

Biostat 203B Homework 5

Due Mar 20 @ 11:59PM

Wenqiang Ge UID:106371961

Table of contents

0.1 Predicting ICU duration	1
---------------------------------------	---

0.1 Predicting ICU duration

Using the ICU cohort `mimiciv_icu_cohort.rds` you built in Homework 4, develop at least three machine learning approaches (logistic regression with enet regularization, random forest, boosting, SVM, MLP, etc) plus a model stacking approach for predicting whether a patient's ICU stay will be longer than 2 days. You should use the `los_long` variable as the outcome. You algorithms can use patient demographic information (gender, age at ICU `intime`, marital status, race), ICU admission information (first care unit), the last lab measurements before the ICU stay, and first vital measurements during ICU stay as features. You are welcome to use any feature engineering techniques you think are appropriate; but make sure to not use features that are not available at an ICU stay's `intime`. For instance, `last_careunit` cannot be used in your algorithms.

1. Data preprocessing and feature engineering.

```
# Load libraries
library(GGally)
```

Loading required package: `ggplot2`

Registered S3 method overwritten by 'GGally':

```
method from
+.gg      ggplot2
library(gtsummary)
library(ranger)
library(tidyverse)
```

```
-- Attaching core tidyverse packages ----- tidyverse 2.0.0 --
v dplyr      1.1.4      v readr      2.1.5
v forcats    1.0.0      v stringr    1.5.1
```

```

v lubridate 1.9.4      v tibble      3.2.1
v purrr      1.0.4      v tidyr      1.3.1

-- Conflicts ----- tidyverse_conflicts() --
x dplyr::filter() masks stats::filter()
x dplyr::lag()     masks stats::lag()
i Use the conflicted package (<http://conflicted.r-lib.org/>) to force all conflicts to become explicit
library(tidymodels)

-- Attaching packages ----- tidymodels 1.3.0 --
v broom      1.0.7      v rsample     1.2.1
v dials      1.4.0      v tune        1.3.0
v infer      1.0.7      v workflows   1.2.0
v modeldata  1.4.0      v workflowsets 1.1.0
v parsnip    1.3.0      v yardstick   1.3.2
v recipes    1.1.1

-- Conflicts ----- tidymodels_conflicts() --
x scales::discard() masks purrr::discard()
x dplyr::filter()   masks stats::filter()
x recipes::fixed()  masks stringr::fixed()
x dplyr::lag()       masks stats::lag()
x yardstick::spec() masks readr::spec()
x recipes::step()    masks stats::step()
library(xgboost)

Attaching package: 'xgboost'

The following object is masked from 'package:dplyr':

    slice
library(doParallel)

Loading required package: foreach

Attaching package: 'foreach'

The following objects are masked from 'package:purrr':

    accumulate, when

Loading required package: iterators
Loading required package: parallel
library(stacks)
library(keras)

```

```
Attaching package: 'keras'
```

```
The following object is masked from 'package:yardstick':
```

```
  get_weights
```

```
library(dplyr)
library(e1071)
```

```
Attaching package: 'e1071'
```

```
The following object is masked from 'package:tune':
```

```
  tune
```

```
The following object is masked from 'package:rsample':
```

```
  permutations
```

```
The following object is masked from 'package:parsnip':
```

```
  tune
```

```
library(rsample)
library(glmnet)
```

```
Loading required package: Matrix
```

```
Attaching package: 'Matrix'
```

```
The following objects are masked from 'package:tidyr':
```

```
  expand, pack, unpack
```

```
Loaded glmnet 4.1-8
```

```
library(vip)
```

```
Attaching package: 'vip'
```

```
The following object is masked from 'package:utils':
```

```
  vi
```

```

#loading data
mimic_icu_cohort <- readRDS("../hw4/mimiciv_shiny/mimic_icu_cohort.rds")
lab <- c("creatinine", "potassium", "sodium", "chloride", "bicarbonate",
        "hematocrit", "glucose", "wbc")

vital <- c("heart_rate", "non-invasive_blood_pressure_systolic",
          "non-invasive_blood_pressure_diastolic",
          "temperature_fahrenheit", "respiratory_rate")

#Select the variables
mimic_icu_cohort <- mimic_icu_cohort |>
  select(c("subject_id", "hadm_id", "stay_id",
          "gender", "first_careunit",
          "insurance", "marital_status",
          "race", "age_intime",
          all_of(lab), all_of(vital),
          "los_long")) |>
  arrange("subject_id", "hadm_id", "stay_id")

#Convert gender to factor variables
mimic_icu_cohort$gender <- ifelse(mimic_icu_cohort$gender == "M", 0, 1)
mimic_icu_cohort$gender <- factor(mimic_icu_cohort$gender)

# replace NA with 0, label other unique values
mimic_icu_cohort[] <- lapply(mimic_icu_cohort, function(x) {
  # Only modify character or factor columns, but skip 'gender'
  if (is.factor(x) || is.character(x)) {
    # Skip the 'gender' column
    if (!"gender" %in% names(mimic_icu_cohort)[
      which(sapply(mimic_icu_cohort, identical, x))
    ]) {
      # Convert to character type if it is not gender
      x <- as.character(x)

      # Replace NA with "0"
      x[is.na(x)] <- "0"

      # Convert back to factor and then numeric
      x <- factor(x)
      return(as.numeric(x)) # Convert factor levels to numeric
    }
  }
  return(x) # Leave other columns unchanged
})

```

```

# Check for outliers in numeric variables
mimic_icu_cohort[] <- lapply(mimic_icu_cohort, function(x) {
  if (is.numeric(x)) {
    # Calculate the Q1 and Q3, and the IQR
    Q1 <- quantile(x, 0.25, na.rm = TRUE)
    Q3 <- quantile(x, 0.75, na.rm = TRUE)
    IQR_value <- Q3 - Q1

    # Define the lower and upper limits for outliers using the IQR method
    lower_limit <- Q1 - 1.5 * IQR_value
    upper_limit <- Q3 + 1.5 * IQR_value

    # Replace values outside the defined limits with NA (outliers)
    x[x < lower_limit | x > upper_limit] <- NA
  }
  return(x)
})

# Replace NA values with the mean of the column
mimic_icu_cohort[] <- lapply(mimic_icu_cohort, function(x) {
  if (is.numeric(x)) {
    x[is.na(x)] <- mean(x, na.rm = TRUE) # Replace NA with the mean value
  }
  return(x)
})

# Convert los_long to a factor
mimic_icu_cohort$los_long <- factor(mimic_icu_cohort$los_long,
                                   levels = c(FALSE, TRUE))

# Remove rows where los_long is NA
mimic_icu_cohort <- mimic_icu_cohort %>%
  drop_na(los_long)

# Convert other variables to factors
mimic_icu_cohort$marital_status <- factor(mimic_icu_cohort$marital_status)
mimic_icu_cohort$race <- factor(mimic_icu_cohort$race)
mimic_icu_cohort$first_careunit <- factor(mimic_icu_cohort$first_careunit)
mimic_icu_cohort$insurance <- factor(mimic_icu_cohort$insurance)

head(mimic_icu_cohort)

# A tibble: 6 x 23
  subject_id hadm_id stay_id gender first_careunit insurance marital_status
    <dbl>    <dbl>   <dbl> <fct>   <fct>         <fct>      <fct>
1  10000032 29079034 39553978 1      2             2          5

```

```

2  10000690 25860671 37081114 1      2      3      5
3  10000980 26913865 39765666 1      2      3      3
4  10001217 24597018 37067082 1      5      6      3
5  10001217 27703517 34592300 1      5      6      3
6  10001725 25563031 31205490 1      3      6      3
# i 16 more variables: race <fct>, age_intime <dbl>, creatinine <dbl>,
#   potassium <dbl>, sodium <dbl>, chloride <dbl>, bicarbonate <dbl>,
#   hematocrit <dbl>, glucose <dbl>, wbc <dbl>, heart_rate <dbl>,
#   `non-invasive_blood_pressure_systolic` <dbl>,
#   `non-invasive_blood_pressure_diastolic` <dbl>,
#   temperature_fahrenheit <dbl>, respiratory_rate <dbl>, los_long <fct>

```

2. Partition data into 50% training set and 50% test set. Stratify partitioning according to `los_long`. For grading purpose, sort the data by `subject_id`, `hadm_id`, and `stay_id` and use the seed 203 for the initial data split. Below is the sample code.

0.1.0.1 Initial split into test and non-test sets

```

set.seed(203)

# sort
mimic_icu_cohort <- mimic_icu_cohort |>
  arrange(subject_id, hadm_id, stay_id)

data_split <- initial_split(
  mimic_icu_cohort,
  # stratify by los_long
  strata = "los_long",
  prop = 0.5
)

mimic_train <- training(data_split) |>
  select(-subject_id, -hadm_id, -stay_id)
dim(mimic_train )

[1] 47221    20

mimic_test <- testing(data_split) |>
  select(-subject_id, -hadm_id, -stay_id)
dim(mimic_test)

[1] 47223    20

```

0.1.0.2 Recipe

```

data_recipe <-
  recipe(

```

```

    los_long ~ .,
    data = mimic_train
  ) |>
  step_dummy(all_nominal_predictors()) |>
  step_zv(all_numeric_predictors()) |>
  step_normalize(all_numeric_predictors())
data_recipe

```

-- Recipe -----

-- Inputs

Number of variables by role

```

outcome:    1
predictor: 19

```

-- Operations

```

* Dummy variables from: all_nominal_predictors()
* Zero variance filter on: all_numeric_predictors()
* Centering and scaling for: all_numeric_predictors()
folds <- vfold_cv(mimic_train, v = 2)

```

0.1.1 logistic regression with enet regularization

0.1.1.1 Model & Workflow

```

# set up the logistic regression model
logit_model <- logistic_reg(
  penalty = tune(),
  mixture = tune()) |>
  set_engine("glmnet",
             standardize = TRUE)

logit_wf <- workflow() |>
  add_recipe(data_recipe) |>
  add_model(logit_model) |>
  print()

```

== Workflow =====

Preprocessor: Recipe

Model: logistic_reg()

```

-- Preprocessor -----
3 Recipe Steps

* step_dummy()
* step_zv()
* step_normalize()

-- Model -----
Logistic Regression Model Specification (classification)

Main Arguments:
  penalty = tune()
  mixture = tune()

Engine-Specific Arguments:
  standardize = TRUE

Computational engine: glmnet

```

0.1.1.2 Tuning grid

```

logit_grid <- grid_regular(
  penalty(range = c(-3, 0)),
  mixture(),
  levels = 3
)

```

0.1.1.3 Cross-validation

```

set.seed(203)

if (file.exists("logit_res.rds")) {
  logit_res <- read_rds("logit_res.rds")
} else {
  logit_res <-
    tune_grid(
      object = logit_wf,
      resamples = folds,
      grid = logit_grid,
      metrics = metric_set(roc_auc, accuracy),
      control = control_stack_grid()
    )
  write_rds(logit_res, "logit_res.rds")
}
logit_res

```



```

# Tuning results
# 2-fold cross-validation
# A tibble: 2 x 5
  splits          id .metrics      .notes      .predictions
  <list>        <chr> <list>      <list>      <list>
1 <split [23610/23611]> Fold1 <tibble [18 x 6]> <tibble [0 x 3]> <tibble>
2 <split [23611/23610]> Fold2 <tibble [18 x 6]> <tibble [0 x 3]> <tibble>

```

visualize the CV results

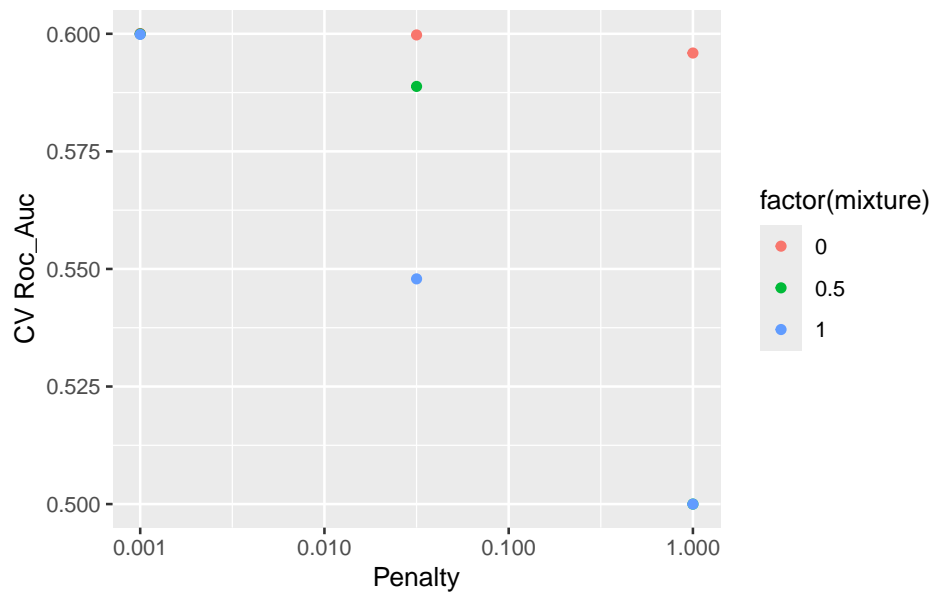
```

logit_res |>
  collect_metrics() |>
  print(width = Inf) |>
  filter(.metric == "roc_auc") |>
  ggplot(mapping = aes(x = penalty, y = mean,
                       color = factor(mixture))) +
  geom_point() +
  labs(x = "Penalty", y = "CV Roc_Auc") +
  scale_x_log10()

```

A tibble: 18 x 8

	penalty	mixture	.metric	.estimator	mean	n	std_err	.config
	<dbl>	<dbl>	<chr>	<chr>	<dbl>	<int>	<dbl>	<chr>
1	0.001	0	accuracy	binary	0.573	2	0.00189	Preprocessor1_Model11
2	0.001	0	roc_auc	binary	0.600	2	0.00194	Preprocessor1_Model11
3	0.0316	0	accuracy	binary	0.572	2	0.00219	Preprocessor1_Model12
4	0.0316	0	roc_auc	binary	0.600	2	0.00201	Preprocessor1_Model12
5	1	0	accuracy	binary	0.563	2	0.00228	Preprocessor1_Model13
6	1	0	roc_auc	binary	0.596	2	0.00165	Preprocessor1_Model13
7	0.001	0.5	accuracy	binary	0.572	2	0.00221	Preprocessor1_Model14
8	0.001	0.5	roc_auc	binary	0.600	2	0.00194	Preprocessor1_Model14
9	0.0316	0.5	accuracy	binary	0.562	2	0.00164	Preprocessor1_Model15
10	0.0316	0.5	roc_auc	binary	0.589	2	0.000653	Preprocessor1_Model15
11	1	0.5	accuracy	binary	0.509	2	0.00111	Preprocessor1_Model16
12	1	0.5	roc_auc	binary	0.5	2	0	Preprocessor1_Model16
13	0.001	1	accuracy	binary	0.572	2	0.00175	Preprocessor1_Model17
14	0.001	1	roc_auc	binary	0.600	2	0.00183	Preprocessor1_Model17
15	0.0316	1	accuracy	binary	0.518	2	0.000329	Preprocessor1_Model18
16	0.0316	1	roc_auc	binary	0.548	2	0.00121	Preprocessor1_Model18
17	1	1	accuracy	binary	0.509	2	0.00111	Preprocessor1_Model19
18	1	1	roc_auc	binary	0.5	2	0	Preprocessor1_Model19



0.1.1.4 Finalize the model

```
#show best models
logit_res |>
  show_best(metric = "roc_auc")
```

A tibble: 5 x 8

	penalty	mixture	.metric	.estimator	mean	n	std_err	.config
	<dbl>	<dbl>	<chr>	<chr>	<dbl>	<int>	<dbl>	<chr>
1	0.001	0	roc_auc	binary	0.600	2	0.00194	Preprocessor1_Model1
2	0.001	0.5	roc_auc	binary	0.600	2	0.00194	Preprocessor1_Model4
3	0.001	1	roc_auc	binary	0.600	2	0.00183	Preprocessor1_Model7
4	0.0316	0	roc_auc	binary	0.600	2	0.00201	Preprocessor1_Model2
5	1	0	roc_auc	binary	0.596	2	0.00165	Preprocessor1_Model3

```
# select the best model
best_logit <- logit_res |>
  select_best(metric = "roc_auc")
best_logit
```

A tibble: 1 x 3

	penalty	mixture	.config
	<dbl>	<dbl>	<chr>
1	0.001	0	Preprocessor1_Model1

```
#Testing
final_logit <- logit_wf |>
  finalize_workflow(best_logit)
```

```

final_logit

== Workflow =====
Preprocessor: Recipe
Model: logistic_reg()

-- Preprocessor -----
3 Recipe Steps

* step_dummy()
* step_zv()
* step_normalize()

-- Model -----
Logistic Regression Model Specification (classification)

Main Arguments:
  penalty = 0.001
  mixture = 0

Engine-Specific Arguments:
  standardize = TRUE

Computational engine: glmnet
# Fit the whole training set, then predict the test cases
set.seed(203)
final_logit_fit <-
  final_logit |>
  last_fit(data_split)
final_logit_fit

# Resampling results
# Manual resampling
# A tibble: 1 x 6
  splits          id      .metrics .notes  .predictions .workflow
  <list>         <chr>    <list>  <list>  <list>       <list>
1 <split [47221/47223]> train/test sp~ <tibble> <tibble> <tibble>    <workflow>

final_logit_fit |>
  collect_metrics()

# A tibble: 3 x 4
  .metric      .estimator .estimate .config
  <chr>        <chr>        <dbl>    <chr>
1 accuracy    binary        0.570    Preprocessor1_Model11
2 roc_auc     binary        0.595    Preprocessor1_Model11

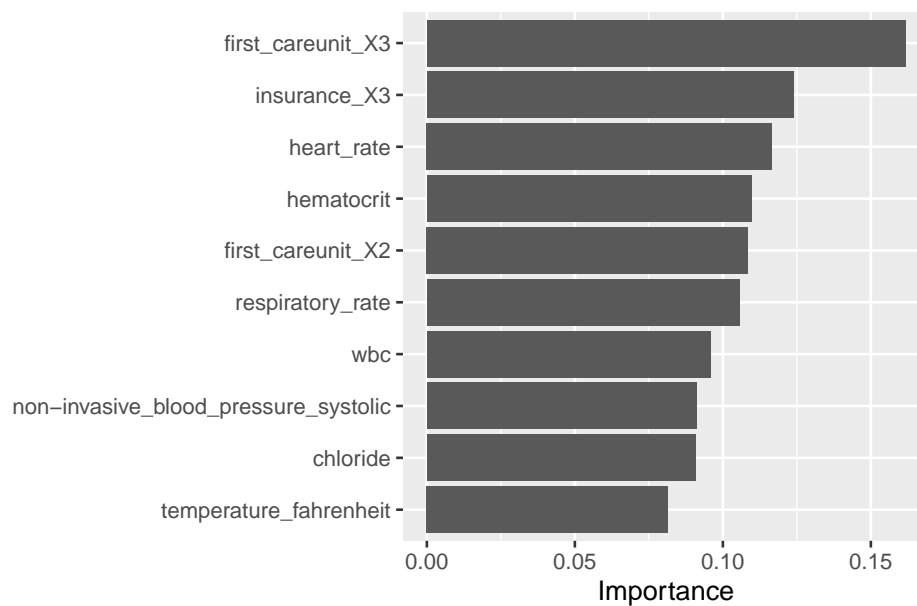
```

3 brier_class binary 0.243 Preprocessor1_Model1

0.1.1.5 Visualize the final model

```
final_logic_tree <- extract_workflow(final_logit_fit)

lg_vip <- final_logic_tree |>
  extract_fit_parsnip() |>
  vip()
lg_vip
```



```
rm(final_logic_tree,
   final_logit, logit_res,
   logit_grid, logit_model,
   logit_wf)
```

0.1.2 Boosting

0.1.2.1 Model & Workflow

```
#Model
bt_mod <-
  boost_tree(
    mode = "classification",
    trees = 100,
    tree_depth = tune(),
```

```

    learn_rate = tune()
  ) %>%
  set_engine("xgboost")
bt_mod

```

Boosted Tree Model Specification (classification)

Main Arguments:

```

trees = 100
tree_depth = tune()
learn_rate = tune()

```

Computational engine: xgboost

```

#Workflow
bt_wf <- workflow() |>
  add_recipe(data_recipe) |>
  add_model(bt_mod)
bt_wf

```

== Workflow =====

Preprocessor: Recipe

Model: boost_tree()

-- Preprocessor -----

3 Recipe Steps

```

* step_dummy()
* step_zv()
* step_normalize()

```

-- Model -----

Boosted Tree Model Specification (classification)

Main Arguments:

```

trees = 100
tree_depth = tune()
learn_rate = tune()

```

Computational engine: xgboost

0.1.2.2 Tuning grid

```

#Tuning
param_grid <- grid_regular(
  tree_depth(range = c(1L, 3L)),
  learn_rate(range = c(-3, 0), trans = log10_trans()),

```

```

    levels = 3
  )
  param_grid

# A tibble: 9 x 2
  tree_depth learn_rate
      <int>      <dbl>
1         1      0.001
2         2      0.001
3         3      0.001
4         1     0.0316
5         2     0.0316
6         3     0.0316
7         1         1
8         2         1
9         3         1

library(doParallel)
registerDoParallel(cores = 4)

```

0.1.2.3 Cross-validation

```

set.seed(203)

if (file.exists("bt_fit.rds")) {
  bt_fit <- read_rds("bt_fit.rds")
} else {
  bt_fit <- bt_wf |>
  tune_grid(
    resamples = folds,
    grid = param_grid,
    metrics = metric_set(yardstick::roc_auc,
                        yardstick::accuracy),
    control = control_stack_grid()
  )
  write_rds(bt_fit, "bt_fit.rds")
}
bt_fit

# Tuning results
# 2-fold cross-validation
# A tibble: 2 x 5
  splits          id   .metrics      .notes      .predictions
  <list>         <chr> <list>      <list>      <list>
1 <split [23610/23611]> Fold1 <tibble [18 x 6]> <tibble [0 x 3]> <tibble>
2 <split [23611/23610]> Fold2 <tibble [18 x 6]> <tibble [0 x 3]> <tibble>

```

```

stopImplicitCluster()
registerDoSEQ()

bt_fit <- readRDS("bt_fit.rds")

bt_fit |>
  collect_metrics() |>
  print(width = Inf) |>
  filter(.metric == "roc_auc") |>
  ggplot(mapping =
    aes(x = learn_rate,
        y = mean,
        color = factor(tree_depth))) +
  geom_point() +
  geom_line() +
  labs(x = "Learning Rate", y = "CV AUC") +
  scale_x_log10()

# A tibble: 18 x 8
  tree_depth learn_rate .metric .estimator mean    n std_err
    <int>      <dbl> <chr>    <chr>    <dbl> <int> <dbl>
1         1    0.001 accuracy binary    0.531     2 0.00207
2         1    0.001 roc_auc  binary    0.545     2 0.0116
3         2    0.001 accuracy binary    0.551     2 0.00134
4         2    0.001 roc_auc  binary    0.560     2 0.00281
5         3    0.001 accuracy binary    0.550     2 0.000266
6         3    0.001 roc_auc  binary    0.570     2 0.00294
7         1    0.0316 accuracy binary    0.564     2 0.00225
8         1    0.0316 roc_auc  binary    0.590     2 0.00350
9         2    0.0316 accuracy binary    0.574     2 0.000542
10        2    0.0316 roc_auc  binary    0.604     2 0.00187
11        3    0.0316 accuracy binary    0.578     2 0.00100
12        3    0.0316 roc_auc  binary    0.612     2 0.00237
13        1     1 accuracy binary    0.580     2 0.00257
14        1     1 roc_auc  binary    0.611     2 0.00346
15        2     1 accuracy binary    0.569     2 0.00259
16        2     1 roc_auc  binary    0.600     2 0.00406
17        3     1 accuracy binary    0.564     2 0.00442
18        3     1 roc_auc  binary    0.587     2 0.00279

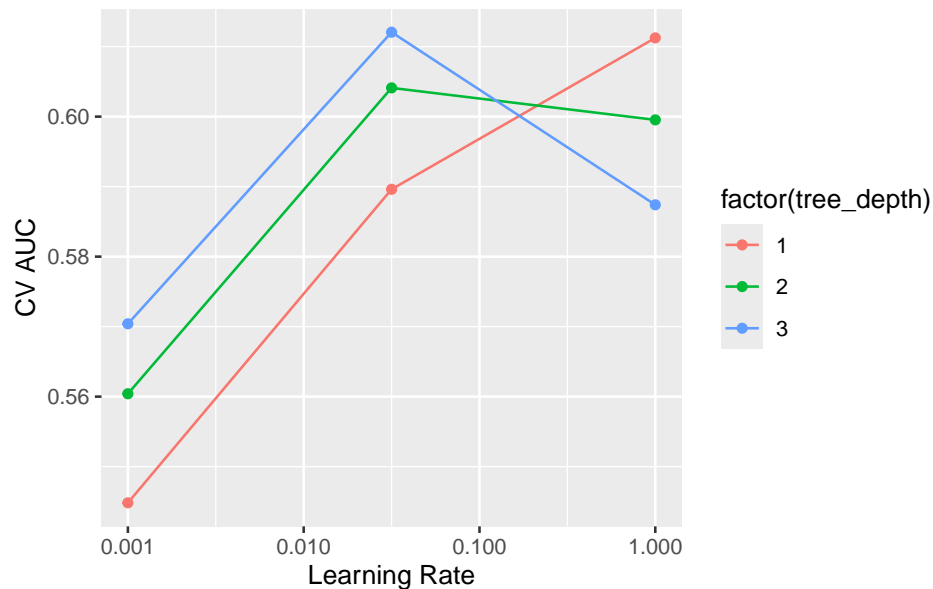
.config
<chr>
1 Preprocessor1_Model1
2 Preprocessor1_Model1
3 Preprocessor1_Model2
4 Preprocessor1_Model2

```

```

5 Preprocessor1_Model3
6 Preprocessor1_Model3
7 Preprocessor1_Model4
8 Preprocessor1_Model4
9 Preprocessor1_Model5
10 Preprocessor1_Model5
11 Preprocessor1_Model6
12 Preprocessor1_Model6
13 Preprocessor1_Model7
14 Preprocessor1_Model7
15 Preprocessor1_Model8
16 Preprocessor1_Model8
17 Preprocessor1_Model9
18 Preprocessor1_Model9

```



0.1.2.4 Finalize the model

```

bt_fit |>
  show_best(metric = "roc_auc")

```

A tibble: 5 x 8

	tree_depth	learn_rate	.metric	.estimator	mean	n	std_err	.config
	<int>	<dbl>	<chr>	<chr>	<dbl>	<int>	<dbl>	<chr>
1	3	0.0316	roc_auc	binary	0.612	2	0.00237	Preprocessor1_Mo~
2	1	1	roc_auc	binary	0.611	2	0.00346	Preprocessor1_Mo~
3	2	0.0316	roc_auc	binary	0.604	2	0.00187	Preprocessor1_Mo~
4	2	1	roc_auc	binary	0.600	2	0.00406	Preprocessor1_Mo~


```

5          1      0.0316 roc_auc binary      0.590      2 0.00350 Preprocessor1_Mo~
best_boost <- bt_fit |>
  select_best(metric = "roc_auc")
best_boost

# A tibble: 1 x 3
  tree_depth learn_rate .config
    <int>      <dbl> <chr>
1         3      0.0316 Preprocessor1_Model6

final_boost_wf <- bt_wf |>
  finalize_workflow(best_boost)
final_boost_wf

== Workflow =====
Preprocessor: Recipe
Model: boost_tree()

-- Preprocessor -----
3 Recipe Steps

* step_dummy()
* step_zv()
* step_normalize()

-- Model -----
Boosted Tree Model Specification (classification)

Main Arguments:
  trees = 100
  tree_depth = 3
  learn_rate = 0.0316227766016838

Computational engine: xgboost

set.seed(203)
final_boost_fit <-
  final_boost_wf |>
  last_fit(data_split)

final_boost_fit |>
  collect_metrics()

# A tibble: 3 x 4
  .metric      .estimator .estimate .config
  <chr>      <chr>      <dbl> <chr>
1 accuracy    binary      0.574 Preprocessor1_Model11

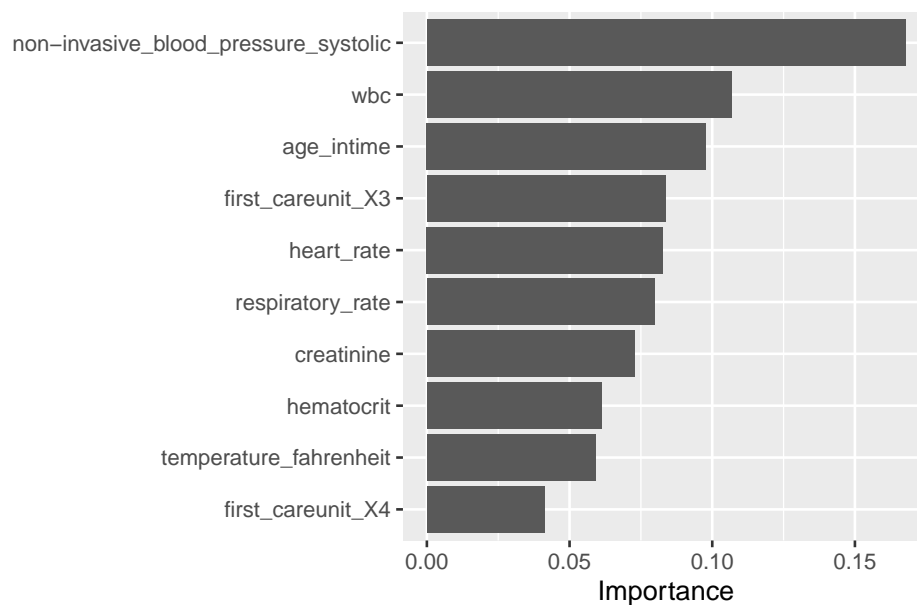
```

```
2 roc_auc      binary      0.609 Preprocessor1_Model11
3 brier_class binary      0.242 Preprocessor1_Model11
```

0.1.2.5 Visualize the final model

```
final_boost_tree <- extract_workflow(final_boost_fit)

bt_vip <- final_boost_tree |>
  extract_fit_parsnip() |>
  vip()
bt_vip
```



```
rm(final_boost_tree,
   final_boost_wf,
   param_grid, bt_mod, bt_wf)
```

0.1.3 Random Forest

0.1.3.1 Model & Workflow

```
rf_mod <-
  rand_forest(
    mode = "classification",
    mtry = tune(),
    trees = tune()
  ) |>
```

```
set_engine("ranger", importance = "impurity")
rf_mod
```

Random Forest Model Specification (classification)

Main Arguments:

```
mtry = tune()
trees = tune()
```

Engine-Specific Arguments:

```
importance = impurity
```

Computational engine: ranger

```
rf_wf <- workflow() |>
  add_recipe(data_recipe) |>
  add_model(rf_mod)
rf_wf
```

== Workflow =====

Preprocessor: Recipe

Model: rand_forest()

-- Preprocessor -----

3 Recipe Steps

```
* step_dummy()
* step_zv()
* step_normalize()
```

-- Model -----

Random Forest Model Specification (classification)

Main Arguments:

```
mtry = tune()
trees = tune()
```

Engine-Specific Arguments:

```
importance = impurity
```

Computational engine: ranger

0.1.3.2 Tuning grid

```
rf_grid <- grid_regular(
  trees(range = c(100L, 300L)),
```

```

    mtry(range = c(3L, 5L)),
    levels = c(2, 2)
  )

rf_grid

```

```

# A tibble: 4 x 2
  trees mtry
  <int> <int>
1   100     3
2   300     3
3   100     5
4   300     5

```

0.1.3.3 Cross-validation

```
set.seed(203)
```

```

if (file.exists("rf_res.rds")) {
  rf_res <- read_rds("rf_res.rds")
} else {
  rf_res <- rf_wf |>
    tune_grid(
      resamples = folds,
      grid = rf_grid,
      metrics = metric_set(roc_auc, accuracy),
      control = control_stack_grid()
    )
  write_rds(rf_res, "rf_res.rds")
}

rf_res

```

```

# Tuning results
# 2-fold cross-validation
# A tibble: 2 x 5
  splits          id   .metrics      .notes      .predictions
  <list>         <chr> <list>      <list>      <list>
1 <split [23610/23611]> Fold1 <tibble [8 x 6]> <tibble [0 x 3]> <tibble>
2 <split [23611/23610]> Fold2 <tibble [8 x 6]> <tibble [0 x 3]> <tibble>

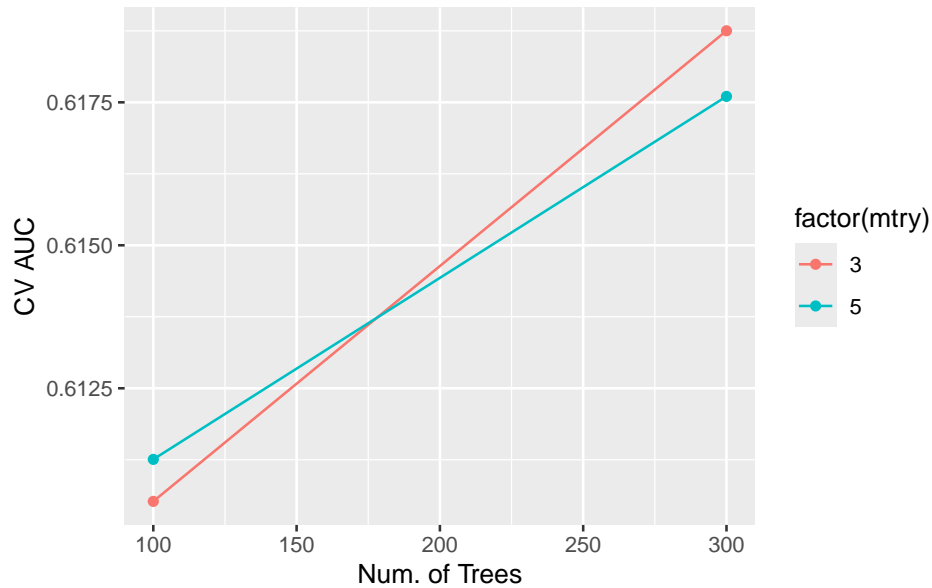
rf_res |>
  collect_metrics() |>
  print(width = Inf) |>
  filter(.metric == "roc_auc") |>
  ggplot(mapping = aes(x = trees, y = mean, color = factor(mtry))) +
  geom_point() +
  geom_line() +

```

```
labs(x = "Num. of Trees", y = "CV AUC")
```

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

	mtry	trees	.metric	.estimator	mean	n	std_err	.config
	<int>	<int>	<chr>	<chr>	<dbl>	<int>	<dbl>	<chr>
1	3	100	accuracy	binary	0.577	2	0.000118	Preprocessor1_Model1
2	3	100	roc_auc	binary	0.611	2	0.000893	Preprocessor1_Model1
3	3	300	accuracy	binary	0.586	2	0.00170	Preprocessor1_Model2
4	3	300	roc_auc	binary	0.619	2	0.00242	Preprocessor1_Model2
5	5	100	accuracy	binary	0.579	2	0.00111	Preprocessor1_Model3
6	5	100	roc_auc	binary	0.611	2	0.00207	Preprocessor1_Model3
7	5	300	accuracy	binary	0.585	2	0.00238	Preprocessor1_Model4
8	5	300	roc_auc	binary	0.618	2	0.00292	Preprocessor1_Model4



0.1.3.4 Finalize the model

```
rf_res |>
  show_best(metric = "roc_auc")
```

```
# A tibble: 4 x 8
```

	mtry	trees	.metric	.estimator	mean	n	std_err	.config
	<int>	<int>	<chr>	<chr>	<dbl>	<int>	<dbl>	<chr>
1	3	300	roc_auc	binary	0.619	2	0.00242	Preprocessor1_Model2
2	5	300	roc_auc	binary	0.618	2	0.00292	Preprocessor1_Model4
3	5	100	roc_auc	binary	0.611	2	0.00207	Preprocessor1_Model3
4	3	100	roc_auc	binary	0.611	2	0.000893	Preprocessor1_Model1

```

#Select best model
best_rf <- rf_res |>
  select_best(metric = "roc_auc")
best_rf

# A tibble: 1 x 3
  mtry trees .config
  <int> <int> <chr>
1     3   300 Preprocessor1_Model2

final_rf <- finalize_workflow(rf_wf, best_rf)

final_rf_fit <- final_rf |>
  last_fit(data_split)

final_rf_fit

# Resampling results
# Manual resampling
# A tibble: 1 x 6
  splits          id          .metrics .notes   .predictions .workflow
  <list>         <chr>         <list>  <list>  <list>       <list>
1 <split [47221/47223]> train/test sp~ <tibble> <tibble> <tibble>    <workflow>

final_rf_fit |>
  collect_metrics()

# A tibble: 3 x 4
  .metric      .estimator .estimate .config
  <chr>       <chr>       <dbl> <chr>
1 accuracy    binary       0.583 Preprocessor1_Model1
2 roc_auc     binary       0.617 Preprocessor1_Model1
3 brier_class binary       0.240 Preprocessor1_Model1

```

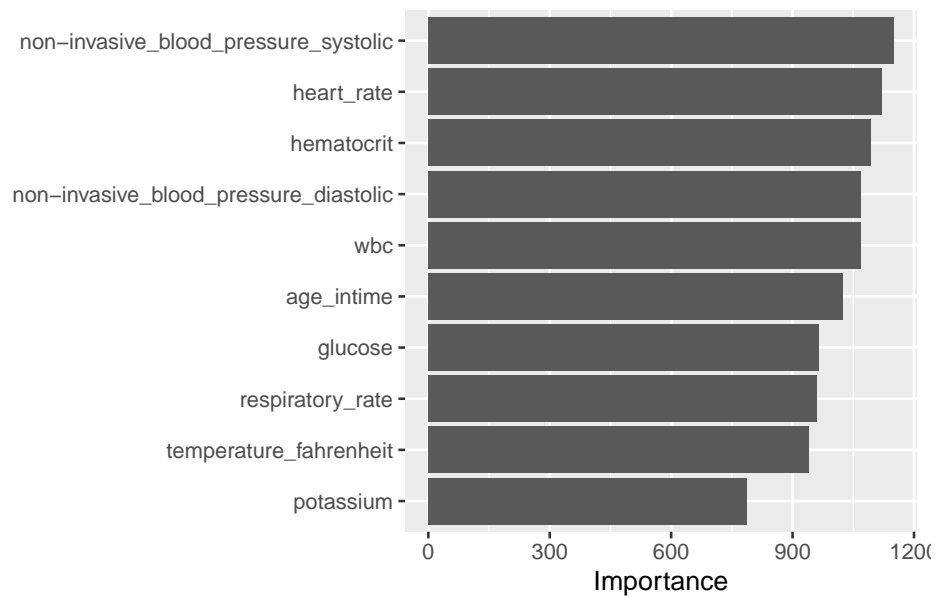
0.1.3.5 Visualize the final model

```

# Extract fitted model
final_rf_tree <- extract_workflow(final_rf_fit)

rf_vip <- final_rf_tree |>
  extract_fit_parsnip() |>
  vip()
rf_vip

```



```
rm(rf_mod, rf_wf, rf_grid,
    folds, rf_res, final_rf,
    final_rf_tree)
```

4. Compare model classification performance on the test set. Report both the area under ROC curve and accuracy for each machine learning algorithm and the model stacking. Interpret the results. What are the most important features in predicting long ICU stays? How do the models compare in terms of performance and interpretability?

See Stacking.qmd file.