# 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 contreto que resolveremos en esta sección son:

> ¿Cómo podemos comparar distintos modelos entre si? ¿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.

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

```
# A tibble: 1,030 × 9
   cement blast_...¹f fly_ash water ...²super ...³coars fine_...   age ...compr
    <dbl>      <dbl>   <dbl> <dbl>    <dbl>    <dbl>    <dbl> <int>   <dbl>
1   540          0       0   162      2.5     1040      676    28    80.0
2   540          0       0   162      2.5     1055      676    28    61.9
3   332.       142.      0   228      0        932      594   270    40.3
# ... with 1,027 more rows, and abbreviated variable names
# ¹  blast_furnace_slag, ²superplasticizer, ³coarse_aggregate,
#   fine_aggregate, compressive_strength
#  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.

```
concrete ←
  concrete %>%
  group_by(across(-compressive_strength)) %>%
  summarize(compressive_strength = mean(compressive_strength),
            .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())
```

Prepararemos nuestras especificaciones de modelos

```
1  library(rules)
2  library(baguette)
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  knn_spec ←
2    nearest_neighbor(neighbors = tune(),
3                     dist_power = tune(),
4                     weight_func = tune()) %>%
5    set_engine("kknn") %>%
6    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(),
8              loss_reduction = tune(),
9              min_n = tune(), sample_size = tune(),
10             trees = tune()) %>%
11   set_engine("xgboost") %>%
12   set_mode("regression")
13
14 cubist_spec ←
15   cubist_rules(committees = tune(), neighbors = tune()) %>%
16   set_engine("Cubist")
```

## 2. SELECCIÓN Y COMPARACIÓN

### 2.1. Flujo de procesamiento

```
1  normalized ←
2    workflow_set(
3      preproc = list(normalized = normalized_rec),
4      models = list(KNN = knn_spec)
5    )
6  normalized
```

```
1  # A workflow set/tibble: 1 × 4
2    wflow_id        info               option      result
3    <chr>           <list>             <list>      <list>
4  1 normalized_KNN <tibble [1 × 4]> <opts[0]> <list [0]>
```

Podemos corroborar que tenemos lo usual

```
1  normalized %>% extract_workflow(id = "normalized_KNN")
```

```
1  == Workflow ========================================================================
2  Preprocessor: Recipe
3  Model: nearest_neighbor()
4
5  -- Preprocessor ----------------------------------------------------------------
6  1 Recipe Step
7
8  - step_normalize()
9
10 -- Model -----------------------------------------------------------------------
11 K-Nearest Neighbor Model Specification (regression)
12
13 Main Arguments:
14   neighbors = tune()
15   weight_func = tune()
16   dist_power = tune()
17
18 Computational engine: kknn
```

ITAM

Para los demás modelos podemos utilizar `dplyr` para definir `respuesta` y `atributos`.

```
model_vars <-
  workflow_variables(outcomes = compressive_strength,
                     predictors = everything())
no_pre_proc <-
  workflow_set(
    preproc = list(simple = model_vars),
    models = list(MARS = mars_spec, CART = cart_spec,
                  CART_bagged = bag_cart_spec,
                  RF = rf_spec, boosting = xgb_spec,
                  Cubist = cubist_spec)
  )
no_pre_proc
```

```
# A workflow set/tibble: 6 × 4
  wflow_id          info            option     result
  <chr>             <list>          <list>     <list>
1 simple_MARS       <tibble [1 × 4]> <opts[0]> <list [0]>
2 simple_CART       <tibble [1 × 4]> <opts[0]> <list [0]>
3 simple_CART_bagged <tibble [1 × 4]> <opts[0]> <list [0]>
4 simple_RF         <tibble [1 × 4]> <opts[0]> <list [0]>
5 simple_boosting   <tibble [1 × 4]> <opts[0]> <list [0]>
6 simple_Cubist     <tibble [1 × 4]> <opts[0]> <list [0]>
```

Agregamos el conjunto de modelos usan términos no lineales e interacciones.

```
with_features <-
  workflow_set(
    preproc = list(full_quad = poly_recipe),
    models = list(linear_reg = linear_reg_spec, KNN = knn_spec)
  )
```

Finalmente, creamos el conjunto completo de procesamiento

```
all_workflows <-
  bind_rows(no_pre_proc, normalized, with_features) %>%
  ## Make the workflow ID's a little more simple:
  mutate(wflow_id = gsub("(simple_)|(normalized_)", "", wflow_id))
all_workflows
```

```
# A workflow set/tibble: 9 × 4
  wflow_id             info            option     result
  <chr>                <list>          <list>     <list>
1 MARS                 <tibble [1 × 4]> <opts[0]> <list [0]>
2 CART                 <tibble [1 × 4]> <opts[0]> <list [0]>
3 CART_bagged          <tibble [1 × 4]> <opts[0]> <list [0]>
4 RF                   <tibble [1 × 4]> <opts[0]> <list [0]>
5 boosting             <tibble [1 × 4]> <opts[0]> <list [0]>
6 Cubist               <tibble [1 × 4]> <opts[0]> <list [0]>
7 KNN                  <tibble [1 × 4]> <opts[0]> <list [0]>
8 full_quad_linear_reg <tibble [1 × 4]> <opts[0]> <list [0]>
9 full_quad_KNN        <tibble [1 × 4]> <opts[0]> <list [0]>
```

## 2.2. Ajuste y evaluación de modelos

Casi todos los modelos tienen parámetros que se tienen que ajustar. Podemos utilizar los métodos de ajuste que ya hemos visto (`tune_grid()`, etc.). Con la función `workflow_map()` se aplica la misma función para **todos** los flujos de entrenamiento.

Usaremos las mismas opciones para cada uno. Es decir, 25 candidatos en cada modelo para validación cruzada, utilizando la misma separación en bloques.

```
grid_ctrl ←
  control_grid(
    save_pred = TRUE,
    parallel_over = "everything",
    save_workflow = TRUE
  )
```

```
all_cores ← parallel::detectCores(logical = TRUE) - 3
library(doParallel)
cl ← makePSOCKcluster(all_cores)
registerDoParallel(cl)
```

```
system.time(
  grid_results ←
    all_workflows %>%
    workflow_map(
      seed = 1503,
      resamples = concrete_folds,
      grid = 25,
      control = grid_ctrl
    )
)
```

```
i Creating pre-processing data to finalize unknown parameter: mtry
    user    system   elapsed
  26.698     5.506 2083.210
```
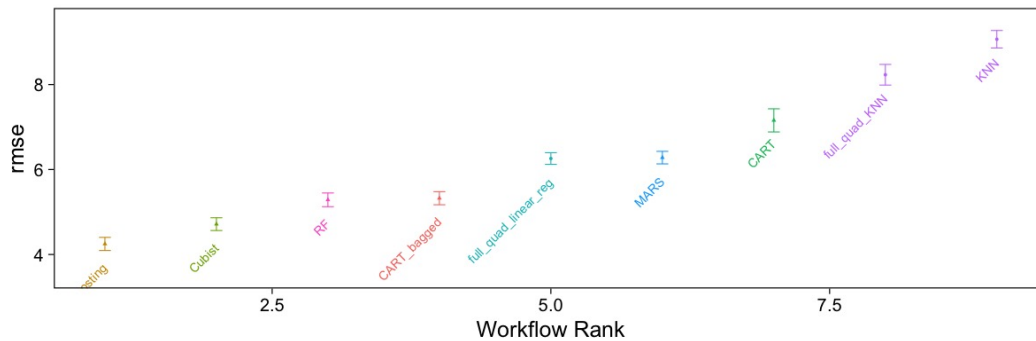
```
grid_results %>%
 rank_results() %>%
 filter(.metric == "rmse") %>%
 select(model, .config, rmse = mean, rank)
```

```
# A tibble: 177 × 4
   model        .config                 rmse   rank
   <chr>        <chr>                  <dbl>  <int>
 1 boost_tree   Preprocessor1_Model04   4.25      1
 2 boost_tree   Preprocessor1_Model06   4.29      2
 3 boost_tree   Preprocessor1_Model13   4.31      3
 4 boost_tree   Preprocessor1_Model14   4.39      4
 5 boost_tree   Preprocessor1_Model16   4.46      5
 6 boost_tree   Preprocessor1_Model03   4.47      6
 7 boost_tree   Preprocessor1_Model15   4.48      7
 8 boost_tree   Preprocessor1_Model05   4.55      8
 9 boost_tree   Preprocessor1_Model20   4.71      9
```

```
13  10 cubist_rules Preprocessor1_Model24   4.71     10
14  # ... with 167 more rows
15  #  Use 'print(n = ...)' to see more rows
```



```
1  library(finetune)
2
3  race_ctrl ←
4    control_race(
5      save_pred = TRUE,
6      parallel_over = "everything",
7      save_workflow = TRUE
8    )
```

```
1  system.time(
2    race_results ←
3      all_workflows %>%
4      workflow_map(
5        "tune_race_anova",
6        seed = 1503,
7        resamples = concrete_folds,
8        grid = 25,
9        control = race_ctrl
10     ))
```

```
1  i Creating pre-processing data to finalize unknown parameter: mtry
2     user   system elapsed
3  157.602    6.237 678.471
```

```
1  race_results
```

```
1  # A workflow set/tibble: 9 × 4
2    wflow_id              info              option       result
3    <chr>                 <list>            <list>       <list>
4  1 MARS                  <tibble [1 × 4]> <opts[3]> <race[+]>
5  2 CART                  <tibble [1 × 4]> <opts[3]> <race[+]>
6  3 CART_bagged           <tibble [1 × 4]> <opts[3]> <rsmp[+]>
7  4 RF                    <tibble [1 × 4]> <opts[3]> <race[+]>
8  5 boosting              <tibble [1 × 4]> <opts[3]> <race[+]>
9  6 Cubist                <tibble [1 × 4]> <opts[3]> <race[+]>
```
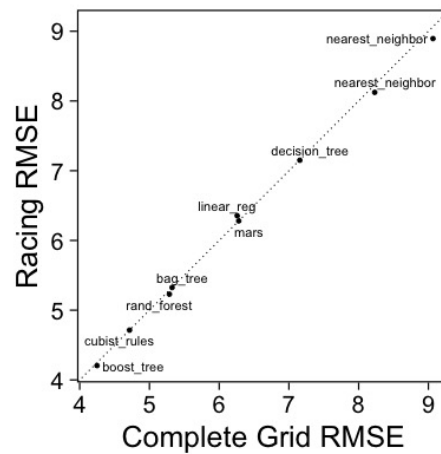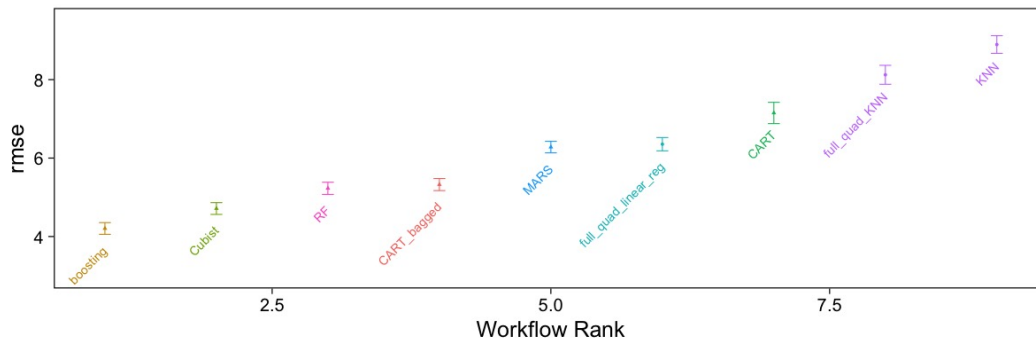
```
10  7 KNN                    <tibble [1 × 4]> <opts[3]> <race[+]>
11  8 full_quad_linear_reg   <tibble [1 × 4]> <opts[3]> <race[+]>
12  9 full_quad_KNN          <tibble [1 × 4]> <opts[3]> <race[+]>
```





## 2.3. Finalizar modelo

```r
1  best_results ←
2    race_results %>%
3    extract_workflow_set_result("boosting") %>%
4    select_best(metric = "rmse")
5  best_results
```

```r
1  # A tibble: 1 × 7
2    trees min_n tree_depth learn_rate loss_reduction sample_size .config
3    <int> <int>      <int>      <dbl>          <dbl>       <dbl> <chr>
4  1  1957     8          7     0.0756    0.000000145       0.679 Preprocessor1_
     Model04
```
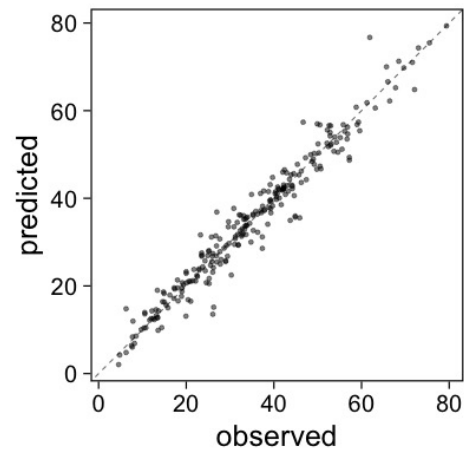
```r
1  boosting_test_results ←
2    race_results %>%
3    extract_workflow("boosting") %>%
4    finalize_workflow(best_results) %>%
5    last_fit(split = concrete_split)
```

```
1  collect_metrics(boosting_test_results)
```

```
1  # A tibble: 2 × 4
2    .metric .estimator .estimate .config
3    <chr>   <chr>          <dbl> <chr>
4  1 rmse    standard        3.43 Preprocessor1_Model1
5  2 rsq     standard       0.953 Preprocessor1_Model1
```



## 3.  MODELO DE ENSAMBLE

## 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