

Joseph Guy Evans Hilaire Jr

**Bootstrap
&
Estimacion intervalo de confianza para la media**

El metodo **Bootstrap** es un metodo que se apoya en la tecnica del muestreo aleatorio sin reemplazo. Con dicho metodo, se pueden estimar las propiedades de un estimador(media, varianza) mediante el muestreo sobre una distribucion aproximada.

En este trabajo, intentaremos estimar intervalos de confianza para la media de una muestra *Normal* usando el metodo Bootstrap.

Algoritmo I

Un algoritmo para estimar la media θ mediante el metodo Bootstrap es el siguiente:

1. Empezamos con una muestra observada x_1, \dots, x_m y $\hat{\theta}$ un estimador *plug-in* de θ .
2. Para $j=1, \dots, N$ generamos:

$$x^{(j)} = \{x_{j1}, \dots, x_{jm}\}$$

donde $j_1, \dots, j_m \sim U\{1, \dots, m\}$

3. Luego para $j=1, \dots, N$ hacemos:

$$\hat{\theta}^{(j)} = \hat{\theta}(x^{(j)})$$

4. Ordenamos la muestra $\{\hat{\theta}^{(1)}, \dots, \hat{\theta}^{(N)}\}$ de menor a mayor
5. Hacemos $l = \lceil \frac{\alpha}{2}N \rceil$ y $u = \lceil (1 - \frac{\alpha}{2})N \rceil$
6. Luego hacemos $t = \hat{\theta}(x_1, \dots, x_m)$
7. Entonces, el intervalo de confianza que buscabamos seria dado por:

$$(2t - \hat{\theta}_{(u)}, 2t - \hat{\theta}_{(l)})$$

A modo de **implementacion** del algoritmo anterior, supongamos que tenemos los datos $X_1, \dots, X_m \sim N(0, 1)$. El intervalo de confianza para la media de esta muestra quedara como sigue:

```

1 import pandas as pd
2 import matplotlib.pyplot as plt
3
4 import random
5 import numpy as np
6 import math
7 import scipy.stats

```

```

1 alpha = 0.05
2 N = 1000
3
4
5 def intervConf_Bootstrap(x, estimador, N):
6     n = int(x.size) ## El tamaño de las muestras con reemplazo
7     theta = np.zeros(N)
8
9     # Generamos las muestras Bootstrap
10    for k in range(0,N):
11        x_bos = np.random.choice(x,n,replace=True) # Muestreo con reemplazo
12        theta[k] = estimador(x_bos)
13
14    # ordenamos los valores de theta
15    theta.sort()
16    thetaOrd = theta
17
18    # Construimos el intervalo de confianza
19    t = estimador(x)
20    l = math.ceil(0.5*alpha*N )
21    u = math.ceil((1-0.5*alpha)*N )
22    return([2*t-thetaOrd[u] ,
23            2*t-thetaOrd[l]])
24
25 # Intervalo de confianza segun su formula exacta:
26 def intervConf(x):
27     n = x.size
28     m = np.mean(x)
29     s = np.std(x)
30
31     # cuantil (1-0.5*alpha) de una T-student
32     cuan = scipy.stats.t.ppf(1-0.5*alpha, n-1)
33
34     return([m-cuan*s/math.sqrt(n) , m+cuan*s/math.sqrt(n)])

```

```

1 n0 = 100
2 IC1 = np.zeros(n0)
3 IC2 = np.zeros(n0)
4
5 II1 = np.zeros(n0)
6 II2 = np.zeros(n0)
7
8 # Generamos los intervalos de confianza
9 # que despues graficaremos
10 for k in range(0,n0):
11     x = np.random.normal(0,1,n0)
12
13     IC1[k] = intervConf_Bootstrap(x,np.mean,N)[0]
14     IC2[k] = intervConf_Bootstrap(x,np.mean,N)[1]
15
16     II1[k] = intervConf(x)[0]
17     II2[k] = intervConf(x)[1]

```

```

1 # Graficamos los intervalos de confianza
2 # bajo el metodo Bootstrap
3 M = np.zeros(n0)
4
5 plt.figure(figsize=(10,5)) ## Determinamos las dimensiones de la grafica
6 plt.plot(M, color='red', label='Media')
7
8 # Graficamos los intervalos obtenidos mediante Bootstrap
9 plt.plot(IC1, 'o', color='royalblue', linestyle='dashed', label='Bootstrap')
10 plt.plot(IC2, 'o', color='royalblue', linestyle='dashed')
11
12 # Graficamos los intervalos obtenidos a partir de su formula exacta
13 plt.plot(II1, 'o', color='grey', linestyle='dashed', label='Metodo Exacto')
14 plt.plot(II2, 'o', color='grey', linestyle='dashed')
15
16
17 plt.legend(bbox_to_anchor=(1.05, 1), loc='upper left')
18 plt.show()

```

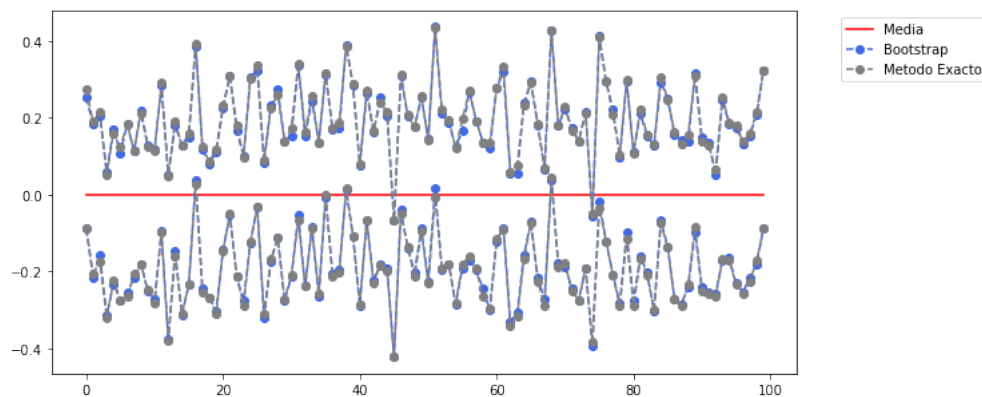


Figure 1: Comparacion Intervalos de Confianza

En la grafica anterior, notamos que los intervalos obtenidos bajo el metodo Bootstrap y el metodo exacto son casi parecidos. Por lo cual, podemos afirmar que el metodo Bootstrap es una buena metodologia para calcular los intervalos de confianza para los parametros de una muestra dada.

Algoritmo II

Con el algoritmo anterior, pudimos obtener una buena aproximacion para los intervalos de confianza de una muestra dada. Tambien, existe otro algoritmo de calculo de intervalos de confianza. Dicho algoritmo se traduce como sigue:

1. Empezamos con una muestra observada x_1, \dots, x_m y $\hat{\theta}$ un estimador *plug-in* de θ .
2. Para $j=1, \dots, N$ generamos:

$$x^{(j)} = \{x_{j1}, \dots, x_{jm}\}$$

donde $j_1, \dots, j_m \sim U\{1, \dots, m\}$

3. Luego para $j=1, \dots, N$ hacemos:

$$\hat{\theta}^{(j)} = \hat{\theta}(x^{(j)})$$

4. Calculamos el siguiente cuantil de una distribucion Normal:

$$\hat{z} = \Phi^{-1}\left(\frac{\#\{j|\hat{\theta}^{(j)} \leq \hat{\theta}\}}{N}\right)$$

5. Despues, para $i=1, \dots, m$, hacemos lo siguiente:

$$\theta_{(i)} = \hat{\theta}_{n-1}(x_1, \dots, x_{i-1}, x_{i+1}, \dots, x_n)$$

en donde $\hat{\theta}_{n-1}(x_1, \dots, x_{i-1}, x_{i+1}, \dots, x_n)$ es la evaluacion de $\hat{\theta}$ en la muestra pero sin el valor x_i , de ahi el subindice $n-1$.

6. Con eso calculamos:

$$\bar{\theta} = \frac{1}{m} \sum_{i=1}^m \theta_{(i)}$$

7. Eso nos sirve para calcular los siguientes valores:

$$\begin{aligned} \hat{a} &= \frac{1}{6} \frac{\sum_{i=1}^m (\bar{\theta} - \theta_{(i)})^3}{\left(\sum_{i=1}^m (\bar{\theta} - \theta_{(i)})^2\right)^{3/2}} \\ q_l &= \Phi\left(\hat{z} + \frac{\hat{z} + \Phi^{-1}(\alpha/2)}{1 - \hat{a}(\hat{z} + \Phi^{-1}(\alpha/2))}\right) \\ q_u &= \Phi\left(\hat{z} + \frac{\hat{z} + \Phi^{-1}(1-\alpha/2)}{1 - \hat{a}(\hat{z} + \Phi^{-1}(1-\alpha/2))}\right) \end{aligned}$$

8. Ordenamos la muestra $\{\hat{\theta}^{(1)}, \dots, \hat{\theta}^{(N)}\}$ de menor a mayor
9. Obtenemos los siguientes valores: $l = \lceil q_l \cdot N \rceil$ y $u = \lceil q_u \cdot N \rceil$
10. Entonces, el intervalo de confianza bajo el metodo **Bias Corrected and Accelerated** seria dado por:

$$(\hat{\theta}_{(l)}, \hat{\theta}_{(u)})$$

Ahora podemos implementar el algoritmo anterior de la siguiente manera:

```

1 def interv.BCa(x, estimador, N):
2     n = x.size
3     theta = np.zeros(N)
4
5     for k in range(0,N):
6         Xb = np.random.choice(x,n,replace=True)
7         theta[k] = estimador(Xb)
8
9     # numero de veces que es menor que el estimador(media)
10    media = np.mean(x)
11    c = np.where(theta <= media)
12    q = len(c[0])/N
13
14    # calculamos el cuantil correspondiente
15    # a partir de una Normal
16    zGorro = scipy.stats.norm.ppf(q)
17
18    # Vector para guardar las medias a continuacion
19    Ti = np.zeros(n)
20
21    for k in range(0,n):
22        # creamos una lista de los indices sin el indice k
23        ind = np.delete(np.arange(n),k)
24        # obtenemos la muestra sin el k-esimo dato
25        xReducido = x[ind]
26        # guardamos las medias
27        Ti[k] = np.mean(xReducido)
28
29    # Sacamos el promedio de las Ti
30    Tmean = np.mean(Ti)
31
32    # Calculamos el valor de (a_gorro)
33    a0 = sum( (Tmean-Ti)**3 )
34    a1 = sum( (Tmean-Ti)**2 )
35    aGorro = a0 / ( 6*a1**(3/2) )
36
37    # Calculamos el valor de (qL)
38    q0 = zGorro + scipy.stats.norm.ppf(alpha/2)
39    q1 = 1- aGorro*q0
40    qL = scipy.stats.norm.cdf(zGorro + q0/q1)
41
42    # Calculamos el valor de (qU)
43    q00 = zGorro + scipy.stats.norm.ppf(1-alpha/2)
44    q11 = 1- aGorro*q00
45    qU = scipy.stats.norm.cdf(zGorro + q00/q11)
46
47    # Ordenamos la muestra
48    theta.sort()
49
50    # Calculamos (l) y (u):
51    l = math.ceil(qL*N)
52    u = math.ceil(qU*N)
53
54    return ([ theta[l], theta[u] ])

```

```

1 BCa1 = np.zeros(n0)
2 BCa2 = np.zeros(n0)
3
4 # Generamos los intervalos de confianza
5 # con el metodo BCa :
6 for k in range(0,n0):
7     x = np.random.normal(0,1,n0)
8
9     BCa1[k] = interv_BCa(x,np.mean,N)[0]
10    BCa2[k] = interv_BCa(x,np.mean,N)[1]
11
12    II1[k] = intervConf(x)[0]
13    II2[k] = intervConf(x)[1]

1 # Graficamos los intervalos de confianza
2 # bajo el metodo Bias Corrected and Accelerated
3 M = np.zeros(n0)
4
5 plt.figure(figsize=(10,5)) ## Determinamos las dimensiones de la grafica
6 plt.plot(M, color='red', label='Media')
7
8 # Graficamos los intervalos obtenidos mediante el metodo BCa
9 plt.plot(BCa1, 'o', color='royalblue', linestyle='dashed', label='metodo BCa')
10 plt.plot(BCa2, 'o', color='royalblue', linestyle='dashed')
11
12 # Graficamos los intervalos obtenidos a partir de su formula exacta
13 plt.plot(II1, 'o', color='grey', linestyle='dashed', label='Metodo Exacto')
14 plt.plot(II2, 'o', color='grey', linestyle='dashed')
15
16
17 plt.legend(bbox_to_anchor=(1.05, 1), loc='upper left')
18 plt.show()

```

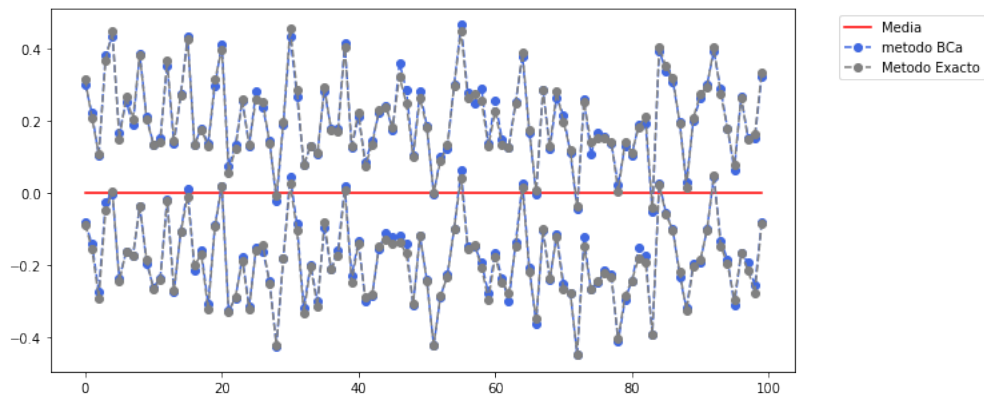


Figure 2: Comparacion Intervalos de Confianza

Con los intervalos de confianza obtenidos bajo las 2 metodologias anteriores, podemos decir que la que estimo mejores intervalos es la que logra capturar en mayor medida el valor de la media de los datos.

Para ello, procedemos a lo siguiente:

```

1 media = 0
2
3 num_01 = [IC1[k]<= media <= IC2[k] for k in range(len(IC1))]
4 num_02 = [BCa1[k]<= media <= BCa2[k] for k in range(len(BCa1))]
5
6 num_01 = np.sum(num_01)
7 num_02 = np.sum(num_02)
8
9 num_01, num_02

```

Obtenemos el siguiente resultado:

Metodologia	Bootstrap	BCa
Cantidad buena	94	87

Entonces podemos ver que hay 94 intervalos *Bootstrap* que logran capturar la media igual a 0 mientras que hay 87 intervalos *BCa* que capturan al valor de 0.

Eso nos dice que en esta iteracion, el intervalo *Bootstrap* fue mejor que la otra metodologia para estimar intervalos de confianza adecuados para la media del conjunto de datos *Normales*.