

EST-25134: Aprendizaje Estadístico

Profesor: Alfredo Garbuno Iñigo — Primavera, 2023 — Flujo de diagnóstico.

Objetivo: Que veremos.

Lectura recomendada: Capítulo 10 de [1] y capítulo 15 de [2].

1. INTRODUCCIÓN

Hemos discutido ya sobre distintos modelos y cómo cada modelo tiene distintas necesidades para pre-procesar los datos antes de realizarse el ajuste. En el capítulo de 10 de Kuhn and Johnson [1] se ajustan varios modelos para predecir la capacidad de compresión de mezclas de concreto en función de los ingredientes que se utiliza para cada mezcla. Las preguntas en concreto que resolveremos en esta sección son:

¿Cómo podemos comparar distintos modelos entre sí? ¿Cómo podemos utilizar un flujo de trabajo que nos ayude a hacerlo de manera eficiente?

Los datos que usaremos para ilustrar estos conceptos son los mismos que usan [1] donde lo que nos interesa es predecir `compressive_strength` y las unidades son kilogramos por metro cúbico.

```
1 library(tidymodels)
2 data(concrete, package = "modeldata")
3 concrete > print(n = 3, width = 70)
```

```
1 # A tibble: 1,030 × 9
2   cement blast_...1f fly_ash water ...2super ...3coars fine_... age ...compr
3   <dbl>      <dbl>   <dbl> <dbl>   <dbl>   <dbl>   <dbl> <int>   <dbl>
4 1    540          0       0   162     2.5    1040    676    28    80.0
5 2    540          0       0   162     2.5    1055    676    28    61.9
6 3   332.      142.       0   228     0      932    594   270    40.3
7 # ... with 1,027 more rows, and abbreviated variable names
8 # 1 blast_furnace_slag, 2superplasticizer, 3coarse_aggregate,
9 #   fine_aggregate, compressive_strength
10 # Use 'print(n = ...)' to see more rows
```

En particular para estos datos tenemos mezclas que se probaron varias veces por lo tanto reduciremos un poco esta multiplicidad.

```
1 concrete <-
2   concrete %>%
3   group_by(across(-compressive_strength)) %>%
4   summarize(compressive_strength = mean(compressive_strength),
5             .groups = "drop")
```

Prepararemos nuestros conjuntos de entrenamiento y prueba

```
1 set.seed(1501)
2 concrete_split <- initial_split(concrete, strata = compressive_strength)
3 concrete_train <- training(concrete_split)
4 concrete_test  <- testing(concrete_split)
5
6 set.seed(1502)
7 concrete_folds <-
8   vfold_cv(concrete_train, strata = compressive_strength, repeats = 5)
```

Usaremos algunas preparaciones de datos, pues hay modelos (no todos) que las requieren

```
1 normalized_rec <-
2   recipe(compressive_strength ~ ., data = concrete_train) %>%
3   step_normalize(all_predictors())
4
5 poly_recipe <-
6   normalized_rec %>%
7   step_poly(all_predictors()) %>%
8   step_interact(~ all_predictors():all_predictors())
```

Preparemos nuestras especificaciones de modelos

```
1 library(rules)
2 library(bagette)
3
4 linear_reg_spec <-
5   linear_reg(penalty = tune(), mixture = tune()) %>%
6   set_engine("glmnet")
7
8 mars_spec <-
9   mars(prod_degree = tune()) %>% #- use GCV to choose terms
10  set_engine("earth") %>%
11  set_mode("regression")
```

```
1 cart_spec <-
2   decision_tree(cost_complexity = tune(), min_n = tune()) %>%
3   set_engine("rpart") %>%
4   set_mode("regression")
5
6 bag_cart_spec <-
7   bag_tree() %>%
8   set_engine("rpart", times = 50L) %>%
9   set_mode("regression")
```

```
1 rf_spec <-
2   rand_forest(mtry = tune(), min_n = tune(), trees = 1000) %>%
3   set_engine("ranger") %>%
4   set_mode("regression")
5
6 xgb_spec <-
7   boost_tree(tree_depth = tune(), learn_rate = tune(), loss_reduction = tune()
8     ,
9     min_n = tune(), sample_size = tune(), trees = tune()) %>%
10  set_engine("xgboost") %>%
```

```
10   set_mode("regression")
11
12   cubist_spec <-
13     cubist_rules(committees = tune(), neighbors = tune()) %>%
14     set_engine("Cubist")
```

2. FLUJO DE PROCESAMIENTO

REFERENCIAS

- [1] M. Kuhn and K. Johnson. *Applied Predictive Modeling*. Springer New York, New York, NY, 2013. ISBN 978-1-4614-6848-6 978-1-4614-6849-3. . [1](#)
- [2] M. Kuhn and J. Silge. *Tidy Modeling with R*. O'Reilly Media, Inc., 2022. [1](#)