

EST-24107: Simulación

Profesor: Alfredo Garbuno Iñigo — Primavera, 2022 — Aplicaciones de *bootstrap*.

Objetivo: Que veremos.

Lectura recomendada: Referencia.

1. BOOTSTRAP PARAMÉTRICO

- Supongamos que tenemos una muestra $X_1, \dots, X_N \stackrel{\text{iid}}{\sim} \mathbb{P}(x; \theta^*)$. Es decir, tenemos un modelo paramétrico que da lugar a nuestros datos.
- En este tipo de problemas de inferencia suponemos la familia paramétrica

$$\mathcal{P}_\Theta = \{\mathbb{P}(\cdot; \theta) : \theta \in \Theta\}, \quad (1)$$

donde Θ denota el **espacio parametral** (los posibles valores de los parámetros de un modelo).

- En esta tarea no conocemos el valor específico de θ^* . Por lo tanto, lo tenemos que estimar. Usualmente a través de resolver un problema de optimización

$$\hat{\theta}_{\text{MLE}} = \arg \max_{\theta \in \Theta} \prod_{i=1}^N \mathbb{P}(X_i; \theta). \quad (2)$$

cuya solución llamamos **estimador de máxima verosimilitud**.

- Adicional, nos encantaría poder establecer una cuantificación de la incertidumbre sobre este valor. En particular, reportar

$$\text{ee}(\hat{\theta}_{\text{MLE}}) = \left(\mathbb{V}(\hat{\theta}_{\text{MLE}}) \right)^{1/2}. \quad (3)$$

- Para algunos modelos es fácil poder estimarlo, utilizando propiedades asintóticas y/o analíticas de nuestros estimadores (lo ven en el curso de Estadística Matemática).
- Sin embargo, ¿qué pasa si nuestro estimador no tiene fórmulas cerradas para el cálculo del error estándar? ¿O si nuestro tamaño de muestra no sugiere que los supuestos del TLC se cumplen?

1.0.1. Definición [Método bootstrap paramétrico]: El error estándar estimado para $\hat{\theta}_{\text{MLE}}$ por medio del *bootstrap* paramétrico se calcula como sigue:

1. Se calcula $\hat{\theta}_{\text{MLE}}$ para la muestra observada.
2. Se simula una muestra iid de tamaño N de $X_1^{(b)}, \dots, X_N^{(b)} \stackrel{\text{iid}}{\sim} \mathbb{P}(x; \hat{\theta}_{\text{MLE}})$ (muestra *bootstrap*).
3. Se recalcula el estimador de máxima verosimilitud para la muestra *bootstrap*, lo cual denotamos por $\hat{\theta}_{\text{MLE}}^{(b)} = s(X_1^{(b)}, \dots, X_N^{(b)})$.
4. Se repiten los pasos 2–3 muchas veces ($B = 1,000 - 10,000$).
5. Se calcula la desviación estándar de los valores $\hat{\theta}_{\text{MLE}}^{(b)}$ obtenidos. Este es el error estándar estimado para el estimador $\hat{\theta}_{\text{MLE}}$.

1.0.2. Observación:

- Nota cómo cambiamos el mecanismo de remuestreo $\hat{\mathbb{P}}_N$ por $\mathbb{P}(x; \hat{\theta}_{\text{MLE}})$.
- En espíritu es lo mismo, pero estamos dispuestos a incorporar mayores supuestos en nuestra tarea de inferencia.

1.1. Ejemplo: Datos normales

Como ejercicio, podemos encontrar los estimadores de máxima verosimilitud cuando tenemos una muestra $X_1, \dots, X_N \stackrel{\text{iid}}{\sim} N(\mu, \sigma^2)$ (puedes derivar e igualar a cero para encontrar el mínimo). También podemos resolver numéricamente.

Supongamos que tenemos la siguiente muestra:

```
1 set.seed(41852)
2 muestra <- rnorm(150, mean = 1, sd = 2)
```

Para la cual podemos calcular los estimadores de máxima verosimilitud de un modelo normal

```
1 mle.obs <- broom::tidy(MASS::fitdistr(muestra, "normal")) >
2   tibble::column_to_rownames("term")
3 mle.obs
```

```
1      estimate std.error
2 mean      1.136    0.1502
3 sd        1.839    0.1062
```

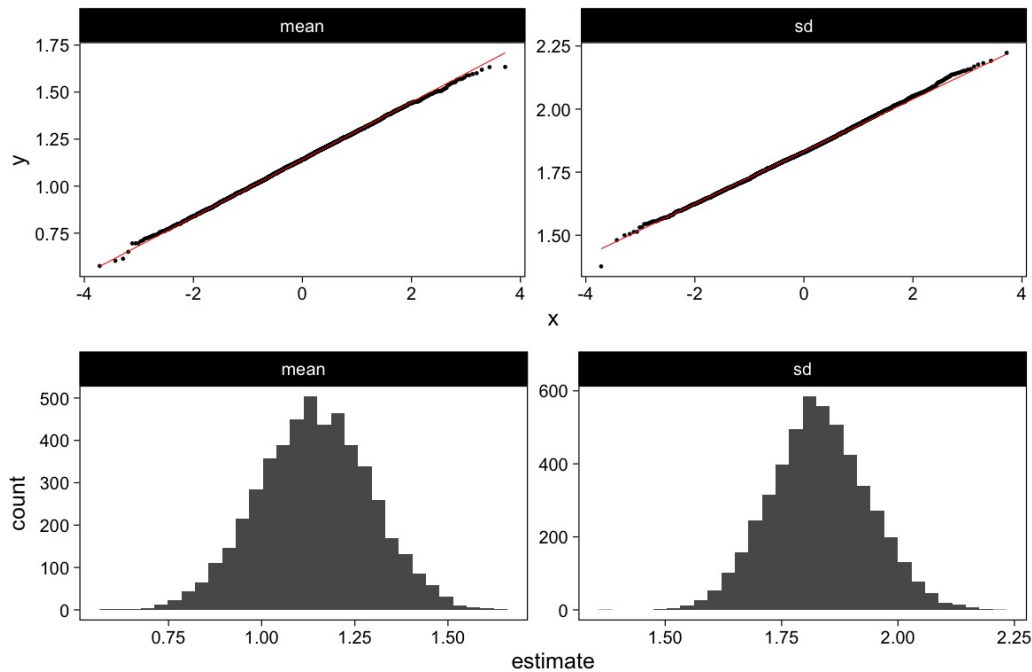
Con esta estimación podemos definir el proceso de remuestreo.

```
1 ## paso 1: define el estimador
2 estimador_mle <- function(datos, modelo = "normal"){
3   datos >
4     MASS::fitdistr(modelo) >
5     broom::tidy() >
6     select(-std.error)
7 }
```

```
1 ## paso 2: define el proceso de remuestreo
2 paramboot_sample <- function(data){
3   rnorm(length(data),
4     mean = mle.obs["mean", "estimate"],
5     sd = mle.obs["sd", "estimate"])
6 }
```

```
1 ## paso 3: define el paso bootstrap
2 paso_bootstrap <- function(id){
3   muestra >
4     paramboot_sample() >
5     estimador_mle()
6 }
```

```
1 ## paso 4: aplica bootstrap parametrico
2 boot_mle <- map_df(1:5000, paso_bootstrap)
```



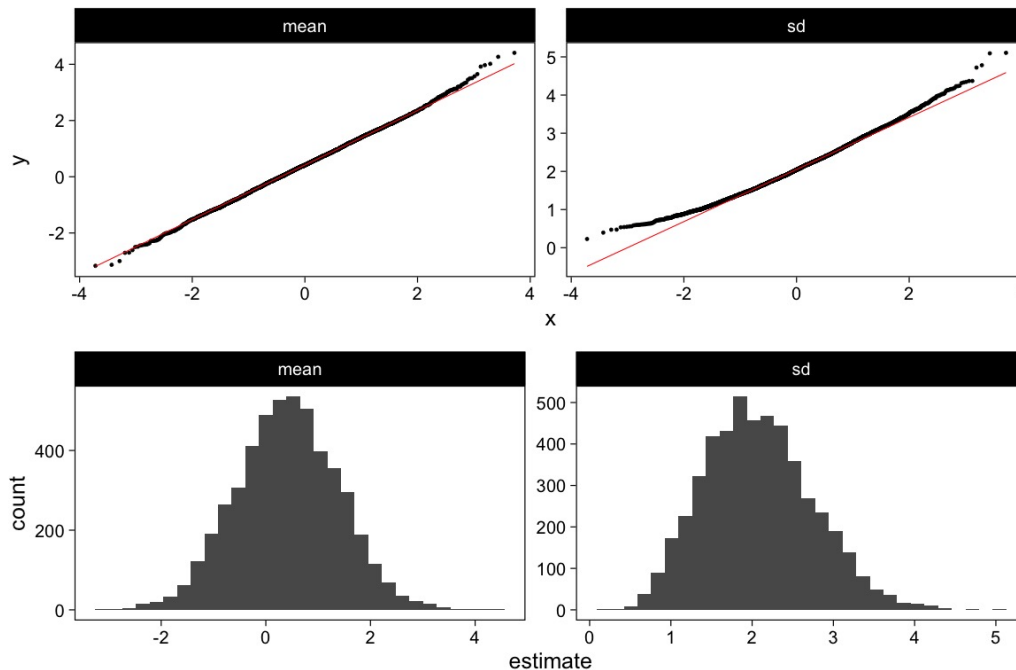
Las distribuciones son aproximadamente normales. Nótese que esto no siempre sucede, especialmente con parámetros de dispersión como σ . (Examina las curvas de nivel del ejemplo de arriba).

Ahora, supongamos que tenemos una muestra más chica. Repasa los pasos para asegurarte que entiendes el procedimiento:

```
1 set.seed(4182)
2 muestra <- rnorm(6, mean = 1, sd = 2)
3 mle.obs <- broom::tidy(MASS::fitdistr(muestra, "normal")) >
4   tibble::column_to_rownames("term")
5 mle.obs
```

```
1      estimate std.error
2 mean    0.3979    0.9794
3 sd      2.3990    0.6925
```

```
1 ## paso 4: aplica bootstrap parametrico
2 boot_mle <- map_df(1:5000, paso_bootstrap)
```



Donde vemos que la distribución de σ tienen sesgo a la derecha, pues en algunos casos obtenemos estimaciones muy cercanas a cero. Podemos usar intervalos de percentiles.

1.2. Comparación bootstrap paramétrico y no paramétrico

```
1 propinas <- read_csv("data/propinas.csv",
2                       progress = FALSE,
3                       show_col_types = FALSE) >
4 mutate(id = 1:244)
```

```
1 ## paso 1: define el estimador
2 estimador <- function(split, ...){
3   muestra <- analysis(split) > group_by(momento)
4   muestra >
5     summarise(estimate = mean(cuenta_total), .groups = 'drop') >
6     mutate(term = momento)
7 }
```

```
1 ## paso 2 y 3: remuestrea y calcula estimador
2 boot_samples <- bootstraps(propinas, strata = momento, 500) >
3 mutate(res_boot = map(splits, estimador))
4 ## paso 4: construye intervalos de confianza
5 intervalos_noparam <- boot_samples >
6 int_pctl(res_boot, alpha = 0.05) >
7 mutate(across(where(is.numeric), round, 2))
8 intervalos_noparam
```

```
1 # A tibble: 2 × 6
2   term   .lower .estimate .upper .alpha .method
3   <chr>   <dbl>   <dbl>   <dbl>   <dbl>   <chr>
```

1 BOOTSTRAP PARAMÉTRICO Comparación bootstrap paramétrico y no paramétrico

```
4 1 Cena      19.7      20.8      22.0      0.1 percentile
5 2 Comida    15.7      17.2      18.7      0.1 percentile
```

```
1 ## paso 1: define estimador
2 estimador_mle_grupos <- function(muestra, modelo = "normal") {
3   muestra >
4     select(momento, cuenta_total) >
5     group_by(momento) >
6     nest(data = cuenta_total) >
7     summarise(mle = map(data, function(x) {
8       nobs <- nrow(x)
9       unlist(x) >
10        estimador_mle(modelo = modelo) >
11        mutate(n = nobs)
12      }))
13 }
```

```
1 mle.obs <- estimador_mle_grupos(propinas, "normal")
2 mle.obs > unnest(mle)
```

```
1 # A tibble: 4 × 4
2   momento term   estimate     n
3   <chr>   <chr>   <dbl> <int>
4 1 Cena   mean     20.8   176
5 2 Cena   sd       9.12   176
6 3 Comida mean     17.2    68
7 4 Comida sd       7.66    68
```

```
1 ## paso 2: define proceso de remuestreo
2 param_boot_grupos <- function(estimadores){
3   estimadores >
4     group_by(momento) >
5     mutate(simulaciones = map(mle, function(m){
6       tibble(cuenta_total = rnorm(m$n[1], m$estimate[1], sd = m$estimate[2]))
7     })) >
8     unnest(simulaciones) >
9     select(-mle) >
10    ungroup()
11 }
```

```
1 ## paso 3: paso bootstrap
2 paso_bootstrap_grupos <- function(id){
3   param_boot_grupos(mle.obs) >
4   estimador_mle_grupos()
5 }
```

```
1 ## paso 4: aplica bootstrap y presenta intervalos
2 intervalos_param <- tibble(id = 1:500) >
3   mutate(estimadores = map(id, paso_bootstrap_grupos)) >
4   unnest(estimadores) >
5   unnest(mle) >
```

```

6 group_by(momento, term) >
7 summarise(.lower = quantile(estimate, 0.025),
8           .estimate = mean(estimate),
9           .upper = quantile(estimate, 0.975),
10          .alpha = .05,
11          .method = "percentile (normal)", .groups = "drop") >
12 filter(term == "mean") > select(-term)
13 intervalos_param

```

```

1 # A tibble: 2 × 6
2 momento .lower .estimate .upper .alpha .method
3 <chr>     <dbl>     <dbl> <dbl> <dbl> <chr>
4 1 Cena      19.6      20.8  22.1   0.1 percentile (normal)
5 2 Comida    15.3      17.1  18.8   0.1 percentile (normal)

```

```

1 # A tibble: 2 × 6
2 term .lower .estimate .upper .alpha .method
3 <chr> <dbl>     <dbl> <dbl> <dbl> <chr>
4 1 Cena      19.7      20.8  22.0   0.1 percentile
5 2 Comida    15.7      17.2  18.7   0.1 percentile

```

```

1 # A tibble: 1 × 6
2 term .lower .estimate .upper .alpha .method
3 <chr> <dbl>     <dbl> <dbl> <dbl> <chr>
4 1 Cena      17.8      20.8  23.9   0.1 percentile (exponential)

```

1.3. Ventajas y desventajas de bootstrap paramétrico

- Ventaja: el *bootstrap* paramétrico puede dar estimadores más precisos e intervalos más angostos y bien calibrados que el no paramétrico, **siempre y cuando el modelo teórico sea razonable**.
- Desventaja: Es necesario decidir el modelo teórico, que tendrá cierto grado de desajuste vs. el proceso generador real de los datos. Si el ajuste es muy malo, los resultados tienen poca utilidad. Para el no paramétrico no es necesario hacer supuestos teóricos.
- Ventaja: el *bootstrap* paramétrico puede ser más escalable que el no paramétrico, pues no es necesario cargar y remuestrear los datos originales, y tenemos mejoras adicionales cuando tenemos expresiones explícitas para los estimadores de máxima verosimilitud (como en el caso normal, donde es innecesario hacer optimización numérica).
- Desventaja: el *bootstrap* paramétrico es conceptualmente más complicado que el no paramétrico, y como vimos arriba, sus supuestos pueden ser más frágiles que los del no paramétrico.

2. ESTRUCTURAS DE DATOS

2.1. Datos agrupados

```

1 departures_raw <- read_csv("https://raw.githubusercontent.com/rfordatascience/
  tidyuesday/master/data/2021/2021-04-27/departures.csv", show_col_types =
  FALSE, progress = FALSE)

```

```

1 departures <- departures_raw %>%
2   filter(departure_code < 9) %>%
3   mutate(involuntary = if_else(departure_code %in% 3:4, "involuntary", "other"
4     )) %>%
5   filter(fyear > 1995, fyear < 2019)
6 departures

```

```

1 df <- departures %>%
2   count(fyear, involuntary) %>%
3   pivot_wider(names_from = involuntary, values_from = n)
4 df

```

```

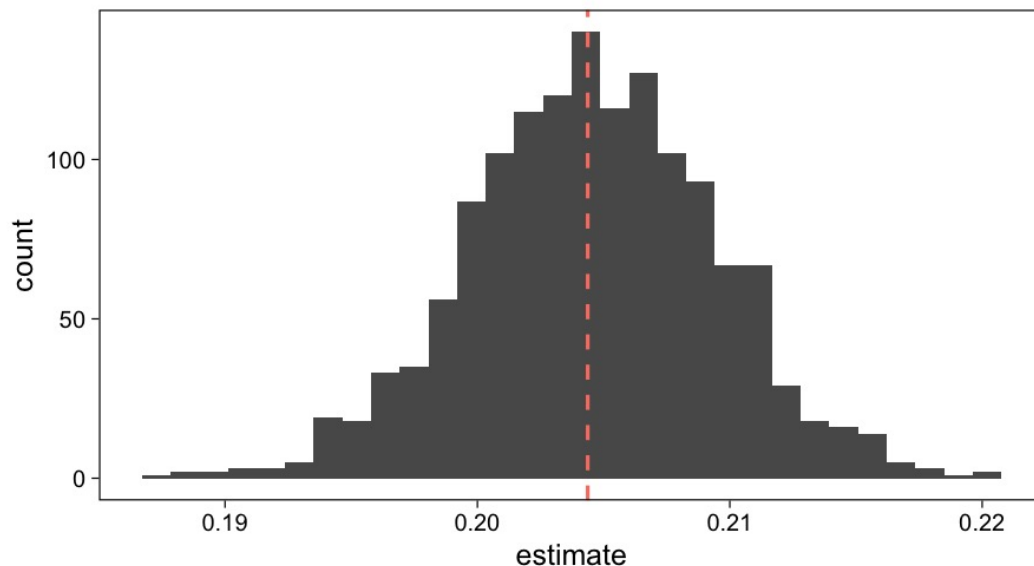
1 ## paso 1: define estimador
2 estimador_razon <- function(split, ...){
3   analysis(split) >
4     count(fyear, involuntary) >
5     pivot_wider(names_from = involuntary, values_from = n) >
6     mutate(prop = involuntary/(involuntary + other)) >
7     summarise(estimate = mean(prop)) >
8     pull(estimate)
9 }

```

```

1 dif_boot <- bootstraps(departures, 1400, apparent = TRUE) >
2   mutate(estimate = map_dbl(splits, estimador_razon))

```



2.2. Muestras independientes

```

1 data_ej <- tibble(respuesta = c(94, 197, 16, 38, 99, 141, 23, 52, 104, 146,
2   10, 51, 30, 40, 27, 46),
3   tipo = c(rep("tratamiento", 7), rep("control", 9)))

```

```

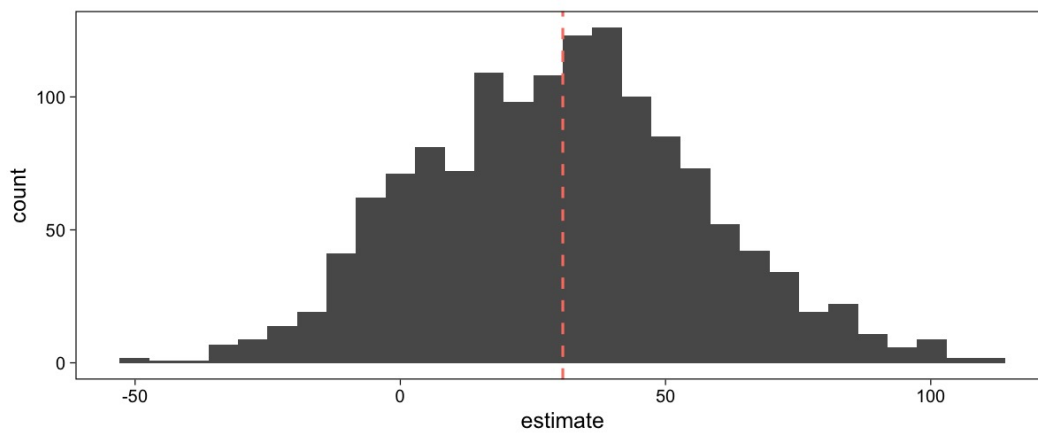
1 diferencia_medias <- function(split, ...){
2   analysis(split) ▷
3   group_by(tipo) ▷
4   summarise(promedio = mean(respuesta)) ▷
5   pivot_wider(names_from = tipo, values_from = promedio) ▷
6   mutate(estimate = tratamiento - control) ▷
7   pull(estimate)
8 }

```

```

1 dif_boot <- bootstraps(data_ej, 1400, strata = tipo, apparent = TRUE) ▷
2   mutate(estimate = map_dbl(splits, diferencia_medias))

```



2.3. Datos pareados

```

1 dormir <- sleep ▷
2   pivot_wider(names_from = group,
3               names_prefix = "medicina_",
4               values_from = extra)
5 dormir

```

```

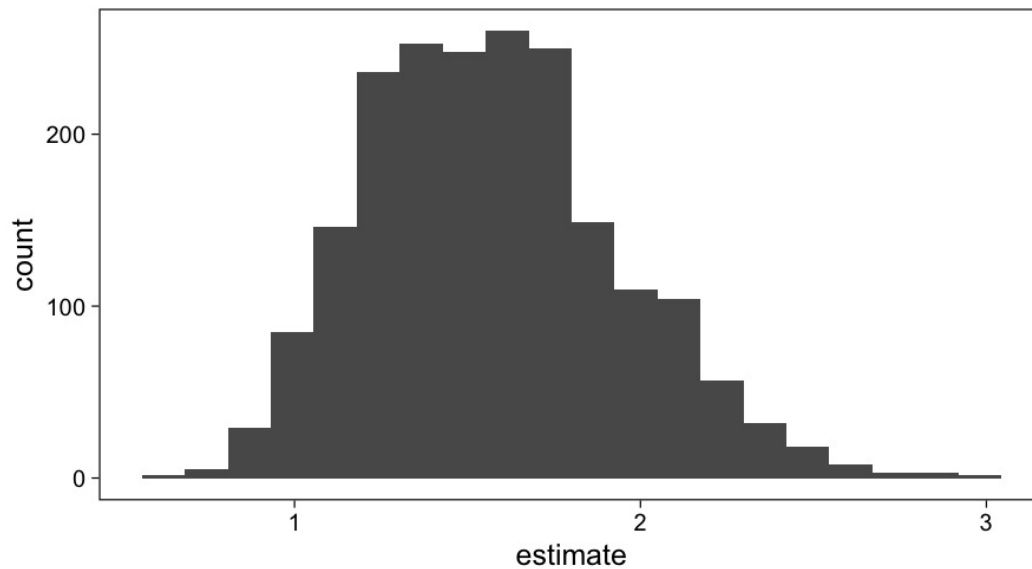
1 estimador_dif_pareados <- function(split, ...){
2   muestra <- analysis(split)
3   muestra %>%
4     mutate(dif_2_menos_1 = medicina_2 - medicina_1) %>%
5     summarise(estimate = mean(dif_2_menos_1), .groups = "drop") %>%
6     mutate(term = "diferencia 2 vs 1")
7 }

```

```

1 pareados_boot <- bootstraps(dormir, 2000, apparent = TRUE) %>%
2   mutate(res_boot = map(splits, estimador_dif_pareados))

```

```
1 pareados_boot %>% int_pctl(res_boot, 0.05)
```

```
1 # A tibble: 1 × 6
2   term          .lower .estimate .upper .alpha .method
3   <chr>         <dbl>   <dbl>  <dbl> <dbl> <chr>
4 1 diferencia 2 vs 1   0.97    1.58   2.39   0.05 percentile
```

2.4. Series temporales

Moving blocks bootstraps

1. origin windows con cumulative apagado
2. Bootstrap sobre los slices

3. MODELOS DE REGRESIÓN

3.1. Observaciones

3.2. Residuales