

# Probabilidad y estadística

## Clase 5



# Estimación no paramétrica

# Función de distribución empírica

**Def:** Sea  $\underline{X}_n$  una m.a. tal que  $X_i \overset{i.i.d.}{\sim} F$ , donde  $F$  es una función de distribución. La **función de distribución empírica (ECDF)** es una función  $\hat{F}_n$  que pone masa  $1/n$  en cada observación  $X_i$ .

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

# Ejercicio 1

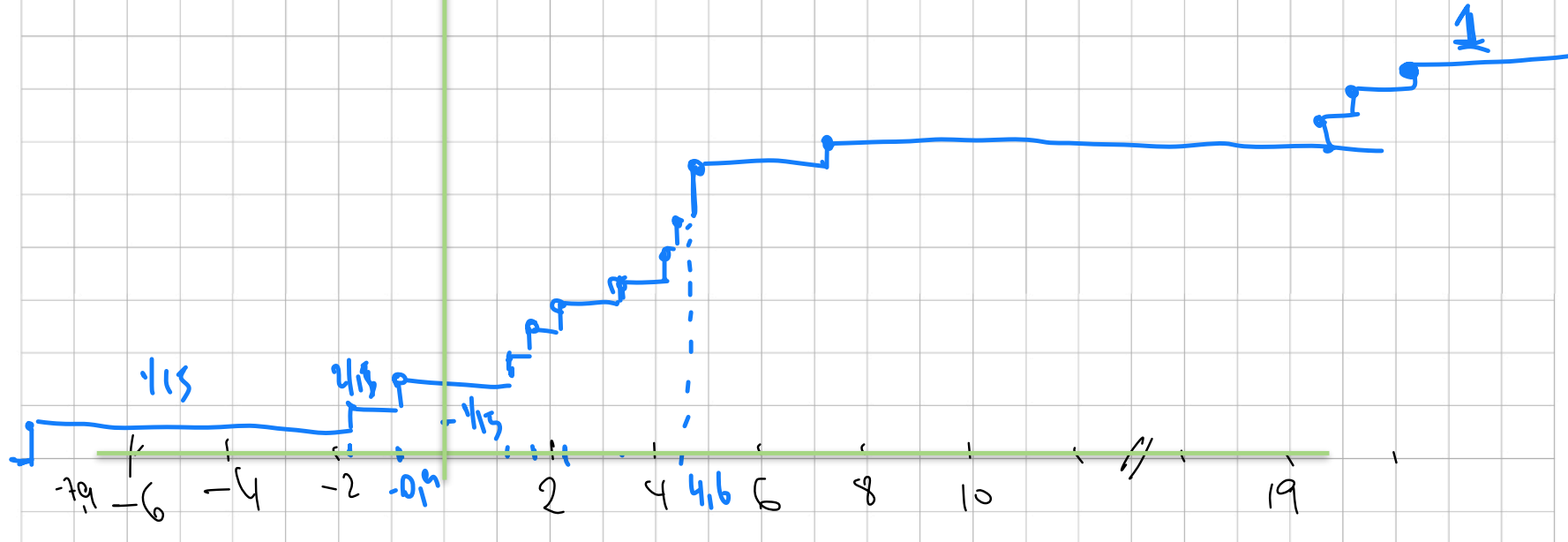
De un experimento en los efectos de un medicamento para la ansiedad, entre otras cosas se midió la diferencia (en segundos) entre el puntaje de un test de memoria antes y después de tomar el medicamento, obteniendo los siguientes resultados:

1.2, 4.6, 4.3, 4.2, -7.9, 7.8, 3.4, 19.8, 25.5, -1.9, 2.1, -0.9, 4.6, 21.1, 1,7

1. Obtener la función de distribución empírica a mano.
2. Utilizar la columna 'Diff' del dataset `Islander_data.csv` y calcular la func. de distribución empírica usando Python.

1.2, 4.6, 4.3, 4.2, 7.9, 7.8, 8.4, 19.8, 25.5, 21.9, 21.1, 4.9, 4.6, 21.1, 4.2

-7.9   -1.9   -0.9   1.2   1.7   2.1   3.4   4.2   4.3  
 4.6   4.6   7.8   19.8   21.1   25.5



# Propiedades de la ECDF

$$\begin{aligned}\mathbb{E} \left( \widehat{F}_n(x) \right) &= F(x), \\ \mathbb{V} \left( \widehat{F}_n(x) \right) &= \frac{F(x)(1 - F(x))}{n}, \\ \text{MSE} &= \frac{F(x)(1 - F(x))}{n} \rightarrow 0, \\ \widehat{F}_n(x) &\xrightarrow{\text{P}} F(x).\end{aligned}$$

# Estimación de densidades (*smoothing*)

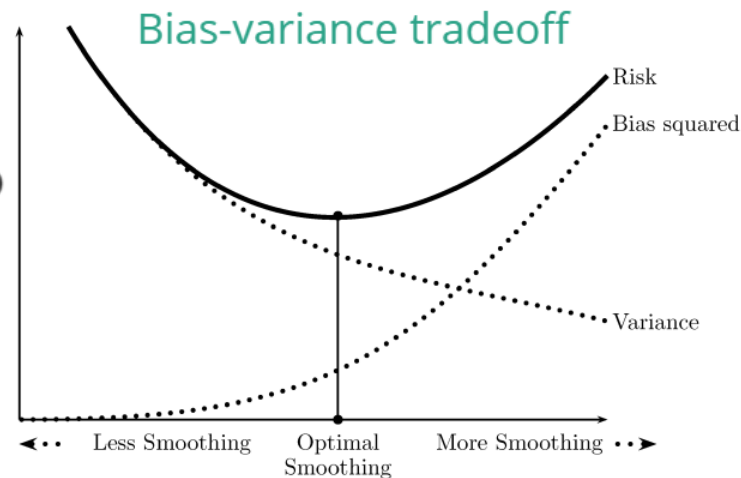
Si deseamos estimar una función de densidad  $f(x)$  o una función de regresión  $\phi(x) = \mathbb{E}[X|Y = y]$ , se deben hacer algunas suposiciones de suavidad.

Sea  $\hat{g}_n$  un estimador de  $g$ .

Definimos el **riesgo** (error cuadrático medio integrado (MISE)) como

$$R(g, \hat{g}_n) = \mathbb{E} \left[ \int (g(u) - \hat{g}_n(u))^2 du \right]$$

All of Statistics, Wasserman





# Histogramas

1. Se selecciona un origen  $x_0$  y se divide la recta real en intervalos de longitud  $h$

$$B_j = [x_0 + (j - 1)h, x_0 + jh], j \in \mathbb{N}$$

Se cuenta cuantas observaciones caen en cada intervalo armando una tabla

1. de frecuencias. Denotamos a la cantidad de observaciones que caen en el intervalo  $j$  como  $n_j$

Para cada intervalo, se divide la frecuencia absoluta por la cantidad total de la muestra  $n$  (para convertirlas en frecuencias relativas, análogo a como se hace con las probabilidades) y por la longitud  $h$  (para asegurarse que el area debajo del histograma sea igual a 1):

Formalmente, el histograma está dado por:

$$\hat{f}_h(x) = \frac{1}{nh} \sum_{i=1}^n \sum_j \mathbf{1}(x_i \in B_j) \mathbf{1}(x \in B_j)$$

Apunte de Histograma - PyE FIUBA

# Ejercicio 2

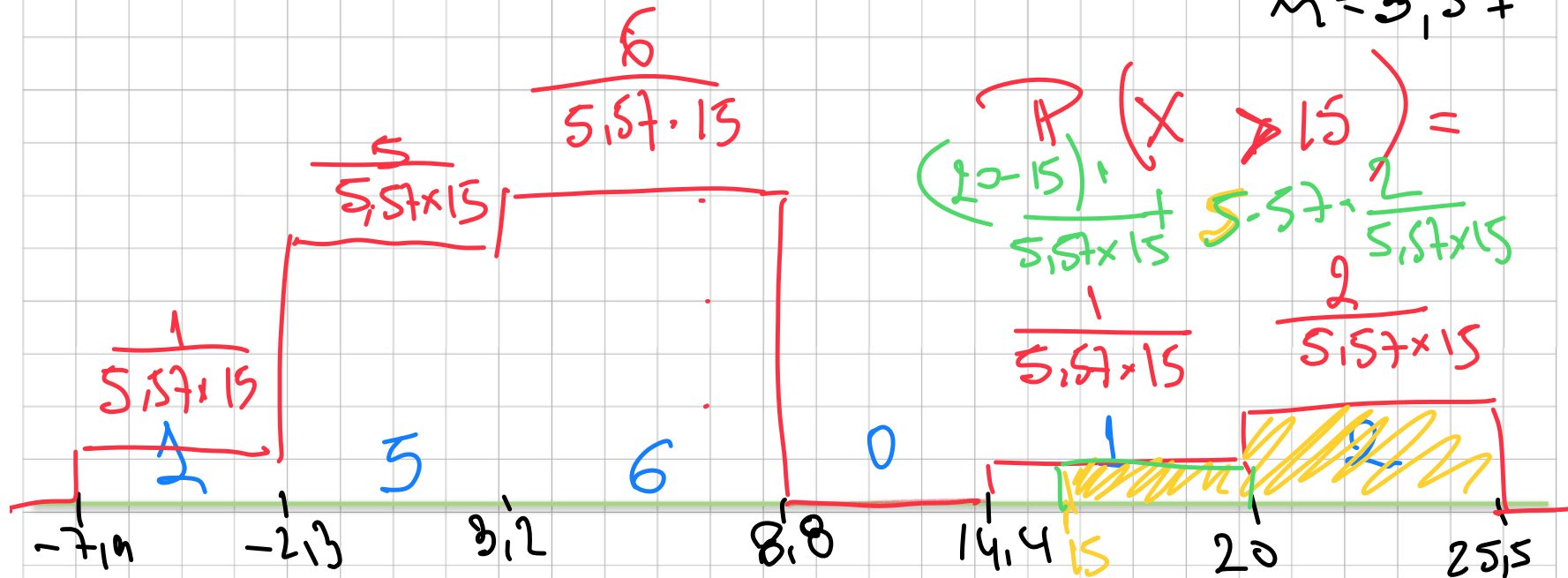
A partir de los datos del ejercicio 1,

1. Calcular a mano, el histograma de 6 bins
2. A partir de los datos del dataset graficar el histograma de la columna 'Diff' utilizando Python

1.2, 4.6, 4.3, 4.2, -7.9, 7.8, 3.4, 19.8, 25.5, -1.9, 2.1, -0.9, 4.6, 21.1, 1.7

~~-7.9~~ -1.9 -0.9 1.2 1.7 2.1 3.4 4.2 4.3  
 4.6 4.6 7.8 19.8 21.1 25.5

$$\mu = 5.57$$

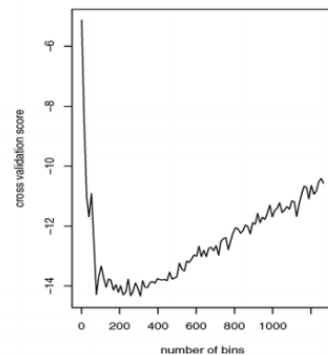
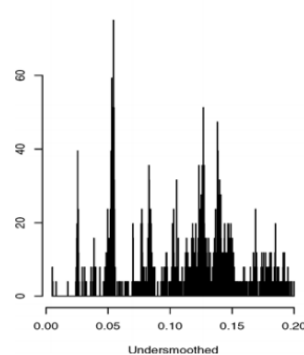
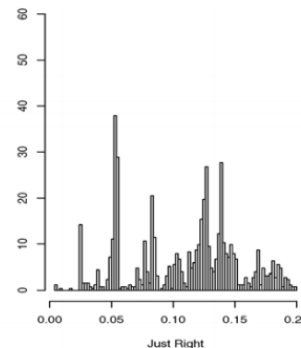
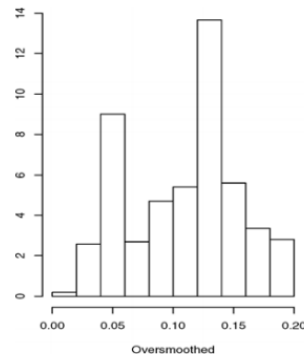


# Propiedades del histograma

**Teorema:** Sea  $x$  y  $m$  fijos, y sea  $B_n$  el bin que contiene a  $x$ , luego

$$\mathbb{E}(\hat{f}_n(x)) = \frac{p_j}{h} \quad \mathbb{V}(\hat{f}_n(x)) = \frac{p_j(1 - p_j)}{nh^2}.$$

**Obs:** Al aumentar la cantidad de bins ( $m$ ), Disminuye el sesgo, pero aumenta la varianza. Acá esta el tradeoff.



# Estimación de densidad por kernel

Los histogramas son discontinuos, los **estimadores de densidad por kernel (KDE)** son una versión más suave y convergen más rápido a la densidad verdadera que el histograma.

# Kernels

Se define un **kernel** como una función  $K$  suave tal que:

$$K(x) \geq 0, \int K(x)dx = 1, \int xK(x)dx=0, \text{ y}$$

$$\sigma_K^2 = \int x^2 K(x)dx > 0.$$

Algunos kernels comunes:

- Epanechnikov:  $K(x) = \begin{cases} \frac{3}{4}(1 - x^2/5)/\sqrt{5}, & |x| < 5 \\ 0 & \text{e. o. c.} \end{cases}$

Es óptima en el sentido de error cuadrático medio

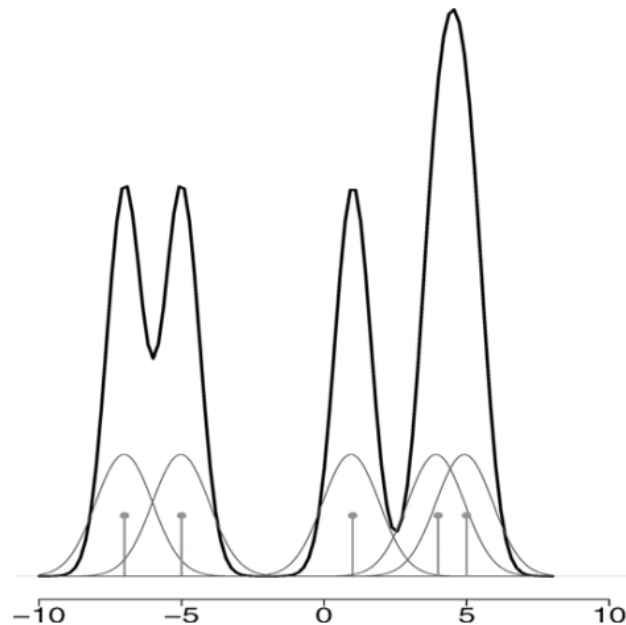
- Gaussiano (simple)

# KDE

**Def:** Dado un kernel  $K$  y un número positivo  $h$ , llamado **ancho de banda**, el **estimador de densidad por kernel** se define como

$$\hat{f}(x) = \frac{1}{n} \sum_{i=1}^n \frac{1}{h} H\left(\frac{x-X_i}{h}\right)$$

Nuevamente el parámetro  $h$  es el que nos controla el tradeoff sesgo-varianza



## Ejercicio 3

A partir de la columna 'Diff' del dataset `Islander_data` estimar la densidad por el método de KDE. Analizar qué ocurre al tomar distintos valores de  $h$ .



# Intervalos de confianza

# Motivación

Hasta ahora habíamos visto estimadores puntuales, que, dada una muestra, nos devuelven un único valor  $\hat{\theta}$  que se aproxima al valor verdadero del parámetro deseado  $\theta$ .

Una forma de obtener información sobre la precisión de la estimación, en el caso de que  $\theta$  sea unidimensional, es proporcionar un intervalo  $[a(X), b(X)]$  de manera que la probabilidad de que dicho intervalo contenga el verdadero valor  $\theta$  sea alta, por ejemplo, 0.95.

# Región de confianza

**Def:** Dada una m.a.  $\underline{X}$  con distribución perteneciente a una familia  $F_\theta(x)$ , con  $\theta \in \Theta$ , una **región de confianza**  $S(\underline{X})$  para  $\theta$  con nivel de confianza  $1 - \alpha$  será un conjunto tal que

$$\mathbb{P}(\theta \in S(\underline{X})) = 1 - \alpha. (*)$$

**Obs:**  $\theta$  **no** es aleatorio, lo aleatorio es  $(*)$  es  $S(\underline{X})$ .

**Obs:** Si  $S(\underline{X}) = (a(\underline{X}), b(\underline{X}))$  diremos que es un **intervalo de confianza**.

Si  $S(\underline{X}) = (\min(\Theta), b(\underline{X}))$  diremos que es una **cota superior**.

Si  $S(\underline{X}) = (a(\underline{X}), \max(\Theta))$  diremos que es una **cota inferior**.

# Juguemos un poquito

Usemos la siguiente api para entender mejor qué es un IC

# Método del pivote

**Teorema:** Sea  $\underline{X}$  una muestra aleatoria con distribución perteneciente a una familia  $F_\theta(x)$ , con  $\theta \in \Theta$ , y sea  $U = g(\underline{X}, \theta)$  una variable cuya distribución **no** depende de  $\theta$ . Sean  $a$  y  $b$  tales que  $\mathbb{P}(a \leq U \leq b) = 1 - \alpha$ . Luego,

$$S(\underline{X}) = \{\theta : a < g(\underline{X}, \theta) \leq b\}$$

es una región de confianza para  $\theta$ . A  $U$  se lo llama **pivote**.

## Ejercicio 4

Sea  $\underline{X} = (X_1, \dots, X_n)$  una muestra aleatoria de tamaño  $n$  de una población con distribución normal de media  $\mu$  y varianza 4. Hallar una cota inferior del 95% para  $\mu$ .

Suponer  $n=20$  y  $\mu=3$ , simular la muestra y obtener el valor de la cota

$$\underline{X} = (X_1, \dots, X_n)$$

$$X_i \stackrel{iid}{\sim} N(\mu, 4)$$

busca ~~con~~  
inferior ~~para~~  $\mu$ .  
 $\mu > d(x)$

$$\hat{\mu} = \frac{1}{n} \sum_{i=1}^n X_i = \bar{X}$$

$$\mu = E[\bar{X}] \Rightarrow$$

LGA dice  
que  $\frac{1}{n} \sum_{i=1}^n X_i \rightarrow E(\bar{X}) = \mu$

$$\underline{\bar{X} \sim N(\mu, \frac{4}{n})}$$

$$\text{var}(\bar{X}) = \text{var}\left(\frac{1}{n} \sum X_i\right) = \frac{1}{n^2} \text{var}\left(\sum X_i\right) \stackrel{\text{indep}}{=} \frac{1}{n^2} \sum \underbrace{\text{var}(X_i)}_4 = \frac{n \cdot 4}{n^2} = \frac{4}{n}$$

~~E~~standarizaje

$$\frac{\bar{X} - \mu}{\sqrt{4/n}}$$

$$= \frac{\bar{X} - 4}{2} \sqrt{n} \sim$$

$$N(0, 1)$$

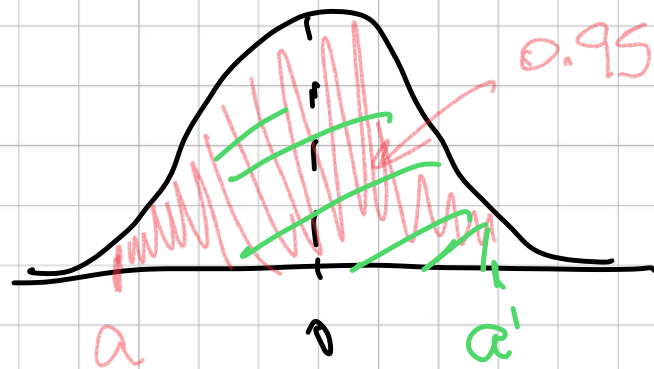
$$U(\bar{X}, \mu)$$

no depende  
de  $\mu$   
Control  
pivot

Busco una cota inferior

$$P(a \leq U(X, \mu)) \geq 0,95$$

$$P(a \leq \frac{\bar{X} - \mu \sqrt{n}}{2}) \geq 0,95$$



$$\hookrightarrow a = z_{0,05} = +1,64$$

0,95

$$P\left(\frac{\bar{X} - \mu \sqrt{n}}{2} \leq a\right)$$

$$\frac{\bar{X} - \mu \sqrt{n}}{2} \geq +1,64$$

Si  $\mu \uparrow \Rightarrow U(X, \mu) \downarrow$

$$\bar{X} - \mu \geq +1,64 \frac{2}{\sqrt{n}} \rightsquigarrow \boxed{\bar{X} \pm 1,64 \frac{2}{\sqrt{n}} \geq \mu}$$

Si  $U(X, \mu)$  es decreciente en  $\mu \Rightarrow$  hay que invertir los signos



# Algunos resultados importantes

**Teorema:** Sea  $\underline{X} = X_1, \dots, X_n$  una m.a. de una distribución  $\mathcal{N}(\mu, \sigma^2)$

$$Z = \sqrt{n} \frac{(\bar{X} - \mu)}{\sigma} \sim \mathcal{N}(0, 1)$$

$$W = \sum_{i=1}^n \frac{(X_i - \bar{X})^2}{\sigma^2} \sim \chi_{n-1}^2$$

$V$  y  $W$  son independientes

$$\text{Si } S^2 = \frac{1}{n-1} \sum_{i=1}^n (X_i - \bar{X})^2, U = \sqrt{n} \frac{(\bar{X} - \mu)}{S} \sim t_{n-1}$$

**Obs:** en general vale que si  $X \sim \mathcal{N}(0, 1)$  y  $Y \sim \chi_n^2$ , con  $X$  e  $Y$  independientes vale que  $\frac{X}{\sqrt{Y/n}} \sim t_n$

# Algunos pivotes para variables normales

Dada  $\underline{X}_n$  una m.a. de una distribución  $\mathcal{N}(\mu, \sigma^2)$  definimos algunos pivotes:

- Para la media con varianza conocida:  $U(\underline{X}, \mu) = \frac{(\bar{X} - \mu)}{\sigma} \sqrt{n} \sim \mathcal{N}(0, 1)$
- Para la media con varianza desconocida:  $U(\underline{X}, \mu) = \frac{(\bar{X} - \mu)}{S} \sqrt{n} \sim t_{n-1}$
- Para el desvío con media conocida:  $U(\underline{X}, \sigma) = \frac{\sum_{i=1}^n (X_i - \mu)^2}{\sigma^2} \sim \chi_n^2$
- Para el desvío con media desconocida:  $U(\underline{X}, \sigma) = \frac{\sum_{i=1}^n (X_i - \bar{X})^2}{\sigma^2} \sim \chi_{n-1}^2$

Dada también  $\underline{Y}_m$  una m.a. de una distribución  $\mathcal{N}(\lambda, \sigma^2)$  y sea :

- Comparación de medias con varianzas conocidas:  $U(\underline{X}, \Delta) = \frac{\bar{X} - \bar{Y} - \Delta}{\sqrt{\frac{\sigma_X^2}{n} + \frac{\sigma_Y^2}{m}}} \sim \mathcal{N}(0, 1)$
- Comparación de medias con varianzas desconocidas e iguales:

$$U(\underline{X}, \Delta) = \frac{\bar{X} - \bar{Y} - \Delta}{S_p \sqrt{\frac{1}{n} + \frac{1}{m}}} \sim t_{n+m-2}, \text{ con } S_p^2 = \frac{(m-1)S_X^2 + (n-1)S_Y^2}{n+m-2}$$

## Ejercicio 5

Dada una muestra aleatoria  $\underline{X} = (X_1, \dots, X_n)$  de una población con distribución normal con media y varianza desconocidas, hallar el intervalo de confianza de nivel 0.99 para la media de la población.

Suponer  $n=50$ ,  $\mu = 2$ ,  $\sigma = 3$ , simular la muestra y calcular el IC resultante de la misma.

- .

$$\underline{X} = (X_1, \dots, X_m) \quad X_i \stackrel{iid}{\sim} N(\mu, \sigma^2)$$

desconocido. Busco IC 0.99

$$U(\underline{X}, \mu) = \frac{\bar{X} - \mu}{S/\sqrt{n}}$$

$$S = \sqrt{\frac{1}{m-1} \sum_{i=1}^m (X_i - \bar{X})^2}$$

x.stat()

$$P\left(a < \frac{\bar{X} - \mu}{S/\sqrt{n}} < b\right) \geq 0.99$$

$$-z_{0.995} < \frac{\bar{X} - \mu}{S/\sqrt{n}} < z_{0.995}$$

$$-z_{0.995} \frac{S}{\sqrt{n}} < \bar{X} - \mu < z_{0.995} \frac{S}{\sqrt{n}}$$

$$+ z_{0.995} \frac{S}{\sqrt{n}} + \bar{X} > \mu > -z_{0.995} \frac{S}{\sqrt{n}} + \bar{X}$$

$$IC: \bar{X} \pm z_{0.995} \frac{S}{\sqrt{n}}$$



$$a = z_{0.005} = -z_{0.995}$$

$$b = z_{0.995}$$

# Bibliografía

- "Notas de Estadística", Graciela Boente y Víctor Yohai, FCEyN, UBA.
- "All of Statistic: A concise Course in Statistical Inference", Larry Wasserman