
```

```
## [1] 0.3000127
```

Minutes to run: 0.005

Global parameters

```
k <- 4 # Number of folds for cross validation
grid_size <- 50 # Number of parameter combination to tune on each model

set.seed(234)
df_folds <- vfold_cv(df_train, v = k,
                      strata = all_of(outcome_column))

max_auc_loss <- 0.01
```

Minutes to run: 0.001

Functions

```
validation = function(model_fit, new_data, plot=TRUE) {
  library(pROC)
  library(caret)

  test_predictions_prob <-
    predict(model_fit, new_data = new_data, type = "prob") %>%
    rename_at(vars(starts_with(".pred_")), ~ str_remove(., ".pred_")) %>%
    .\$`1` 

  pROC_obj <- roc(
    new_data[[outcome_column]],
    test_predictions_prob,
    direction = "<",
    levels = c(0, 1),
    smoothed = TRUE,
    ci = TRUE,
    ci.alpha = 0.9,
    stratified = FALSE,
    plot = plot,
    auc.polygon = TRUE,
    max.auc.polygon = TRUE,
    grid = TRUE,
    print.auc = TRUE,
    show.thres = TRUE
  )

  test_predictions_class <-
    predict(model_fit, new_data = new_data, type = "class") %>%
    rename_at(vars(starts_with(".pred_")), ~ str_remove(., ".pred_")) %>%
    .\$class

  conf_matrix <- table(test_predictions_class, new_data[[outcome_column]])

  if (plot) {
    sens.ci <- ci.se(pROC_obj)
    plot(sens.ci, type = "shape", col = "lightblue")
    plot(sens.ci, type = "bars")

    confusionMatrix(conf_matrix) %>% print
```

```

}

return(pROC_obj)
}

```

Minutes to run: 0

Feature Selection

```

model_fit_wf <- function(features, outcome_column, hyperparameters){
  model_recipe <-
    recipe(formula = sprintf("%s ~ .", outcome_column) %>% as.formula,
           data = df_train %>% select(all_of(c(features, outcome_column)))) %>%
    step_novel(all_nominal_predictors()) %>%
    step_unknown(all_nominal_predictors()) %>%
    step_other(all_nominal_predictors(), threshold = 0.05, other=".merged") %>%
    step_impute_mean(all_numeric_predictors()) %>%
    step_zv(all_predictors())

  model_spec <-
    do.call(boost_tree, hyperparameters) %>%
    set_engine("lightgbm") %>%
    set_mode("classification")

  model_workflow <-
    workflow() %>%
    add_recipe(model_recipe) %>%
    add_model(model_spec)

  model_fit_rs <- model_workflow %>%
    fit_resamples(df_folds)

  model_fit <- model_workflow %>%
    fit(df_train)

  model_auc <- validation(model_fit, df_test, plot=F)

  raw_model <- parsnip::extract_fit_engine(model_fit)

  feature_importance <- lgb.importance(raw_model, percentage = TRUE)

  return(list(cv_auc = collect_metrics(model_fit_rs) %>% filter(.metric == 'roc_auc') %>% .$mean,
             importance = feature_importance,
             auc = as.numeric(model_auc$auc),
             auc_lower = model_auc$ci[1],
             auc_upper = model_auc$ci[3]))
}

}

```

Minutes to run: 0

```

hyperparameters <- readRDS(
  sprintf(
    "../EDA/auxiliar/hyperparameters/model_selection/lightgbm_parameters_%s.rds",
    outcome_column
  )
)

full_model <- model_fit_wf(features, outcome_column, hyperparameters)

sprintf('Full Model CV Train AUC: %.3f' ,full_model$cv_auc)

```

```

## [1] "Full Model CV Train AUC: 0.709"
sprintf('Full Model Test AUC: %.3f' ,full_model$auc)

## [1] "Full Model Test AUC: 0.668"
Minutes to run: 0.376

Features with zero importance on the initial model:
unimportant_features <- setdiff(features, full_model$importance$Feature)

unimportant_features %>%
  gluedown::md_order()

## 1. underlying_heart_disease
## 2. heart_disease
## 3. transplant
## 4. hemodialysis
## 5. disch_outcomes_t0
## 6. cir_toracica
## 7. angioplastia
## 8. angio_rm
## 9. arteriografia

Minutes to run: 0

trimmed_features <- full_model$importance$Feature

trimmed_model <- model_fit_wf(trimmed_features,
                               outcome_column, hyperparameters)

sprintf('Trimmed Model CV Train AUC: %.3f' ,trimmed_model$cv_auc)

## [1] "Trimmed Model CV Train AUC: 0.709"
sprintf('Trimmed Model Test AUC: %.3f' ,trimmed_model$auc)

## [1] "Trimmed Model Test AUC: 0.669"

Minutes to run: 0.368

current_features <- trimmed_features
current_model <- trimmed_model
current_least_important <- tail(trimmed_model$importance$Feature, 1)
current_auc_loss <- full_model$cv_auc - trimmed_model$cv_auc

selection_results <- tibble::tribble(
  ~`Number of Features`, ~`AUC Loss`, ~`Least Important Feature`,
  length(features), 0, tail(full_model$importance$Feature, 1),
  length(trimmed_features), current_auc_loss, tail(trimmed_model$importance$Feature, 1)
)

while (current_auc_loss < max_auc_loss){
  last_feature_dropped <- current_least_important

  current_features <- setdiff(current_features, current_least_important)
  hyperparameters$mtry = min(hyperparameters$mtry, length(current_features))
  current_model <- model_fit_wf(current_features, outcome_column, hyperparameters)
  current_least_important <- tail(current_model$importance$Feature, 1)

  current_auc_loss <- full_model$cv_auc - current_model$cv_auc

  selection_results <- selection_results %>%
    add_row(`Number of Features` = length(current_features),
           `AUC Loss` = current_auc_loss,
           `Least Important Feature` = current_least_important)
}

```

```

    print(c(length(current_features), current_auc_loss))
}

## [1] 74.0000000000 -0.0004080554
## [1] 7.300000e+01 5.272114e-04
## [1] 72.0000000000 -0.0002879659
## [1] 7.100000e+01 6.522527e-04
## [1] 70.0000000000 -0.0008374449
## [1] 6.900000e+01 8.210115e-04
## [1] 68.0000000000 0.001876782
## [1] 6.700000e+01 4.342318e-04
## [1] 6.60000e+01 6.97658e-06
## [1] 6.50000e+01 9.14889e-04
## [1] 64.0000000000 0.001345526
## [1] 63.0000000000 0.001071606
## [1] 62.0000000000 0.002434281
## [1] 61.0000000000 0.002920854
## [1] 60.0000000000 0.001210984
## [1] 59.0000000000 0.002204165
## [1] 58.0000000000 0.002171603
## [1] 57.0000000000 0.002239099
## [1] 56.0000000000 0.003367496
## [1] 55.0000000000 0.002466279
## [1] 54.0000000000 0.002258522
## [1] 53.0000000000 0.002611354
## [1] 52.0000000000 0.002523445
## [1] 51.0000000000 0.003778555
## [1] 50.0000000000 0.005056016
## [1] 49.0000000000 0.004610443
## [1] 48.0000000000 0.005081828
## [1] 47.0000000000 0.004731453
## [1] 46.0000000000 0.004219859
## [1] 45.0000000000 0.004796832
## [1] 44.0000000000 0.005467358
## [1] 43.00000000 0.00658657
## [1] 42.0000000000 0.006686176
## [1] 41.0000000000 0.007075376
## [1] 40.00000000 0.00553965
## [1] 39.0000000000 0.005999608
## [1] 38.0000000000 0.005854617
## [1] 37.0000000000 0.007453258
## [1] 36.0000000000 0.009395409
## [1] 35.00000000 0.01139753

selection_results

## # A tibble: 42 x 3
##   `Number of Features` `AUC Loss` `Least Important Feature`
##   <int>        <dbl> <chr>
## 1 84          0     icp
## 2 75          0.000263 icp
## 3 74         -0.000408 antiplaquetario_ev
## 4 73          0.000527 dialysis_t0
## 5 72         -0.000288 espiro_ergoespiro
## 6 71          0.000652 education_level
## 7 70         -0.000837 valvopathy
## 8 69          0.000821 transfusao
## 9 68          0.00188  cec
## 10 67         0.000434 cve_desf
## # ... with 32 more rows
## # i Use `print(n = ...)` to see more rows

```

Minutes to run: 13.816

```
selected_features <- c(current_features, last_feature_dropped)

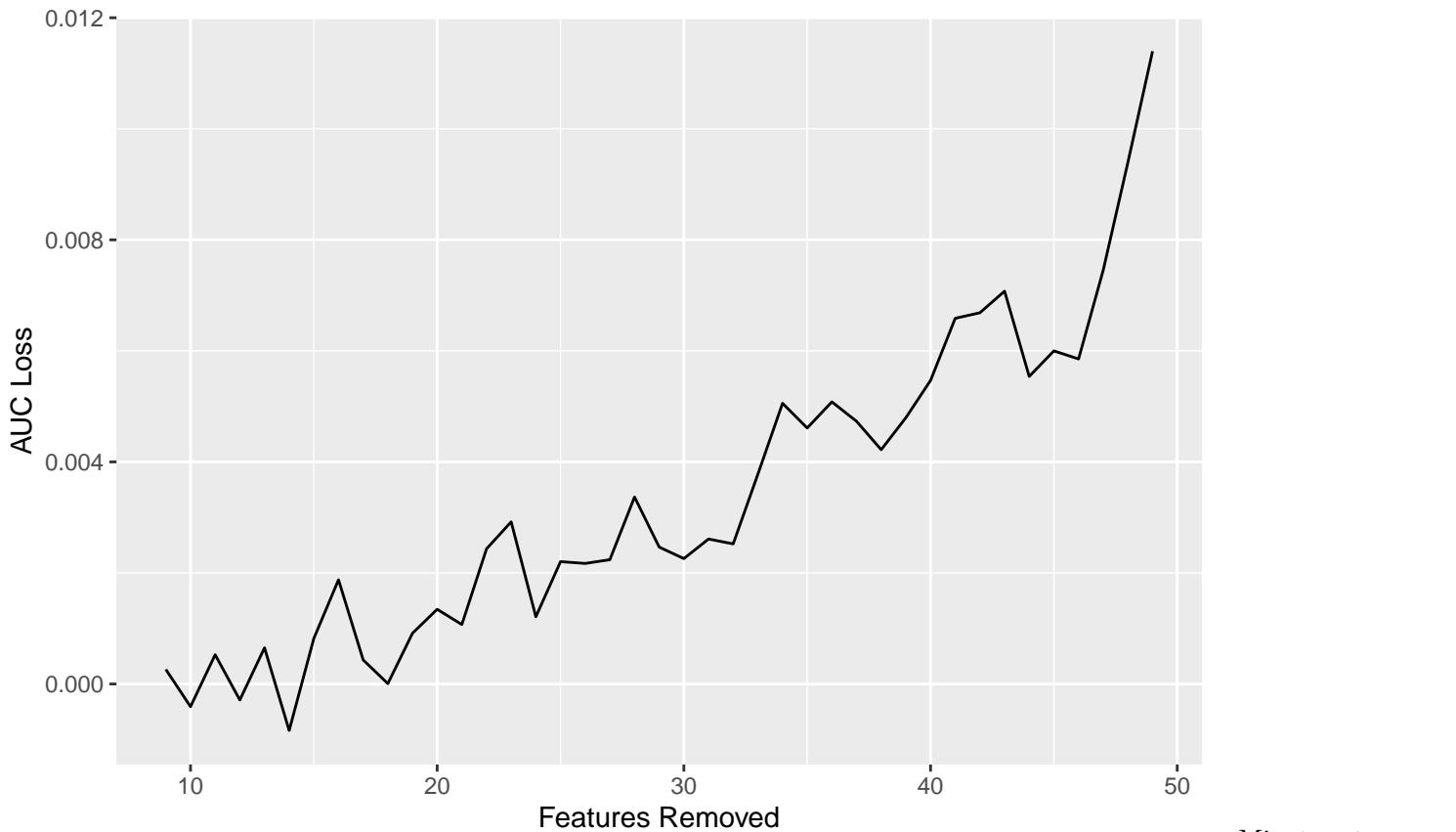
feature_selected_model <- model_fit_wf(selected_features,
                                         outcome_column, hyperparameters)

sprintf('Trimmed Model CV Train AUC: %.3f', feature_selected_model$cv_auc)

## [1] "Trimmed Model CV Train AUC: 0.700"
sprintf('Trimmed Model Test AUC: %.3f', feature_selected_model$auc)

## [1] "Trimmed Model Test AUC: 0.662"

selection_results %>%
  filter(`Number of Features` < length(features)) %>%
  mutate(`Features Removed` = length(features) - `Number of Features`) %>%
  ggplot(aes(x = `Features Removed`, y = `AUC Loss`)) +
  geom_line()
```



Minutes to run:

0.322

Hyperparameter tuning

```
lightgbm_recipe <-
  recipe(formula = sprintf("%s ~ .", outcome_column) %>% as.formula,
         data = df_train %>% dplyr::select(all_of(c(selected_features, outcome_column)))) %>%
  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(), one_hot = TRUE) %>%
  step_impute_mean(all_numeric_predictors()) %>%
  step_zv(all_predictors())
```

```

lightgbm_spec <- boost_tree(
  mtry = tune(),
  trees = tune(),
  min_n = tune(),
  tree_depth = tune(),
  learn_rate = tune(),
  loss_reduction = tune()
) %>%
  set_engine("lightgbm") %>%
  set_mode("classification")

lightgbm_grid <- grid_latin_hypercube(
  finalize(mtry()),
  df_train %>% dplyr::select(all_of(c(selected_features, outcome_column))), 
  dials::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)

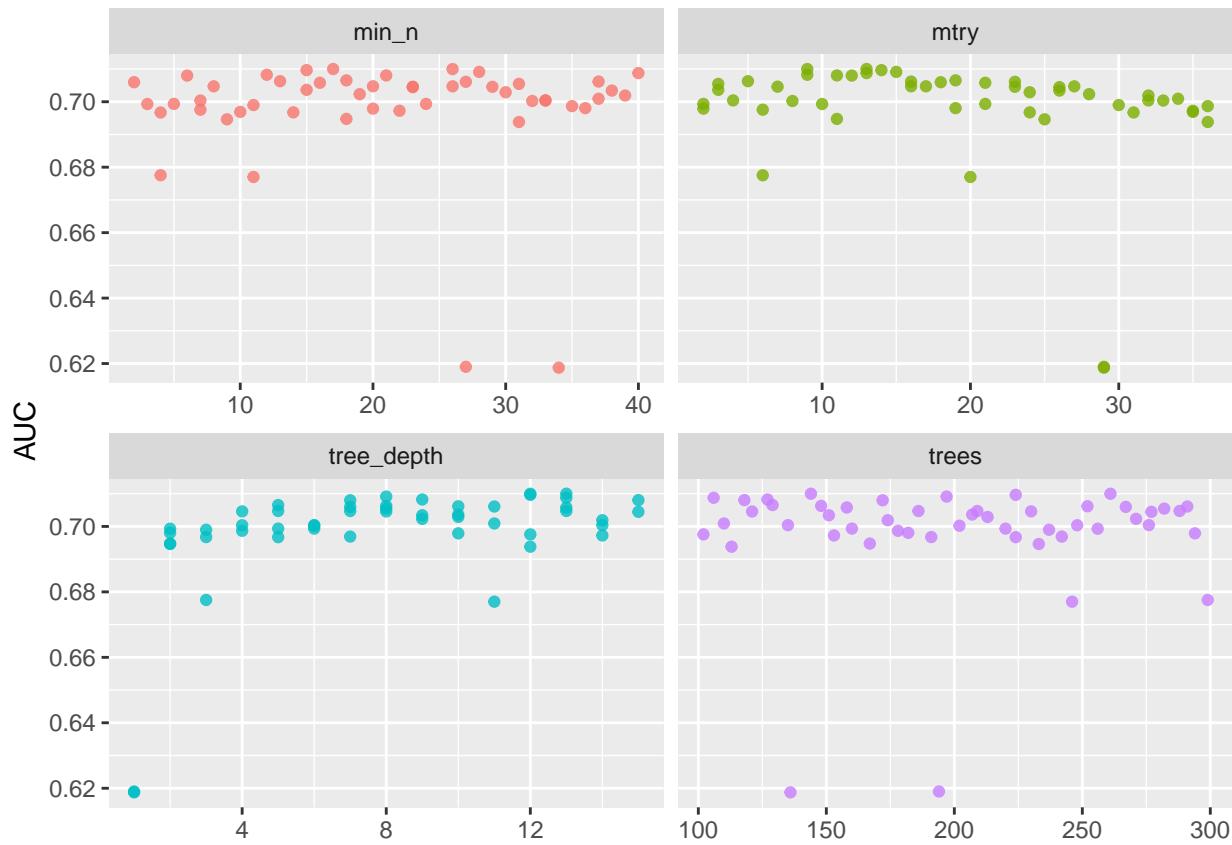
lightgbm_tune %>%
  show_best("roc_auc")

## # A tibble: 5 x 12
##   mtry trees min_n tree_depth learn_rate loss_~1 .metric .esti~2  mean      n std_err .config
##   <int> <int> <int>     <int>       <dbl>    <dbl> <chr>  <chr>    <dbl> <int> <dbl> <chr>
## 1    13    144    17        13    2.22e-4 4.91e-3 roc_auc binary  0.710     4  0.0171 Prepro-
## 2     9    261    26        12    1.58e-9 8.51e-6 roc_auc binary  0.710     4  0.0170 Prepro-
## 3    14    224    15        12    2.71e-5 3.62e-3 roc_auc binary  0.710     4  0.0164 Prepro-
## 4    15    197    28         8    5.17e-5 6.20e-4 roc_auc binary  0.709     4  0.0159 Prepro-
## 5    13    106    40        13    5.82e-4 6.43e-8 roc_auc binary  0.709     4  0.0160 Prepro-
## # ... with abbreviated variable names 1: loss_reduction, 2: .estimator

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

lightgbm_tune %>%
  collect_metrics() %>%
  filter(.metric == "roc_auc") %>%
  select(mean, mtry:tree_depth) %>%
  pivot_longer(mtry:tree_depth,
               values_to = "value",
               names_to = "parameter"
  ) %>%
  ggplot(aes(value, mean, color = parameter)) +
  geom_point(alpha = 0.8, show.legend = FALSE) +
  facet_wrap(~parameter, scales = "free_x") +
  labs(x = NULL, y = "AUC")

```



```

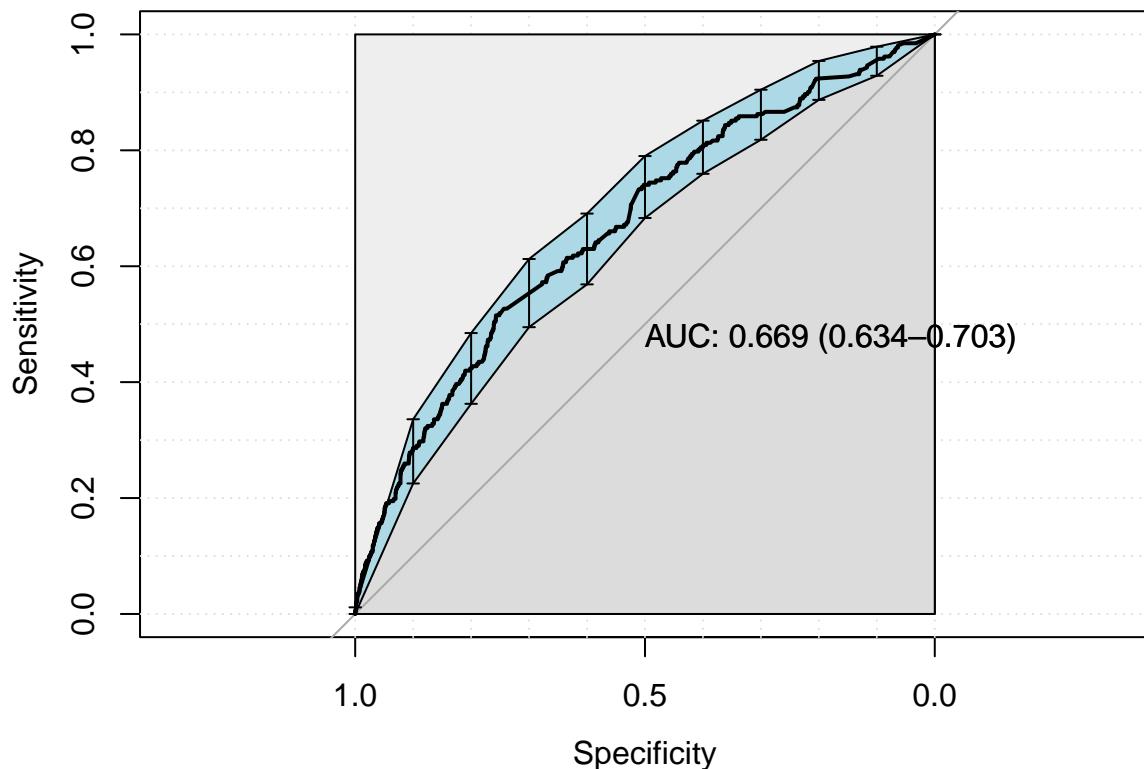
final_lightgbm_workflow <-
  lightgbm_workflow %>%
  finalize_workflow(best_lightgbm)

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)

```



```

## Confusion Matrix and Statistics
##
## test_predictions_class      0      1
##                      0 4468 262
##                      1     0     0
##
##          Accuracy : 0.9446
## 95% CI : (0.9377, 0.951)
## No Information Rate : 0.9446
## P-Value [Acc > NIR] : 0.5164
##
##          Kappa : 0
##
## Mcnemar's Test P-Value : <2e-16
##
##          Sensitivity : 1.0000
##          Specificity : 0.0000
## Pos Pred Value : 0.9446
## Neg Pred Value :    NaN
## Prevalence : 0.9446
## Detection Rate : 0.9446
## Detection Prevalence : 1.0000
## Balanced Accuracy : 0.5000
##
## 'Positive' Class : 0
##

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

```

Minutes to run: 5.605

```

lightgbm_model <- parsnip::extract_fit_engine(final_lightgbm_fit)

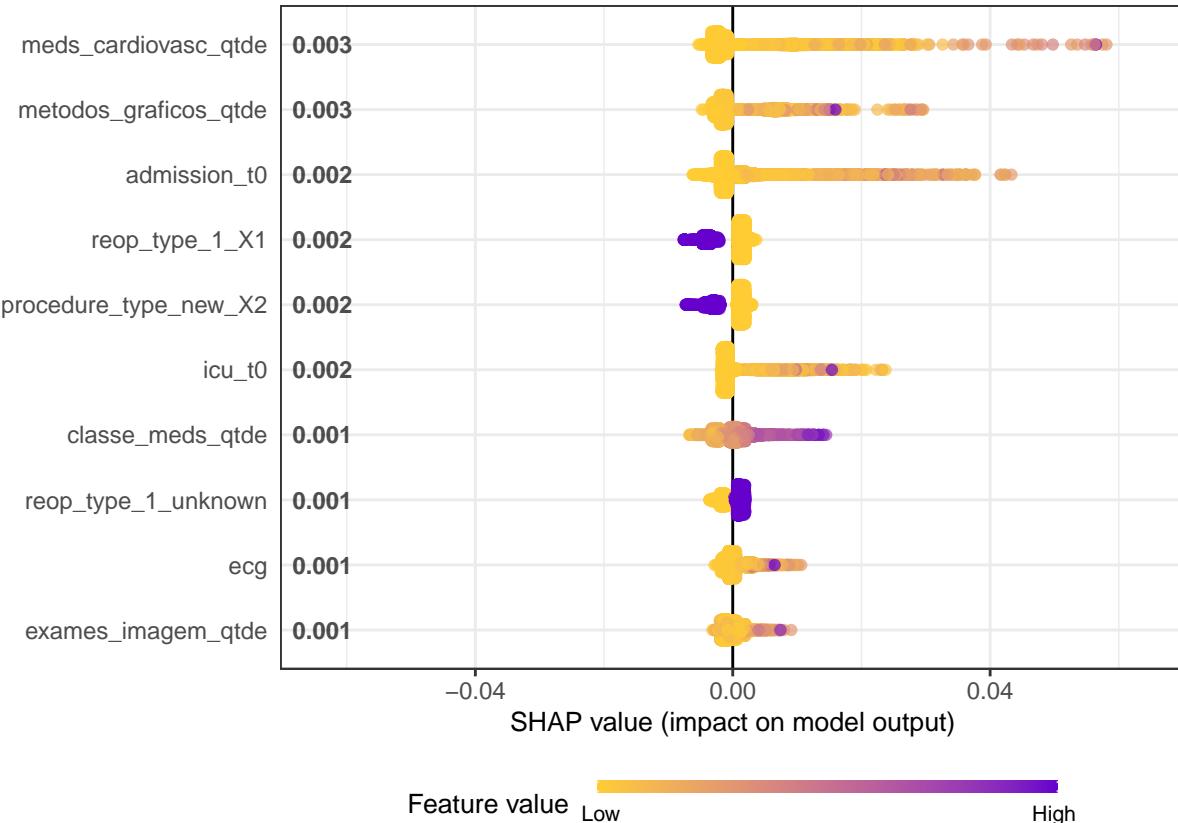
trained_rec <- prep(lightgbm_recipe, training = df_train)

df_train_baked <- bake(trained_rec, new_data = df_train)
df_test_baked <- bake(trained_rec, new_data = df_test)

matrix_train <- as.matrix(df_train_baked %>% select(-all_of(outcome_column)))
matrix_test <- as.matrix(df_test_baked %>% select(-all_of(outcome_column)))

shap.plot.summary.wrap1(model = lightgbm_model, X = matrix_train, top_n = 10, dilute = F)

```



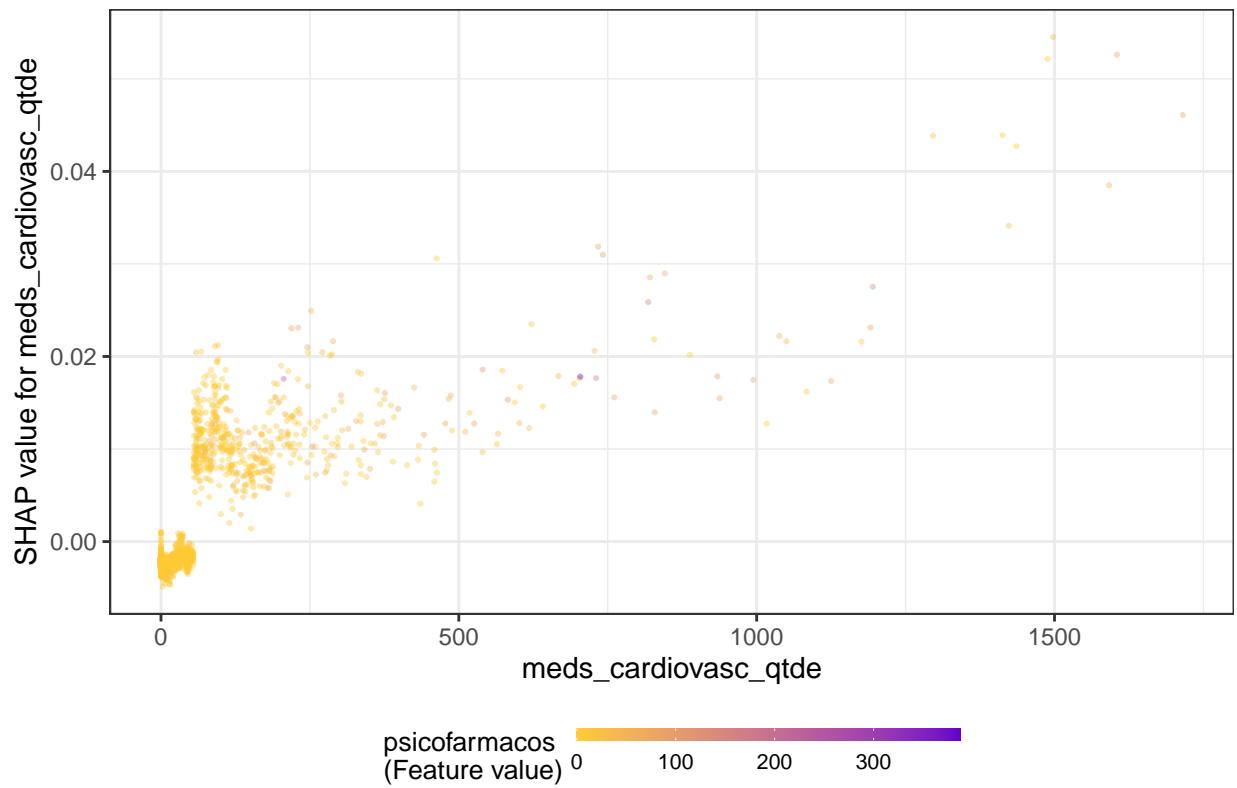
```

# Crunch SHAP values
shap <- shap.prep(lightgbm_model, X_train = matrix_test)

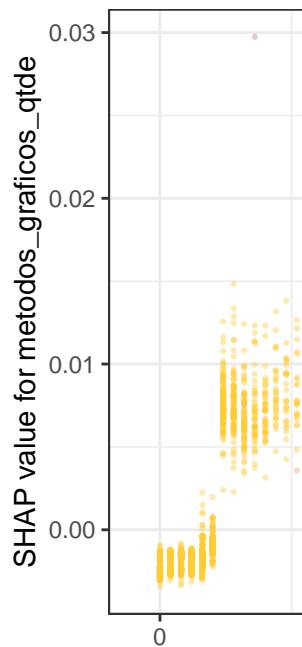
for (x in shap.importance(shap, names_only = TRUE)[1:5]) {
  p <- shap.plot.dependence(
    shap,
    x = x,
    color_feature = "auto",
    smooth = FALSE,
    jitter_width = 0.01,
    alpha = 0.4
  ) +
  ggtitle(x)
  print(p)
}

```

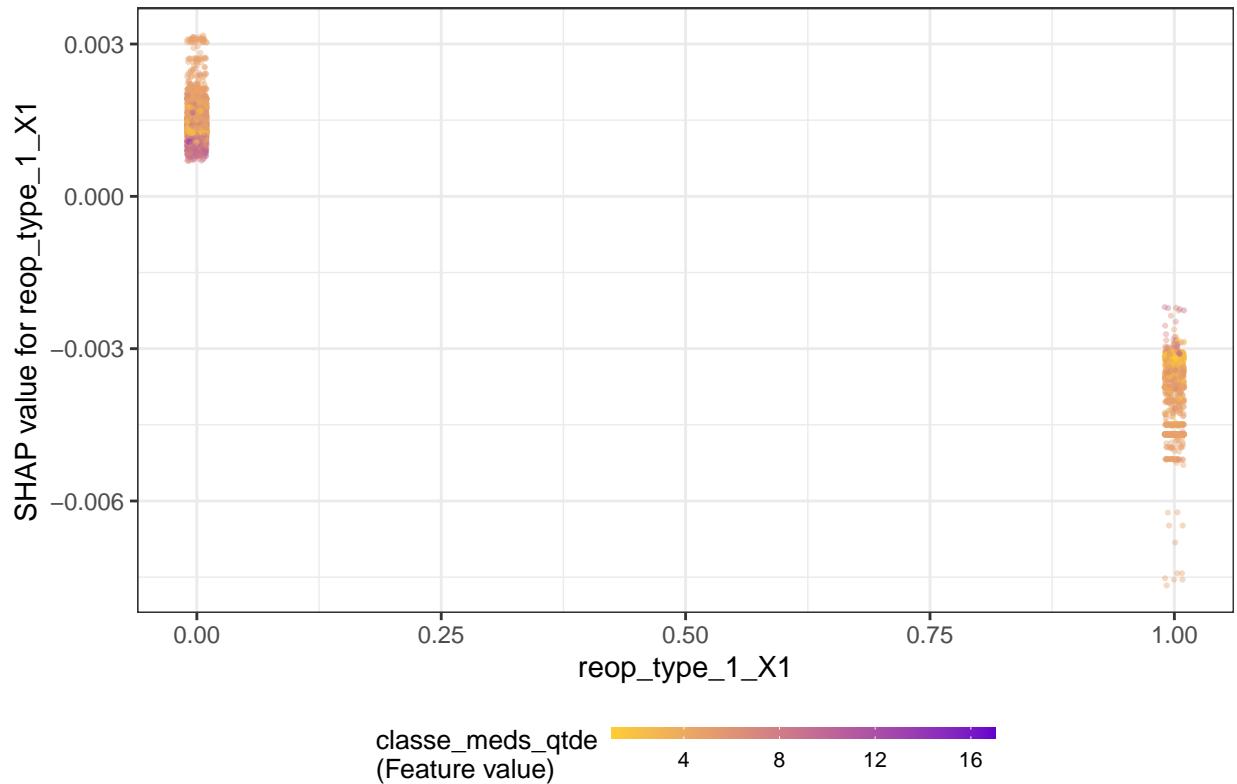
meds_cardiovasc_qtde



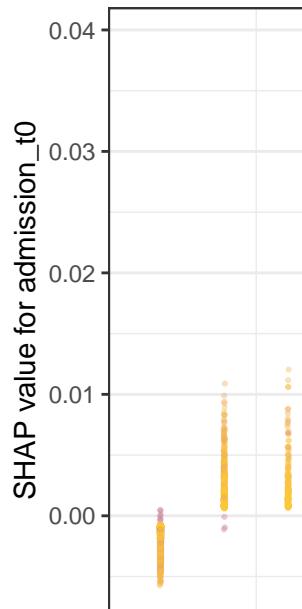
metodos_gra



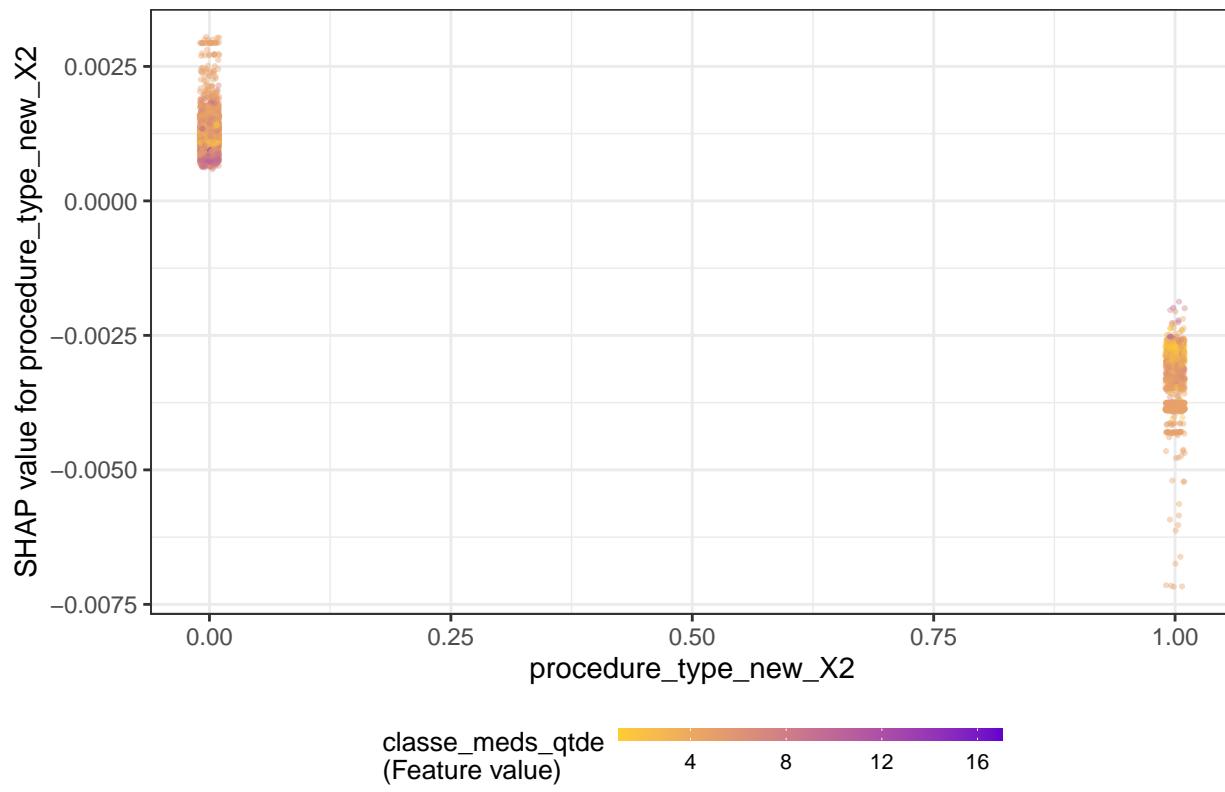
reop_type_1_X1



admission_t0



procedure_type_new_X2

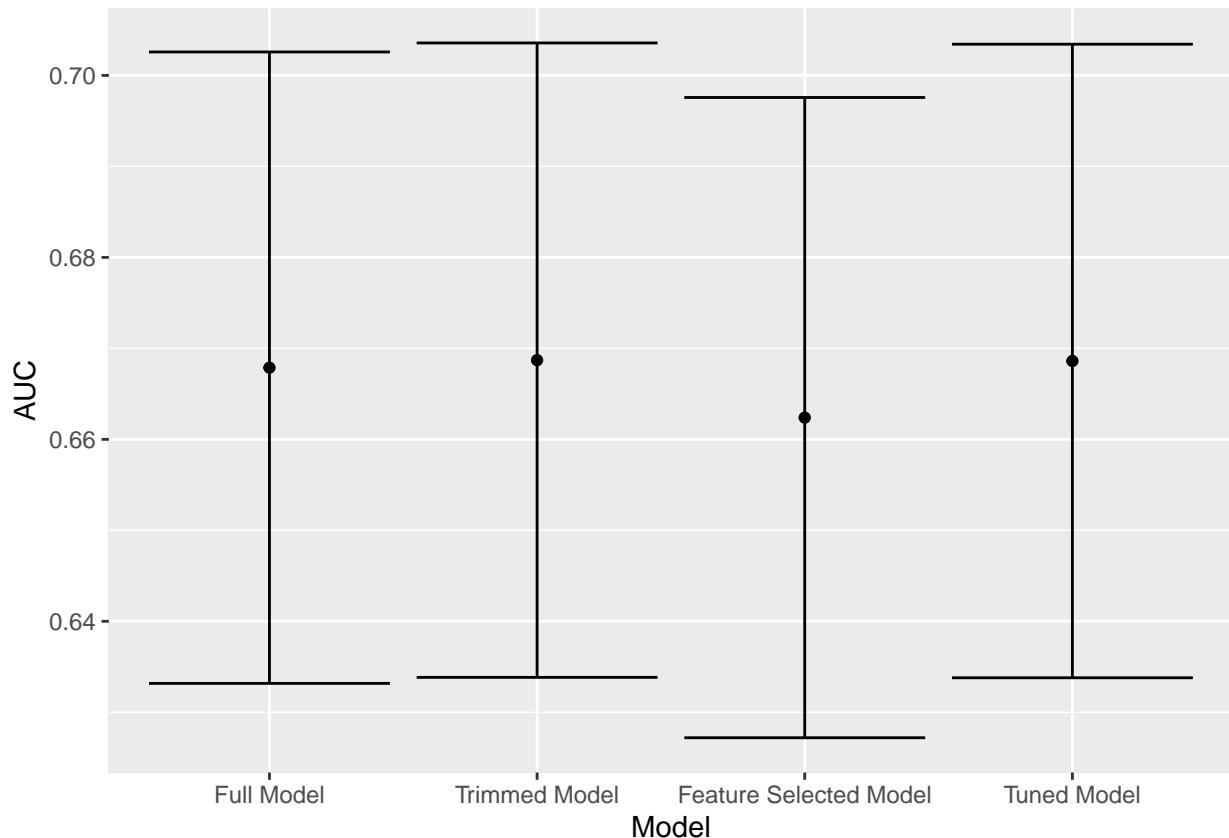


Minutes to run: 0

Models Comparison

```
df_auc <- tibble::tribble(
  ~Model, ~`AUC`, ~`Lower Limit`, ~`Upper Limit`,
  'Full Model', full_model$auc, full_model$auc_lower, full_model$auc_upper,
  'Trimmed Model', trimmed_model$auc, trimmed_model$auc_lower, trimmed_model$auc_upper,
  'Feature Selected Model', feature_selected_model$auc, feature_selected_model$auc_lower, feature_selected_model$auc_upper,
  'Tuned Model', as.numeric(lightgbm_auc$auc), lightgbm_auc$ci[1], lightgbm_auc$ci[3]
) %>%
  mutate(Target = outcome_column,
    Model = factor(Model,
      levels = c('Full Model', 'Trimmed Model',
      'Feature Selected Model', 'Tuned Model')))

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



```
saveRDS(df_auc, sprintf("../EDA/auxiliar/performance/tuning/%s_auc_result.RData", outcome_column))
```

Minutes to run: 0.002