

Aprendizaje Estadístico Supervisado

Bruno Tancredi

2024

Laboratorio remuestreo

- Este laboratorio vamos a ver como aplicar distintos métodos de remuestreo.
- Usaremos varios paquetes, `rsample` para crear conjuntos de datos basados en remuestreo y `yardstick` para calcular métricas de performance.
- Usaremos `rsample` y `dials` para la selección de los valores de los parámetros de tuneo.

Laboratorio remuestreo

Vamos a seguir un ejemplo del blog de Julia Silge [link blog](#)

¿Conocen a Julia?

Importante

Julia Silge va a participar en esta edición de LatinR!

Paquetes que usaremos

```
1 library(tidymodels) #Recordar que al importar tidymodels estamos im  
2 library(schrute) #Contiene los datos que vamos a utilizar, es necesari  
3 library(vip) #Importancia de variables  
4 set.seed(1234)
```

Datos

Qué son los datos

- Utilizaremos el conjunto de datos **theoffice** del paquete **schrute**. Cada fila del conjunto de datos es una línea de diálogo.
- Queremos predecir la calificación del episodio obtenida en la plataforma **IMDB** en base a predictores que están presentes en el conjunto de datos y que iremos construyendo.

.

Datos de calificación

```
1 office_info <- theoffice |>
2   select(season, episode_name, director, writer,
3   character, text, imdb_rating)
4 office_info |> head()
```

A tibble: 6 × 7

	season	episode_name	director	writer	character	text	imdb_rating
	<int>	<chr>	<chr>	<chr>	<chr>	<chr>	<dbl>
1	1	Pilot	Ken Kwapis	Ricky Gervais;Step...	Michael	All ...	7.6
2	1	Pilot	Ken Kwapis	Ricky Gervais;Step...	Jim	Oh, ...	7.6
3	1	Pilot	Ken Kwapis	Ricky Gervais;Step...	Michael	So y...	7.6
4	1	Pilot	Ken Kwapis	Ricky Gervais;Step...	Jim	Actu...	7.6
5	1	Pilot	Ken Kwapis	Ricky Gervais;Step...	Michael	All ...	7.6
6	1	Pilot	Ken Kwapis	Ricky Gervais;Step...	Michael	Yes,...	7.6

Creación de atributos

- Cantidad de lineas que tiene cada personaje en el episodio.
- Guionistas/Directores (no diferenciaremos el rol) que participan en el episodio.

Lineas en el capítulo por personaje.

```
1 characters <- office_info |>
2   count(episode_name, character) |> # Cantidad de lineas por cap
3   add_count(character, wt = n, name = "character_count") |> #Ca
4   filter(character_count > 800) |> #Obtenemos los personajes que
5   select(-character_count) |>
6   pivot_wider(
7     names_from = character,
8     values_from = n,
9     values_fill = list(n = 0)
10  ) #Obtenemos una fila sola por capitulo, ponemos 0 si el person
```

Directores/Esritores del capítulo.

```
1 creators <- office_info |>
2   distinct(episode_name, director, writer) |>
3   pivot_longer(director:writer, names_to = "role", values_to = "person") |>
4   separate_rows(person, sep = ";") |>
5   add_count(person) |>
6   mutate(person = case_when(
7     n <= 10 ~ 'Guest',
8     n > 10 ~ person
9   )) |>
10  distinct(episode_name, person) |>
11  mutate(person_value = 1) |>
12  pivot_wider(
13    names_from = person,
14    values_from = person_value,
15    values_fill = list(person_value = 0)
16  )
```

Conjunto final

```
1 office <- office_info |>
2   distinct(season, episode_name, imdb_rating) |>
3   inner_join(characters) |>
4   inner_join(creators) |>
5   mutate_at("season", as.factor)
```

Exploración

```
1 office |>  
2   ggplot(aes(season, imdb_rating, fill = as.factor(season))) +  
3   geom_boxplot(show.legend = FALSE)
```

División de datos

- Dividimos los datos en train/test. Tratando de conservar la proporción de datos por temporada.

```
1 office_split <- initial_split(office,  
2                               strata = season,  
3                               prop = 3/4)  
4 office_train <- training(office_split)  
5 office_test  <- testing(office_split)
```

Modelo

recipe

Permite especificar una serie de pasos de preprocesamiento de datos, facilitando la creación de flujos de trabajo reproducibles.

```
1 office_rec <- recipe(imdb_rating ~ ., data = office_train) |>  
2   update_role(episode_name, new_role = "ID") |> #Nos permite c  
3   step_dummy(season) |> #codificamos categorias a columnas bi  
4   step_normalize(all_numeric(), -all_outcomes()) #Normalizamos
```

Más info de [step_functions](#) **link**

parsnip

Proporciona una interface ordenada y unificada de modelos que puede ser usada para aplicar una amplia variedad de modelos sin preocuparnos por las variaciones en la sintaxis de los distintos paquetes.

- mixture=1 para Lasso, mixture=0 para Ridge.
- La implementación de glmnet en parsnip tiene ciertas peculiaridades, consultar **documentación**
- Con `show_engines("linear_reg")` vemos todas los posibles paquetes para implementar el tipo de modelo de `parsnip linear_reg`

```
1 tune_spec <- linear_reg(penalty = tune(), mixture = 1) |> #Con tu
2   set_engine("glmnet")
```


workflow

Proporciona una manera de combinar una especificación de modelo, una receta de preparación de datos y un conjunto de parámetros de ajuste del modelo en un solo objeto. Esto facilita la gestión de todo el proceso de modelado.

- Cuando usamos workflow la llamada a `fit()` realiza el preprocesamiento de los datos y el ajuste del modelo.

```
1 tune_wf <- workflow() |>  
2   add_recipe(office_rec) |>  
3   add_model(tune_spec)
```

Ajuste de hiperparámetros

- Si queremos encontrar el mejor valor para λ implica el tuneo de este hiperparámetro y podemos usar k-Fold Cross-Validation.

Para poder hacer validación cruzada necesitamos usar el paquete **tune** y tres cosas para que funcione:

- Un objeto **parsnip** / **workflow** con uno o más argumentos indicados para ser tuneados.
- Un objeto **vfold_cv** **rsample** con las particiones de validación cruzada.
- Una **tibble** que denote los valores de los hiperparámetros a ser explorados.

Ajuste de hiperparámetros

```
1 office_cv <- vfold_cv(office_train, v = 5) #5 particiones debido a c
2 lambda_grid <- grid_regular(penalty(c(-10,-1)), levels = 50) #Def
3
4
5 lasso_tune<- tune_grid( #Realizamos el ajuste, pensar que estar
6   tune_wf,
7   resamples = office_cv,
8   grid = lambda_grid
9 )
```

Ajuste de hiperparámetros

```
1 lasso_tune |>  
2 collect_metrics()
```

A tibble: 100 × 7

	penalty	.metric	.estimator	mean	n	std_err	.config
	<dbl>	<chr>	<chr>	<dbl>	<int>	<dbl>	<chr>
1	1	e-10	rmse	standard	0.482	5	0.0315 Preprocessor1_Model01
2	1	e-10	rsq	standard	0.341	5	0.0400 Preprocessor1_Model01
3	1.53e-10	rmse	standard	0.482	5	0.0315 Preprocessor1_Model02	
4	1.53e-10	rsq	standard	0.341	5	0.0400 Preprocessor1_Model02	
5	2.33e-10	rmse	standard	0.482	5	0.0315 Preprocessor1_Model03	
6	2.33e-10	rsq	standard	0.341	5	0.0400 Preprocessor1_Model03	
7	3.56e-10	rmse	standard	0.482	5	0.0315 Preprocessor1_Model04	
8	3.56e-10	rsq	standard	0.341	5	0.0400 Preprocessor1_Model04	
9	5.43e-10	rmse	standard	0.482	5	0.0315 Preprocessor1_Model05	
10	5.43e-10	rsq	standard	0.341	5	0.0400 Preprocessor1_Model05	

i 90 more rows

Ajuste de hiperparámetros

```
1 lasso_tune |>
2   collect_metrics() |>
3   ggplot(aes(penalty, mean, color = .metric)) +
4   geom_errorbar(aes(
5     ymin = mean - std_err,
6     ymax = mean + std_err
7   ),
8   alpha = 0.5
9   ) +
10  geom_line(size = 1.5) +
11  facet_wrap(~.metric, scales = "free", nrow = 2) +
12  scale_x_log10() +
13  theme(legend.position = "none")
```

Ajuste de hiperparámetros

Ajuste de hiperparámetros

Seleccionamos el modelo cuyo RMSE es más bajo. Podríamos haber utilizado otra métrica (**métricas disponibles**). Existen otros criterios, para la selección como `select_by_one_std_err()` que selecciona el modelo más simple cuyo error este dentro de un desvío estándar del modelo de mejor resultado.

```
1 lowest_rmse <- lasso_tune |>  
2   select_best(metric = "rmse")  
3  
4 final_lasso <- finalize_workflow(tune_wf, lowest_rmse)
```

Resultados de test.

```
1 last_fit(  
2   final_lasso,  
3   office_split  
4 ) |>  
5 collect_metrics()
```

A tibble: 2 × 4

	.metric	.estimator	.estimate	.config
	<chr>	<chr>	<dbl>	<chr>
1	rmse	standard	0.393	Preprocessor1_Model1
2	rsq	standard	0.376	Preprocessor1_Model1

Podríamos haber hecho lo mismo con bootstrap!

```
1 office_boot <- bootstraps(office_train, strata = season)
2 lambda_grid <- grid_regular(penalty(c(-10,-1)), levels = 50) #Defi
3
4 lasso_tune <- tune_grid(
5   tune_wf,
6   resamples = office_boot,
7   grid = lambda_grid
8 )
```

Resumen del flujo de trabajo.

.

Diagnóstico de modelo

Importancia de variable

Queremos ver que variables son más importantes a la hora de predecir la calificación. `vi` toma al valor absoluto del coeficiente asociado a la variable (recordar que las variables están normalizadas) como indicador de la importancia.

- Para calcular la importancia de variable fue necesario extraer el modelo y trabajar con el objeto `glmnet`.

```
1 final_lasso |>  
2   fit(office_train) |>  
3   extract_fit_engine() |>  
4   vi(lambda = lowest_rmse$penalty) |> #Es muy importante marca  
5   ggplot(aes(x = Importance, y = reorder(Variable, Importance), fil  
6   geom_col() +  
7   scale_x_continuous(expand = c(0, 0)) +  
8   labs(y = NULL)
```

Importancia de variable

Tu turno!

Haremos el ejercicio 9 del capítulo 6 del libro **ISL**. Predicaremos el número de aplicaciones recibidas usando las otras variables del conjunto de datos **College** como predictores.

1. Cargar el conjunto de datos **College**, del paquete **ISLR2**. Dividir estos datos en entrenamiento/testeo.
2. Ajustar un modelo lineal utilizando mínimos cuadrados. Reportar el error obtenido en el conjunto de testeo.
3. Ajustar un modelo lineal utilizando Ridge. Escoger λ con validación cruzada. Reportar el error obtenido en el conjunto de testeo.

Tu turno!

4. Ajustar un modelo lineal utilizando Lasso. Escoger λ con validación cruzada. Reportar el error obtenido en el conjunto de testeo.
5. Ajustar un modelo de vecinos más cercanos. Escoger k con validación cruzada. Reportar el error obtenido en el conjunto de testeo.
6. Comentar los resultados obtenidos. ¿Con cuánta precisión podemos predecir el número de aplicaciones?. ¿Hay un mejor modelo? ¿Hay mucha diferencia en el error de testeo de los modelos?