

Final Model - readmission_180d

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 (968 bytes)
## . max_auc_loss = double 1= 0.01
## . outcome_column = character 1= readmission_180d
## . 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. sex
2. age
3. education_level
4. patient_state
5. underlying_heart_disease
6. heart_disease
7. nyha_basal

8. prior_mi
9. heart_failure
10. af
11. cardiac_arrest
12. transplant
13. valvopathy
14. endocardites
15. diabetes
16. renal_failure
17. hemodialysis
18. copd
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. icu_t0
28. dialysis_t0
29. n_procedure_t0
30. admission_t0_emergency
31. aco
32. antiarritmico
33. betabloqueador
34. ieca_bra
35. dva
36. digoxina
37. estatina
38. diuretico
39. vasodilatador
40. insuf_cardiaca
41. espironolactona
42. bloq_calcio
43. antiplaquetario_ev
44. insulin
45. anticonvulsivante
46. psicofarmacos
47. antifungico
48. antiviral
49. antiretroviral
50. classe_meds_qtde
51. meds_cardiovasc_qtde
52. meds_antimicrobianos
53. vni
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. arteriografia
 84. cintilografia
 85. ecocardiograma
 86. endoscopia
 87. flebografia
 88. pet_ct
 89. ultrassom
 90. tomografia
 91. ressonancia
 92. exames_imagem_qtde
 93. bic
 94. mpp
 95. 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: 0.713"
sprintf('Full Model Test AUC: %.3f' ,full_model$auc)

## [1] "Full Model Test AUC: 0.693"

```

Features with zero importance on the initial model:

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

unimportant_features %>%
  gluedown::md_order()

1. prior_mi
2. transplant
3. valvopathy
4. endocardites
5. renal_failure
6. hemodialysis
7. dialysis_t0
8. n_procedure_t0
9. antiplaquetario_ev
10. antifungico
11. antiviral
12. antiretroviral
13. vni
14. cec
15. transplante_cardiaco
16. cir_toracica
17. icp
18. intervencao_cv
19. angioplastia
20. cateterismo
21. cve_desf
22. transfusao
23. teste_esforco
24. espiro_ergoespiro
25. tilt_teste
26. citologia
27. biopsia
28. angio_rm
29. angio_tc
30. arteriografia
31. flebografia
32. pet_ct
33. mpp

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: 0.714"
sprintf('Trimmed Model Test AUC: %.3f' ,trimmed_model$auc)

## [1] "Trimmed Model Test AUC: 0.693"
selection_results <- tibble::tribble(
  ~`Tested Feature`, ~`Dropped`, ~`Number of Features`, ~`CV AUC`, ~`CV AUC Std Error`, ~`Total AUC Loss`, ~`Ins
  '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] "61"           "-5e-04"      "0"          "reop_type_1"
## [1] "60"           "-7e-04"      "-2e-04"    "insulina"
## [1] "59"           "0"          "6e-04"     "bic"
## [1] "58"           ""           "-5e-04"    "-4e-04"
## [4] "outros_proced_cirurgicos"
## [1] "57"           ""           "-4e-04"    "1e-04"      "eletrofisiologia"
## [1] "56"           ""           "-8e-04"    "-4e-04"    "tomografia"
## [1] "55"           ""           "-9e-04"    "-2e-04"    "af"
## [1] "54"           ""           "-0.0014"   "-4e-04"    "ventilacao_mecanica"
## [1] "53"           ""           "-0.0013"   "0"          "diabetes"
## [1] "52"           ""           "-0.0024"   "-0.0011"   "aco"
## [1] "51"           ""           "-0.0012"   "0.0013"
## [4] "admission_t0_emergency"
## [1] "50"           ""           "-0.0016"   "-5e-04"    "anticonvulsivante"
## [1] "49"           ""           "-0.0014"   "2e-04"     "cintilografia"
## [1] "48"           ""           "-0.0017"   "-4e-04"    "ressonancia"
## [1] "47"           ""           "-0.001"    "8e-04"     "cultura"
## [1] "46"           ""           "-0.0013"   "-4e-04"
## [4] "cateter Venoso_Central"
## [1] "45"           ""           "-0.0014"   "0"          "digoxina"
## [1] "44"           ""           "-0.0012"   "2e-04"     "cardiac_arrest"
## [1] "43"           ""           "-0.0019"   "-7e-04"    "heart_disease"
## [1] "42"           ""           "-0.0019"   "0"          "copd"
## [1] "41"           ""           "-0.0027"   "-8e-04"    "heart_failure"
## [1] "40"           ""           "-9e-04"    "0.0018"    "endoscopia"
## [1] "39"           ""           "-6e-04"    "4e-04"     "holter"
## [1] "38"           ""           "-9e-04"    "-4e-04"    "ultrassom"
## [1] "37"           ""           "-0.0011"   "-1e-04"    "cied_final_group_1"
## [1] "36"           ""           "-0.0014"   "-3e-04"    "bloq_calcio"
## [1] "35"           ""           "-0.0012"   "2e-04"
## [4] "analises_clinicas_qtde"
## [1] "34"           ""           "-0.0017"   "-6e-04"    "interconsulta"
## [1] "33"           ""           "-0.0018"   "-1e-04"
## [4] "underlying_heart_disease"
## [1] "32"           ""           "-0.0011"   "7e-04"     "patient_state"
## [1] "31"           "6e-04"     "0.0017"   "sex"
## [1] "30"           ""           "-2e-04"    "-8e-04"    "procedure_type_new"
## [1] "29"           ""           "-3e-04"    "-1e-04"    "betabloqueador"
## [1] "28"           ""           "-6e-04"    "-3e-04"    "exames_imagem_qtde"
## [1] "27"           ""           "1e-04"     "6e-04"     "insuf_cardiaca"
## [1] "26"           ""           "0.002"    "0.0019"
## [4] "proced_invasivos_qtde"
## [1] "25"           "6e-04"     "-0.0014"   "estatina"
## [1] "24"           "6e-04"     "0"          "nyha_basal"
## [1] "23"           "6e-04"     "0"          "diuretico"
## [1] "22"           ""           "0.0013"    "6e-04"     "histopatologia_qtde"
## [1] "21"           ""           "0.0011"    "-2e-04"    "espironolactona"
## [1] "20"           ""           "0.0016"    "5e-04"     "comorbidities_count"
## [1] "19"           "0.0024"   "7e-04"     "dva"
## [1] "18"           ""           "0.0032"    "8e-04"
## [4] "metodos_graficos_qtde"
## [1] "18"           ""           "0.0032"    "0.0036"    "procedure_type_1"
## [1] "17"           ""           "0.0018"    "-0.0014"
## [4] "admission_pre_t0_180d"
## [1] "16"           "0.001"    "-8e-04"    "ecocardiograma"
## [1] "16"           "0.001"    "0.0028"    "equipe_multiprof"
## [1] "15"           ""           "0.002"     "0.001"     "meds_antimicrobianos"
## [1] "15"           "0.002"    "0.0027"    "cied_final_1"
## [1] "14"           "0.0031"   "0.0011"    "psicofarmacos"
## [1] "13"           "0.0043"   "0.0012"    "education_level"
## [1] "12"           "0.0045"   "3e-04"     "ieca_bra"

```

```

## [1] "11"           "0.0065"        "0.002"          "antiarritmico"
## [1] "10"           "0.0065" "1e-04"       "icu_t0"
## [1] "10"           "0.0065"        "0.0032"         "vasodilatador"
## [1] "9"            "0.0044"        "-0.0022"        "laboratorio"
## [1] "9"            "0.0044"        "0.0025"         "classe_meds_qtde"
## [1] "9"            "0.0044"        "0.014"          ""
## [4] "admission_pre_t0_count"
## [1] "8"            "0.0063"        "0.0019"         "meds_cardiovasc_qtde"
## [1] "7"            "0.0055" "8e-04"       "age"
## [1] "7"            "0.0055"        "0.0125"        "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	95	0.7131	0.0064	0.0000	0.0000
All unimportant	TRUE	62	0.7136	0.0063	-0.0005	-0.0005
reop_type_1	TRUE	61	0.7136	0.0063	-0.0005	0.0000
insulina	TRUE	60	0.7138	0.0063	-0.0007	-0.0002
bic	TRUE	59	0.7131	0.0064	0.0000	0.0006
outros_proced_cirurgicos	TRUE	58	0.7136	0.0063	-0.0005	-0.0004
eletrofisiologia	TRUE	57	0.7134	0.0063	-0.0004	0.0001
tomografia	TRUE	56	0.7138	0.0064	-0.0008	-0.0004
af	TRUE	55	0.7140	0.0067	-0.0009	-0.0002
ventilacao_mecanica	TRUE	54	0.7144	0.0063	-0.0014	-0.0004
diabetes	TRUE	53	0.7144	0.0066	-0.0013	0.0000
aco	TRUE	52	0.7155	0.0064	-0.0024	-0.0011
admission_t0_emergency	TRUE	51	0.7142	0.0063	-0.0012	0.0013
anticonvulsivante	TRUE	50	0.7147	0.0066	-0.0016	-0.0005
cintilografia	TRUE	49	0.7144	0.0061	-0.0014	0.0002
ressonancia	TRUE	48	0.7148	0.0063	-0.0017	-0.0004
cultura	TRUE	47	0.7140	0.0064	-0.0010	0.0008
cateter Venoso Central	TRUE	46	0.7144	0.0063	-0.0013	-0.0004
digoxina	TRUE	45	0.7145	0.0065	-0.0014	0.0000
cardiac_arrest	TRUE	44	0.7143	0.0066	-0.0012	0.0002
heart_disease	TRUE	43	0.7150	0.0062	-0.0019	-0.0007
copd	TRUE	42	0.7149	0.0062	-0.0019	0.0000
heart_failure	TRUE	41	0.7158	0.0062	-0.0027	-0.0008
endoscopia	TRUE	40	0.7140	0.0063	-0.0009	0.0018
holter	TRUE	39	0.7136	0.0065	-0.0006	0.0004
ultrassom	TRUE	38	0.7140	0.0062	-0.0009	-0.0004
cied_final_group_1	TRUE	37	0.7141	0.0057	-0.0011	-0.0001
bloq_calcio	TRUE	36	0.7145	0.0056	-0.0014	-0.0003
analises_clinicas_qtde	TRUE	35	0.7142	0.0060	-0.0012	0.0002
interconsulta	TRUE	34	0.7148	0.0060	-0.0017	-0.0006
underlying_heart_disease	TRUE	33	0.7149	0.0059	-0.0018	-0.0001
patient_state	TRUE	32	0.7142	0.0056	-0.0011	0.0007
sex	TRUE	31	0.7125	0.0058	0.0006	0.0017
procedure_type_new	TRUE	30	0.7133	0.0061	-0.0002	-0.0008
betabloqueador	TRUE	29	0.7133	0.0062	-0.0003	-0.0001
exames_imagem_qtde	TRUE	28	0.7137	0.0058	-0.0006	-0.0003
insuf_cardiaca	TRUE	27	0.7130	0.0059	0.0001	0.0006
proced_invasivos_qtde	TRUE	26	0.7111	0.0056	0.0020	0.0019
estatina	TRUE	25	0.7125	0.0059	0.0006	-0.0014
nyha_basal	TRUE	24	0.7125	0.0055	0.0006	0.0000

Table 1: (continued)

Tested Feature	Dropped	Features	CV AUC	CV AUC Std Error	Total AUC Loss	Instant AUC Loss
diuretico	TRUE	23	0.7124	0.0055	0.0006	0.0000
histopatologia_qtde	TRUE	22	0.7118	0.0054	0.0013	0.0006
espironolactona	TRUE	21	0.7120	0.0056	0.0011	-0.0002
comorbidities_count	TRUE	20	0.7115	0.0056	0.0016	0.0005
dva	TRUE	19	0.7107	0.0056	0.0024	0.0007
metodos_graficos_qtde	TRUE	18	0.7099	0.0058	0.0032	0.0008
procedure_type_1	FALSE	17	0.7063	0.0067	0.0032	0.0036
admission_pre_t0_180d	TRUE	17	0.7113	0.0058	0.0018	-0.0014
ecocardiograma	TRUE	16	0.7121	0.0059	0.0010	-0.0008
equipe_multiprof	FALSE	15	0.7093	0.0059	0.0010	0.0028
meds_antimicrobianos	TRUE	15	0.7111	0.0059	0.0020	0.0010
cied_final_1	FALSE	14	0.7084	0.0056	0.0020	0.0027
psicofarmacos	TRUE	14	0.7100	0.0060	0.0031	0.0011
education_level	TRUE	13	0.7088	0.0063	0.0043	0.0012
ieca_bra	TRUE	12	0.7086	0.0062	0.0045	0.0003
antiarritmico	TRUE	11	0.7066	0.0062	0.0065	0.0020
icu_t0	TRUE	10	0.7065	0.0061	0.0065	0.0001
vasodilatador	FALSE	9	0.7034	0.0061	0.0065	0.0032
laboratorio	TRUE	9	0.7087	0.0062	0.0044	-0.0022
classe_meds_qtde	FALSE	8	0.7062	0.0061	0.0044	0.0025
admission_pre_t0_count	FALSE	8	0.6947	0.0062	0.0044	0.0140
meds_cardiovasc_qtde	TRUE	8	0.7067	0.0061	0.0063	0.0019
age	TRUE	7	0.7076	0.0063	0.0055	-0.0008
hospital_stay	FALSE	6	0.6951	0.0056	0.0055	0.0125

```

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)

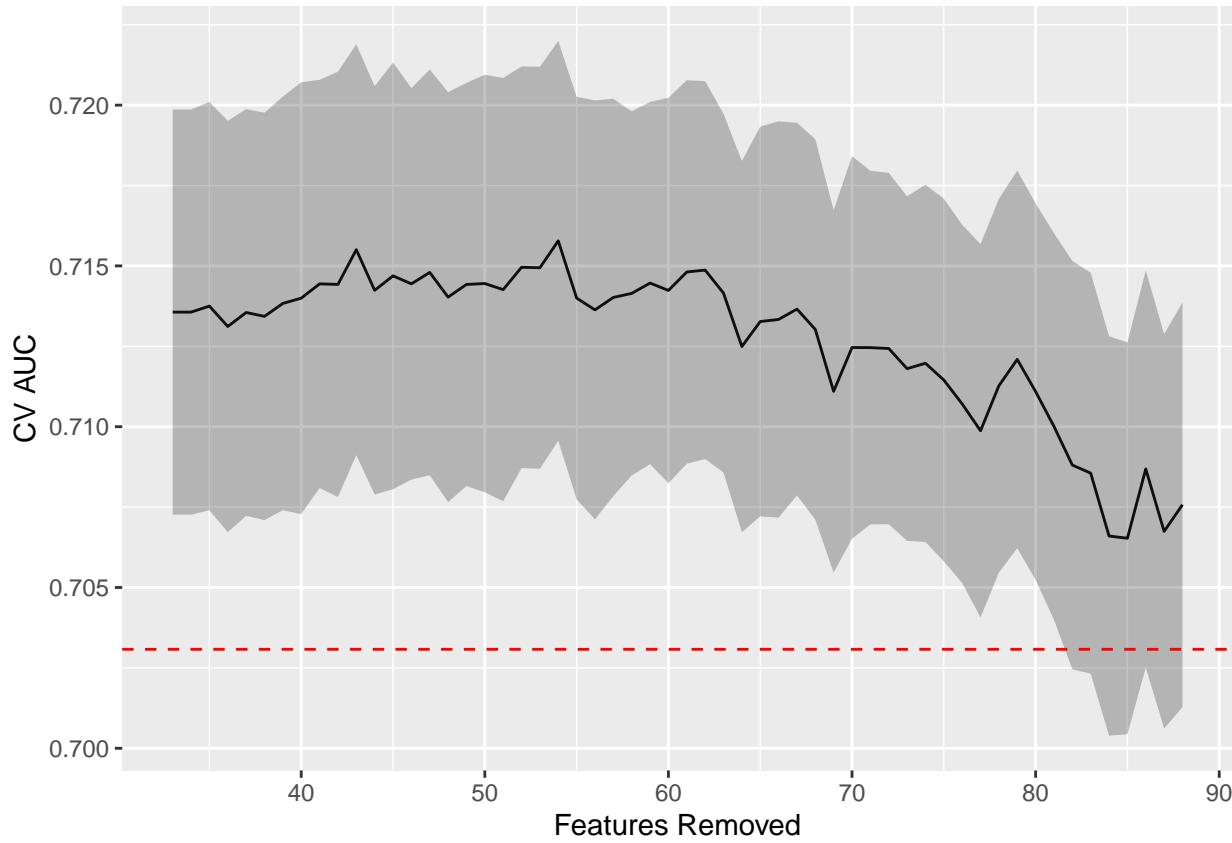
## [1] "Selected Model CV Train AUC: 0.708"
sprintf('Selected Model Test AUC: %.3f', feature_selected_model$auc)

## [1] "Selected Model Test AUC: 0.679"

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. vasodilatador
4. classe_meds_qtde
5. cied_final_1
6. equipe_multiprof
7. procedure_type_1

Standard

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

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)

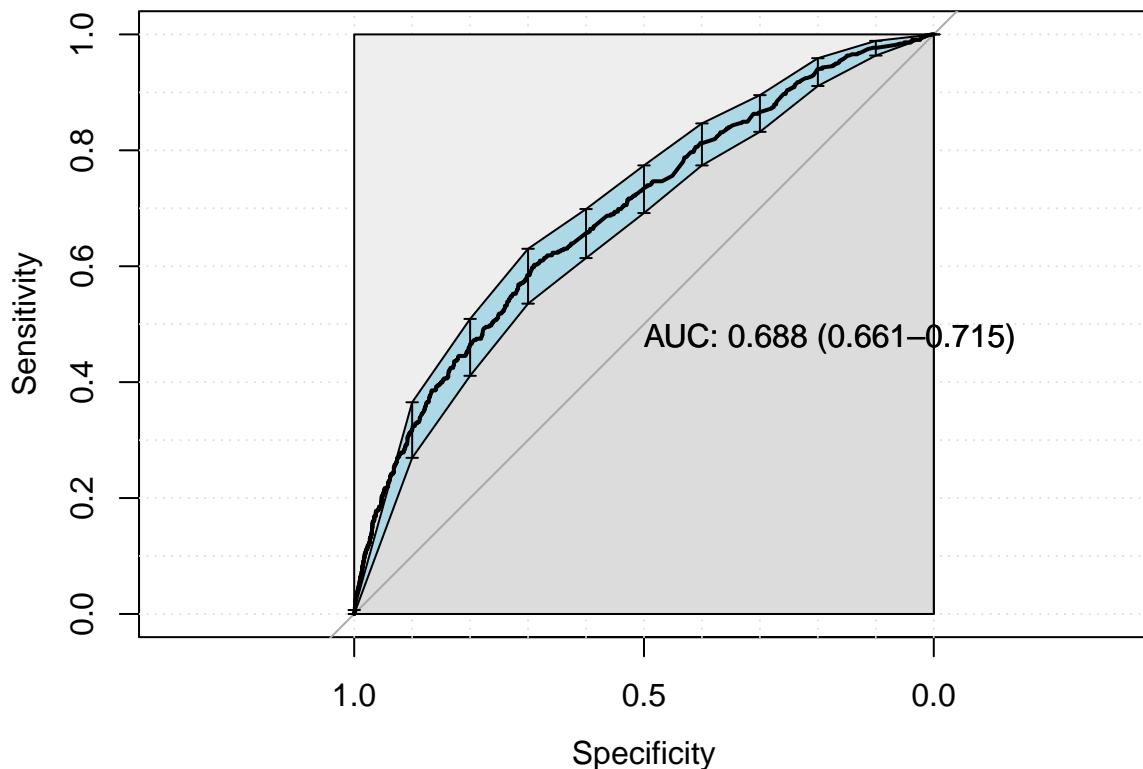
lightgbm_auc <- validation(final_lightgbm_fit, df_test)

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: 0.09"
## Confusion Matrix and Statistics
##
##      reference
## data      0      1
##   0 2958  174
##   1 1334  264
##
##                  Accuracy : 0.6812
##                  95% CI : (0.6677, 0.6945)
##      No Information Rate : 0.9074
##      P-Value [Acc > NIR] : 1
##
##                  Kappa : 0.1334
##
## Mcnemar's Test P-Value : <2e-16
##
##      Sensitivity : 0.6892
##      Specificity : 0.6027
##      Pos Pred Value : 0.9444
##      Neg Pred Value : 0.1652
##      Prevalence : 0.9074
##      Detection Rate : 0.6254
##      Detection Prevalence : 0.6622
##      Balanced Accuracy : 0.6460
##
##      '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)
}

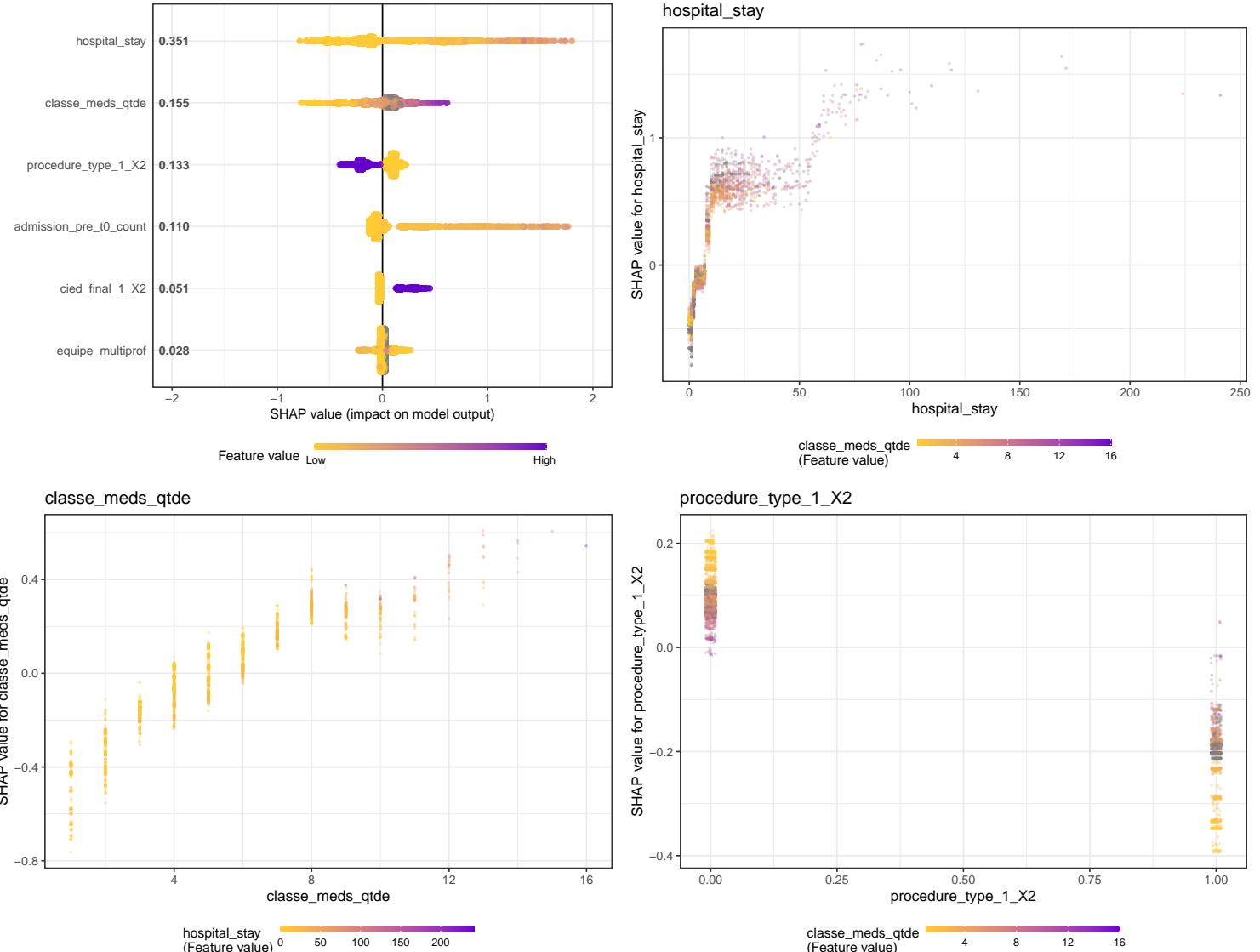
## Saving 6.5 x 5 in image
## Warning: Removed 1472 rows containing missing values (geom_point).
## Saving 6.5 x 5 in image

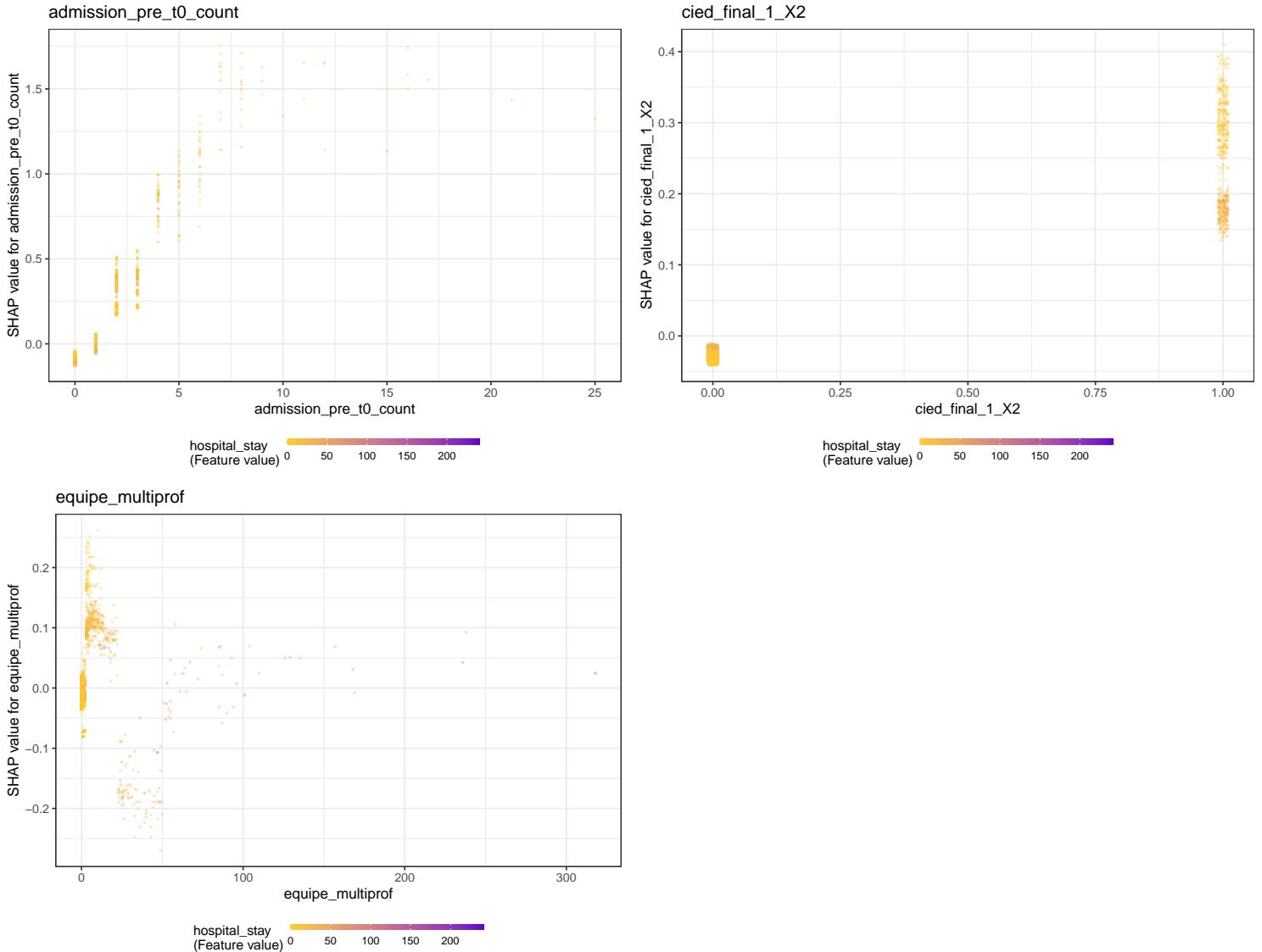
```

```

## Warning: Removed 1472 rows containing missing values (geom_point).
## Saving 6.5 x 5 in image
## Saving 6.5 x 5 in image
## Saving 6.5 x 5 in image
## Warning: Removed 822 rows containing missing values (geom_point).
## Saving 6.5 x 5 in image
## Warning: Removed 822 rows containing missing values (geom_point).
## Saving 6.5 x 5 in image
## Warning: Removed 1064 rows containing missing values (geom_point).
## Saving 6.5 x 5 in image
## Saving 6.5 x 5 in image
## Saving 6.5 x 5 in image

```





```

## $num_iterations
## [1] 86
##
## $learning_rate
## [1] 0.04322432
##
## $max_depth
## [1] 3
##
## $feature_fraction
## [1] 1
##
## $min_data_in_leaf
## [1] 80
##
## $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] 36014
##
## $deterministic
## [1] TRUE
##
## $verbose
## [1] -1
##
## $metric
## list()
##
## $interaction_constraints
## list()
##
## $feature_pre_filter
## [1] FALSE

```

Models Comparison

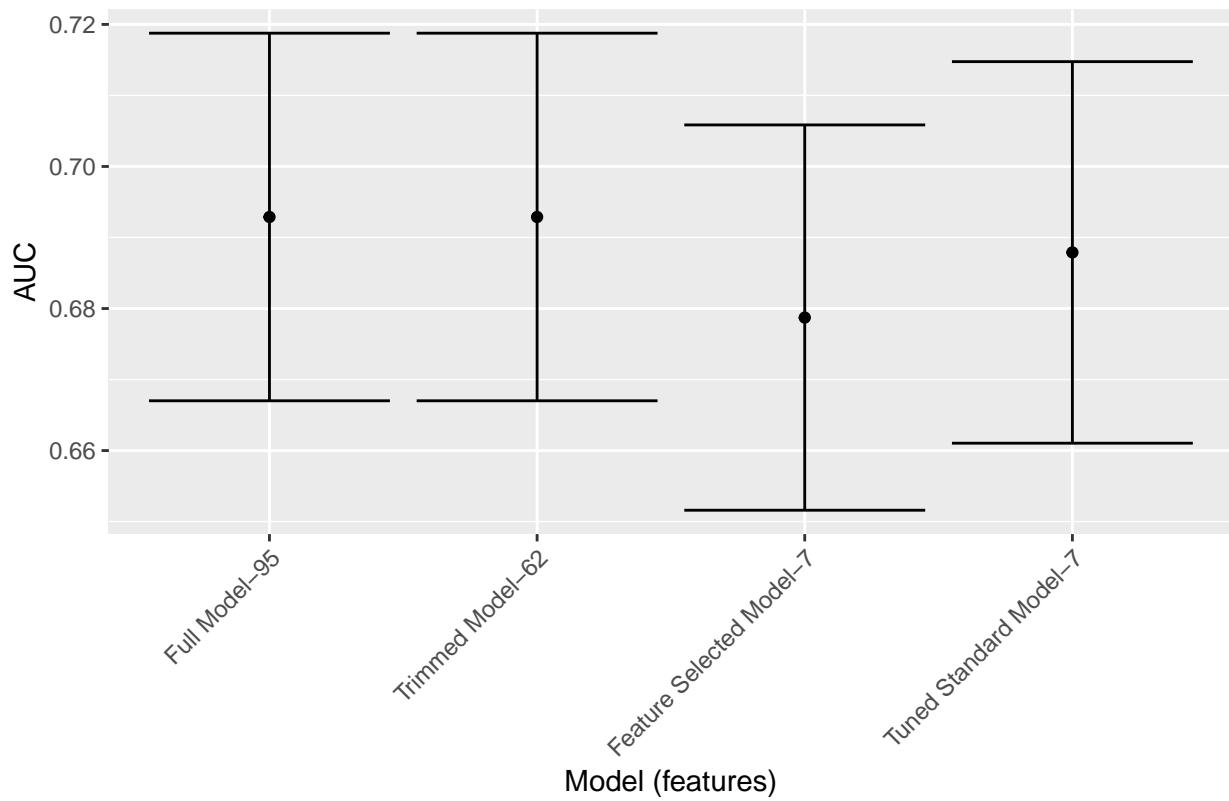
```

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$auc_upper,
  'Tuned Standard Model', standard_results$auc, standard_results$auc_lower, standard_results$auc_upper, length(standard_features)
) %>%
  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))

```

readmission_180d



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