

EST-25134: Aprendizaje Estadístico

Profesor: Alfredo Garbuno Iñigo — Primavera, 2023 — Ensamble de modelos.

Objetivo: En competencias de modelado predictivo es usual observar que combinaciones de modelos logran los mejores resultados. En esta sección del curso explicaremos por qué y presentaremos una implementación de esto de la librería `tidymodels`.

Lectura recomendada: Capítulo 10 de [2] y capítulo 15 de [3]. La última sección del libro de [1] menciona las garantías estadísticas que presentamos en esta clase.

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 [2] se ajustan varios modelos para predecir la capacidad de compresión de mezclas de concreto en función de los ingredientes que se utilizan para cada mezcla. Las preguntas 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 [2] 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, prueba y agendamos nuestro presupuesto de datos en un proceso de validación cruzada.

```

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 instrucciones de pre-procesamiento 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 knn_spec <-
2   nearest_neighbor(neighbors = tune(),
3                     dist_power = tune(),
4                     weight_func = tune()) ▷
5   set_engine("kknn") ▷
6   set_mode("regression")

```

```

1 linear_reg_spec <-
2   linear_reg(penalty = tune(), mixture = tune()) ▷
3   set_engine("glmnet")

```

```

1 library(bagette)
2 cart_spec <-
3   decision_tree(cost_complexity = tune(), min_n = tune()) ▷
4   set_engine("rpart") ▷
5   set_mode("regression")
6
7 bag_cart_spec <-
8   bag_tree() ▷
9   set_engine("rpart", times = 50L) ▷
10  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")
```

2. SELECCIÓN Y COMPARACIÓN

La idea de poder comparar distintos modelos es poder tener una idea de qué estructura de modelos tienen un mejor desempeño. Para después hacer más optimizaciones en las estructuras que funcionan.

Usaremos la infraestructura de `tidymodels` para poder utilizar una interfase para ejecutar un flujo de trabajo ordenado.

2.1. Flujo de procesamiento

Empezamos combinando la *receta* estandarizadora con el modelo adecuado.

```
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 -- Preprocessor -----
5 1 Recipe Step
6 - step_normalize()
7 -- Model -----
8 K-Nearest Neighbor Model Specification (regression)
9 Main Arguments:
10   neighbors = tune()
11   weight_func = tune()
12   dist_power = tune()
13 Computational engine: kkn
```

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

```
1 model_vars <- workflow_variables(
2   outcomes = compressive_strength,
3   predictors = everything()
4 )
```

```

1 no_pre_proc <- workflow_set(
2   preproc = list(simple = model_vars),
3   models = list(CART = cart_spec,
4                 CART_bagged = bag_cart_spec,
5                 RF = rf_spec,
6                 boosting = xgb_spec)
7 )
8 no_pre_proc

```

```

1 # A workflow set/tibble: 4 × 4
2   wflow_id      info      option      result
3   <chr>         <list>    <list>    <list>
4 1 simple_CART   <tibble [1 × 4]> <opts[0]> <list [0]>
5 2 simple_CART_bagged <tibble [1 × 4]> <opts[0]> <list [0]>
6 3 simple_RF     <tibble [1 × 4]> <opts[0]> <list [0]>
7 4 simple_boosting <tibble [1 × 4]> <opts[0]> <list [0]>

```

Agregamos otro conjunto de modelos que usen términos no lineales e interacciones.

```

1 with_features <-
2   workflow_set(
3     preproc = list(fullquad = poly_recipe),
4     models = list(linear_reg = linear_reg_spec, KNN = knn_spec)
5   )

```

Finalmente, creamos el conjunto completo de procesamiento (preparación, entrenamiento, evaluación)

```

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

```

```

1 # A workflow set/tibble: 7 × 4
2   wflow_id      info      option      result
3   <chr>         <list>    <list>    <list>
4 1 CART         <tibble [1 × 4]> <opts[0]> <list [0]>
5 2 CART_bagged  <tibble [1 × 4]> <opts[0]> <list [0]>
6 3 RF           <tibble [1 × 4]> <opts[0]> <list [0]>
7 4 boosting     <tibble [1 × 4]> <opts[0]> <list [0]>
8 5 KNN          <tibble [1 × 4]> <opts[0]> <list [0]>
9 6 fullquad_linear_reg <tibble [1 × 4]> <opts[0]> <list [0]>
10 7 fullquad_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.

La idea de este proceso es ilustrar un mecanismo para condensar en una misma ejecución lo que hemos visto a lo largo de todo el curso. Si, cada modelo tiene distintos hiperparámetros pero de momento nos concentraremos en explorar capacidades predictivas.

```
1 grid_ctrl ←
2   control_grid(
3     save_pred = TRUE,
4     parallel_over = "everything",
5     save_workflow = TRUE
6   )
```

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

```
1 system.time(
2   grid_results ← all_workflows ▷
3     workflow_map(
4       seed = 1503,
5       resamples = concrete_folds,
6       grid = 25,
7       control = grid_ctrl
8     )
9 )
```

```
1 i Creating pre-processing data to finalize unknown parameter: mtry
2   user      system elapsed
3 16.785      2.903 1418.576
```

El tibble del flujo se actualiza con las leyendas en `option` y `results`. Los indicadores `tune[+]` y `rsmp[+]` significan que no hubo problemas para procesar ese modelo.

```
1 grid_results
```

```
1 # A workflow set/tibble: 7 × 4
2   wflow_id      info      option      result
3   <chr>         <list>    <list>    <list>
4 1 CART         <tibble [1 × 4]> <opts [3]> <tune [+]>
5 2 CART_bagged  <tibble [1 × 4]> <opts [3]> <rsmp [+]>
6 3 RF           <tibble [1 × 4]> <opts [3]> <tune [+]>
7 4 boosting     <tibble [1 × 4]> <opts [3]> <tune [+]>
8 5 KNN          <tibble [1 × 4]> <opts [3]> <tune [+]>
9 6 fullquad_linear_reg <tibble [1 × 4]> <opts [3]> <tune [+]>
10 7 fullquad_KNN  <tibble [1 × 4]> <opts [3]> <tune [+]>
```

Por último, con la función `rank_results()` ordenamos los modelos de acuerdo a su capacidad predictiva. La opción `select.best` nos muestra dentro de cada familia el mejor modelo para ordenar por capacidad predictiva.

```

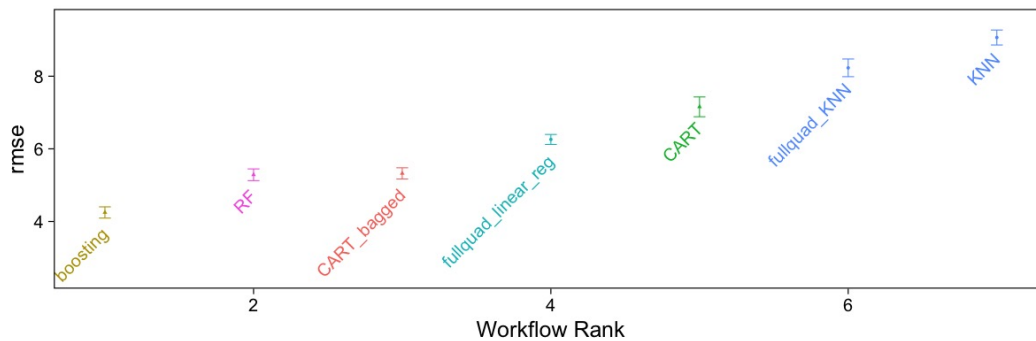
1 grid_results >
2   rank_results(select_best = TRUE) >
3   filter(.metric == "rmse") >
4   select(model, .config, rmse = mean, rank)

```

```

1 # A tibble: 7 × 4
2   model          .config          rmse  rank
3   <chr>         <chr>         <dbl> <int>
4 1 boost_tree   Preprocessor1_Model104  4.25    1
5 2 rand_forest Preprocessor1_Model118  5.29    2
6 3 bag_tree     Preprocessor1_Model11  5.32    3
7 4 linear_reg   Preprocessor1_Model116  6.26    4
8 5 decision_tree Preprocessor1_Model119  7.16    5
9 6 nearest_neighbor Preprocessor1_Model118  8.23    6
10 7 nearest_neighbor Preprocessor1_Model116  9.07    7

```



2.3. Ajuste y comparación eficiente

Utilizaremos el mismo proceso eficiente de comparación de modelos para determinar la mejor configuración dentro de cada uno.

```

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 96.917   1.464 374.666

```

El método ajusta 11 modelos de los 151 posibles. Es decir, sólo requiere ajustar el 7.3% .

```

1 race_results

```

```

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

```

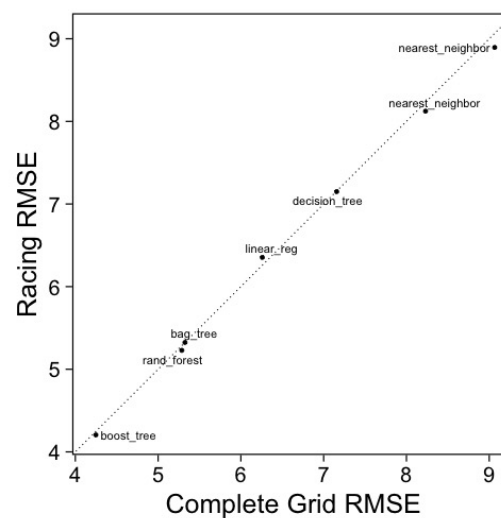
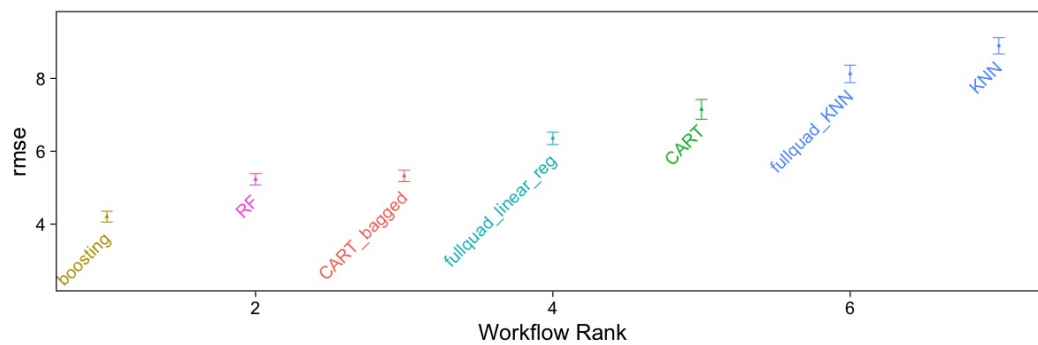


FIGURA 1. Comparación de procedimiento completo contra paro acelerado.

2.4. Finalizar modelo

```

1 best_results <-
2   race_results >
3   extract_workflow_set_result("boosting") >
4   select_best(metric = "rmse")
5 best_results

```

```

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

```

```

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      3.41 Preprocessor1_Model1
5 2 rsq      standard      0.954 Preprocessor1_Model1

```

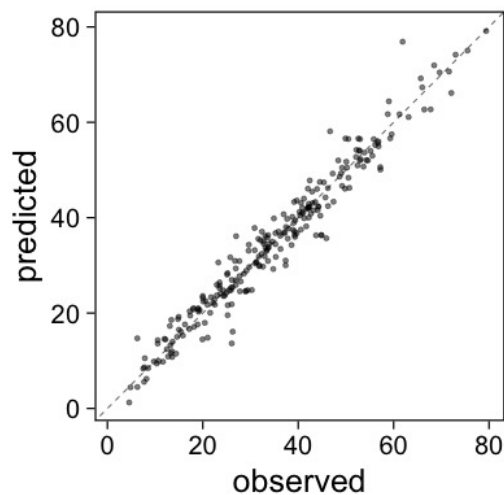


FIGURA 2. Predicciones contra datos reales.

3. ENSAMBLE DE MODELOS

Consideremos que tenemos una colección de modelos entrenados con los cuales podemos obtener predicciones $\hat{f}_1(x), \dots, \hat{f}_M(x)$. Considerando pérdida cuadrática podemos buscar pesos $\omega = (\omega_1, \dots, \omega_M)$ tales que

$$\hat{\omega} = \arg \min_{\omega \in \mathbb{R}_+^M} \mathbb{E}_{\mathcal{D}_n, Y} \left(Y - \sum_{m=1}^M \omega_m \hat{f}_m(x) \right)^2, \quad (1)$$

donde el valor esperado se calcula considerando la variabilidad en conjuntos de entrenamiento y la aleatoriedad en la respuesta para un atributo fijo.

La solución de este problema será la solución **poblacional** considerando el vector que acopla las predicciones de una familia de modelos $\hat{F}^\top = [\hat{f}_1(x), \dots, \hat{f}_M(x)]$:

$$\hat{\omega} = \mathbb{E}[\hat{F}(x) \hat{F}^\top(x)]^{-1} \mathbb{E}[\hat{F}(x) Y]. \quad (2)$$

La combinación lineal, sabemos, tiene mejor capacidad predictiva que cualquier modelo

$$\mathbb{E} \left(Y - \sum_{m=1}^M \hat{\omega}_m \hat{f}_m(x) \right)^2 \leq \mathbb{E} \left(Y - \hat{f}_m(x) \right)^2 \quad \forall m. \quad (3)$$

El problema es que no tenemos acceso a la solución poblacional del problema de regresión. Así que podemos usar el concepto de **stacking** para construir un estimador basado en predicciones fuera de muestra.

3.1. Stacking de modelos

Una manera que tenemos para darle la vuelta al problema de antes es utilizar el mecanismo de **validación cruzada** para ajustar los pesos de la siguiente manera.

Denotemos por $\hat{f}_m^{-i}(x)$ la salida del m -ésimo modelo en x el cual fue entrenado sin la observación i -ésima. Los coeficientes se ajustan por medio de mínimos cuadrados utilizando la regresión de respuestas y_i con atributos $\hat{f}_m^{-i}(x)$. De tal forma que tenemos

$$\hat{\omega}^{\text{stack}} = \arg \min_{\omega \in \mathbb{R}^M} \sum_{i=1}^N \left(y_i - \sum_{m=1}^M \omega_m \hat{f}_m^{-i}(x) \right)^2. \quad (4)$$

Las predicciones se realizan por medio de

$$\sum_{m=1}^M \hat{\omega}_m^{\text{stack}} \hat{f}_m^{-i}(x). \quad (5)$$

El resultado es una combinación lineal de predicciones. En la práctica se obtienen los mejores resultados si se restringen los pesos de la combinación lineal a ser no-negativos y que sumen 1.

4. ILUSTRACIÓN NUMÉRICA

Retomaremos nuestra colección de modelos que hemos ajustado. La intención es poder crear una combinación de éstos para mejorar nuestra capacidad predictiva. Para estos usaremos **stacks** de **tidymodels**.

```
1 library(stacks)
```

4.0.1. Para pensar:

1. En el contexto de validación cruzada con K bloques: ¿cuántas predicciones fuera de muestra tenemos para cada observación con el conjunto de entrenamiento?
2. En el contexto de validación cruzada con K bloques: ¿cuántas predicciones fuera de muestra tenemos para cada observación con el conjunto de entrenamiento si repetimos validación cruzada B veces?
3. En el contexto de los modelos que hemos usado: ¿cuántas predicciones fuera de muestra tenemos para *bagging*?

4.1. Construyendo nuestra colección de predicciones

Le pasamos nuestros modelos entrenados a un *stack* vacío. La función `add_candidates()` se encarga de filtrar modelos tienen **todo** el perfil predictivo para cada observación en el conjunto de entrenamiento.

```
1 concrete_stack ←
2   stacks() ▷
3   add_candidates(race_results)
4
5 concrete_stack
```

```
1 # A data stack with 7 model definitions and 13 candidate members:
2 #   CART: 1 model configuration
3 #   CART_bagged: 1 model configuration
4 #   RF: 1 model configuration
5 #   boosting: 1 model configuration
6 #   KNN: 3 model configurations
7 #   full_quad_linear_reg: 5 model configurations
8 #   full_quad_KNN: 1 model configuration
9 # Outcome: compressive_strength (numeric)
```

4.1.1. Para pensar: ¿Qué pasaría si en lugar de pasar resultados de ANOVA usamos los resultados de la función `tune_grid()`?

4.2. Mezcla de predicciones

Entrenamos nuestro *meta*-modelo utilizando las predicciones fuera de muestra. En esta situación debemos de considerar:

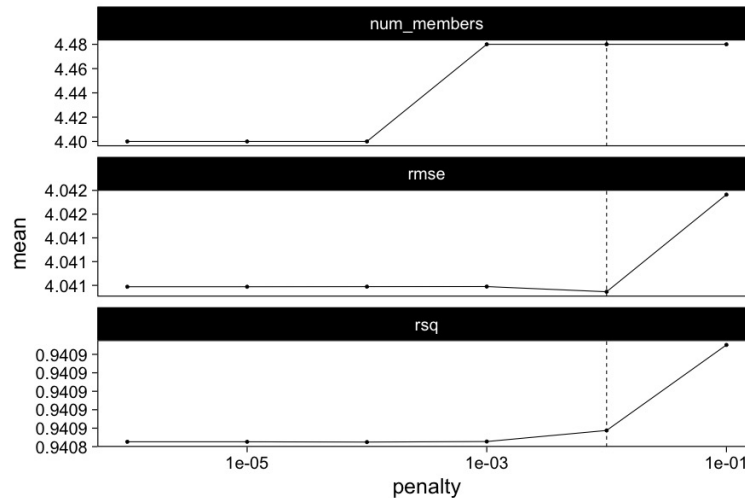
1. Las predicciones entre familias de predictores pueden estar altamente correlacionadas.
2. Habrá predictores que no son necesarios si ya hay algún elemento de la misma familia.

4.2.1. Para pensar: ¿Qué estrategia de regularización hemos visto para resolver estos problemas?

4.2.2. Implementación La función `stacks::blend_predictions()` nos permite ajustar un modelo lineal regularizado que es capaz de evitar estos problemas.

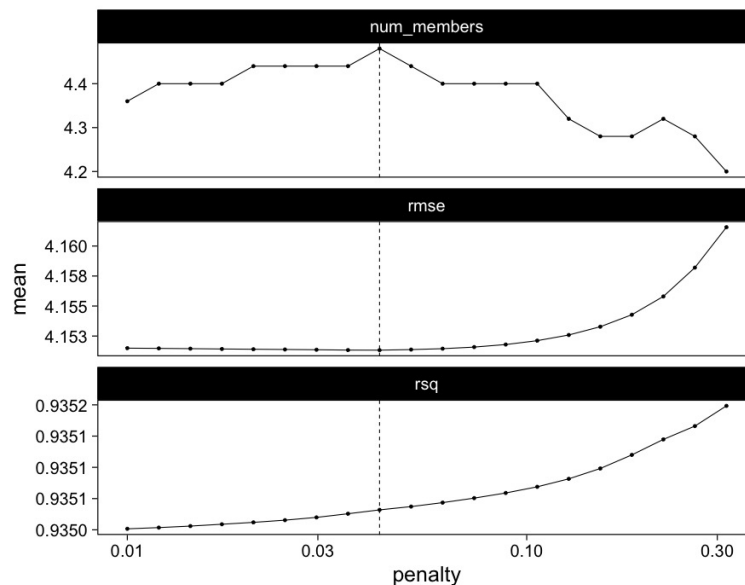
```
1 set.seed(2001)
2 ens ← blend_predictions(concrete_stack)
```

El procedimiento realiza un ajuste interno con remuestreo que permite construir trazar la curva de error predictivo como función de factor de penalización.



Al igual que en validación cruzada los resultados obtenidos nos pueden ayudar a concentrar nuestros esfuerzos computacionales en zonas de mayor interés.

```
1 set.seed(2002)
2 ens <- blend_predictions(concrete_stack, penalty = 10^seq(-2, -0.5, length =
  20))
```



La combinación final queda construida de la siguiente manera. x

```
1 ens
```

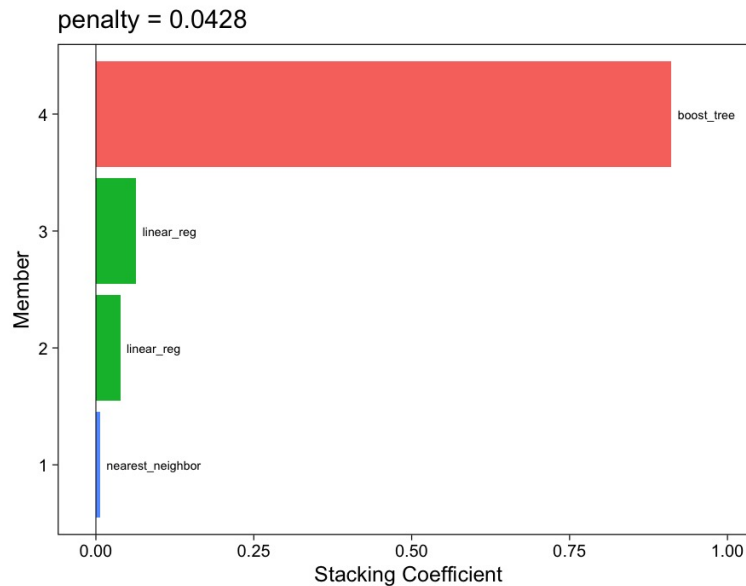
```
1 -- A stacked ensemble model -----
2 Out of 13 possible candidate members, the ensemble retained 4.
3 Penalty: 0.0428133239871939.
4 Mixture: 1.
5
6 The 4 highest weighted members are:
7 # A tibble: 4 × 3
```

```

8 member          type          weight
9 <chr>            <chr>         <dbl>
10 1 boosting_1_04    boost_tree    0.911
11 2 fullquad_linear_reg_1_17 linear_reg    0.0638
12 3 fullquad_linear_reg_1_16 linear_reg    0.0387
13 4 KNN_1_12         nearest_neighbor 0.00704
14
15 Members have not yet been fitted with 'fit_members()'.

```

Las contribuciones de cada modelo se pueden resumir de manera gráfica con la función de `autoplot()` [†].



4.2.3. *Para pensar:* Por último hay que ajustar los modelos finales a todo el conjunto de entrenamiento. ¿En nuestro ejemplo cuántos modelos se entrenan? ¿Cómo cambia esto cuando lo comparamos con el contexto de validación cruzada?

```
1 ens <- fit_members(ens)
```

```

1 reg_metrics <- metric_set(rmse, rsq)
2 ens_test_pred <-
3   predict(ens, concrete_test) >
4   bind_cols(concrete_test)
5
6 ens_test_pred >
7   reg_metrics(compressive_strength, .pred)

```

```

1 # A tibble: 2 × 3
2   .metric .estimator .estimate
3   <chr>    <chr>         <dbl>
4 1 rmse    standard         3.36
5 2 rsq     standard         0.956

```

5. CONCLUSIONES

- Las mejoras pueden ser marginales comparadas contra las del mejor modelo individual.
- En la práctica la mezcla de modelos vuelve el resultado muy poco interpretable.
- Para las tareas de interés predictivo (sin explicaciones y sin restricciones computacionales) son lo mejor.
- Una combinación lineal de predictores es sólo una forma de combinar.

REFERENCIAS

- [1] T. Hastie, R. Tibshirani, and J. Friedman. *The Elements of Statistical Learning*. Springer Series in Statistics. Springer New York, New York, NY, 2009. ISBN 978-0-387-84857-0 978-0-387-84858-7. . [1](#)
- [2] 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](#)
- [3] M. Kuhn and J. Silge. *Tidy Modeling with R*. O'Reilly Media, Inc., 2022. [1](#)