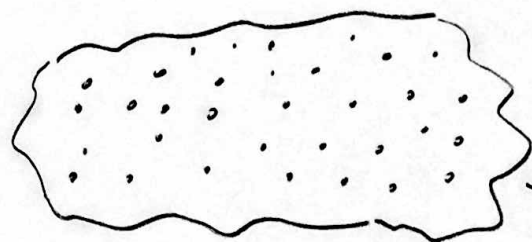


INTERVALOS DE CONFIANZA

1

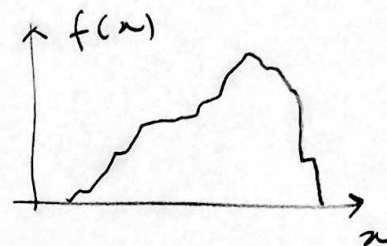


UNIVERSO
(ej. Estudiantes universitarios)

→ Cada muestra puede ser medida con una variable x (ej. Edad)

x : variable aleatoria con función de densidad de probabilidad $f(x)$, i.e.,

$$\int_{-\infty}^{+\infty} f(x) dx = 1$$



$f(x)$ tiene una media $\mu = \int_{-\infty}^{+\infty} x f(x) dx$

y una varianza $\sigma^2 = \int_{-\infty}^{+\infty} (x - \mu)^2 f(x) dx$

¿Cómo estimar μ ?

¿Qué tan buena es su estimación?

TEOREMA DEL LIMITE CENTRAL

Si se toman N muestras del universo: x_1, x_2, \dots, x_N

se puede calcular $\hat{\mu} = \frac{1}{N} \sum_{i=1}^N x_i$, Si este proceso lo

repetimos M veces (sacando N muestras aleatorias y promediándolas cada vez) entonces tendremos una nueva variable aleatoria $\hat{\mu}_1, \hat{\mu}_2, \dots, \hat{\mu}_M$ cuya media y varianza:

$$\tilde{\mu} = \frac{1}{M} \sum_{k=1}^M \hat{\mu}_k$$

$$\tilde{\sigma}^2 = \frac{1}{M-1} \sum_{k=1}^M (\hat{\mu}_k - \tilde{\mu})^2$$

Si M es suficientemente grande, la variable aleatoria tiene una distribución normal con media $\tilde{\mu} = \mu$ y varianza $\tilde{\sigma}^2 = \frac{\sigma^2}{N}$. \rightarrow PATOS - Teorema Central.m

(Si M es pequeño, e.g. $M < 20$, la distribución es t-Student con $M-1$ grados de libertad).

Ejemplo



En una bolsa hay millones de "1" y "0" escritos en bolas. ¿Cuál es la probabilidad de sacar x "1" al sacar N bolas?

Supongamos que la probabilidad de sacar un "1" es p (y la de sacar un "0" es $q = 1-p$).

bolas: $\left\{ \begin{array}{cccccc} N & N-1 & \dots & 2 & 1 \\ 0 & 0 & & 0 & 0 \\ 0 & 0 & & 0 & 1 \\ 0 & 0 & & 1 & 0 \\ 0 & 0 & & 1 & 1 \\ & & & & & \vdots \\ 1 & 1 & & 1 & 1 \end{array} \right.$

2^N posibilidades

existen $\frac{N!}{x!(N-x)!}$ posibilidades.

des (en las 2^N posibles combinaciones) de sacar exactamente x "1".

$$p(x) = \frac{N!}{x!(N-x)!} p^x q^{N-x}$$

Distribución de Bernoulli
con $\mu = p$ $\sigma^2 = pq$

Ejemplo MATLAB:

$$M = 1000000 \quad p = 0.75 \quad q = 1-p$$

$$T = \text{rand}(M, 1) \geq q$$

$$\text{mean}(T) \approx 0.75 (p)$$

$$\text{std}(T) \approx \sqrt{pq}$$

Si repetimos el experimento M veces, i.e. sacamos M veces N bolas

$M \left\{ \begin{array}{cccccc} 0 & 1 & 0 & \dots & 0 & 1 \rightarrow \hat{\mu}_1 \\ 0 & 0 & 1 & \dots & 1 & 1 \rightarrow \hat{\mu}_2 \\ & & & & & \vdots \\ 1 & 0 & 1 & \dots & 0 & 0 \rightarrow \hat{\mu}_M \end{array} \right.$

Por el teorema central $\tilde{\mu} = \frac{1}{M} \sum \hat{\mu}_k \approx \mu = p$

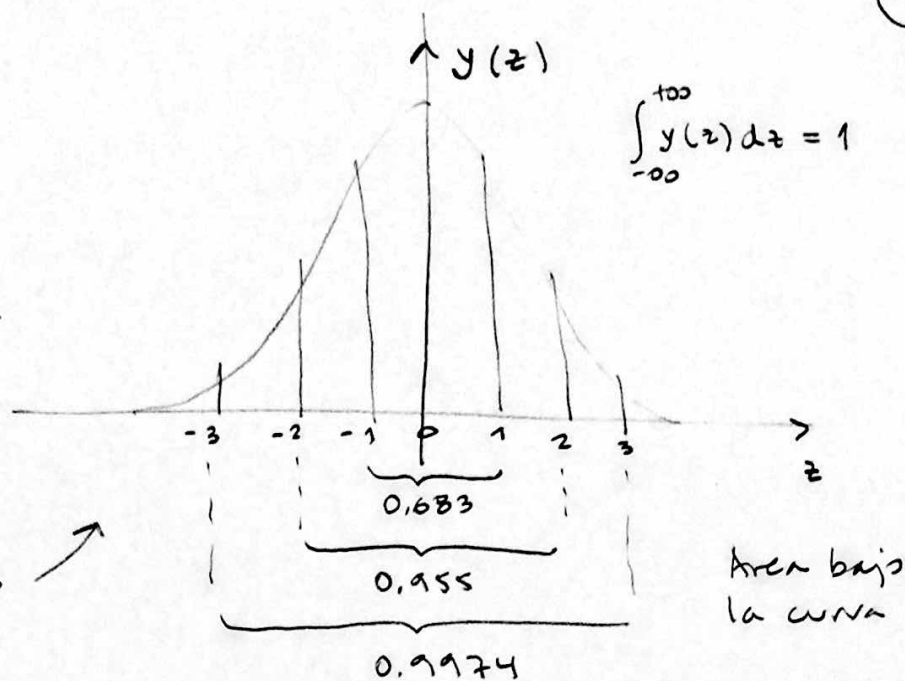
$$\tilde{\sigma}^2 \approx \frac{\sigma^2}{N} = \frac{pq}{N}$$

\rightarrow PATOS - Bernoulli.m

DISTRIBUCIÓN NORMAL

Si tenemos una distribución normal con media $\tilde{\mu}$ y varianza $\tilde{\sigma}^2$ podemos hacer un cambio de coordenadas

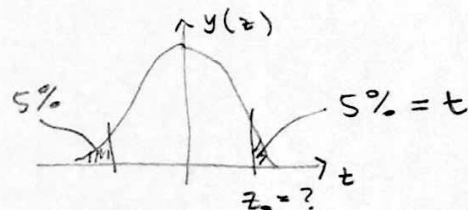
$$z = \frac{x - \tilde{\mu}}{\tilde{\sigma}} \text{ y obtenemos}$$



Si se tiene una variable x de un universo (ej. edad de estudiantes) y calculamos μ_k como el promedio de N muestras aleatorias, y repetimos este experimento M veces, tendremos una estimación de $\tilde{\mu}$ que tiene una varianza $\tilde{\sigma}$. Sabemos entonces que con un 95.5% de probabilidad esta estimación de $\hat{\mu}$ está entre $z = -2$ y $z = 2$, es decir entre $\mu \pm 2\sigma$. Este es el intervalo de confianza de la estimación.

Ejemplo: ¿Cómo calcular el intervalo de confianza para una probabilidad del 90%?

1) Se define $t = \frac{1-c}{2}$ ($c = 0.90$)



2) Se busca la inversa de la normal $z_0 = \text{norminv}(1-t)^*$

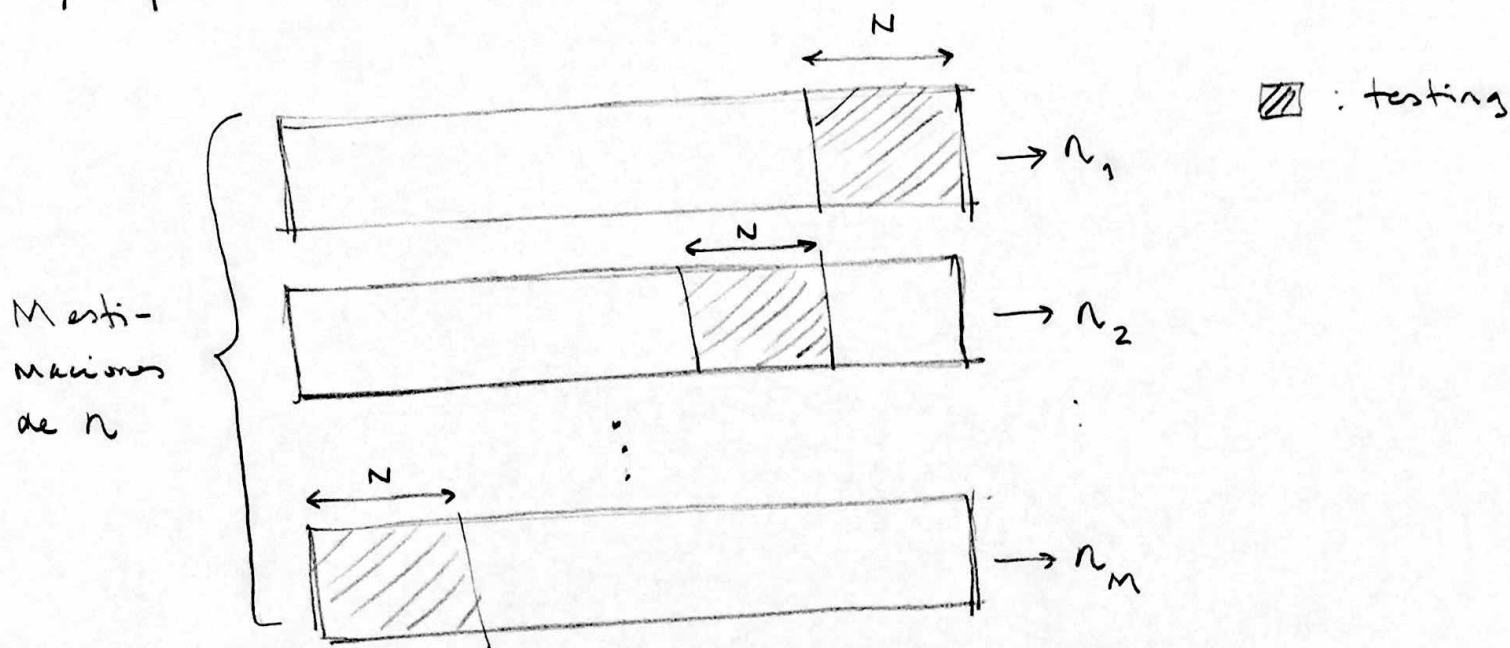
3) Intervalo de confianza:

$$\tilde{\mu} - \tilde{\sigma} z_0 \leq \tilde{\mu} \leq \tilde{\mu} + \tilde{\sigma} z_0$$

→ ojo: $\tilde{\sigma} = \frac{\sigma}{\sqrt{N}}$

* Para t-student $z_0 = t_{inv}(1-t, M-1)$

Ejemplo de Intervalos de Confianza en Validación Cruzada:



N muestras de testing en cada experimento

Diagram showing the testing results for M experiments, each with N samples.

	N				
M	0	1	0		$\rightarrow n_1$
	0	0	1		$\rightarrow n_2$
	1	1	1		\vdots
	1	0	1		$\rightarrow n_M$

"1" bien clasificado

"0" mal clasificado

\Rightarrow Distribución de Bernoulli $\rightarrow \hat{p} = \frac{1}{M} \sum_{k=1}^M n_k \approx p$

$\hat{\sigma}^2 = \frac{\hat{\sigma}^2}{N} \approx \frac{p(1-p)}{N}$

¿Cómo se calculan los intervalos de confianza?

Ø. Definir la probabilidad del intervalo de confianza c , calcular $t = (1-c)/2$.

1. Realizar M experimentos con N datos q/n (ver arriba)

2. Estimar n_k for $k=1, \dots, M$

3. Calcular $p = \frac{1}{M} \sum n_k$; $q = 1-p$; $\hat{\sigma} = \sqrt{pq}/N$

4. If $M > 20$ $z_0 = \text{norminv}(1-t)$ else $z_0 = t_{\text{inv}}(1-t, M-1)$

5. Intervalo de confianza $\hat{n} - \Delta n \leq \hat{n} = p \leq \hat{n} + \Delta n$
con $\Delta n = z_0 \hat{\sigma}$.

DISTRIBUCIONES:NORMAL:

$$y(z) = \frac{1}{\sqrt{2\pi}} e^{-z^2/2}$$

t-Student:

$$y(z) = \frac{\Gamma\left(\frac{v+1}{2}\right)}{\sqrt{v\pi} \Gamma\left(\frac{v}{2}\right)} \left(1 + \frac{z^2}{v}\right)^{-(v+1)/2}$$

 v : grados de libertad

$$\Gamma(v) = (v-1)!$$

→ PADS-normal-vs-t.m