

# EST-24107: Simulación

**Profesor:** Alfredo Garbuno Iñigo — Primavera, 2022 — *Bootstrap* paramétrico.

**Objetivo:** Que veremos.

**Lectura recomendada:** Capítulo 6 de Efron and Tibshirani [2]. Sección 13.2 de Chihara and Hesterberg [1].

## 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  $\Theta$  denota el **espacio parametral** (los posibles valores de los parámetros de un modelo).

- En esta tarea no conocemos el valor específico de  $\theta^*$ . 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).
- Consideremos ejemplos:
  1. Modelo Poisson,  $X_1, \dots, X_N \stackrel{\text{iid}}{\sim} \text{Poisson}(\lambda)$ .
  2. Modelo Bernoulli,  $X_1, \dots, X_N \stackrel{\text{iid}}{\sim} \text{Bernoulli}(\theta)$ .
  3. Modelo uniforme,  $X_1, \dots, X_N \stackrel{\text{iid}}{\sim} \text{U}(0, \theta)$ .
  4. Modelo normal,  $X_1, \dots, X_N \stackrel{\text{iid}}{\sim} \text{N}(\mu, \sigma^2)$ .
- 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.

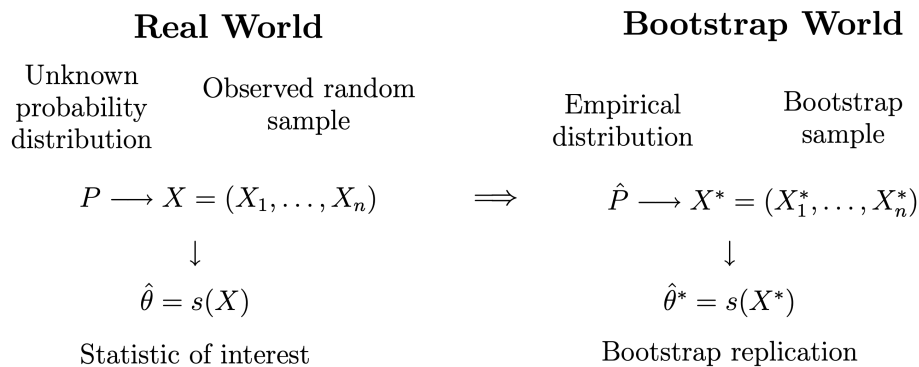


FIGURA 1. Imagen tomada del material del curso de *Cómputo Estadístico* de Michael Eichler.

## 1.1. Ejemplo: Datos normales

Como ejercicio, podemos encontrar los estimadores de máxima verosimilitud cuando tenemos una muestra  $X_1, \dots, X_N \stackrel{\text{iid}}{\sim} \text{N}(\mu, \sigma^2)$  (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 }

```

```

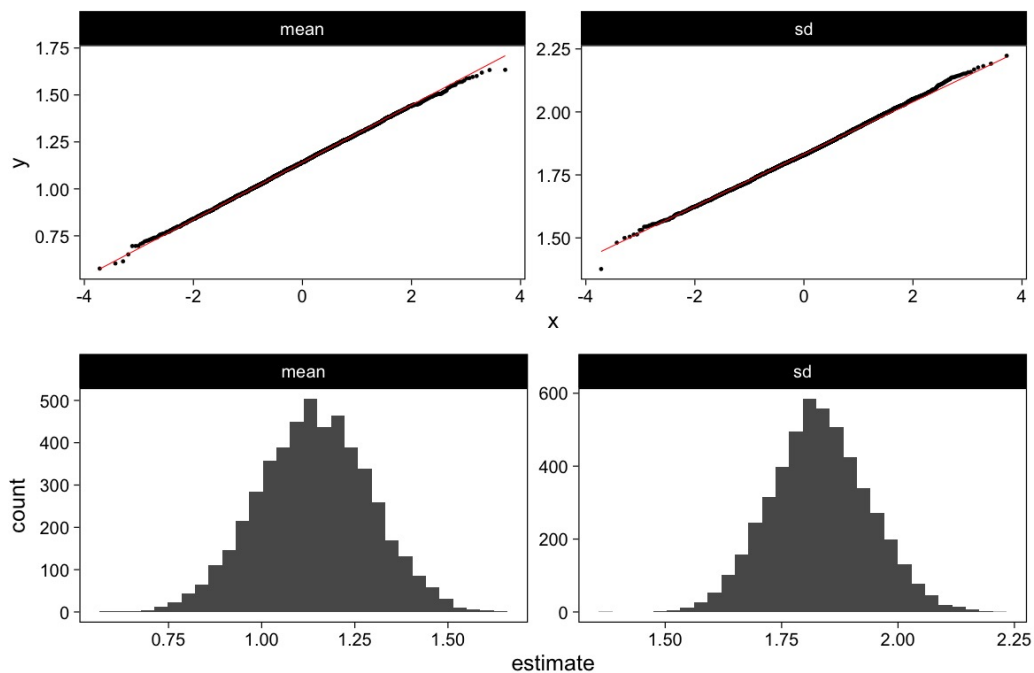
1 ## paso 3: define el paso bootstrap
2 paso_bootstrap <- function(id){
3   muestra >
4   paramboot_sample() >
5   estimador_mle()
6 }

```

```

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

```



## 1.2 Comparación bootstrap paramétrico y no paramétrico

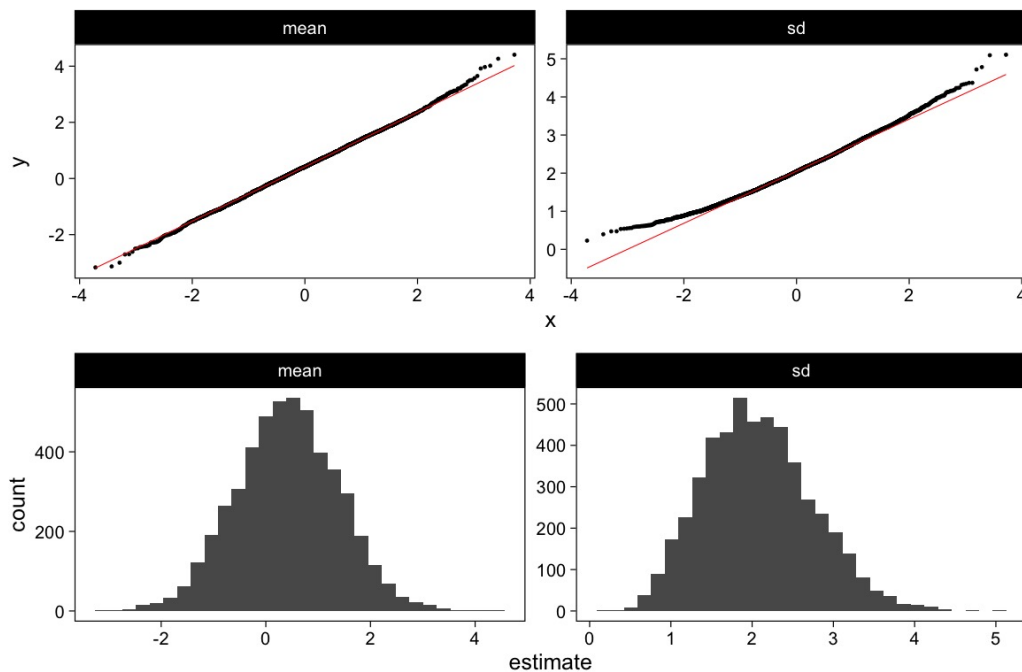
Las distribuciones son aproximadamente normales. Nótese que esto no siempre sucede, especialmente con parámetros de dispersión como  $\sigma$ . (Examina las curvas de nivel del ejemplo de arriba).

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 paramétrico
2 boot_mle <- map_df(1:5000, paso_bootstrap)
```



Donde vemos que la distribución de  $\sigma$  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>
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 >

```

## 1.2 Comparación bootstrap paramétrico y no paramétrico BOOTSTRAP PARAMÉTRICO

```
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)
```

El modelo exponencial nos da intervalos mas anchos (maypr incertidumbre) lo cual ilustra que si el modelo paramétrico no es el adecuado, los supuestos adicionales sirven poco para mejorar la estimación de incertidumbre.

### 1.3. Ejemplo: Datos de viento

Consideremos los siguientes datos que corresponden a datos de producción energética por medio de una turbina de viento. En este caso nos interesa estimar el percentil 10 % pues es lo que esperaríamos que la turbina genere el 90 % de las veces.

```
1 library(resampledData)
2 data(Turbine)
3 Turbine > tibble()
```

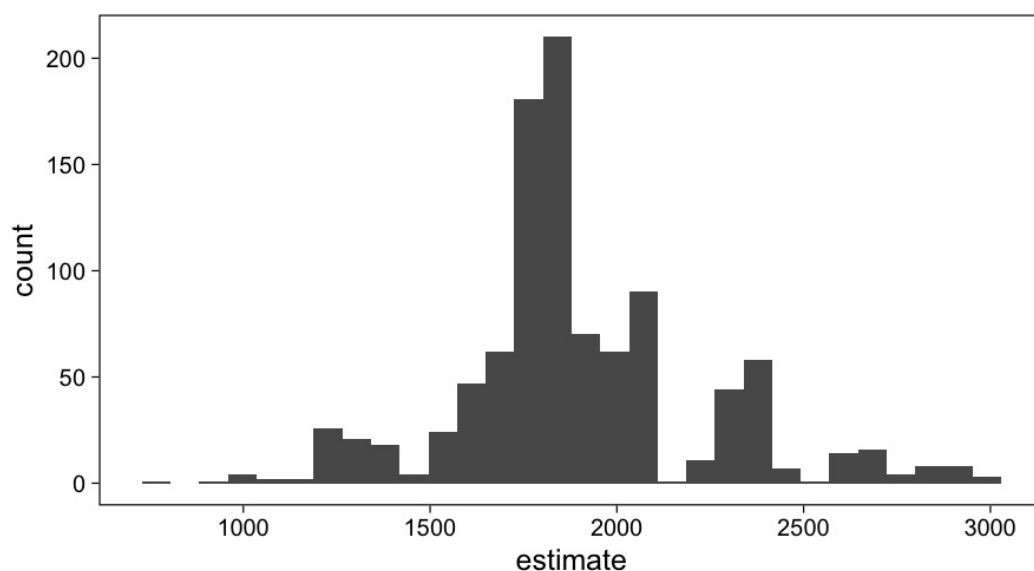
Esperamos los problemas usuales de nuestro estimador si utilizáramos el **bootstrap no paramétrico**.

```
1 Turbine >
2   summarise(estimate = quantile(Production, probs = .1))
```

```
1   estimate
2 1      1817
```

```
1 ## paso 1: define el estimador
2 calcula_percentil <- function(split, ...) {
3   split >
4     analysis() >
5     summarise(estimate = quantile(Production, probs = .1)) >
6     mutate(term = "Percentil")
7 }
```

```
1 nonparam_boot <- bootstraps(Turbine, 1000) >
2   mutate(resultados = map(splits, calcula_percentil))
```



Si asumimos modelo Weibull( $k, \lambda$ ) para los datos y estimamos los parámetros obtenemos lo siguiente. Revisa los pasos para asegurarte que queda claro el procedimiento.

```

1  ## paso 1: define el estimador
2  ajusta_weibull <- function(data){
3    tibble(data) ▷
4    filter(Production > 0) ▷
5    pull(Production) ▷
6    MASS::fitdistr("weibull") ▷
7    broom::tidy() ▷
8    select(-std.error) ▷
9    tibble::column_to_rownames("term")
10 }
11
12 mle.weibull <- ajusta_weibull(Turbine)
13 mle.weibull

```

```

1      estimate
2 shape      1.283
3 scale 11795.041

```

```

1  ## paso 2: define el proceso de remuestreo
2  paramboot_sample <- function(data){
3    tibble(Production = rweibull(nrow(data),
4                                scale = mle.weibull["scale", "estimate"],
5                                shape = mle.weibull["shape", "estimate"])
6  )
7  }

```

```

1  ## paso 1.5: complementa el estimador
2  extrae_cuantil <- function(params){
3    qweibull(scale = params["scale", "estimate"],
4             shape = params["shape", "estimate"],
5             p = .10) %>%
6    tibble(estimate = .)
7  }

```

```

1  ## paso 3: define el paso bootstrap
2  paso_bootstrap <- function(id){
3    Turbine ▷
4    paramboot_sample() ▷
5    ajusta_weibull() ▷
6    extrae_cuantil()
7  }

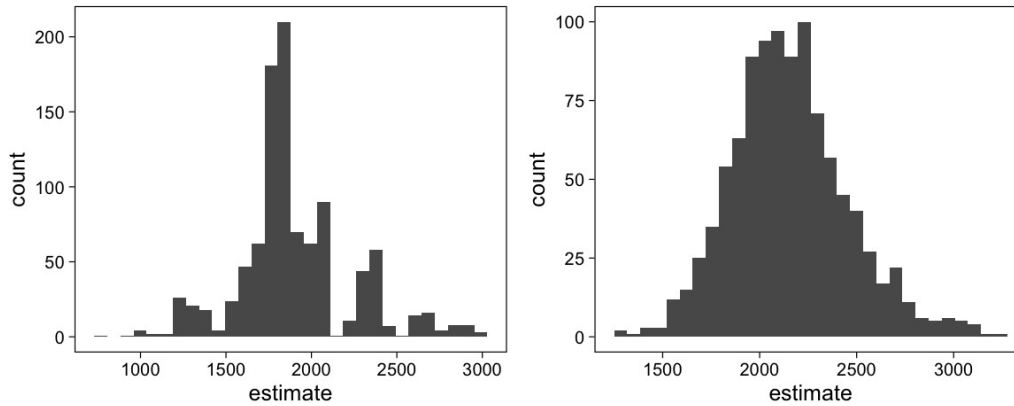
```

```

1  ## paso 4: aplica bootstrap parametrico
2  param_boot <- map_df(1:1000, paso_bootstrap)

```





```

1 # A tibble: 1 × 7
2   term      .lower .estimate .upper .alpha .method      .length
3   <chr>      <dbl>      <dbl>  <dbl> <dbl> <chr>      <dbl>
4 1 Percentil 1261.      1898.  2715.  0.05 percentile 1454.

```

```

1 # A tibble: 1 × 7
2   term      .lower .estimate .upper .alpha .method      .length
3   <chr>      <dbl>      <dbl>  <dbl> <dbl> <chr>      <dbl>
4 1 Percentil 1699.      2155.  2687.  0.05 percentile (param) 988.

```

## 1.4. El método de momentos

Utilizar máxima verosimilitud **no** es la única manera de poder realizar *bootstrap* paramétrico. Podemos utilizar el **método de momentos**, el cual es otra aplicación directa de la ley de los grandes números.

**1.4.1. Definición [método de momentos]:** Supongamos que queremos estimar  $k$  parámetros de un modelo paramétrico  $X \sim \mathbb{P}(\cdot; \theta)$ . Es decir, queremos realizar inferencia sobre  $\theta \in \Theta \subseteq \mathbb{R}^k$ . Supongamos que podemos escribir el siguiente sistema de ecuaciones

$$\begin{aligned}
 \mu_1 &= \mathbb{E}[X] = g_1(\theta_1, \dots, \theta_k), \\
 \mu_2 &= \mathbb{E}[X^2] = g_2(\theta_1, \dots, \theta_k), \\
 &\vdots \\
 \mu_k &= \mathbb{E}[X^k] = g_k(\theta_1, \dots, \theta_k).
 \end{aligned}$$

Sea  $X_1, \dots, X_N \stackrel{\text{iid}}{\sim} \mathbb{P}(\cdot; \theta)$  una muestra del modelo probabilístico y denotemos por

$$\hat{\mu}_k = \frac{1}{N} \sum_{n=1}^N x_n^k, \quad (4)$$

los promedios basados en la muestra. Entonces, el **estimador de momentos** del vector  $\theta \in$

$\Theta \subseteq \mathbb{R}^k$  está dado por la solución del sistema de ecuaciones

$$\begin{aligned}\hat{\mu}_1 &= g_1(\hat{\theta}_1, \dots, \hat{\theta}_k), \\ \hat{\mu}_2 &= g_2(\hat{\theta}_1, \dots, \hat{\theta}_k), \\ &\vdots \\ \hat{\mu}_k &= g_k(\hat{\theta}_1, \dots, \hat{\theta}_k).\end{aligned}$$

1.4.2. *Ejemplo:* Consideremos los datos  $X_1, \dots, X_N \stackrel{\text{iid}}{\sim} \text{Gamma}(\alpha, \beta)$  donde tenemos el siguiente sistema de ecuaciones

$$\alpha = \frac{\mathbb{E}(X)^2}{\mathbb{V}(X)}, \quad \beta = \frac{\mathbb{V}(X)}{\mathbb{E}(X)}.$$

Los cuales podemos estimar utilizando las aproximaciones

$$\mathbb{E}(X^k) \approx \frac{1}{N} \sum_{n=1}^N x_n^k. \quad (5)$$

### 1.5. 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.

### REFERENCIAS

- [1] L. M. Chihara and T. C. Hesterberg. *Mathematical Statistics with Resampling and R*. John Wiley & Sons, Inc., Hoboken, NJ, USA, aug 2018. ISBN 978-1-119-50596-9 978-1-119-41654-8. . 1
- [2] B. Efron and R. J. Tibshirani. *An Introduction to the Bootstrap*. Springer US, Boston, MA, 1993. ISBN 978-0-412-04231-7 978-1-4899-4541-9. . 1