

# Final Model - readmission\_60d

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
max_auc_loss <- params$max_auc_loss # Max accepted loss of AUC for reducing num of features
repeats <- params$repeats
Hmisc::list.tree(params)

##  params = list 5 (952 bytes)
## .  max_auc_loss = double 1= 0.01
## .  outcome_column = character 1= readmission_60d
## .  k = double 1= 10
## .  grid_size = double 1= 50
## .  repeats = double 1= 2
```

## Imports

```
library(tidyverse)
library(yaml)
library(tidymodels)
library(usemodels)
library(vip)
library(kableExtra)
library(SHAPforxgboost)
library(xgboost)
library(Matrix)
library(mltools)
library(bonsai)
library(lightgbm)
library(pROC)
library(caret)
library(themis)

source("aux_functions.R")

select <- dplyr::select
```

## 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[columns_list$outcome_columns] <- lapply(df[columns_list$outcome_columns], factor)
```

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

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

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

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

## 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)

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)
```

Starting features:

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

1. age
2. education\_level
3. underlying\_heart\_disease
4. heart\_disease
5. nyha\_basal
6. prior\_mi
7. heart\_failure

8. af
9. 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\_pre\_t0\_count
21. admission\_pre\_t0\_180d
22. icu\_t0
23. dialysis\_t0
24. admission\_t0\_emergency
25. aco
26. antiarritmico
27. betabloqueador
28. ieca\_bra
29. dva
30. digoxina
31. estatina
32. diuretico
33. vasodilatador
34. insuf\_cardiaca
35. espirolactona
36. bloq\_calcio
37. antiplaquetario\_ev
38. insulina
39. anticonvulsivante
40. psicofarmacos
41. antifungico
42. antiviral
43. classe\_meds\_qtde
44. meds\_cardiovasc\_qtde
45. meds\_antimicrobianos
46. ventilacao\_mecanica
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. interconsulta
60. equipe\_multiprof
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. ressonancia
82. exames\_imagem\_qtde
83. bic
84. hospital\_stay

## 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) %>% select(all_of(c(features, outcome_column)))
df_test <- testing(df_split) %>% select(all_of(c(features, outcome_column)))

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

## Feature Selection

```
custom_dummy_names <- function(var, lvl, ordinal = FALSE) {
  dummy_names(var, lvl, ordinal = FALSE, sep = "__")
}

model_fit_wf <- function(df_train, 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_dummy(all_nominal_predictors(), naming = custom_dummy_names)

  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) %>%
  separate(Feature, c("Feature", "value"), "__", fill = 'right') %>%
  group_by(Feature) %>%
  summarise(Gain = sum(Gain),
            Cover = sum(Cover),
            Frequency = sum(Frequency)) %>%
  ungroup() %>%
  arrange(desc(Gain))

cv_results <- collect_metrics(model_fit_rs) %>% filter(.metric == 'roc_auc')

return(
  list(
    cv_auc = cv_results$mean,
    cv_auc_std_err = cv_results$std_err,
    importance = feature_importance,
    auc = as.numeric(model_auc$auc),
    auc_lower = model_auc$ci[1],
    auc_upper = model_auc$ci[3]
  )
)
}

```

```

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

hyperparameters$sample_size <- 1

full_model <- model_fit_wf(df_train, features, outcome_column, hyperparameters)

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

## [1] "Full Model CV Train AUC: 0.711"

sprintf('Full Model Test AUC: %.3f' ,full_model$auc)

```

```
## [1] "Full Model Test AUC: 0.672"
```

Features with zero importance on the initial model:

```
unimportant_features <- setdiff(features, full_model$importance$Feature)
```

```
unimportant_features %>%
  gluedown::md_order()
```

1. transplant
2. procedure\_type\_1
3. dialysis\_t0
4. cec

5. transplante\_cardiaco
6. cir\_toracica
7. teste\_esforco
8. espiro\_ergoespiro
9. tilt\_teste
10. biopsia
11. arteriografia
12. pet\_ct
13. ressonancia

```

trimmed_features <- full_model$importance$Feature
trimmed_model <- model_fit_wf(df_train, trimmed_features,
                             outcome_column, hyperparameters)

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

## [1] "Trimmed Model CV Train AUC: 0.712"

sprintf('Trimmed Model Test AUC: %.3f' ,trimmed_model$auc)

## [1] "Trimmed Model Test AUC: 0.672"

selection_results <- tibble::tribble(
  ~`Tested Feature`, ~`Dropped`, ~`Number of Features`, ~`CV AUC`, ~`CV AUC Std Error`, ~`Total AUC Loss`, ~`Instant AUC Loss`
  'None', TRUE, length(features), full_model$cv_auc, full_model$cv_auc_std_err, 0, 0
)

whitelist <- c()

if (full_model$cv_auc - trimmed_model$cv_auc < max_auc_loss) {
  current_features <- trimmed_features
  current_model <- trimmed_model
  current_auc_loss <- full_model$cv_auc - current_model$cv_auc
  instant_auc_loss <- full_model$cv_auc - current_model$cv_auc

  selection_results <- selection_results %>%
    add_row(`Tested Feature` = 'All unimportant',
            `Dropped` = TRUE,
            `Number of Features` = length(trimmed_features),
            `CV AUC` = current_model$cv_auc,
            `CV AUC Std Error` = current_model$cv_auc_std_err,
            `Total AUC Loss` = current_auc_loss,
            `Instant AUC Loss` = instant_auc_loss)
} else {
  current_features <- features
  current_model <- full_model
  current_auc_loss <- 0
}

while (current_auc_loss < max_auc_loss & mean(current_features %in% whitelist) < 1) {
  zero_importance_features <-
    setdiff(current_features, current_model$importance$Feature) %>%
    setdiff(whitelist)
  if (length(zero_importance_features) > 0) {
    current_least_important <- zero_importance_features[1]
  } else {
    current_least_important <-
      tail(setdiff(current_model$importance$Feature, whitelist), 1)
  }
  test_features <-
    setdiff(current_features, current_least_important)
  current_model <-
    model_fit_wf(df_train, test_features, outcome_column, hyperparameters)
}

```

```

instant_auc_loss <-
  tail(selection_results %>% filter(Dropped) %>% .$`CV AUC`, n = 1) - current_model$cv_auc

if (instant_auc_loss < max_auc_loss / 5 &
    current_auc_loss < max_auc_loss) {
  dropped <- TRUE
  current_features <- test_features
  current_auc_loss <- full_model$cv_auc - current_model$cv_auc
} else {
  dropped <- FALSE
  whitelist <- c(whitelist, current_least_important)
}

selection_results <- selection_results %>%
  add_row(
    `Tested Feature` = current_least_important,
    `Dropped` = dropped,
    `Number of Features` = length(test_features),
    `CV AUC` = current_model$cv_auc,
    `CV AUC Std Error` = current_model$cv_auc_std_err,
    `Total AUC Loss` = current_auc_loss,
    `Instant AUC Loss` = instant_auc_loss
  )

print(c(
  length(current_features),
  round(current_auc_loss, 4),
  round(instant_auc_loss, 4),
  current_least_important
))
}

## [1] "70"          "-0.0011"      "2e-04"        "reop_type_1"
## [1] "69"          "-7e-04"       "4e-04"        "angioplastia"
## [1] "68"          "-0.001"       "-3e-04"       "angio_tc"
## [1] "67"          "-0.0011"      "-1e-04"       "cateter_venoso_central"
## [1] "66"          "-9e-04"       "2e-04"        "valvopathy"
## [1] "65"          "-7e-04"       "2e-04"        "prior_mi"
## [1] "64"          "-0.001"       "-3e-04"       "antiplaquetario_ev"
## [1] "63"          "-0.0012"      "-2e-04"       "insulina"
## [1] "62"          "-6e-04"       "6e-04"        "holter"
## [1] "61"          "0"            "6e-04"        "heart_failure"
## [1] "60"          "3e-04"        "2e-04"        "cardiac_arrest"
## [1] "59"          "4e-04"        "1e-04"        "analises_clinicas_qtde"
## [1] "58"          "-1e-04"       "-4e-04"       "endoscopia"
## [1] "57"          "-9e-04"       "-8e-04"       "betabloqueador"
## [1] "56"          "-8e-04"       "1e-04"        "cve_desf"
## [1] "55"          "-3e-04"       "5e-04"        "hemodialysis"
## [1] "54"          "-1e-04"       "2e-04"        "antiviral"
## [1] "53"          "1e-04"        "2e-04"        "tomografia"
## [1] "52"          "-4e-04"       "-5e-04"       "antifungico"
## [1] "51"          "-2e-04"       "2e-04"        "cultura"
## [1] "50"          "7e-04"        "9e-04"        "ventilacao_mecanica"
## [1] "49"          "-6e-04"       "-0.0012"      "ultrassom"
## [1] "48"          "-7e-04"       "-1e-04"       "insuf_cardiaca"
## [1] "47"          "-0.001"       "-3e-04"       "cateterismo"
## [1] "46"          "-9e-04"       "1e-04"        "procedure_type_new"
## [1] "45"          "-4e-04"       "5e-04"        "ecocardiograma"
## [1] "44"          "-0.0022"      "-0.0018"      "digoxina"
## [1] "43"          "-0.0018"      "4e-04"        "interconsulta"
## [1] "42"          "-0.0012"      "6e-04"        "histopatologia_qtde"

```

```

## [1] "41"      "-0.0017" "-5e-04" "icp"
## [1] "40"      "-7e-04" "0.001" "af"
## [1] "39"      "-0.001" "-3e-04" "eletrofisiologia"
## [1] "38"      "-0.0012" "-2e-04" "cintilografia"
## [1] "37"      "-0.0012" "0" "citologia"
## [1] "36"      "-0.001" "2e-04" "transfusao"
## [1] "35"      "-0.0014" "-4e-04"
## [4] "outros_proced_cirurgicos"
## [1] "34"      "-3e-04" "0.0011" "cied_final_group_1"
## [1] "33"      "-4e-04" "-1e-04" "aco"
## [1] "32"      "-0.0015" "-0.0012" "heart_disease"
## [1] "31"      "-3e-04" "0.0012" "nyha_basal"
## [1] "30"      "-8e-04" "-5e-04" "bic"
## [1] "29"      "-0.0015" "-6e-04" "espironolactona"
## [1] "28"      "-0.0027" "-0.0012" "angio_rm"
## [1] "27"      "-0.0022" "5e-04" "admission_pre_t0_180d"
## [1] "26"      "-0.0019" "3e-04"
## [4] "underlying_heart_disease"
## [1] "25"      "-0.0025" "-5e-04" "diabetes"
## [1] "24"      "-0.002" "5e-04" "bloq_calcio"
## [1] "23"      "-0.0022" "-2e-04" "equipe_multiprof"
## [1] "22"      "-0.0029" "-7e-04" "diuretico"
## [1] "21"      "-0.0029" "-1e-04" "anticonvulsivante"
## [1] "20"      "-0.0036" "-7e-04" "psicofarmacos"
## [1] "19"      "-0.0044" "-8e-04" "admission_t0_emergency"
## [1] "18"      "-0.0026" "0.0018" "cied_final_1"
## [1] "17"      "-0.0025" "1e-04" "antiarritmico"
## [1] "16"      "-0.0035" "-0.001" "comorbidities_count"
## [1] "15"      "-0.0039" "-4e-04" "exames_imagem_qtde"
## [1] "14"      "-0.0051" "-0.0012" "dva"
## [1] "13"      "-0.0051" "0" "meds_cardiovasc_qtde"
## [1] "12"      "-0.0057" "-6e-04" "meds_antimicrobianos"
## [1] "12"      "-0.0057" "0.0026" "vasodilatador"
## [1] "11"      "-0.0049" "7e-04" "education_level"
## [1] "10"      "-0.0037" "0.0013" "proced_invasivos_qtde"
## [1] "9"      "-0.003" "7e-04" "estatina"
## [1] "8"      "-0.0053" "-0.0024" "icu_t0"
## [1] "7"      "-0.0056" "-2e-04" "metodos_graficos_qtde"
## [1] "6"      "-0.0038" "0.0017" "ieca_bra"
## [1] "6"      "-0.0038" "0.0101" "classe_meds_qtde"
## [1] "5"      "-0.0051" "-0.0013" "age"
## [1] "5"      "-0.0051" "0.0152" "admission_pre_t0_count"
## [1] "4"      "-0.0041" "0.001" "laboratorio"
## [1] "4"      "-0.0041" "0.0676" "hospital_stay"

```

```

selection_results %>%
  rename(Features = `Number of Features`) %>%
  niceFormatting(digits = 4, label = 1)

```

Table 1:

| Tested Feature         | Dropped | Features | CV AUC | CV AUC Std Error | Total AUC Loss | Instant AUC Loss |
|------------------------|---------|----------|--------|------------------|----------------|------------------|
| None                   | TRUE    | 84       | 0.7105 | 0.0088           | 0.0000         | 0.0000           |
| All unimportant        | TRUE    | 71       | 0.7118 | 0.0086           | -0.0013        | -0.0013          |
| reop_type_1            | TRUE    | 70       | 0.7116 | 0.0085           | -0.0011        | 0.0002           |
| angioplastia           | TRUE    | 69       | 0.7112 | 0.0084           | -0.0007        | 0.0004           |
| angio_tc               | TRUE    | 68       | 0.7115 | 0.0084           | -0.0010        | -0.0003          |
| cateter_venoso_central | TRUE    | 67       | 0.7116 | 0.0084           | -0.0011        | -0.0001          |
| valvopathy             | TRUE    | 66       | 0.7114 | 0.0084           | -0.0009        | 0.0002           |
| prior_mi               | TRUE    | 65       | 0.7112 | 0.0084           | -0.0007        | 0.0002           |
| antiplaquetario_ev     | TRUE    | 64       | 0.7115 | 0.0084           | -0.0010        | -0.0003          |



Table 1: (continued)

| Tested Feature           | Dropped | Features | CV AUC | CV AUC Std Error | Total AUC Loss | Instant AUC Loss |
|--------------------------|---------|----------|--------|------------------|----------------|------------------|
| insulina                 | TRUE    | 63       | 0.7118 | 0.0083           | -0.0012        | -0.0002          |
| holter                   | TRUE    | 62       | 0.7111 | 0.0084           | -0.0006        | 0.0006           |
| heart_failure            | TRUE    | 61       | 0.7105 | 0.0084           | 0.0000         | 0.0006           |
| cardiac_arrest           | TRUE    | 60       | 0.7103 | 0.0084           | 0.0003         | 0.0002           |
| analises_clinicas_qtde   | TRUE    | 59       | 0.7102 | 0.0084           | 0.0004         | 0.0001           |
| endoscopia               | TRUE    | 58       | 0.7106 | 0.0086           | -0.0001        | -0.0004          |
| betabloqueador           | TRUE    | 57       | 0.7114 | 0.0085           | -0.0009        | -0.0008          |
| cve_desf                 | TRUE    | 56       | 0.7113 | 0.0085           | -0.0008        | 0.0001           |
| hemodialysis             | TRUE    | 55       | 0.7108 | 0.0084           | -0.0003        | 0.0005           |
| antiviral                | TRUE    | 54       | 0.7106 | 0.0084           | -0.0001        | 0.0002           |
| tomografia               | TRUE    | 53       | 0.7104 | 0.0085           | 0.0001         | 0.0002           |
| antifungico              | TRUE    | 52       | 0.7109 | 0.0084           | -0.0004        | -0.0005          |
| cultura                  | TRUE    | 51       | 0.7107 | 0.0085           | -0.0002        | 0.0002           |
| ventilacao_mecanica      | TRUE    | 50       | 0.7098 | 0.0084           | 0.0007         | 0.0009           |
| ultrassom                | TRUE    | 49       | 0.7111 | 0.0084           | -0.0006        | -0.0012          |
| insuf_cardiaca           | TRUE    | 48       | 0.7112 | 0.0086           | -0.0007        | -0.0001          |
| cateterismo              | TRUE    | 47       | 0.7115 | 0.0085           | -0.0010        | -0.0003          |
| procedure_type_new       | TRUE    | 46       | 0.7114 | 0.0087           | -0.0009        | 0.0001           |
| ecocardiograma           | TRUE    | 45       | 0.7109 | 0.0084           | -0.0004        | 0.0005           |
| digoxina                 | TRUE    | 44       | 0.7127 | 0.0085           | -0.0022        | -0.0018          |
| interconsulta            | TRUE    | 43       | 0.7123 | 0.0085           | -0.0018        | 0.0004           |
| histopatologia_qtde      | TRUE    | 42       | 0.7117 | 0.0085           | -0.0012        | 0.0006           |
| icp                      | TRUE    | 41       | 0.7122 | 0.0085           | -0.0017        | -0.0005          |
| af                       | TRUE    | 40       | 0.7112 | 0.0085           | -0.0007        | 0.0010           |
| eletrofisiologia         | TRUE    | 39       | 0.7115 | 0.0084           | -0.0010        | -0.0003          |
| cintilografia            | TRUE    | 38       | 0.7117 | 0.0087           | -0.0012        | -0.0002          |
| citologia                | TRUE    | 37       | 0.7117 | 0.0086           | -0.0012        | 0.0000           |
| transfusao               | TRUE    | 36       | 0.7116 | 0.0087           | -0.0010        | 0.0002           |
| outros_proced_cirurgicos | TRUE    | 35       | 0.7119 | 0.0088           | -0.0014        | -0.0004          |
| cied_final_group_1       | TRUE    | 34       | 0.7108 | 0.0084           | -0.0003        | 0.0011           |
| aco                      | TRUE    | 33       | 0.7109 | 0.0085           | -0.0004        | -0.0001          |
| heart_disease            | TRUE    | 32       | 0.7120 | 0.0085           | -0.0015        | -0.0012          |
| nyha_basal               | TRUE    | 31       | 0.7108 | 0.0085           | -0.0003        | 0.0012           |
| bic                      | TRUE    | 30       | 0.7114 | 0.0086           | -0.0008        | -0.0005          |
| espironolactona          | TRUE    | 29       | 0.7120 | 0.0086           | -0.0015        | -0.0006          |
| angio_rm                 | TRUE    | 28       | 0.7132 | 0.0084           | -0.0027        | -0.0012          |
| admission_pre_t0_180d    | TRUE    | 27       | 0.7127 | 0.0090           | -0.0022        | 0.0005           |
| underlying_heart_disease | TRUE    | 26       | 0.7125 | 0.0087           | -0.0019        | 0.0003           |
| diabetes                 | TRUE    | 25       | 0.7130 | 0.0088           | -0.0025        | -0.0005          |
| bloq_calcio              | TRUE    | 24       | 0.7125 | 0.0088           | -0.0020        | 0.0005           |
| equipe_multiprof         | TRUE    | 23       | 0.7127 | 0.0091           | -0.0022        | -0.0002          |
| diuretico                | TRUE    | 22       | 0.7134 | 0.0090           | -0.0029        | -0.0007          |
| anticonvulsivante        | TRUE    | 21       | 0.7134 | 0.0090           | -0.0029        | -0.0001          |
| psicofarmacos            | TRUE    | 20       | 0.7141 | 0.0090           | -0.0036        | -0.0007          |
| admission_t0_emergency   | TRUE    | 19       | 0.7149 | 0.0089           | -0.0044        | -0.0008          |
| cied_final_1             | TRUE    | 18       | 0.7131 | 0.0092           | -0.0026        | 0.0018           |
| antiarritmico            | TRUE    | 17       | 0.7130 | 0.0094           | -0.0025        | 0.0001           |
| comorbidities_count      | TRUE    | 16       | 0.7140 | 0.0095           | -0.0035        | -0.0010          |
| exames_imagem_qtde       | TRUE    | 15       | 0.7144 | 0.0094           | -0.0039        | -0.0004          |
| dva                      | TRUE    | 14       | 0.7156 | 0.0095           | -0.0051        | -0.0012          |
| meds_cardiovasc_qtde     | TRUE    | 13       | 0.7156 | 0.0094           | -0.0051        | 0.0000           |
| meds_antimicrobianos     | TRUE    | 12       | 0.7162 | 0.0096           | -0.0057        | -0.0006          |
| vasodilatador            | FALSE   | 11       | 0.7136 | 0.0093           | -0.0057        | 0.0026           |

Table 1: (continued)

| Tested Feature         | Dropped | Features | CV AUC | CV AUC Std Error | Total AUC Loss | Instant AUC Loss |
|------------------------|---------|----------|--------|------------------|----------------|------------------|
| education_level        | TRUE    | 11       | 0.7154 | 0.0091           | -0.0049        | 0.0007           |
| proced_invasivos_qtde  | TRUE    | 10       | 0.7142 | 0.0092           | -0.0037        | 0.0013           |
| estatina               | TRUE    | 9        | 0.7135 | 0.0090           | -0.0030        | 0.0007           |
| icu_t0                 | TRUE    | 8        | 0.7159 | 0.0091           | -0.0053        | -0.0024          |
| metodos_graficos_qtde  | TRUE    | 7        | 0.7161 | 0.0090           | -0.0056        | -0.0002          |
| ieca_bra               | TRUE    | 6        | 0.7144 | 0.0098           | -0.0038        | 0.0017           |
| classe_meds_qtde       | FALSE   | 5        | 0.7042 | 0.0095           | -0.0038        | 0.0101           |
| age                    | TRUE    | 5        | 0.7156 | 0.0097           | -0.0051        | -0.0013          |
| admission_pre_t0_count | FALSE   | 4        | 0.7004 | 0.0093           | -0.0051        | 0.0152           |
| laboratorio            | TRUE    | 4        | 0.7146 | 0.0095           | -0.0041        | 0.0010           |
| hospital_stay          | FALSE   | 3        | 0.6470 | 0.0091           | -0.0041        | 0.0676           |

```

selected_features <- current_features

con <- file(sprintf('./auxiliar/final_model/selected_features/%s.yaml', outcome_column), "w")
write_yaml(selected_features, con)
close(con)

feature_selected_model <- model_fit_wf(df_train, selected_features,
                                       outcome_column, hyperparameters)

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

## [1] "Selected Model CV Train AUC: 0.715"

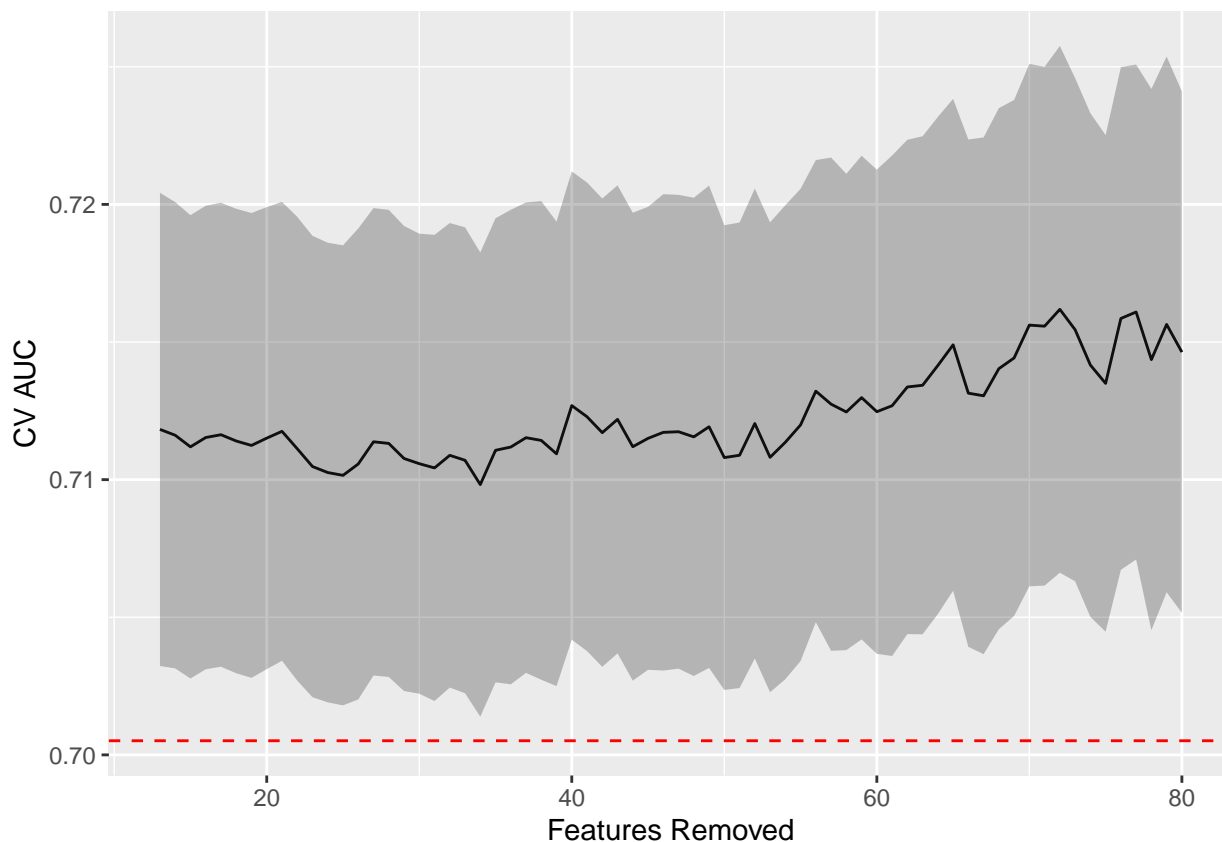
sprintf('Selected Model Test AUC: %.3f', feature_selected_model$auc)

## [1] "Selected Model Test AUC: 0.670"

selection_results <- selection_results %>%
  filter(`Number of Features` < length(features)) %>%
  mutate(`Features Removed` = length(features) - `Number of Features`,
         `CV AUC Low` = `CV AUC` - `CV AUC Std Error`,
         `CV AUC High` = `CV AUC` + `CV AUC Std Error`)

selection_results %>%
  filter(Dropped) %>%
  ggplot(aes(x = `Features Removed`, y = `CV AUC`,
             ymin = `CV AUC Low`, ymax = `CV AUC High`)) +
  geom_line() +
  geom_ribbon(alpha = .3) +
  geom_hline(yintercept = full_model$cv_auc - max_auc_loss,
            linetype = "dashed", color = "red")

```



## Hyperparameter tuning

Selected features:

```
gluedown::md_order(selected_features, seq = TRUE, pad = TRUE)
```

1. hospital\_stay
2. admission\_pre\_t0\_count
3. classe\_meds\_qtde
4. vasodilatador

## Standard

```
lightgbm_recipe <-
  recipe(formula = sprintf("%s ~ .", outcome_column) %>% as.formula,
    data = df_train %>% select(all_of(c(selected_features, outcome_column)))) %>%
  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())

lightgbm_tuning <- function(recipe) {

  lightgbm_spec <- boost_tree(
    trees = tune(),
    min_n = tune(),
    tree_depth = tune(),
    learn_rate = tune(),
    sample_size = 1.0
  ) %>%
  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(2L, 15L)),
  learn_rate(range = c(-3, -1), trans = log10_trans()),
  size = grid_size
)

lightgbm_workflow <-
  workflow() %>%
  add_recipe(recipe) %>%
  add_model(lightgbm_spec)

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

lightgbm_tune %>%
  show_best("roc_auc") %>%
  niceFormatting(digits = 5, label = 4)

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)

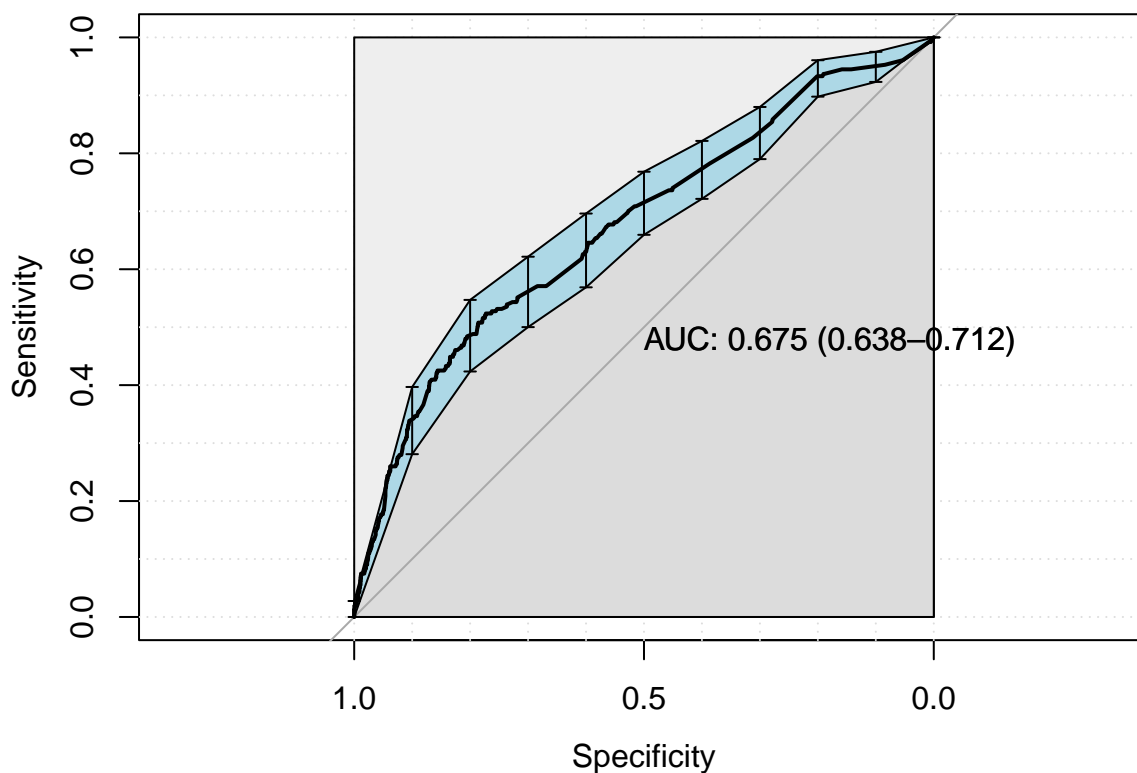
lightgbm_auc <- validation(final_lightgbm_fit, df_test)

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

return(list(auc = as.numeric(lightgbm_auc$auc),
            auc_lower = lightgbm_auc$ci[1],
            auc_upper = lightgbm_auc$ci[3],
            parameters = lightgbm_parameters,
            fit = final_lightgbm_fit))
}

standard_results <- lightgbm_tuning(lightgbm_recipe)

```



```
## [1] "Optimal Threshold: 0.07"
## Confusion Matrix and Statistics
##
##      reference
## data    0    1
##    0 3460 121
##    1 1016 133
##
##              Accuracy : 0.7596
##              95% CI   : (0.7472, 0.7717)
##    No Information Rate : 0.9463
##    P-Value [Acc > NIR] : 1
##
##              Kappa   : 0.1114
##
##  McNemar's Test P-Value : <2e-16
##
##              Sensitivity : 0.7730
##              Specificity : 0.5236
##              Pos Pred Value : 0.9662
##              Neg Pred Value : 0.1158
##              Prevalence : 0.9463
##              Detection Rate : 0.7315
##              Detection Prevalence : 0.7571
##              Balanced Accuracy : 0.6483
##
##              'Positive' Class : 0
##
final_lightgbm_fit <- standard_results$fit
lightgbm_parameters <- standard_results$parameters

saveRDS(
  lightgbm_parameters,
```

```

file = sprintf(
  "./auxiliar/final_model/hyperparameters/lightgbm_%s.rds",
  outcome_column
)
)

# Save the final model. We need it for the calculator
lgb.save(
  parsnip::extract_fit_engine(final_lightgbm_fit),
  sprintf("./results/%s/final_model.txt", outcome_column)
)
saveRDS(final_lightgbm_fit,
  sprintf("./results/%s/final_model_wf.rds", outcome_column))

```

## SHAP values

```

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)))

n_plots <- min(6, length(selected_features))
plotted <- 0

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

shap <- shap.prep(lightgbm_model, X_train = matrix_test)

for (x in shap.importance(shap, names_only = TRUE)) {
  p <- shap.plot.dependence(
    shap,
    x = x,
    color_feature = "auto",
    smooth = FALSE,
    jitter_width = 0.01,
    alpha = 0.3
  ) +
    labs(title = x)

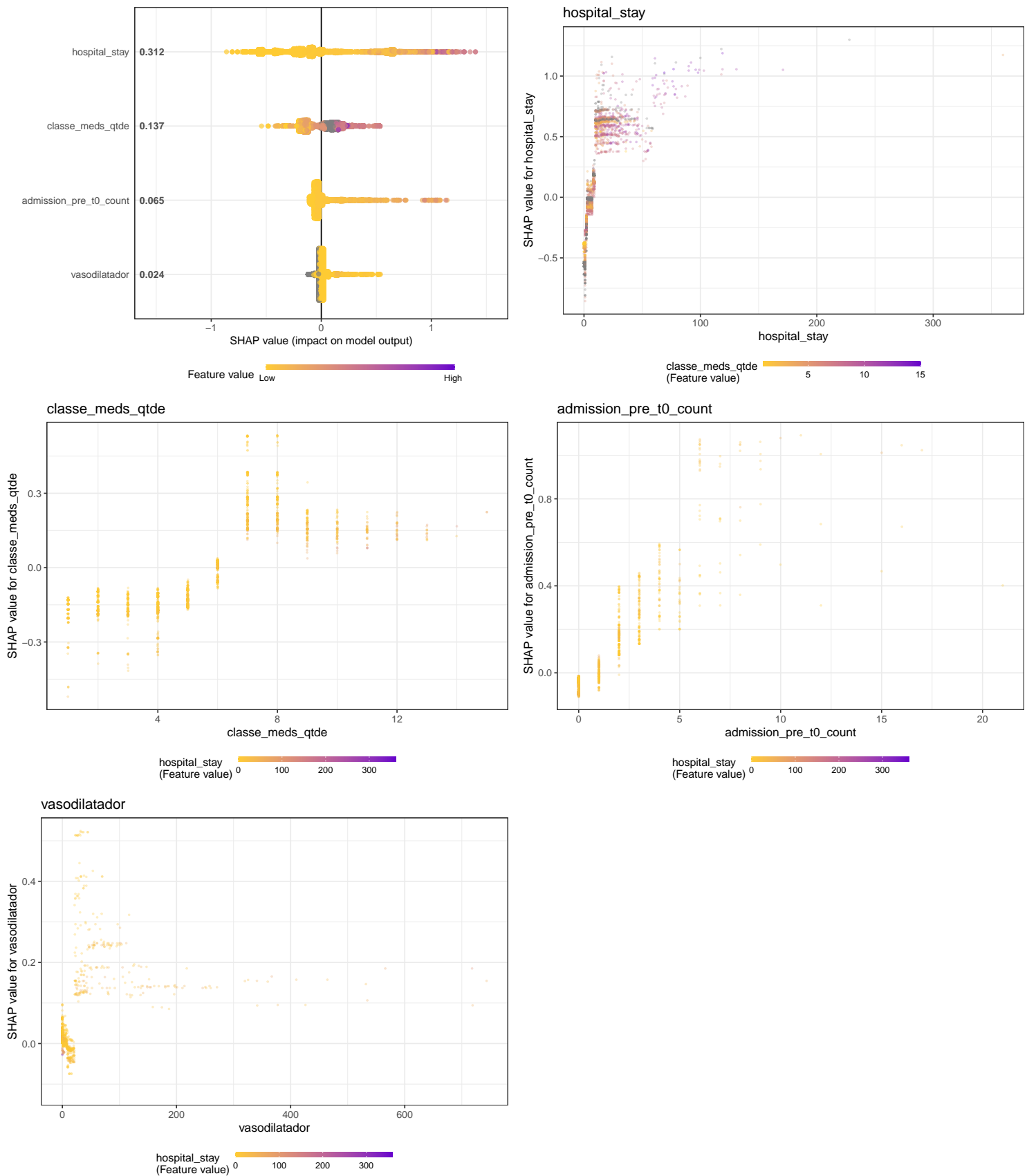
  if (plotted < n_plots) {
    print(p)
    plotted <- plotted + 1
  }

  ggsave(sprintf("./auxiliar/final_model/shap_plots/%s.png", x),
    plot = p,
    dpi = 300)
}

## Saving 6.5 x 5 in image
## Warning: Removed 1455 rows containing missing values (geom_point).
## Saving 6.5 x 5 in image

```

## Warning: Removed 1455 rows containing missing values (geom\_point).  
## Saving 6.5 x 5 in image  
## Warning: Removed 1041 rows containing missing values (geom\_point).  
## Saving 6.5 x 5 in image  
## Warning: Removed 1041 rows containing missing values (geom\_point).



```
## $num_iterations
## [1] 135
##
## $learning_rate
## [1] 0.008090184
##
## $max_depth
## [1] 4
##
## $feature_fraction
## [1] 1
##
## $min_data_in_leaf
## [1] 33
##
## $min_gain_to_split
## [1] 0
##
## $bagging_fraction
## [1] 1
##
## $num_class
## [1] 1
##
## $objective
## [1] "binary"
##
## $num_threads
## $num_threads$num_threads
## [1] 0
##
##
## $nthread
## [1] 8
##
## $seed
## [1] 1591
##
## $deterministic
## [1] TRUE
##
## $verbose
## [1] -1
##
## $metric
## list()
##
## $interaction_constraints
## list()
##
## $feature_pre_filter
## [1] FALSE
```

## Models Comparison

```
df_auc <- tibble::tribble(
  ~Model, ~`AUC`, ~`Lower Limit`, ~`Upper Limit`, ~`Features`,
  'Full Model', full_model$auc, full_model$auc_lower, full_model$auc_upper, length(features),
  'Trimmed Model', trimmed_model$auc, trimmed_model$auc_lower, trimmed_model$auc_upper, length(trimmed_features),
  'Feature Selected Model', feature_selected_model$auc, feature_selected_model$auc_lower, feature_selected_model$auc_upper, length(feature_selected_features)
```

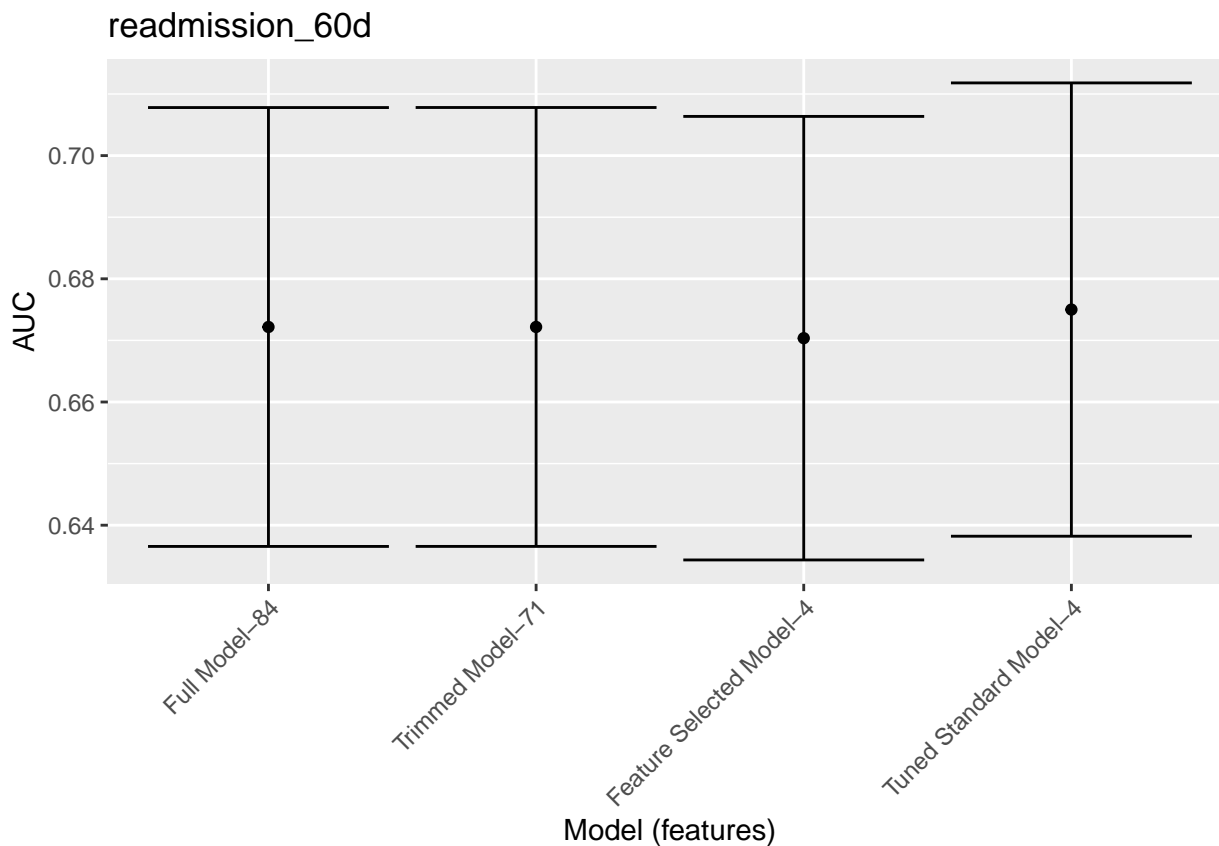


```

'Tuned Standard Model', standard_results$auc, standard_results$auc_lower, standard_results$auc_upper, length(s
) %>%
  mutate(Target = outcome_column,
         `Model (features)` = fct_reorder(paste0(Model, "-", Features), -Features))

df_auc %>%
  ggplot(aes(
    x = `Model (features)`,
    y = AUC,
    ymin = `Lower Limit`,
    ymax = `Upper Limit`
  )) +
  geom_point() +
  geom_errorbar() +
  labs(title = outcome_column) +
  theme(axis.text.x = element_text(angle = 45, hjust = 1))

```



```

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

```