

Final Model - death_30days

Eduardo Yuki Yada

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)
```

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))
```

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_t0', # com admission_pre_t0_count
                      '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)
}

```

Starting features:

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

1. sex
2. age
3. education_level
4. underlying_heart_disease
5. heart_disease
6. nyha_basal
7. hypertension
8. prior_mi
9. heart_failure
10. af
11. valvopathy
12. diabetes
13. renal_failure
14. hemodialysis
15. cancer
16. comorbidities_count
17. procedure_type_1
18. reop_type_1
19. procedure_type_new
20. cied_final_1
21. cied_final_group_1
22. admission_pre_t0_count
23. admission_pre_t0_180d
24. year_adm_t0
25. icu_t0
26. antiarritmico
27. antihipertensivo
28. betabloqueador
29. dva
30. diuretico
31. vasodilatador
32. espironolactona
33. antiplaquetario_ev
34. insulina
35. psicofarmacos
36. antifungico
37. classe_meds_qtde
38. meds_cardiovasc_qtde
39. meds_antimicrobianos
40. vni
41. intervencao_cv
42. cateter_venoso_central
43. proced_invasivos_qtde
44. transfusao
45. interconsulta
46. equipe_multiprof
47. ecg
48. holter
49. metodos_graficos_qtde
50. laboratorio
51. cultura

```

52. analises_clinicas_qtde
53. citologia
54. histopatologia_qtde
55. angio_tc
56. angiografia
57. cintilografia
58. ecocardiograma
59. flebografia
60. ultrassom
61. tomografia
62. radiografia
63. ressonancia
64. exames_imagem_qtde
65. bic

```

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)))

```

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

```

Functions

```

niceFormatting = function(df, caption="", digits = 2, font_size = NULL){
  df %>%
    kbl(booktabs = T, longtable = T, caption = caption, digits = digits, format = "latex") %>%
    kable_styling(font_size = font_size,
                  latex_options = c("striped", "HOLD_position", "repeat_header"))
}

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)
}

```

Feature Selection

```

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_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]))
}

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

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: 0.761"
sprintf('Full Model Test AUC: %.3f' ,full_model$auc)

## [1] "Full Model Test AUC: 0.789"

Features with zero importance on the initial model:

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

unimportant_features %>%
  gluedown::md_order()

  1. education_level
  2. underlying_heart_disease
  3. heart_disease
  4. nyha_basal
  5. hemodialysis
  6. procedure_type_new
  7. cied_final_1
  8. vni

trimmed_features <- full_model$importance$Feature
hyperparameters$mtry = min(hyperparameters$mtry, length(trimmed_features))
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: 0.757"
sprintf('Trimmed Model Test AUC: %.3f' ,trimmed_model$auc)

## [1] "Trimmed Model Test AUC: 0.782"

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

```

```

)

if (full_model$cv_auc - trimmed_model$cv_auc < max_auc_loss) {
  current_features <- trimmed_features
  current_model <- trimmed_model
  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(trimmed_features),
            `AUC Loss` = current_auc_loss,
            `Least Important Feature` = current_least_important)
} else {
  current_features <- features
  current_model <- full_model
  current_least_important <- tail(current_model$importance$Feature, 1)
  current_auc_loss <- 0
}

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(df_train, 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))
}

selection_results %>% niceFormatting(digits = 4)

```

Table 1:

Number of Features	AUC Loss	Least Important Feature
65	0.0000	valvopathy
57	0.0039	intervencao_cv
56	-0.0124	reop_type_1
55	0.0091	antiplaquetario_ev
54	0.0226	cancer

```

if (exists('last_feature_dropped')) {
  selected_features <- c(current_features, last_feature_dropped)
} else {
  selected_features <- current_features
}

con <- file(sprintf('../EDA/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('Trimmed Model CV Train AUC: %.3f', feature_selected_model$cv_auc)

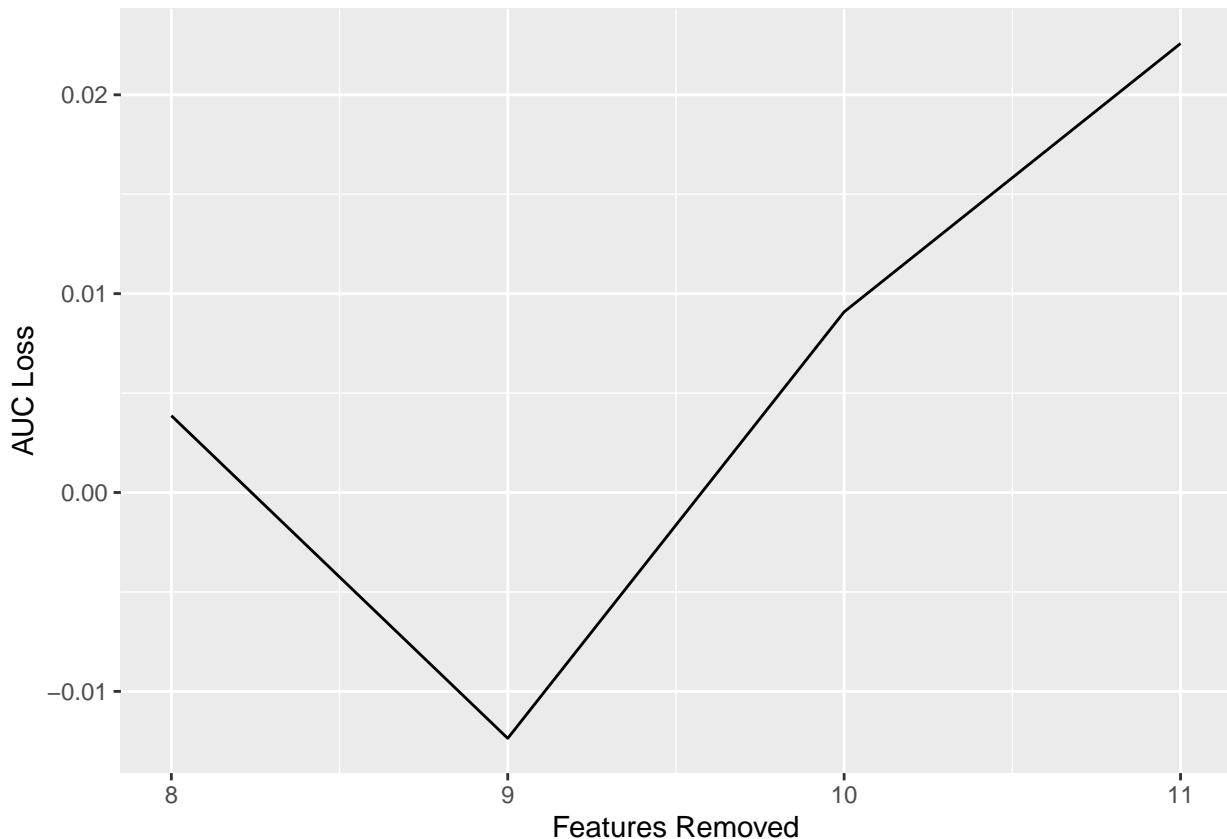
## [1] "Trimmed Model CV Train AUC: 0.756"

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

## [1] "Trimmed Model Test AUC: 0.789"

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()

```



Hyperparameter tuning

Selected features:

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

1. icu_t0
2. age
3. vasodilatador
4. ecocardiograma
5. psicofarmacos
6. admission_pre_t0_count
7. bic
8. laboratorio
9. meds_antimicrobianos
10. year_adm_t0
11. hypertension
12. analises_clinicas_qtde
13. angiografia
14. citologia

```

15. comorbidities_count
16. dva
17. exames_imagem_qtde
18. metodos_graficos_qtde
19. espironolactona
20. ecg
21. meds_cardiovasc_qtde
22. cultura
23. holter
24. equipe_multiprof
25. diuretico
26. admission_pre_t0_180d
27. antiaritmico
28. proced_invasivos_qtde
29. ressonancia
30. radiografia
31. tomografia
32. interconsulta
33. cateter_venoso_central
34. af
35. ultrassom
36. renal_failure
37. classe_meds_qtde
38. cied_final_group_1
39. histopatologia_qtde
40. antihipertensivo
41. betabloqueador
42. prior_mi
43. antifungico
44. angio_tc
45. diabetes
46. flebografia
47. heart_failure
48. cintilografia
49. insulina
50. sex
51. transfusao
52. cancer
53. procedure_type_1
54. valvopathy

```

```

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") %>%
  niceFormatting(digits = 5)

```

Table 2:

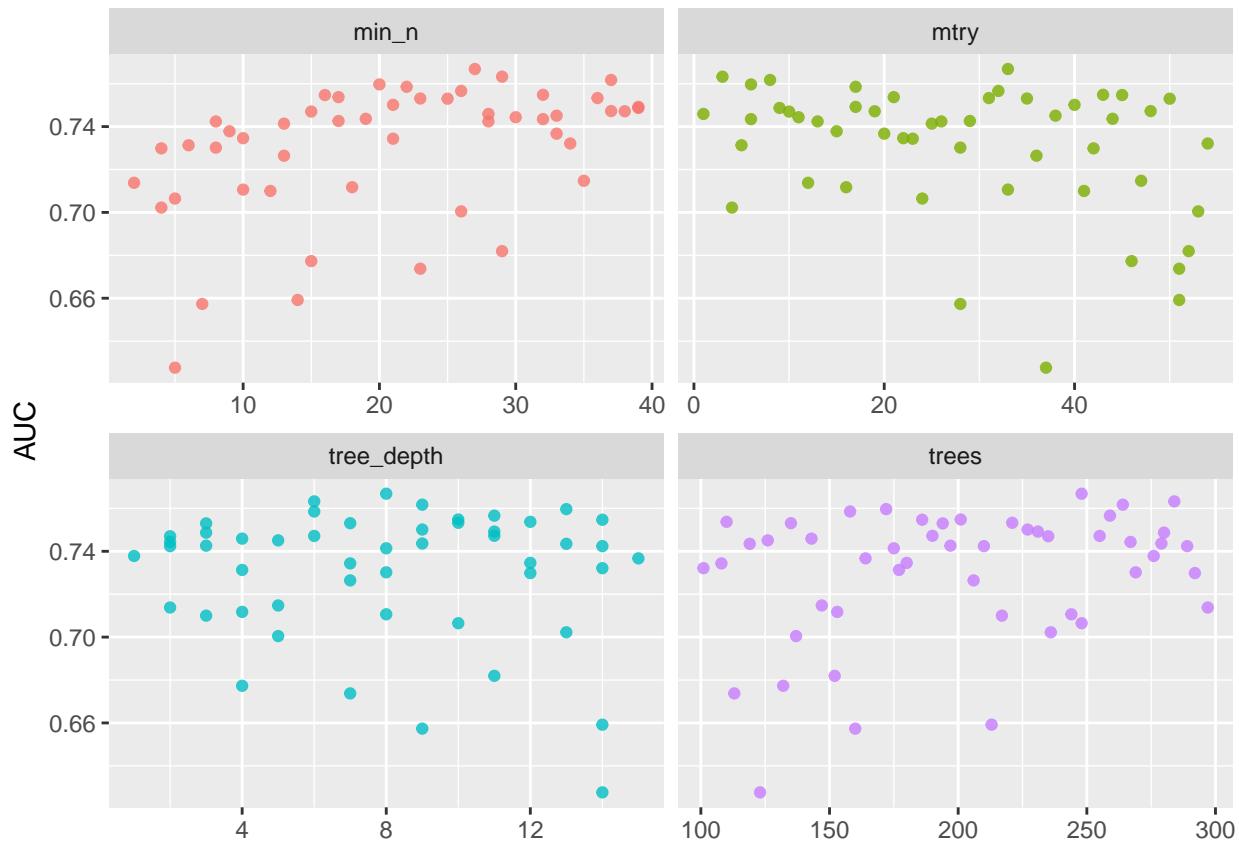
mtry	trees	min_n	tree_depth	learn_rate	loss_reduction	.metric	.estimator	mean	n	std_err	.config
33	248	27	8	0.00016	0.00000	roc_auc	binary	0.76683	4	0.02305	Preprocessor1
3	284	29	6	0.00000	0.78980	roc_auc	binary	0.76326	4	0.03543	Preprocessor1
8	264	37	9	0.00000	0.00000	roc_auc	binary	0.76173	4	0.04433	Preprocessor1
6	172	20	13	0.00000	0.00000	roc_auc	binary	0.75966	4	0.03258	Preprocessor1
17	158	22	6	0.00284	0.02648	roc_auc	binary	0.75854	4	0.02457	Preprocessor1

```

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")

```



```

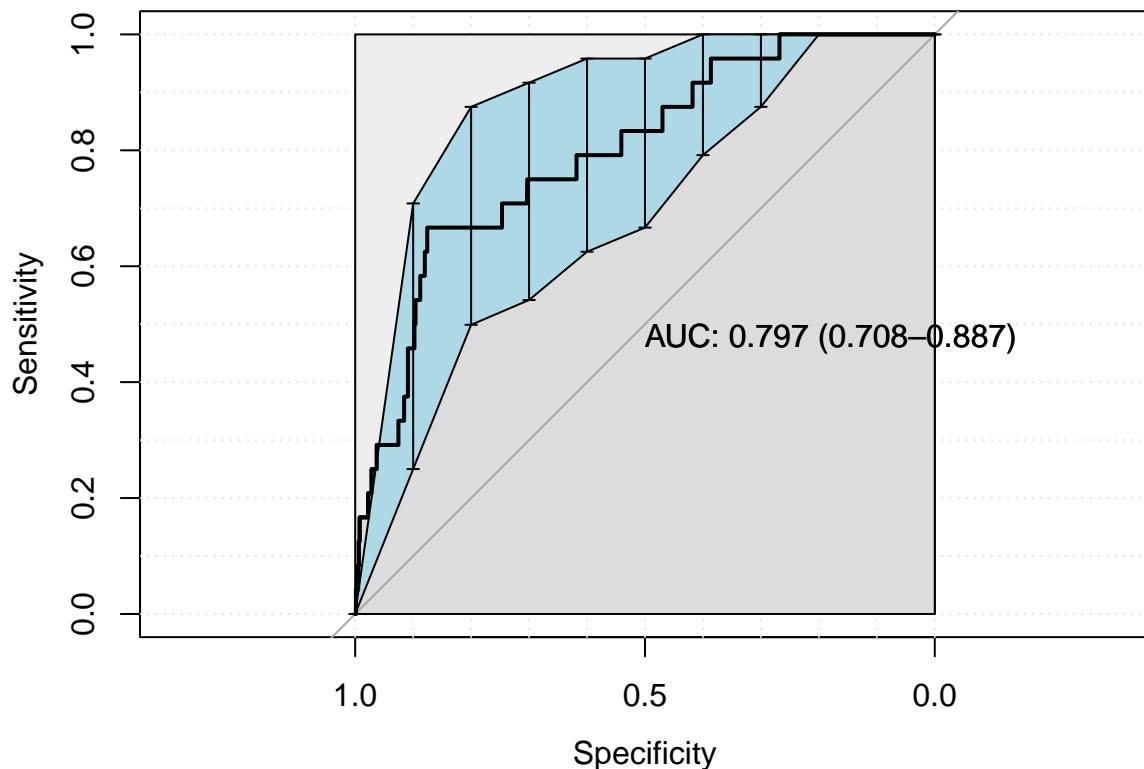
final_lightgbm_workflow <-
  lightgbm_workflow %>%
  finalize_workflow(best_lightgbm)

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 4706   24
##                      1     0     0
##
##          Accuracy : 0.9949
## 95% CI : (0.9925, 0.9967)
## No Information Rate : 0.9949
## P-Value [Acc > NIR] : 0.554
##
##          Kappa : 0
##
## McNemar's Test P-Value : 2.668e-06
##
##          Sensitivity : 1.0000
##          Specificity  : 0.0000
## Pos Pred Value  : 0.9949
## Neg Pred Value  : NaN
## Prevalence       : 0.9949
## Detection Rate  : 0.9949
## 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

saveRDS(

```

```

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

```

SHAP values

```

lightgbm_model <- parsnip::extract_fit_engine(final_lightgbm_fit)

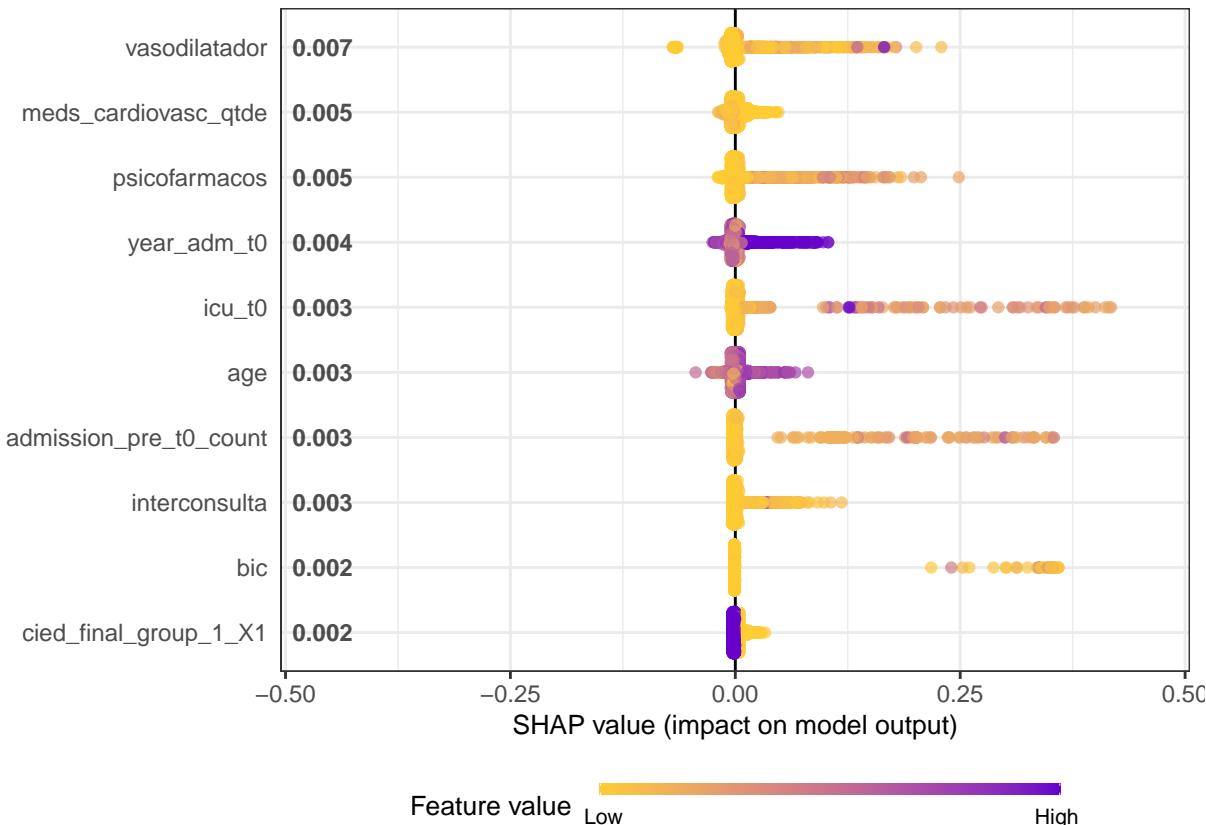
trained_rec <- prep(lightgbm_recipe, training = df_train)

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)))

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 = TRUE,
    jitter_width = 0.01,
    alpha = 0.3
  )
}

```

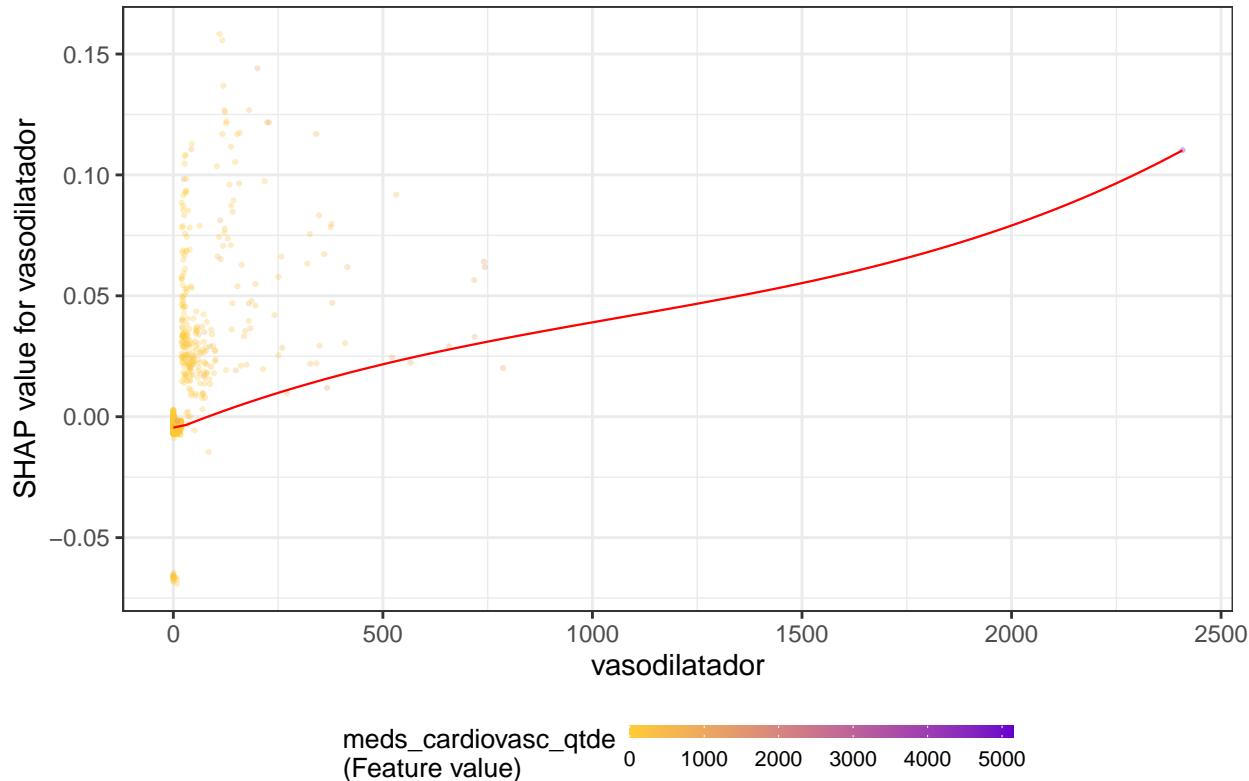
```

) +
  labs(title = x)
print(p)
}

## `geom_smooth()` using formula 'y ~ x'

```

vasodilatador

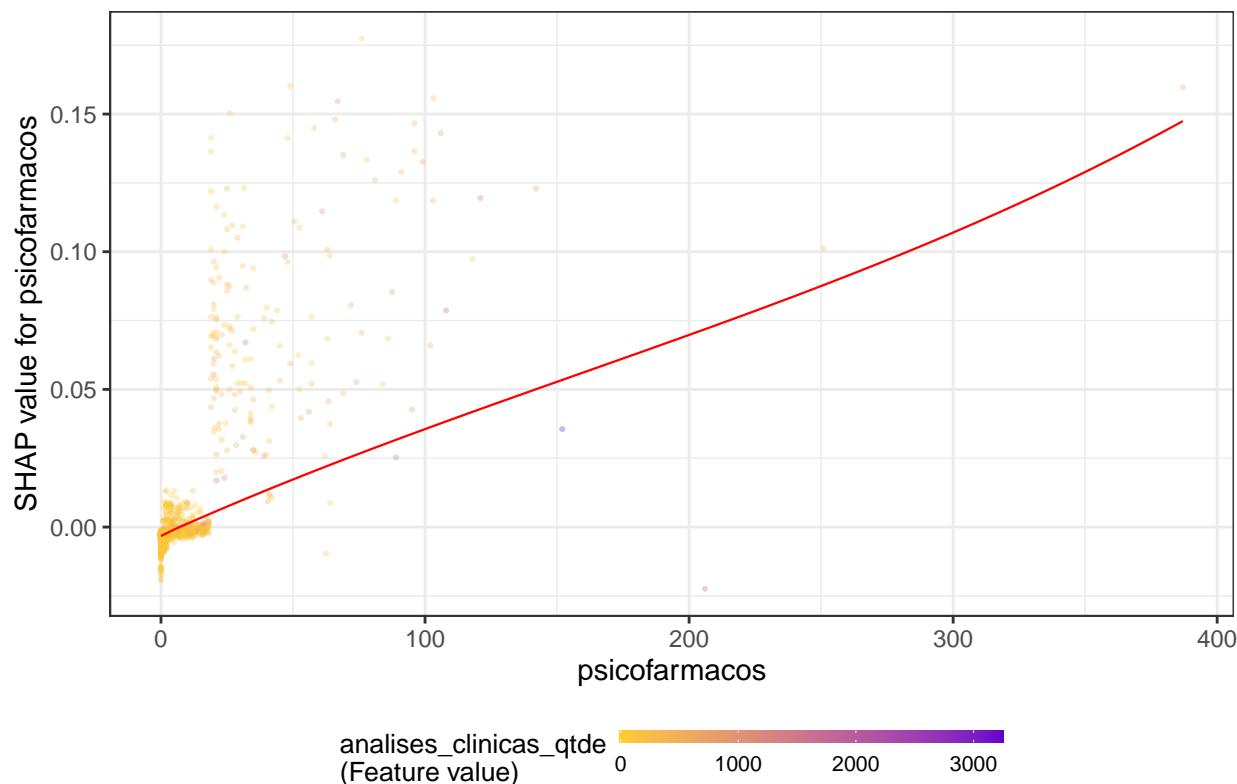


```

## `geom_smooth()` using formula 'y ~ x'

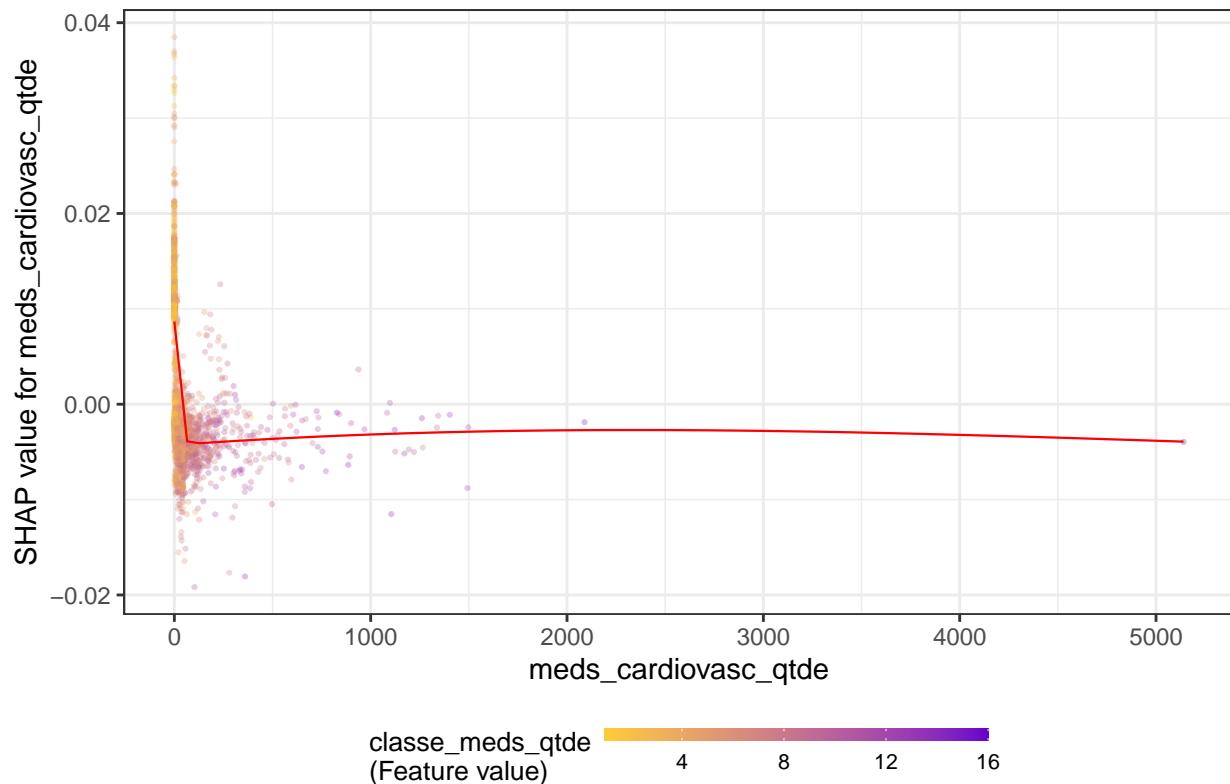
```

psicofarmacos



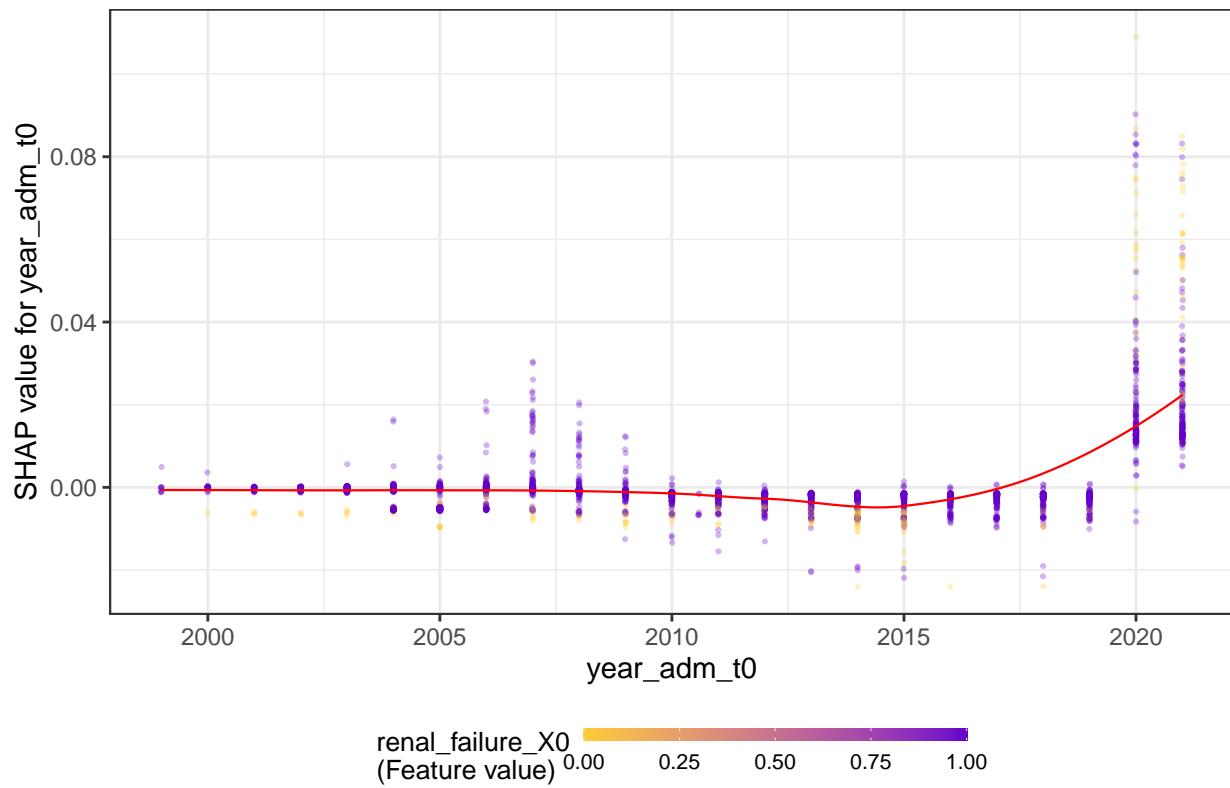
```
## `geom_smooth()` using formula 'y ~ x'
```

meds_cardiovasc_qtde



```
## `geom_smooth()` using formula 'y ~ x'
```

year_adm_t0



```

## `geom_smooth()` using formula 'y ~ x'

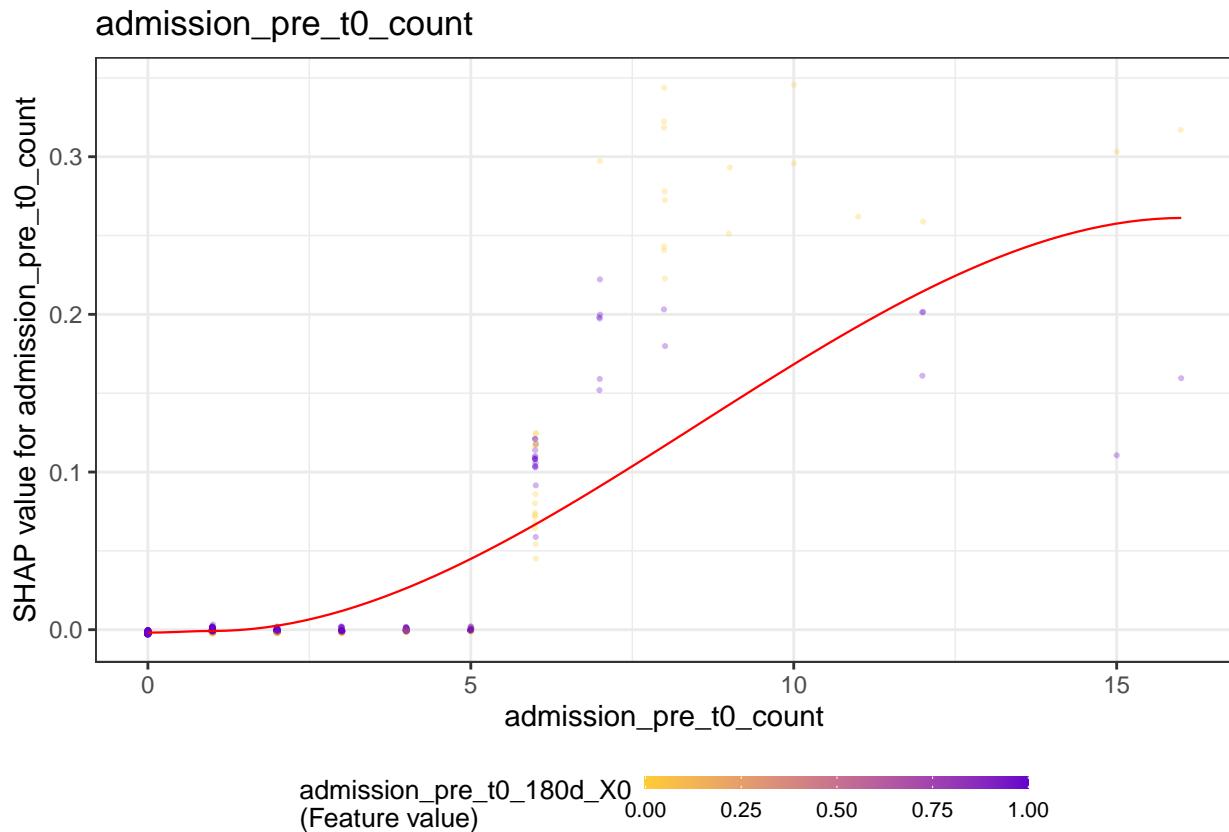
## Warning in simpleLoess(y, x, w, span, degree = degree, parametric = parametric, :
##   pseudoinverse
## used at -0.08

## Warning in simpleLoess(y, x, w, span, degree = degree, parametric = parametric, :
##   neighborhood
##   radius 1.08

## Warning in simpleLoess(y, x, w, span, degree = degree, parametric = parametric, :
##   reciprocal
##   condition number 2.7567e-28

## Warning in simpleLoess(y, x, w, span, degree = degree, parametric = parametric, :
##   There are
##   other near singularities as well. 1

```

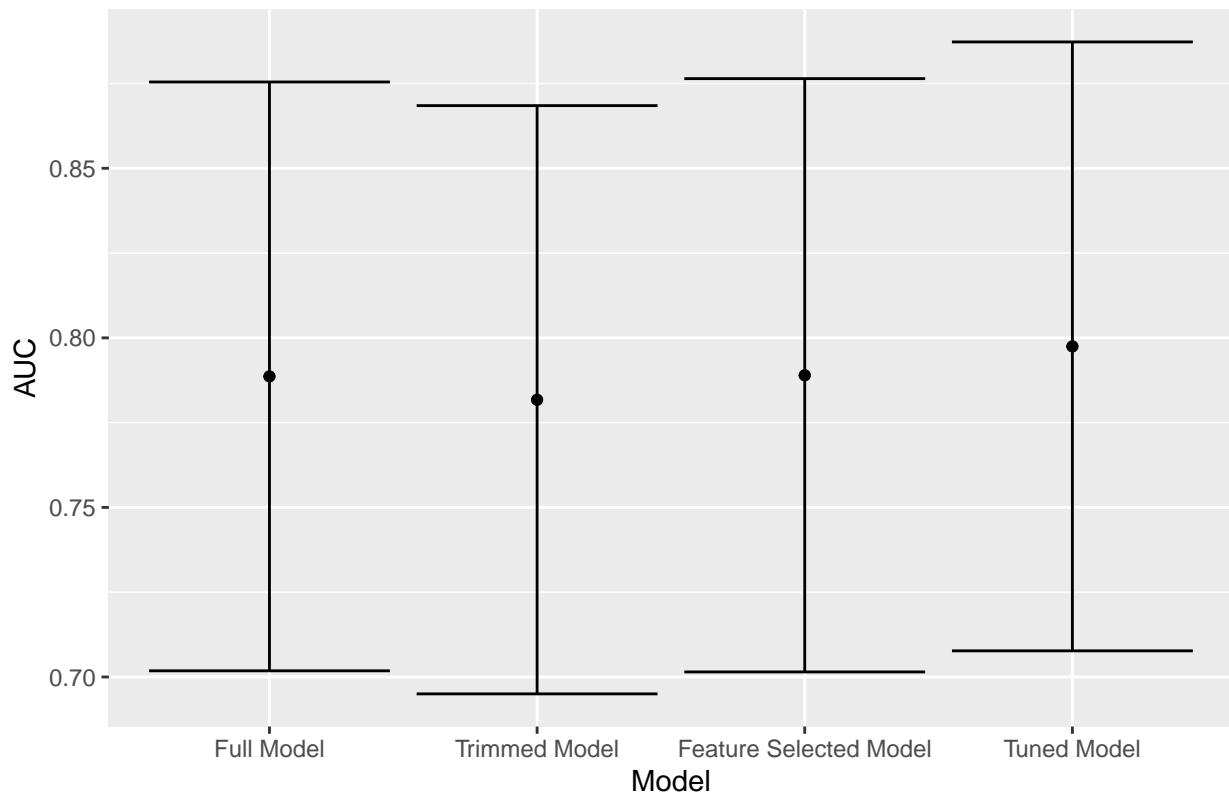


Models Comparison

```
df_auc <- 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],
  # 'Oversampled Model', as.numeric(oversampled_model$auc), oversampled_model$auc_lower, oversampled_model$auc_upper,
  # 'Undersampled Model', as.numeric(undersampled_model$auc), undersampled_model$auc_lower, undersampled_model$auc_upper
) %>%
  mutate(Target = outcome_column,
    Model = factor(Model,
      levels = c('Full Model', 'Trimmed Model',
      'Feature Selected Model', 'Tuned Model',
      'Oversampled Model', 'Undersampled Model')))

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

death_30days



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