

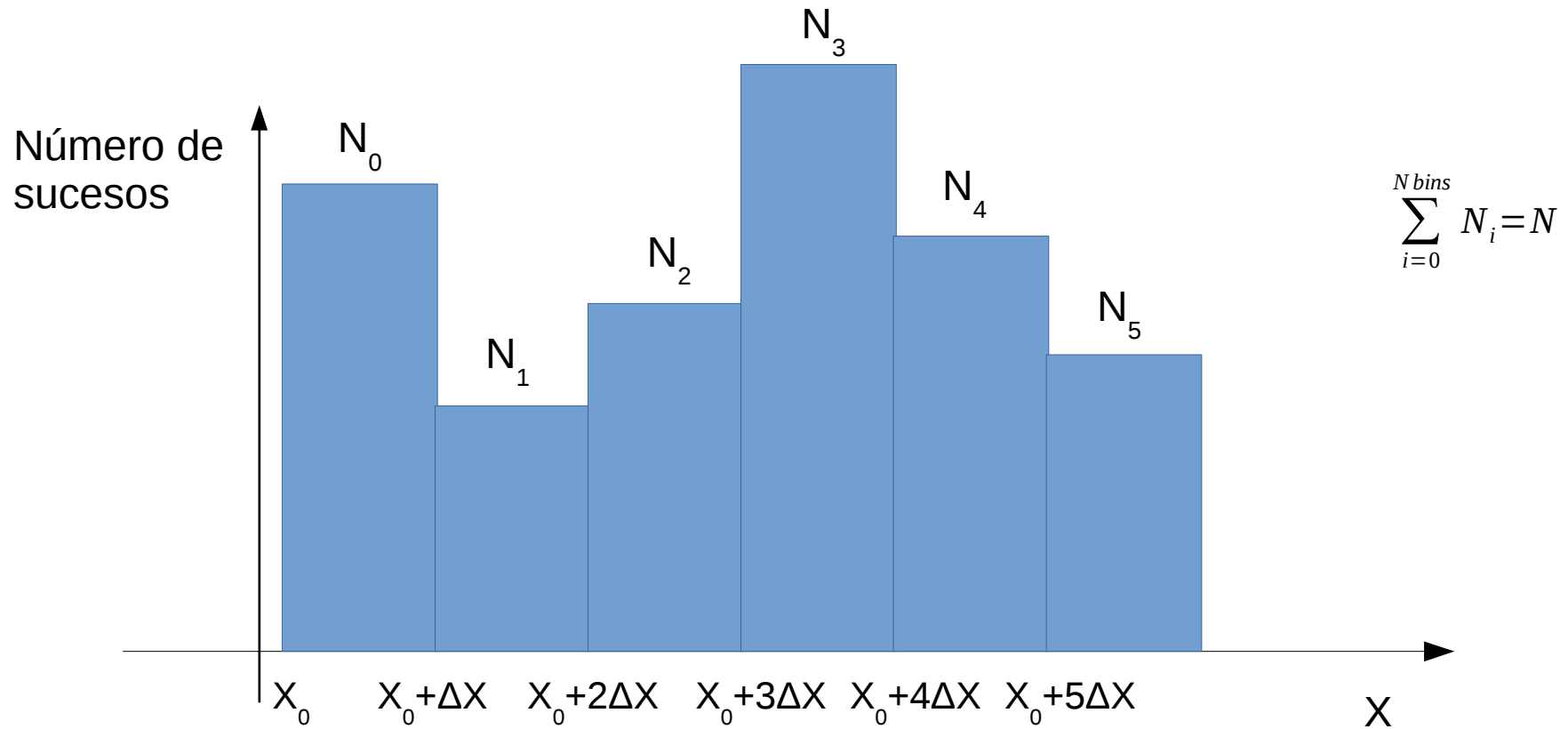
Estadística [continuación]

Santander, 2017-2018

Técnicas de remuestreo

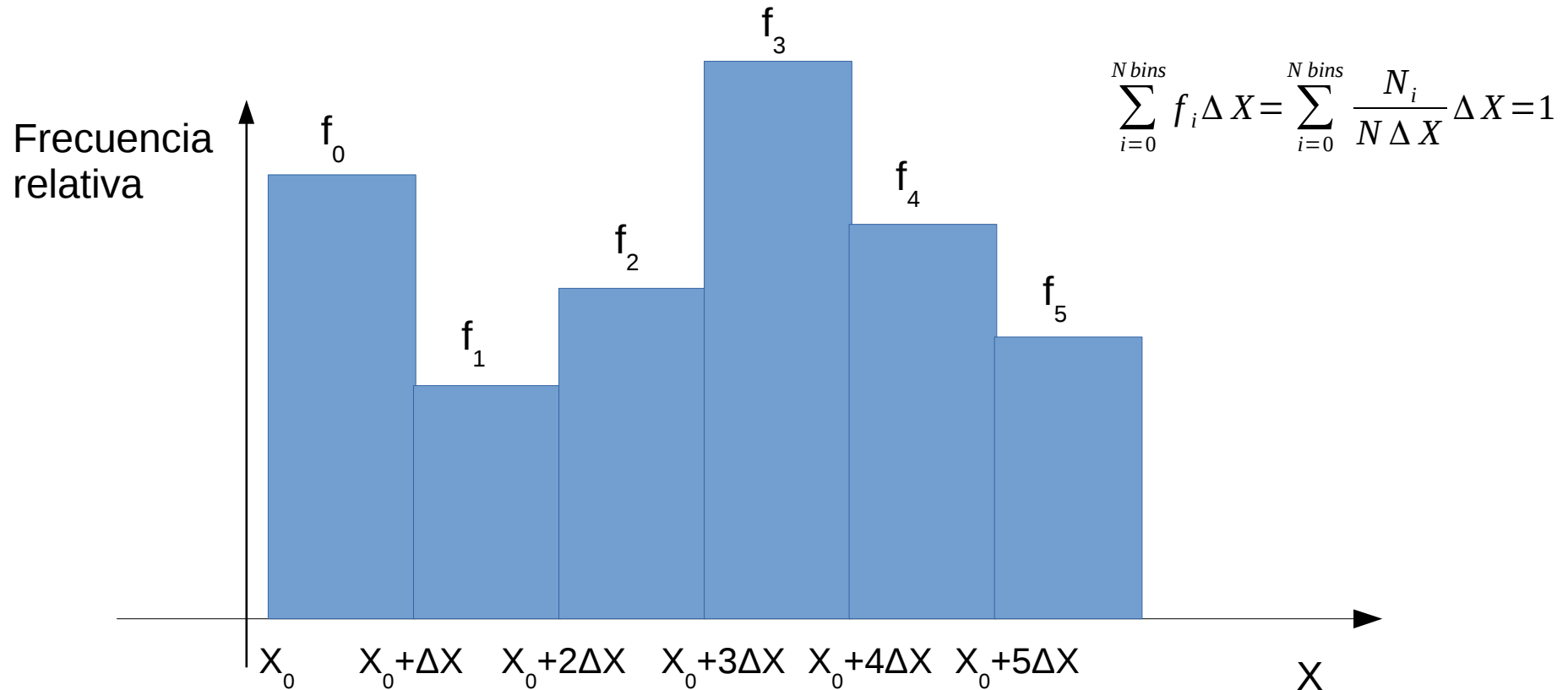
- En muchas ocasiones es posible definir estrategias **sobre la muestra obtenida**, que permiten estimar o mejorar los sesgos de nuestros estimadores (o procedimientos).
- **No conocemos de forma exacta la pdf del proceso** → la muestra **“aproxima” esta pdf.**
- Consideremos una variable aleatoria X con una densidad de probabilidad $F(X)$.
- Consideremos también una muestra de “ N ” medidas $\{X_1, X_2, \dots, X_N\}$.
- Podemos colocar estas N medidas en un histograma de frecuencias.

Relación entre pdf e histograma de frecuencias (I)



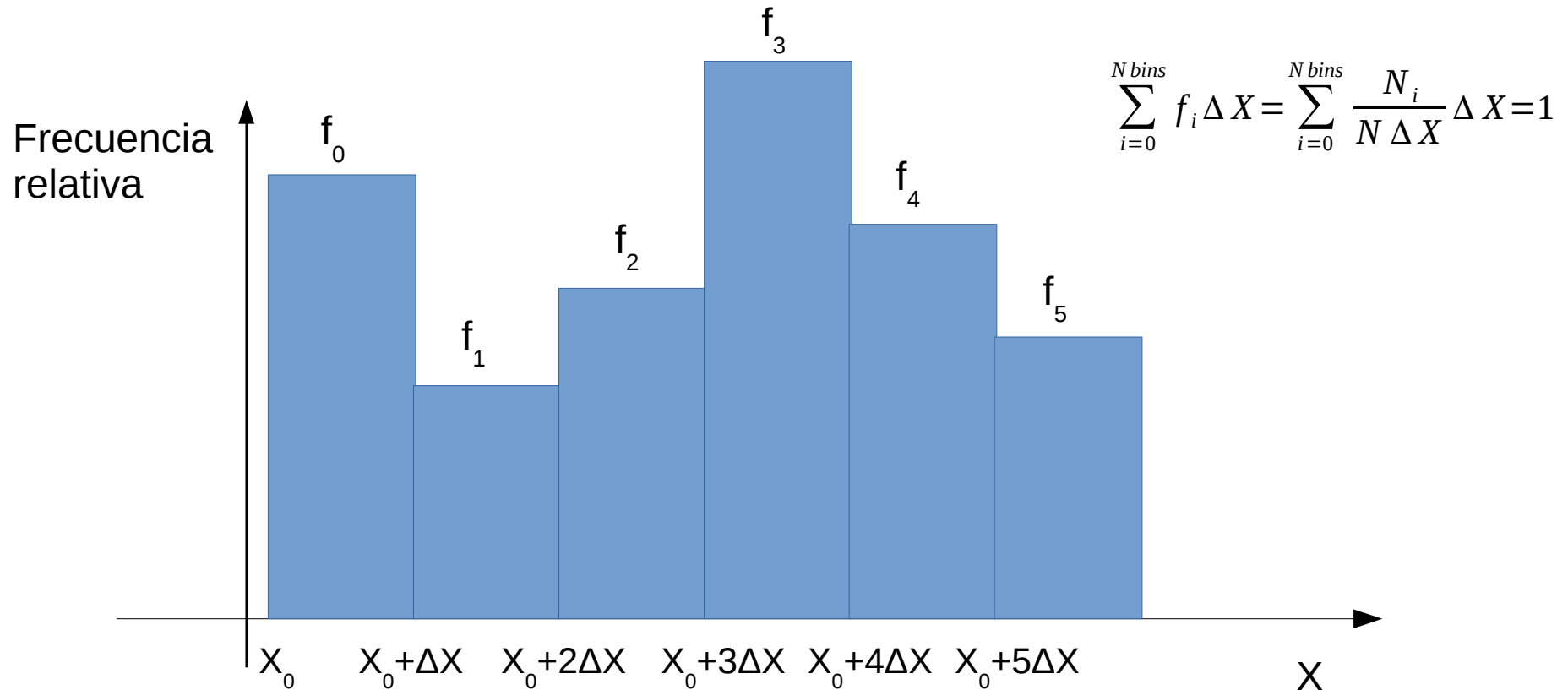
- En este histograma ΔX es la anchura del bin y N_i es el número de medidas en cada bin.
- Cada bin nos dice cuántas medidas están dentro del intervalo $X_i+\Delta X$.

Relación entre pdf e histograma de frecuencias (II)



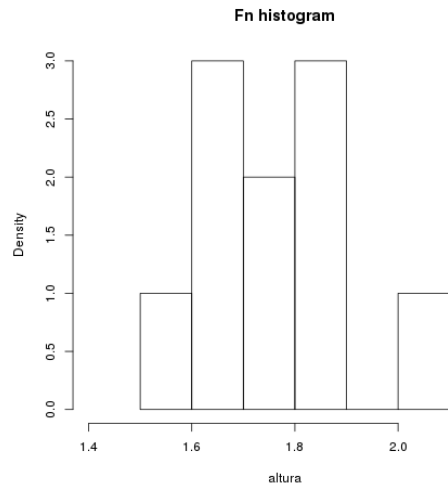
- Podemos obtener la densidad de frecuencia relativa en cada bin dividiendo por el número total y ΔX .
- Esta distribución comienza a guardar similitud con una pdf: en particular su integral es 1.

Relación entre pdf e histograma de frecuencias (III)

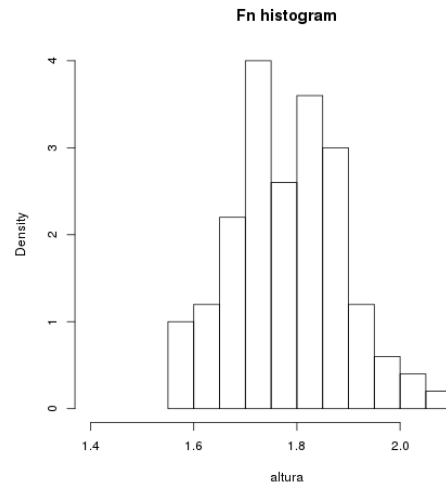


- Este histograma nos dice la densidad de sucesos que cayeron en un bin determinado.
- Si el tamaño del bin fuese cero y hubiese infinitas medidas este histograma sería la pdf.

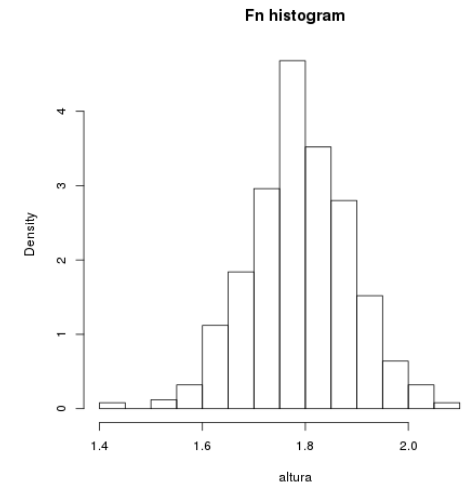
Relación entre pdf e histograma de frecuencias (IV)



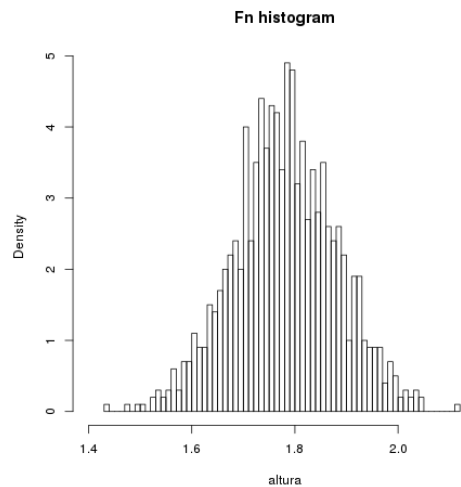
N=10, Bin=0.15m



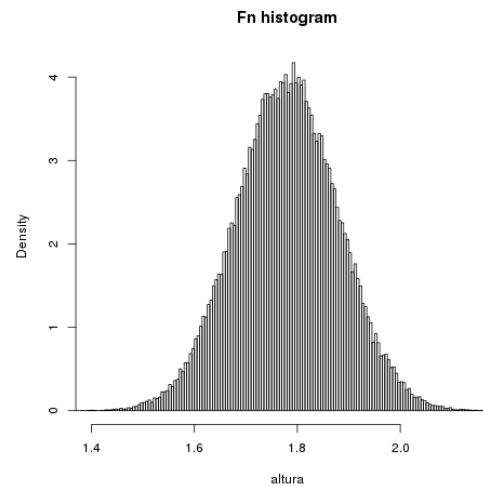
N=100, Bin=0.08m



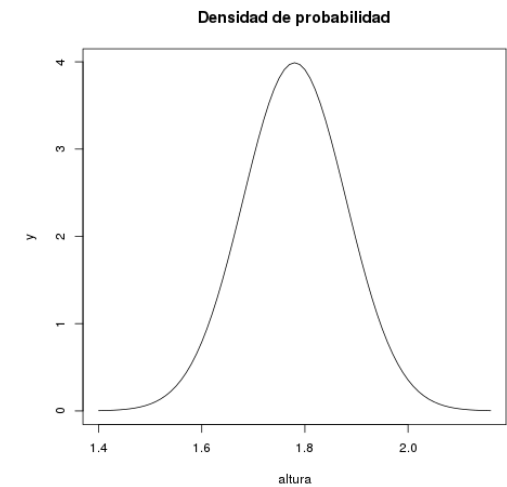
N=500, Bin=0.08m



N=1000, Bin=0.02m



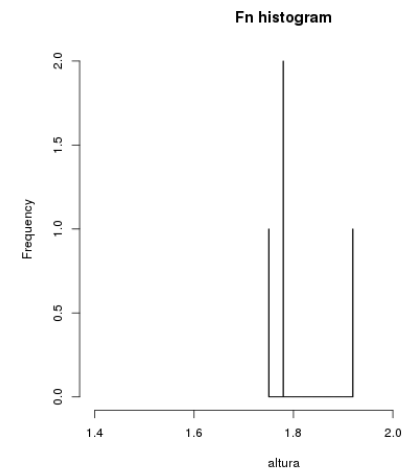
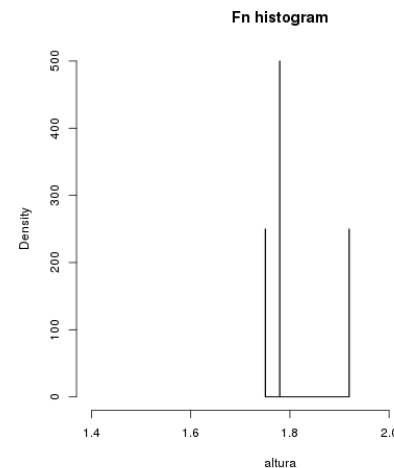
N=100000, Bin=0.004m



Relación entre pdf e histograma de frecuencias (IV)

- Es posible aproximar la pdf de un proceso por el histograma de densidad de frecuencias.
- Para ello es necesario construir F_n con un tamaño de bin infinitesimalmente pequeño.
- Cuando hacemos esto básicamente cada valor medido diferente $\{X_1, X_2, \dots, X_N\}$ ocupa un bin.
- Ejemplo: volviendo al ejemplo de las alturas en donde teníamos $\{1.78, 1.78, 1.92, 1.75\}$.
- F_n no es una aproximación perfecta de la pdf → pero permite simular pseudo-muestras (toys)

Es posible aplicar técnicas de MC sobre el histograma de frecuencias/fracción de frecuencias para generar pseudo-muestras.



Técnica bootstrap

- Supongamos un parámetro Θ medido sobre una muestra obteniéndose un estadístico $\bar{\Theta}$
- Si fuésemos capaces de conseguir M nuevas muestras y computar en cada una $\bar{\Theta}_i$, entonces

$$\bar{\bar{\Theta}} = \frac{1}{M} \sum_{i=1}^M \bar{\Theta}_i$$

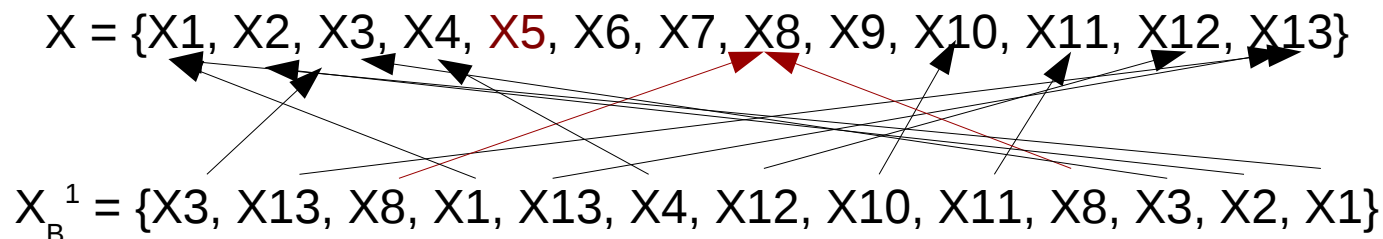
- El sesgo de esta variable es obviamente inferior al sesgo de $\bar{\Theta}$ que sólo involucra una muestra
- Tomar M nuevas muestras resulta costoso, sin embargo obtenido Fn de la muestra tomada, es posible generar M nuevas pseudomuestras (**bootstrap**) sampleando dicha distribución.
- Siendo posible demostrar que la distribución de la variable de bootstrap generada así se distribuye de modo que

$$\bar{\bar{\Theta}} - \bar{\Theta} \approx \bar{\Theta} - \Theta$$

Efron, B. (1979). "Bootstrap methods: Another look at the jackknife". *The Annals of Statistics*. 7 (1): 1–26. doi:10.1214/aos/1176344552.

¿Cómo hacer el resampling en bootstrap?

- Puesto que nuestra distribución de frecuencias F_n tiene un bin para cada valor diferente...
- ...generar una nueva pseudomuestra consiste en sustituir cada uno de los elementos de la muestra por cualquier otro elegido aleatoriamente de la propia muestra hasta completar N



- Cada una de las nuevas pseudomuestras tiene el mismo número de elementos que la original
- No todos los elementos de la muestra original tienen por qué aparecer en una pseudomuestra
- Algunos elementos de la muestra original pueden aparecer repetidos en una pseudomuestra

Aplicación 1: Aproximación del error standard de un estadístico

- Supongamos un parámetro Θ del que hemos obtenido el estimador $\bar{\Theta}$ y queremos conocer su desviación standard

$$SD = \sqrt{\frac{1}{N} \sum_{i=1}^N (\bar{\Theta}_i - \Theta)^2}$$

- Cada uno de los $\bar{\Theta}_i$ se correspondería con el valor obtenido en un experimento diferente.
- Hacer diferentes experimentos resulta costoso, así que podemos estimarlo usando bootstrapping ya que:

$$\bar{\Theta}_B - \bar{\Theta} \approx \bar{\Theta} - \Theta$$

- Utilizando la Fn de la muestra que ha dado lugar a $\bar{\Theta}_i$ podemos generar M muestras de tipo bootstrap.
- Y podemos sustituir en la fórmula $\bar{\Theta}_i$ por cada una de las muestras bootstrap y Θ por $\bar{\Theta}$

$$SD = \sqrt{\frac{1}{N} \sum_{i=1}^N (\bar{\Theta}_i^B - \bar{\Theta})^2}$$

Aplicación 2: Aproximación del sesgo de un estadístico

→ Tal y como hemos definido previamente el sesgo de un estimador viene dada por la cantidad

$$\text{Sesgo}(\bar{\Theta}) = E[\bar{\Theta} - \Theta] = E[\bar{\Theta}] - \Theta$$

→ Es posible utilizar una vez más la igualdad vista anteriormente entre las distribuciones anteriores

$$\bar{\Theta}_B - \bar{\Theta} \approx \bar{\Theta} - \Theta$$

→ Para aproximar el sesgo sustituyendo por las variables de bootstrap y dando el siguiente resultado

$$\text{Sesgo}(\bar{\Theta}) = E[\bar{\Theta} - \Theta] \approx \frac{1}{N} \sum_{i=0}^M \Theta_i^B - \bar{\Theta}$$

El método jackknife

- La técnica de bootstrapping es la más extendida aunque existió otro algoritmo muy popular previo.
- Se trata del **jackknife**, muy utilizado cuando los recursos de cálculo eran limitados.
- Las ideas detrás de esta técnica son las mismas que las asociadas al bootstrap.
- La diferencia consiste en la manera de construir los nuevo datasets.
- En Jackknife cada uno se corresponde simplemente con la muestra formada tras quitar 1 elemento.
- Jackknife es una opción menos potente que Bootstrap y en general se aconseja el otro.

{X1, X2, X3, X4, X5, X6, X7, X8}

Las muestras tienen tamaño N-1

{X1, X2, X3, X4, X5, X6, X7}

{X1, X2, X3, X4, X5, X6, X8}

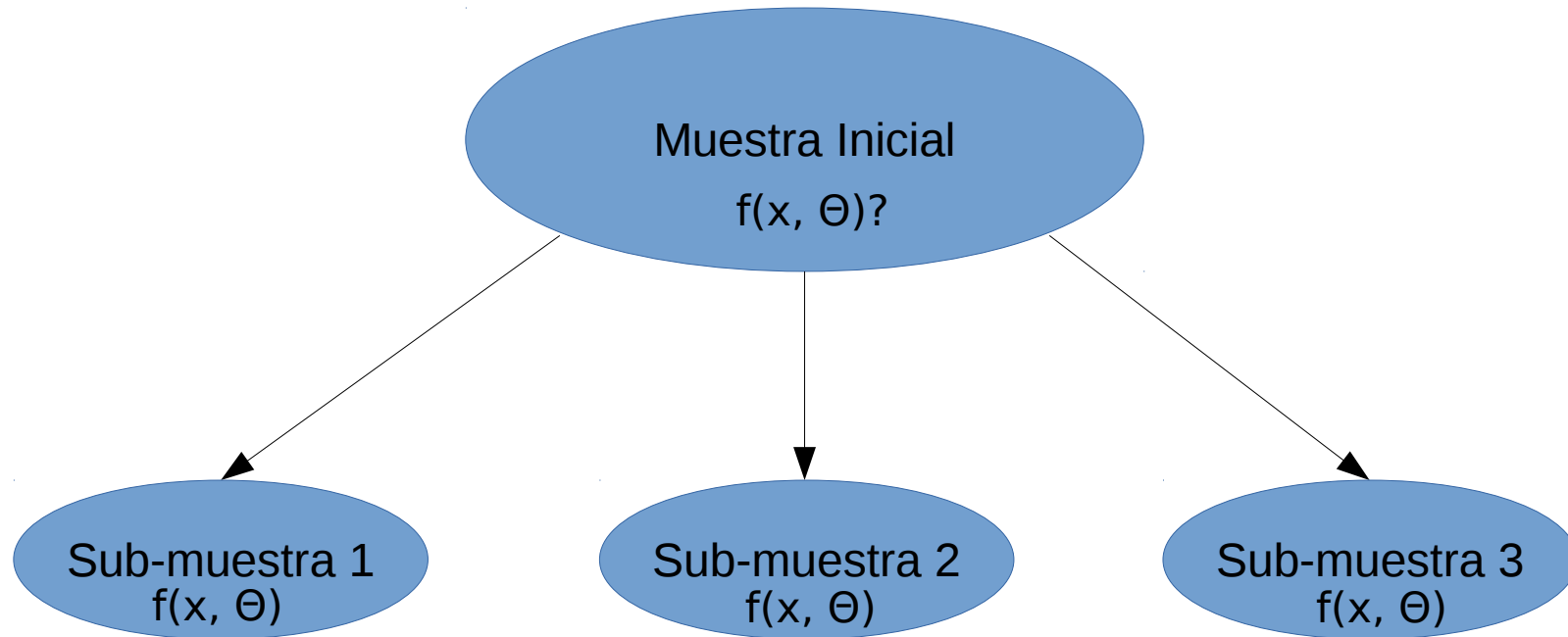
{X1, X2, X3, X4, X5, X7, X8}

{X1, X2, X3, X4, X6, X7, X8}

Cross validation

- **Cross validation** es una técnica que consiste en la partición de una muestra en **N bloques**...
- ...si la muestra se rige por una pdf determinada cada uno de los N bloques también lo hace.
- Y esto permite **comprobar los resultados obtenidos en un bloque con los de los otros bloques**.
- Supongamos que “creemos” que una muestra se distribuye según una pdf $f(x, \Theta)$.
- Supongamos también que contamos con algún tipo de estimador del parámetro Θ .
- **¿Cómo estar seguros de que el estimador tiene sentido si en realidad no conocemos f ?**
- Podemos estimarlo en el bloque 1, para luego compararlo con el resto de bloques.
- **Si las hipótesis son correctas, los resultados deberían ser consistentes entre bloques.**

Cross validation



- El concepto de cross-validation se utiliza en **problemas de modelado** (siguiente tema).
- Normalmente una de las muestras es llamada de “**training**” ya que sobre ella se estiman los parámetros, mientras que otra se suele llamar de “**test**” ya que sobre ella se comprueba la calidad del modelado. En ocasiones se habla también de muestra de “**validación**”.

Incidiremos sobre este concepto.

Ejercicio 2

- Escribe una función de R que reciba un vector de números x (la muestra), y genere **una** muestra bootstrap de ese vector.
- Utilizando la función anterior, escribe una función que reciba un vector de números x (la muestra original) y un número natural N , y que genere una matrix que tenga N columnas, siendo cada una una de las muestras de bootstrap.
- Escribe una función de R que reciba un vector de números x (la muestra) y genere una matriz que contenga TODAS las muestras jackknife con el mismo formato del ejercicio anterior.
- Utilizando las funciones anteriores considera la desviación estándar de la media muestral para una muestra de $N=10000$ que se distribuya como en el ejercicio 1 (gaussiana centrada en $1.70m$ y $\sigma=1.7$). Compara la desviación estándar obtenida, con la obtenida con bootstrap.
- Repite el ejercicio anterior utilizando la técnica jackknife. ¿Cuál da mejor resultado?