

Introducción a la metodología bootstrap

Mètodes no paramètrics i de remostratge

Grau Interuniversitari en Estadística UB - UPC

Departament d'Estadística

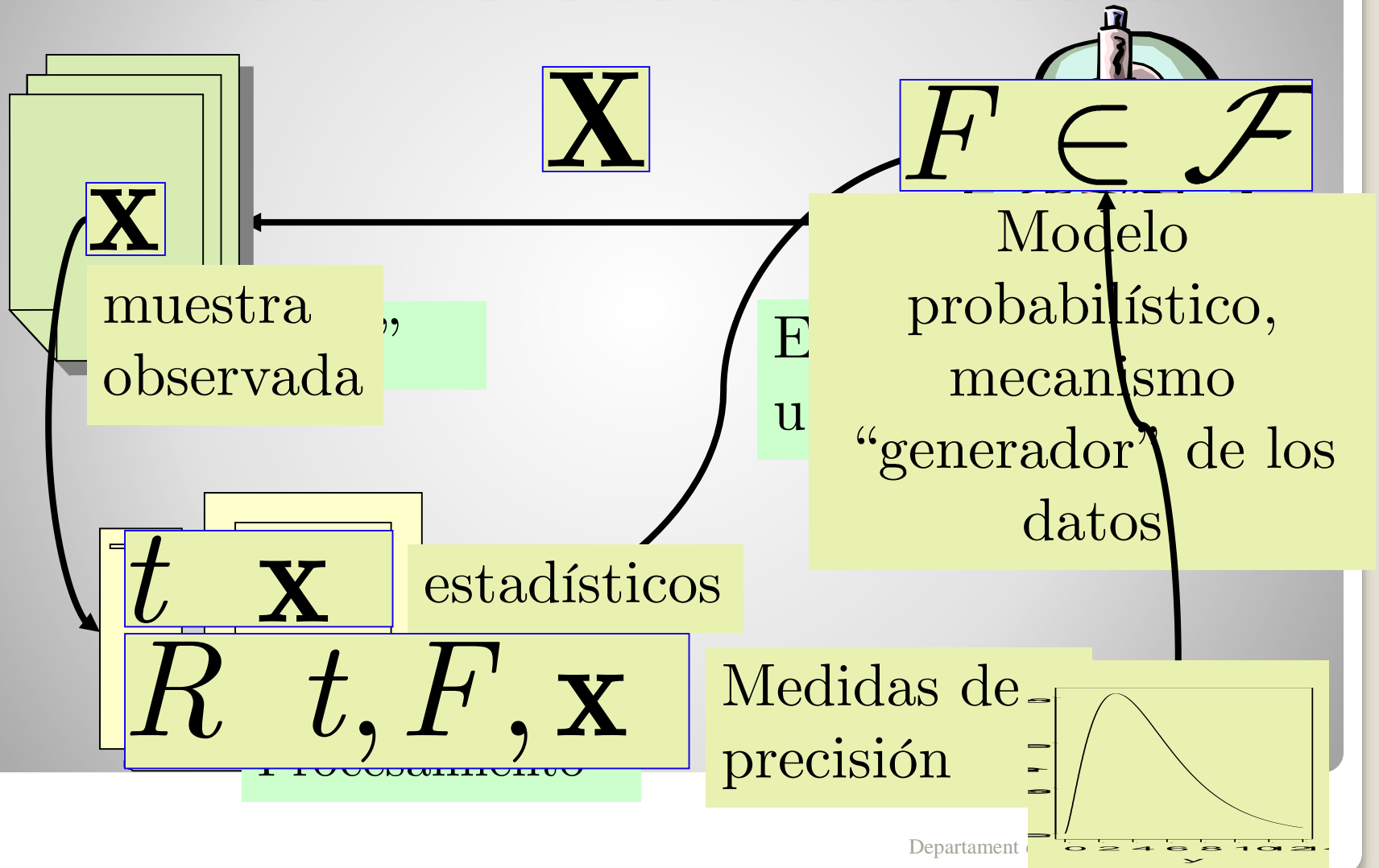
Universitat de Barcelona

Jordi Ocaña Rebull

Puntos a tratar

- Elementos de un problema de inferencia estadística
- Determinación de la distribución muestral (o de alguna de sus características)
- Principio “plug-in” y bootstrap
- Principio de Montecarlo y bootstrap
- Necesaria correspondencia entre “mundo real” y “mundo bootstrap”
- Ejemplos

Elementos de un problema de inferencia estadística



Elementos de un problema de I.E. Ejemplo introductorio

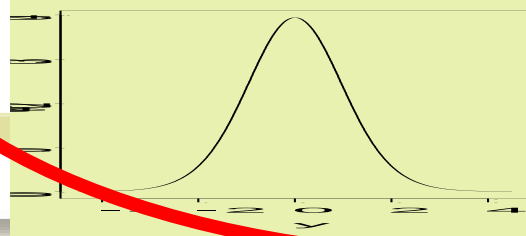
$$\mathbf{x} = x_1, \dots, x_n$$

1. muestra
aleatoria
simple de
tamaño n

$$f(\mathbf{x}; \mu, \sigma^2) = \prod_{i=1}^n \frac{1}{\sigma\sqrt{2\pi}} e^{-\frac{(x_i - \mu)^2}{2\sigma^2}}$$

Normal de media y
varianza
desconocidas

t \mathbf{x}



Ayuda

Opciones...

- Cálculo de p-valores, valores críticos...: indispensable conocer la distribución muestral del estadístico de interés

$$G(t, n, F(x; \theta), \dots)$$

- Cómo determinarla:
 - Paramétrico exacto
 - Asintótico
 - Rangos...

Distribución muestral de un estadístico

Principio “plug-in” y bootstrap (en sentido amplio)

- Fijémonos en el paso $G = G(F; \theta, \dots)$
- ◆ Si \hat{F} es buena estimación de F a partir de los datos, razonable aproximar G mediante:
$$G(\hat{F}, \dots)$$

→ Principio “plug-in”
- ◆ Metodología bootstrap \equiv inferencia basada en el Principio “plug-in”

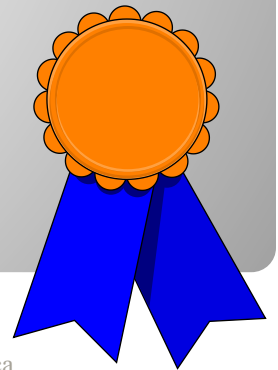
Ejemplo: aplicación automática del Principio "plug-in"

- ◆ A menudo \hat{F} es la distribución empírica, F_n , discreta, que asigna probabilidad $1/n$ a cada valor muestral y 0 a cualquier otro
- ◆ Si interessa característica concreta como

$$\text{var}_F(\bar{X}) = \frac{\text{var}_F(X)}{n}$$

... Según Principio "plug-in":

$$\text{var}_{F_n}(\bar{X}) = \frac{\text{var}_{F_n}(X)}{n} = \frac{s^2}{n}$$



Detalles del cálculo anterior

$$\text{var}_{F_{Fn}}(\bar{X}^*) = \frac{E_{F_m} \left(\left(X^* - E_{F_{Fn}}(X^*) \right)^2 \right)}{n}$$

$$E_{F_{Fn}}(X^*) = \sum_{i=1}^m x_{ij} \frac{1}{m} = \bar{X} \left(= E_{F_m}(\bar{X}) \right)$$

$$E_{F_{Fn}} \left((X^* - \bar{X})^2 \right) = \sum_{i=1}^m (x_{ij} - \bar{X})^2 \frac{1}{m} = s^2$$

- ◆ Conveniencia de notación X^* en lugar de X : no es la misma v.a

Dificultades en la aplicación del Principio “plug-in”

- No tan (o a veces nada) clara su aplicación en situaciones más complejas:
 - otras características de la distribución muestral, incluso para estadísticos sencillos como la media muestral (p.e. un cuantil, ...)
 - otros estadísticos que no sean medias ni funciones sencillas de medias
 - determinación de la distribución muestral completa

$$G(; \hat{F})$$

El método de Montecarlo

F

Modelo probabilístico,
completamente especificado

p.e. n réplicas $N \mu, \sigma^2$ iid

$$\mathbf{x}_1 = (x_{11}, \dots, x_{1n}) \mapsto U(\mathbf{x}_1) = u_1$$

$$\mathbf{x}_2 = (x_{21}, \dots, x_{2n}) \mapsto U(\mathbf{x}_2) = u_2$$

\vdots

$$\mathbf{x}_m = (x_{m1}, \dots, x_{mn}) \mapsto U(\mathbf{x}_m) = u_m$$

(gran)
muestra de
 m valores
del
estadístico
“Leyes de
los grandes
números”

Generación de m
muestras independientes
(o no) según F

$$\frac{1}{m-1} \sum_{j=1}^m (u_j - \bar{u})^2 \cong \text{var}_F U$$

$$\hat{G} \cong G ; F, \text{ etc.}$$

Departament d'estadística

Bootstrap y Montecarlo

$$\hat{F}$$

estimación del Modelo

probabilístico,
p.e. $P^*[X^* = x^*] = \begin{cases} \frac{1}{n} & \text{si } x^* \in x_1, \dots, x_n \\ 0 & \text{en caso contrario} \end{cases}$

$$\mathbf{x}_1^* = (x_{11}^*, \dots, x_{1n}^*) \mapsto U(\mathbf{x}_1^*) = u_1^*$$

$$\mathbf{x}_2^* = (x_{21}^*, \dots, x_{2n}^*) \mapsto U(\mathbf{x}_2^*) = u_2^*$$

\vdots

$$\mathbf{x}_B^* = (x_{B1}^*, \dots, x_{Bn}^*) \mapsto U(\mathbf{x}_B^*) = u_B^*$$

muestra de
 B valores
del
estadístico

“Leyes de
los grandes
números”

Generación de B

“remuestras” de tamaño n

(muestras aleatorias con
reemplazo de los elementos
de \mathbf{x})

$$\frac{1}{B-1} \sum_{b=1}^B (u_b^* - \bar{u}^*)^2 \cong \text{var}_{\hat{F}} U^*$$

$$\hat{G}^* \cong G ; \hat{F} , \text{ etc.}$$

Departament d'estadística

Qué estimamos a partir del Montecarlo bootstrap?

Montecarlo bootstrap

$$\hat{G}^* = \hat{G}(u_1^*, \dots, u_B^*)$$

$$\bar{u}_* = \frac{1}{B} \sum_{b=1}^B u_b^*$$

$$\text{var}_*(U^*) = \frac{1}{B-1} \sum_{b=1}^B (u_b^* - \bar{u}_*)^2$$

$$\hat{P}_* [U^* \geq U(\mathbf{x})] = \frac{\#\{u_b^* \geq U(\mathbf{x})\}}{B}$$

\approx

\approx

\approx

\approx

\approx

Plug-in

$$G(; \hat{F})$$

$$E_{\hat{F}}(U^*)$$

$$\text{var}_{\hat{F}}(U^*)$$

$$P_{\hat{F}}[U^* \geq U(\mathbf{x})]$$

\approx

\approx

\approx

\approx

\approx

"Verdadero" valor del funcional

$$G(; F)$$

$$E_F(U)$$

$$\text{var}_F(U)$$

$$P_F[U \geq U(\mathbf{x})]$$

Error de aproximación de Montecarlo

Problema "clásico" de precisión estadística

- Resultado general (pero no muy útil):
 - Según Leyes de los grandes números, $F_n(x)$ tiende (en diversos sentidos) hacia $F(x)$. Extensible a funciones suficientemente "suaves"
- Validez: resultado sobre funcionales, funciones globales de F_n (u otras estimaciones) y de F : teoremas límite sobre distancias entre distribuciones
- Más interés práctico: comparación entre aproximación bootstrap y otras, para n finito

Validez de la aproximación bootstrap

Características generales de los ejemplos

- Modelo probabilístico subyacente conocido
 - Normal $\mu = 15$, $\sigma = 3$, o bien
 - Exponencial $\alpha = 1/\mu = 1/15$(\Rightarrow **distribución muestral** conocida)
- Análisis de única muestra (pequeña, $n = 10$), generada según uno u otro modelo.
 - caso normal: 15.54, 21.06, 16.52, 13.62, 16.14, 10.98, 13.53, 16.02, 16.79, 15.90
 - caso exponencial: 8.51, 8.71, 69.19, 10.05, 23.64, 8.67, 1.51, 20.36, 1.23, 5.27

Características generales de los ejemplos

- Estadístico bajo estudio: t
- aproximaciones: normal, bootstrap no paramétrico y bootstrap paramétrico
- aproximaciones bootstrap: estima “kernel” a partir de $B = 10000$ valores del estadístico (media o t , según el caso)
- Cada uno de estos valores calculado sobre una remuestra de tamaño $n = 10$

Estadístico t , caso normal: $n = 10, \mu = 15, \sigma = 3$

Verdadera distribución: $t \sim t_{n-1} = 9$

Aproximación normal: $t \approx N(0,1)$

Bootstrap: 1000 valores $t^* = t(\mathbf{x}^*)$

para remuestras $\mathbf{x}^* = (x_1^*, \dots, x_n^*)$

no paramétrico: cada x_i^* escogido con probabilidad $1/n$ entre los de la muestra original

paramétrico: cada x_i^* generado según $N(15.62, 2.63)$

Detalle y justificación del proceso de remuestreo

"Mundo real"

$$\mu = E(X, F)$$

F



$$\mathbf{x} = (x_1, \dots, x_n)$$



$$\bar{x} = \bar{X}(\mathbf{x})$$

$$\hat{s} = \hat{S}(\mathbf{x}) = \sqrt{\frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})^2}$$



$$t = \frac{\sqrt{n}(\bar{x} - \mu)}{\hat{s}}$$

"Mundo bootstrap"

$$\bar{x} = \hat{\mu} = E(X^*, F_n)$$

F_n



$$\mathbf{x}^* = (x_1^*, \dots, x_n^*)$$



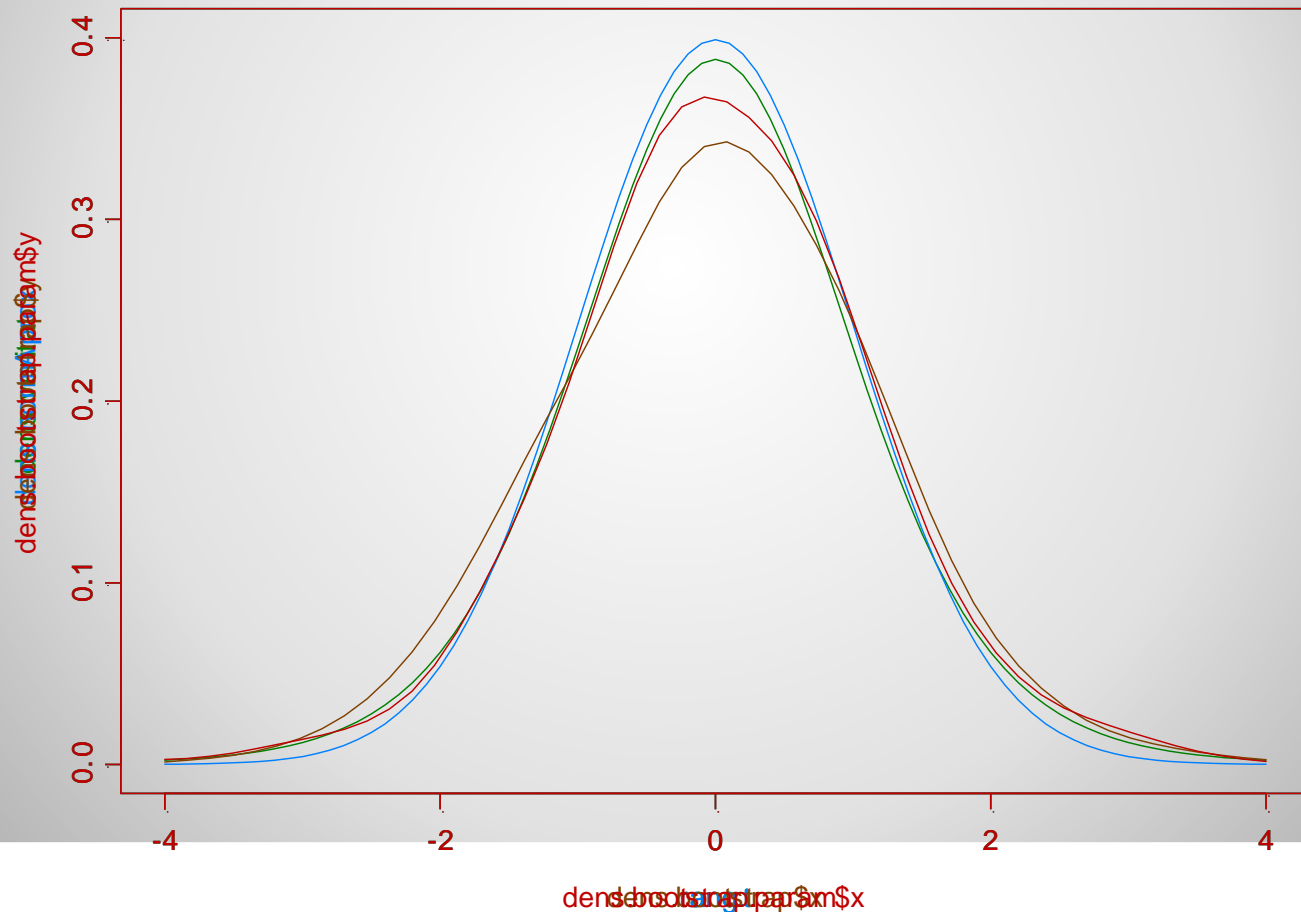
$$\bar{x}^* = \bar{X}(\mathbf{x}^*)$$

$$\hat{s}^* = \hat{S}(\mathbf{x}^*) = \sqrt{\frac{1}{n-1} \sum_{i=1}^n (x_i^* - \bar{x}^*)^2}$$



$$t^* = \frac{\sqrt{n}(\bar{x}^* - \bar{x})}{\hat{s}^*}$$

Estadístico t , normal: verdadera densidad, aprox normal, bootstrap no paramétrico y paramétrico



Estadístico t , exponencial: $n = 10, \alpha = 1/\mu = 1/15$

Verdadera distribución:
estimada por simulación

Aproximación normal: $t \approx N(0,1)$

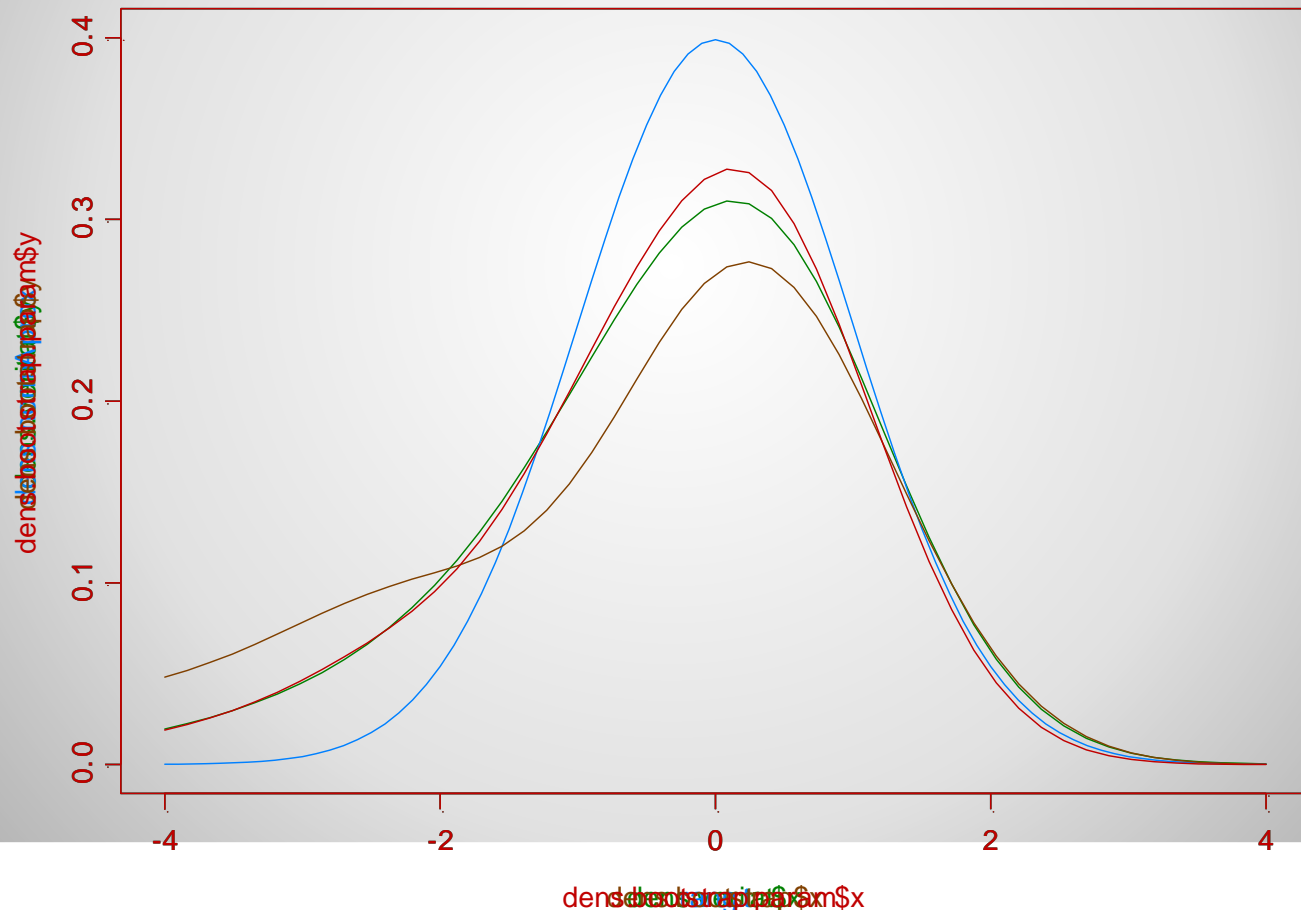
Bootstrap: 1000 valores $t^* = t(\mathbf{x}^*)$

para remuestras $\mathbf{x}^* = (x_1^*, \dots, x_n^*)$

no paramétrico: cada x_i^* escogido con probabilidad $1/n$ entre los de la muestra original

paramétrico: cada x_i^* generado según $Exp(1/15)$

Estadístico t , exponencial: verdadera dens, aprox normal, boot no paramétrico y paramétrico



Caso exponencial, $t, n = 40$

