

2.2. Bootstrap

Este método es un poco más sencillo de implementar que Jackknife y es igualmente de eficaz propuesto por [4].

Primero recordemos que estamos estimando un estadístico a partir de una muestra de modo que $T_n = g(X_1, \dots, X_n)$ donde g es cualquier función (media, varianza, quantiles, etc).

Supongamos que conocemos la distribución real de los X 's, llamada $F(x)$. Si uno quisiera estimar la varianza de X basta con hacer

2 en

$$\text{Var}_F(T_n) = \frac{\sigma^2}{n} = \frac{\int x^2 dF(x) - \left(\int x dF(x)\right)^2}{n}$$

donde $\sigma^2 = \text{Var}(X)$ y el subíndice F es solo para indicar la dependencia con la distribución real.

Ahora dado que no tenemos la distribución real $F(x)$, una opción es encontrar un estimador de esta llamado \hat{F}_n .

La técnica de bootstrap se basa en extraer muchas muestras iid de la distribución \hat{F}_n de modo que se pueda conocer su varianza.

En simple pasos la técnica es

1. Seleccione $X_1^*, \dots, X_n^* \sim \hat{F}_n$

2. Estime $T_n^* = g(X_1^*, \dots, X_n^*)$

3. Repita los Pasos 1 y 2 B veces para obtener $T_{n,1}^*, \dots, T_{n,B}^*$

$X_1, \dots, X_n \sim F(x)$

\hat{F}_n

Muestras

4. Estime

$$v_{\text{boot}} = \frac{1}{B} \sum_{b=1}^B \left(T_{n,b}^* - \frac{1}{B} \sum_{r=1}^B T_{n,r}^* \right)^2$$

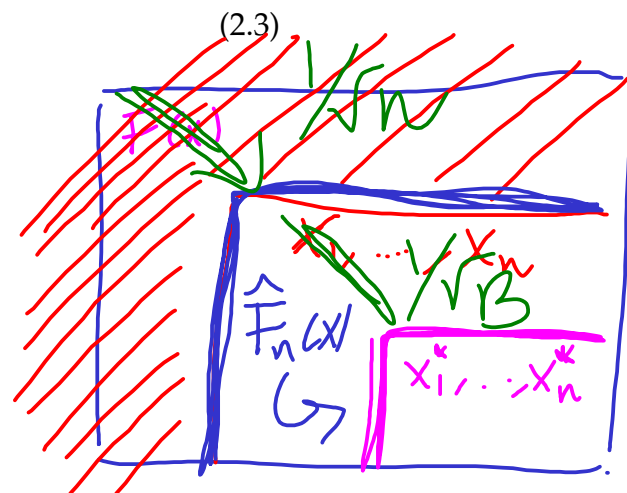
Por la ley de los grandes números tenemos que

$$v_{\text{boot}} \xrightarrow{\text{a.s.}} \mathbb{V}_{\hat{F}_n}(T_n), \quad \text{si } B \rightarrow \infty.$$

además llamaremos,

$$\hat{\text{se}}_{\text{boot}} = \sqrt{v_{\text{boot}}}$$

En pocas palabras lo que tenemos es que



$$\text{Mundo Real: } F \implies X_1, \dots, X_n \implies T_n = g(X_1, \dots, X_n)$$

$$\text{Mundo Bootstrap: } \hat{F}_n \implies X_1^*, \dots, X_n^* \implies T_n^* = g(X_1^*, \dots, X_n^*)$$

En términos de convergencia lo que se tiene es que

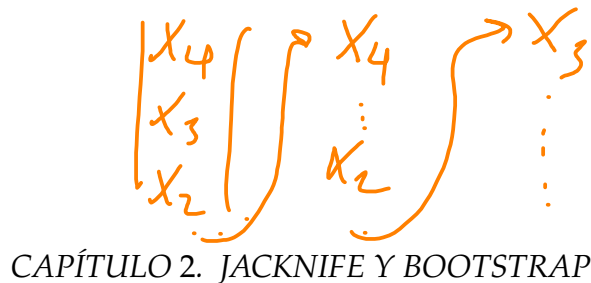
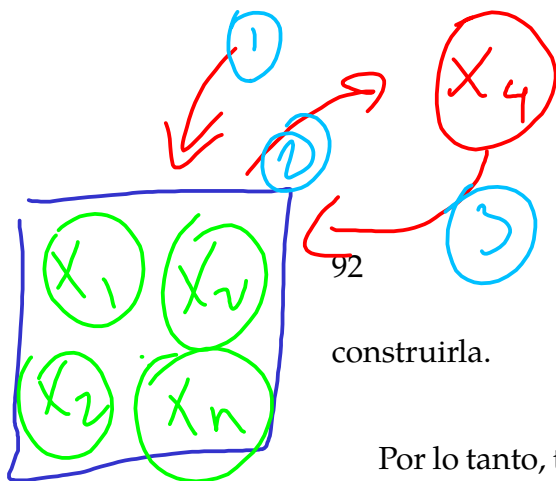
$$\text{Var}_F(T_n) \approx \underbrace{O(1/\sqrt{n})} \approx \text{Var}_{\hat{F}_n}(T_n) \approx \underbrace{O(1/\sqrt{B})} \approx v_{\text{boot}}$$

Pregunta 2.2.1

¿Cómo extraemos una muestra de \hat{F}_n ?

Recuerden que \hat{F}_n asigna la probabilidad de $\frac{1}{n}$ a cada valor usado para

$$\hat{F}_n(x) = \frac{1}{n} \sum_{i=1}^n \mathbb{I}(X_i \leq x)$$



CAPÍTULO 2. JACKKNIFE Y BOOTSTRAP

construirla.

Por lo tanto, todos los puntos originales X_1, \dots, X_n tienen probabilidad $\frac{1}{n}$ de ser escogidos, que resulta ser equivalente a un muestreo con remplazo n -veces.

Así que basta cambiar el punto 1. del algoritmo mencionando anteriormente con

1. Seleccione una muestra con remplazo X_1^*, \dots, X_n^* de X_1, \dots, X_n .

Laboratorio 2.2.2

En este ejemplo podemos tomar $B = 1000$ y construir esa cantidad de veces nuestro estimador.

```
B <- 1000  
Tboot_b <- NULL
```

```
for(b in 1:B) {
```

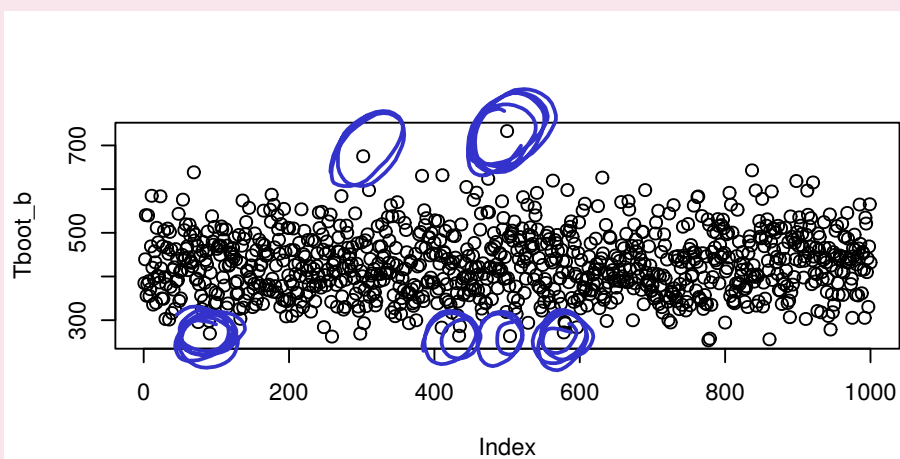
```
  xb <- sample(x, size = n, replace = TRUE)
```

```
  Tboot_b[b] <- var(xb)  
}
```

```
Tboot_b[1:10]
```

```
## [1] 384.8959 439.3992 540.4429 356.0368 382.5001 539.7398 389.2075 355.4033  
## [9] 371.0053 469.0120
```

```
plot(Tboot_b)
```



Por supuesto podemos encontrar los estadísticos usuales para esta nueva muestra

```
(Tboot <- mean(Tboot_b))
```

```
## [1] 428.3197
```

```
(Vboot <- var(Tboot_b))
```

```
## [1] 5345.401
```

```
(sdboot <- sqrt(Vboot))
```

```
## [1] 73.11225
```

2.2.1. Intervalos de confianza

Intervalo Normal

Este es el más sencillo y se escribe como

$$T_n \pm z_{\alpha/2} \hat{S}e_{\text{boot}} \quad (2.4)$$

Cuidado 2.2.3

Este intervalo solo funciona si la distribución de T_n es normal.

Laboratorio 2.2.4

El cálculo de este intervalo es

$c(Tn - z * sdboot,$

$Tn + z * sdboot)$

[1] 285.9510 572.5458

Intervalo pivotal

Sea $\theta = T(F)$ y $\hat{\theta}_n = T(\hat{F}_n)$ y defina la cantidad pivotal $R_n = \hat{\theta}_n - \theta$.

Sea $H(r)$ la función de distribución del pivote:

$$H(r) = \mathbb{P}_F(R_n \leq r).$$

Además considere $C_n^* = (a, b)$ donde

$$a = \hat{\theta}_n - H^{-1}\left(1 - \frac{\alpha}{2}\right) \quad \text{y} \quad b = \hat{\theta}_n - H^{-1}\left(\frac{\alpha}{2}\right).$$

$$\hat{\theta}_n - a = H^{-1}\left(1 - \frac{\alpha}{2}\right)$$

Se sigue que

$$\begin{aligned}
 \mathbb{P}(a \leq \theta \leq b) &= \mathbb{P}(\hat{\theta}_n - b \leq R_n \leq \hat{\theta}_n - a) \\
 &= H(\hat{\theta}_n - a) - H(\hat{\theta}_n - b) \\
 &= H\left(H^{-1}\left(1 - \frac{\alpha}{2}\right)\right) - H\left(H^{-1}\left(\frac{\alpha}{2}\right)\right) \\
 &= 1 - \frac{\alpha}{2} - \frac{\alpha}{2} = 1 - \alpha
 \end{aligned}$$

$$\begin{aligned}
 \theta_n - b &= H^{-1}\left(\frac{\alpha}{2}\right)
 \end{aligned}$$

Nota 2.2.5

$C_n^* = (a, b)$ es un intervalo de confianza al $1 - \alpha$ de confianza.

El problema es que este intervalo depende de H desconocido.

Para resolver este problema, se puede construir una versión bootstrap de H usando lo que sabemos hasta ahora.

Bootstrap

$$\hat{H}(r) = \frac{1}{B} \sum_{b=1}^B I(R_{n,b}^* \leq r)$$

donde $R_{n,b}^* = \hat{\theta}_{n,b}^* - \hat{\theta}_n$

→ DATOS ORIGINALS

Sea r_β^* el cuantil muestral de tamaño β de $(R_{n,1}^*, \dots, R_{n,B}^*)$ y sea θ_β^* el cuantil muestral de tamaño β de $(\theta_{n,1}^*, \dots, \theta_{n,B}^*)$.

Nota 2.2.6

Según la notación anterior note que

$$r_\beta^* = \theta_\beta^* - \hat{\theta}_n$$

Con estas observaciones It follows that an approximate $1 - \alpha$ confidence interval is $C_n = (\hat{a}, \hat{b})$ where

$$\begin{aligned}\hat{a} &= \hat{\theta}_n - \hat{H}^{-1}\left(1 - \frac{\alpha}{2}\right) = \hat{\theta}_n - r_{1-\alpha/2}^* = \hat{\theta}_n - (\theta_{1-\alpha/2}^* - \hat{\theta}_n) = 2\hat{\theta}_n - \theta_{1-\alpha/2}^* \\ \hat{b} &= \hat{\theta}_n - \hat{H}^{-1}\left(\frac{\alpha}{2}\right) = \hat{\theta}_n - r_{\alpha/2}^* = \hat{\theta}_n - (\theta_{\alpha/2}^* - \hat{\theta}_n) = 2\hat{\theta}_n - \theta_{\alpha/2}^*\end{aligned}$$

Nota 2.2.7

El intervalo de confianza pivotal de tamaño $1 - \alpha$ es

$$C_n = \left(2\hat{\theta}_n - \hat{\theta}_{((1-\alpha/2)B)}^*, 2\hat{\theta}_n - \hat{\theta}_{((\alpha/2)B)}^*\right)$$

Laboratorio 2.2.8

El intervalo anterior para un nivel de 95 % se estima de la siguiente forma

```
c(2 * Tn - quantile(Tboot_b, 1 - 0.05 / 2) ,
  2 * Tn - quantile(Tboot_b, 0.05 / 2))

##      97.5%      2.5%
## 275.3768 558.6741
```

Handwritten notes:

- 200, 250, 300 } $\alpha/2$
- 350, 360 } $1 - \alpha/2$

Intervalo pivotal studentizado

Una mejora del intervalo anterior sería normalizar los estimadores previamente

$$Z_n = \frac{T_n - \theta}{\widehat{se}_{boot}}$$

Como θ es desconocido, entonces la versión a estimar es

$$Z_{n,b}^* = \frac{T_{n,b}^* - T_n}{\widehat{se}_b^*}$$

donde \widehat{se}_b^* es un estimador del error estándar de $T_{n,b}^*$ no de T_n .

Cuidado 2.2.9

Esto requerirá estimar la varianza de $T_{n,b}^*$ para cada b .

Δ 2 veces Bootstrap.

Con esto se puede obtener cantidades $Z_{n,1}^*, \dots, Z_{n,B}^*$ que debería ser próximos a Z_n .

Sea z_α^* del α cuantil de $Z_{n,1}^*, \dots, Z_{n,B}^*$, entonces $\mathbb{P}(Z_n \leq z_\alpha^*) \approx \alpha$.

Define el intervalo

$$C_n = (T_n - z_{1-\alpha/2}^* \widehat{se}_{boot}, T_n - z_{\alpha/2}^* \widehat{se}_{boot})$$

Justificado por el siguiente cálculo:

$$\begin{aligned}\mathbb{P}(\theta \in C_n) &= \mathbb{P}\left(T_n - z_{1-\alpha/2}^* \widehat{\text{se}}_{\text{boot}} \leq \theta \leq T_n - z_{\alpha/2}^* \widehat{\text{se}}_{\text{boot}}\right) \\ &= \mathbb{P}\left(z_{\alpha/2}^* \leq \frac{T_n - \theta}{\widehat{\text{se}}_{\text{boot}}} \leq z_{1-\alpha/2}^*\right) \\ &= \mathbb{P}\left(z_{\alpha/2}^* \leq Z_n \leq z_{1-\alpha/2}^*\right) \\ &\approx 1 - \alpha\end{aligned}$$

Laboratorio 2.2.10

Note que para este caso tenemos que hacer bootstrap para cada estimador bootstrap calculado.