



# Inferència estadística

Bloc 4 – Probabilitat i Estadística

Novembre 2019

# Índex

1. Introducció a la inferència estadística
2. Distribucions provinents de la Normal
3. Estimació puntual
4. Estimació per interval
  - a. Intervals de Confiança (IC) de  $\mu$ ,  $\pi$  i  $\sigma$
5. Proves d'hipòtesis (PH)
  - a. Mecànica de les PH
  - b. P-valor
  - c. PH per  $\mu$ ,  $\pi$  i  $\sigma$
6. Annexes: Premissa de Normalitat

# Inferència estadística. Guió

Guió de la part d'Estadística de PE:

- **B4:** Tècnica general de la inferència estadística
  - Estimar un paràmetre (*Intervals de Confiança*)
  - Refutar un paràmetre (*Proves d'Hipòtesis*)
- **B5:** Aplicació (I): Avaluació de millores
  - *Disseny d'experiments*: comparació de dues poblacions.
- **B6:** Aplicació (II): Predicció
  - *Previsió* d'una variable resposta en funció d'una variable explicativa.

# Inferència estadística. Introducció

- **La ciència vol ser refutable:** El criteri de refutabilitat sosté que per ser col·locats en el rang de científics, els enunciats han de poder entrar en conflicte amb observacions possibles. [Ex: “Els marcians existeixen” : no hi ha mitjans per contradir-ho → no és científic]
- **La tècnica vol ser documentable:** S’ha de aportar evidència basada en dades [Ex: “El meu programa funciona”: si no aportes proves/dades → no és tècnic]
- La metodologia estadística permet la inducció: **inferir les característiques de la població a partir de les observacions d’una mostra** [Ex: Per conèixer la velocitat mitjana real de connexió a tota la població amb un determinat proveïdor recullo una mostra de velocitats]



# Inferència estadística. Introducció

- La Inferència Estadística defineix i **quantifica els riscos** d'aquest procés [Ex: No es pot conèixer la mitjana de la vel. de connexió a tota la població a no ser que es tingui dades de tota la població, però la estadística em permet quantificar l'error en l'estimació en una mostra concreta]
- Mètode científic i tècnic (estadístic):
  - per **deducció** → disseny de la recollida de dades (Població → Mostra)
  - per **inducció** → inferir (estimar) resultats (Mostra → Població)

## Exemples:

- Vull dir “El meu programa funciona bé”
  - 1) recollida de dades (*proves o evidència*)
  - 2) anàlisi: estimar una mesura (ex: mitjana del rendiment)
- “El meu programa millora els resultats de ...”
  - 1) recollida de dades (*proves o evidència*)
  - 2) anàlisi: poder refutar la igualtat d'una mesura (ex: mitjana de rendiments)

# Inferència estadística. Mostra Aleatòria Simple (MAS)

Sigui la VA:

$$X: \Omega \rightarrow \mathbb{R}$$

$$\omega_i \rightarrow X(\omega_i) = x_i$$

Direm que

M.A.S. de grandària  $n$  de la v.a.  $X$

és una funció vectorial  $M = (X_1, X_2, \dots, X_n)$  tal que

$$M: \Omega^n \rightarrow \mathbb{R}^n$$

$$\omega = (\omega_1, \omega_2, \dots, \omega_n) \rightarrow M(\omega) = (X_1, X_2, \dots, X_n)$$

Direm que és una MAS si i només si es compleixen les dues condicions següents:

- (1) **Tots els elements** de la població tenen la **mateixa probabilitat** de pertànyer a la mostra.
- (2) **Qualsevol combinació** de  $n$  elements té la **mateixa probabilitat** de pertànyer a la mostra.

La informació aportada per les diferents unitats ha de ser **independent** entre sí:

- les  $X_i$  han de ser VA independents i idènticament distribuïdes (i.i.d.)

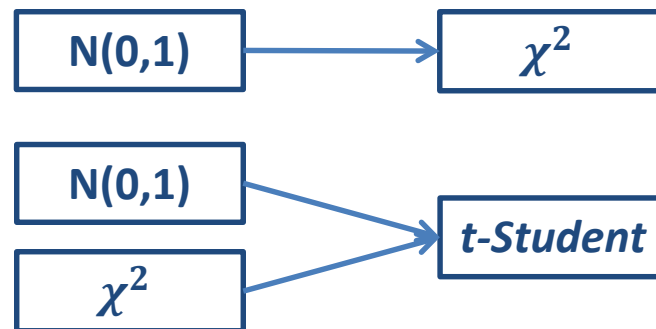
# Inferència estadística. Conceptes bàsics

- **Paràmetre**: Indicador de la població que estem interessats en conèixer o estimar. [Ex: La  $\mu$  (esperança) de les alçades dels estudiants de la FIB]
- **Estadístic**: qualsevol indicador que s'obtingui com a funció de les dades d'una mostra. [Ex: La suma de les alçades dels estudiants recollits en una mostra]
- **Estimador**: estadístic d'una mostra que s'utilitza per conèixer el valor d'un paràmetre de la població. [Ex: La mitjana de les alçades en una mostra d'alumnes de la FIB és una estimador de la  $\mu$  (esperança) de les alçades dels estudiants de la FIB]

**Atenció:** La paraula *mitjana* pot voler dir *paràmetre esperança* quan parlem del centre de gravetat d'una distribució poblacional, o *estadístic mitjana* quan ens referim al promig d'una sèrie de valors obtinguts d'una mostra.

# Models derivats de la Normal: $\chi^2$ i *t-Student*

- Hi ha un parell de distribucions noves que serviran per abordar el contingut d'aquest tema:  $\chi^2$  i **t-Student**
- Aquestes distribucions provenen de fer operacions amb VA provinents d'altres distribucions, entre elles la Normal estàndard.



- A diferència de les distribucions vistes prèviament NO modelen fenòmens de la vida real, sinó el comportament dels estadístics entre les possibles mostres.



# Distribució $\chi^2$ (chi-quadrat)

- Definició:** Siguin  $X_i \sim N(0,1)$ . Llavors:

$$X_1^2 + X_2^2 + \dots + X_n^2 \sim \chi_n^2$$

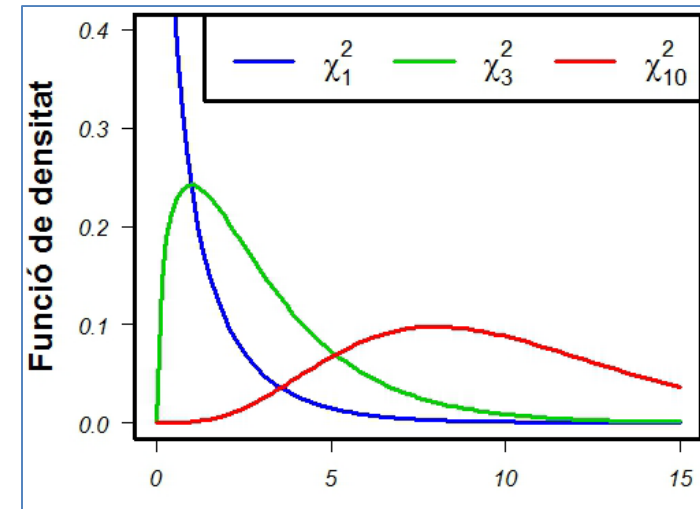
[ Concretament, per  $n = 1 \rightarrow X_1^2 \sim \chi_1^2$  ]

- Notació:**  $X \sim \chi_n^2$
- Paràmetres:**  $n$  (graus de llibertat)
- Funció de probabilitat i distribució:**

$$f(x) = \frac{x^{k/2-1} \cdot e^{-x/2}}{2^{k/2} \cdot \Gamma(k/2)} \quad \text{per } x > 0$$

$$F(x) = \frac{\gamma(k/2, x/2)}{\Gamma(k/2)} \quad \text{per } x > 0$$

$\Gamma$ : funció Gamma  
 $\gamma$ : funció Gamma incompleta  
 $n$ : graus de llibertat



**R:** dchisq, pchisq, qchisq

*Script per veure que la suma de Normals estàndard al quadrat és una  $\chi^2$*

```
M = 500
n = 7
sample = array(rnorm(M*n),dim=c(M,n))
sample2 = sample*sample
sum = apply(sample2, 1, sum)
hist(sum, breaks="Scott", freq=FALSE)
curve(dchisq(x, n),add=TRUE,col=2,lwd=2)
quantile(sum,c(0.25,0.50,0.75))
qchisq(c(0.25, 0.50, 0.75),n)
```

# Mostres de normals  
 # Graus de llibertat  
 # n mostres de N(0,1)  
 # n mostres de (N(0,1))^2  
 # Suma de les mostres al^2  
 # Distribució empírica sumant Normals  
 # Distribució teòrica de la chi-quadrat  
 # Q1, Mediana i Q3 de la suma de Normals  
 # Q1, Mediana i Q3 de la chi-quadrat

# Distribució t-Student

- Definició:** Siguin dues VA independents,  $Z \sim N(0,1)$  i  $Y_n \sim \chi_n^2$ . Llavors:

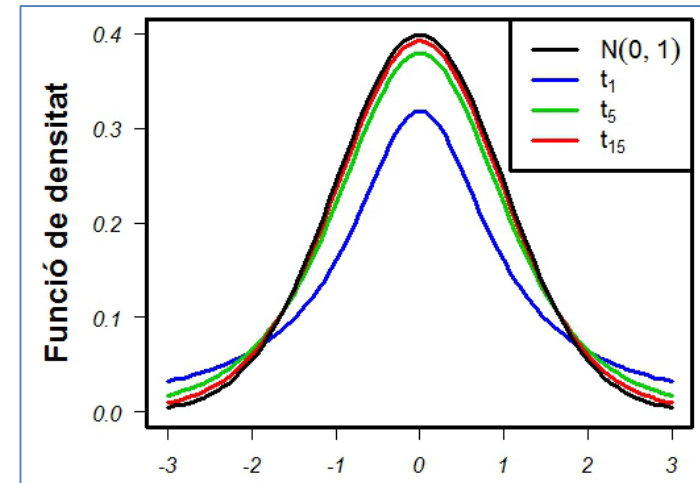
$$\frac{Z}{\sqrt{Y_n/n}} \sim t_n$$

[Quan  $n \rightarrow \infty$  ( $n > 30$ ), llavors  $t_n \rightarrow N(0,1)$ ]

- Notació:**  $X \sim t_n$
- Paràmetres:**  $n$  (graus de llibertat)
- Funció de probabilitat i distribució:**

$$f(x) = \frac{\Gamma\left(\frac{n+1}{2}\right)}{\sqrt{n\pi} \cdot \Gamma(n/2)} \left(1 + \frac{x^2}{n}\right)^{-\frac{n+1}{2}} \text{ per } x > 0$$

$$F(x) = \frac{\Gamma\left(\frac{n+1}{2}\right)}{\sqrt{n} \cdot B(1/2, n/2)} \left(1 + \frac{x^2}{n}\right)^{-\frac{n+1}{2}} \text{ per } x > 0$$



$\Gamma$  : funció Gamma  
 $B$ : funció Beta  
 $n$ : graus de llibertat

R: dt, pt, qt

*Script per  
 veure que  
 a partir de  
 una Z i una  
 $Y_n$  s'obté  
 una t*

```
M = 500; n = 7
samplez = rnorm(M, 0, 1)
samplechi2 = rchisq(M,n)
samplechi2n = sqrt(samplechi2/n)
t = samplez / samplechi2n
hist(t, breaks="Scott", freq=FALSE)
curve(dt(x, n),add=TRUE,col=2,lwd=2)
quantile(t, c(0.25, 0.50, 0.75))
qt(c(0.25, 0.50, 0.75),n)
```

# Número de mostres i graus de llibertat  
 # Mostra de normals  
 # Mostra de chi-quadrats  
 # Càlcul dels denominadors  
 # Càlcul de la t-student  
 # Distribució empírica  
 # Distribució teòrica d'una t-Student  
 # Q1, Mediana i Q3 de  $Z/\sqrt{Y_n/n}$   
 # Q1, Mediana i Q3 de la chi-quadrat

# Estimació puntual

- Un estimador  $\hat{\theta}$  del paràmetre desconegut  $\theta$ , a partir de la mostra  $M(\omega_i) (X_1, X_2, \dots, X_n)$  és una funció de les VA:

$$\hat{\theta} = f(X_1, X_2, \dots, X_n)$$

- Estimació puntual:** valor que l'estimador  $\hat{\theta}$  pren en una mostra concreta.

Ex:  $\bar{x} = \frac{\sum x_i}{n}$  és la mitjana mostral i és una estimació puntual de  $\mu$

**Nota:** Distingiu entre el valor  $\bar{x}$  i la variable aleatòria mitjana mostral  $\bar{X}$  que no s'extreu de cap mostra

- Error tipus o error estàndard:** variabilitat de l'estimador. Ex: en el cas anterior de la MITJANA, l'**error tipus (o estàndard) de la mitjana** (o *mean standard error* o *se*) és:

$$se = \sqrt{V(\bar{X}_n)} = \sqrt{E[(\bar{X}_n - \mu)^2]} = \frac{\sigma}{\sqrt{n}}$$

**Nota:** Generalment, la  $\sigma$  serà desconeguda i l'error tipus l'haurem d'aproximar emprant l'estimador

pertinent ( $\hat{\sigma}$ ) amb les dades de la mostra:  $\widehat{se} = \frac{\hat{\sigma}}{\sqrt{n}} = \frac{s}{\sqrt{n}} = \sqrt{\sum \frac{(x_i - \bar{x})^2}{n-1}} \cdot \frac{1}{\sqrt{n}}$

# Estimació puntual. Exemples

**Nota:** per als paràmetres s'utilitzen lletres de l'alfabet grec

Paràmetre ( $\theta$ ) ( <b>POBLACIÓ</b> )	Estimador ( $\hat{\theta}$ ) ( <b>MOSTRA</b> )
$\mu$ (esperança, mitjana poblacional)	$\bar{x}$ (mitjana mostral)
$\sigma^2$ (variància poblacional) $\sigma$ (desviació tipus poblacional)	$s^2$ (variància mostral) $s$ (desviació tipus mostral)
$\pi$ (probabilitat)	$p$ (proporció)

Exemple:

En 9 dies consecutius s'ha observat el nombre de terminals en una Universitat connectats a internet: 587, 470, 676, 451, 436, 672, 584, 697 i 408  $[\sum x = 4981 ; \sum x^2 = 2860855]$

```
nterm <- c(587, 470, 676, 451, 436, 672, 584, 697, 408)
```

Una estimació puntual del nombre esperat ( $\mu$ ) de terminals diaris connectats és:

`mean(nterm)`  $\rightarrow \bar{x} = 553.44$     o     $\bar{x} = (\sum x_i)/n = 553.44$

Una estimació puntual de la desviació tipus ( $\sigma$ ) del nombre de terminals connectats és:

`sd(nterm)`  $\rightarrow s = 114.0988$     o     $s = \sqrt{(\sum (x_i - \bar{x})^2)/(n - 1)} = 114.0988$

L'estimació de l'error tipus o variabilitat de la mitjana és:

`sd(nterm)/sqrt(length(nterm))`  $\rightarrow se = 38.03 = \sqrt{(\sum (x_i - \bar{x})^2)/(n - 1)} \cdot 1/\sqrt{n} = 38.03$

# Estimació puntual. Propietats dels estimadors

- Inevitablement, les estimacions puntuals **fallen** o, millor dit, com depenen de la mostra que “ens ha tocat”, **fluctuen** (encara que usualment tan sols observem un valor)
- Les 2 obsessions de l'Estadística són:
  - **quantificar** els errors d'estimació
  - **minimitzar** aquests errors
- L'error tipus o típic informa de l'**error esperat** a l'equiparar el valor de l'estimador obtingut en l'estudi amb el valor del paràmetre poblacional.
- Com l'estimador és “qualsevol” estadístic que s'utilitzi amb fins inferencials, hem de definir les propietats que permeten definir els “millors”.

**Nota:** *l'error exacte en una mostra concreta roman desconegut, podent ser inferior o superior que l'error típic o esperat.*

# Estimació puntual. Propietats dels estimadors

## Propietats desitjables

- **No tenir biaix** (= *sesgo, bias*)  $\rightarrow \text{Biaix} = 0$

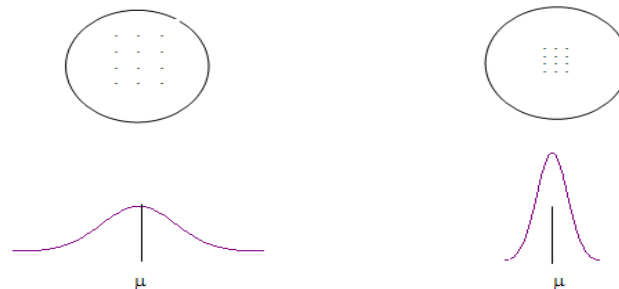
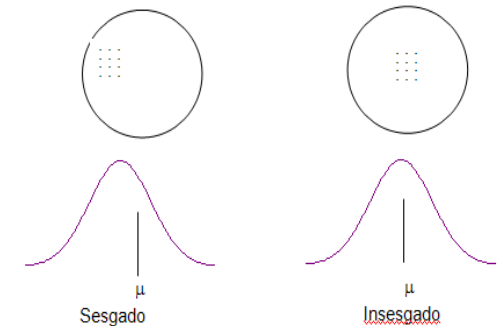
- Biaix és la diferència entre el centre de la distribució del estimador  $[E(\hat{\theta})]$  i el valor del paràmetre a estimar  $[\theta]$

$$\text{Biaix} = E(\hat{\theta}) - \theta$$

- Un estimador  $\hat{\theta}$  del paràmetre  $\theta$  és NO esbiaixat si  $\text{Biaix} = 0$

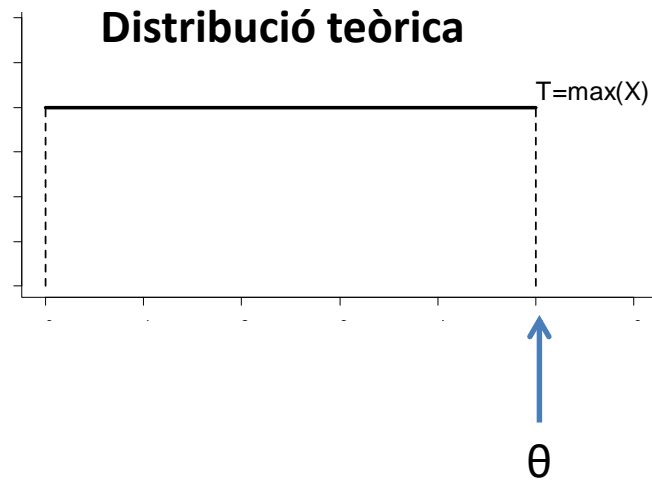
- **Ser Eficient**  $\rightarrow V(\hat{\theta}) \downarrow$

- Entre dos estimadors NO esbiaixats, es diu que és més eficient el que té una variància menor.



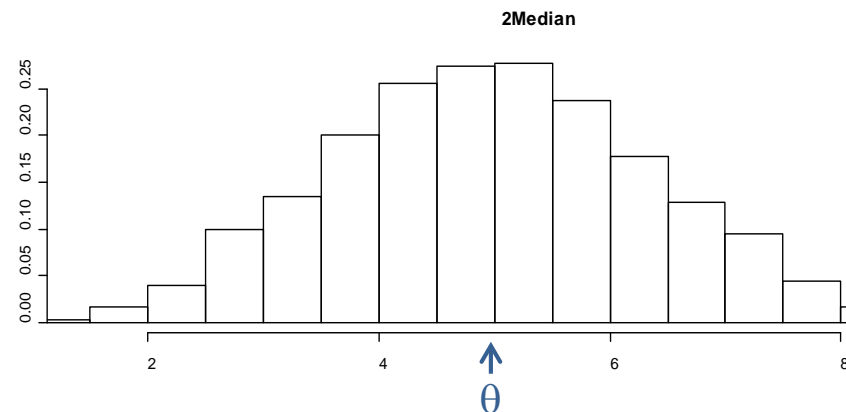
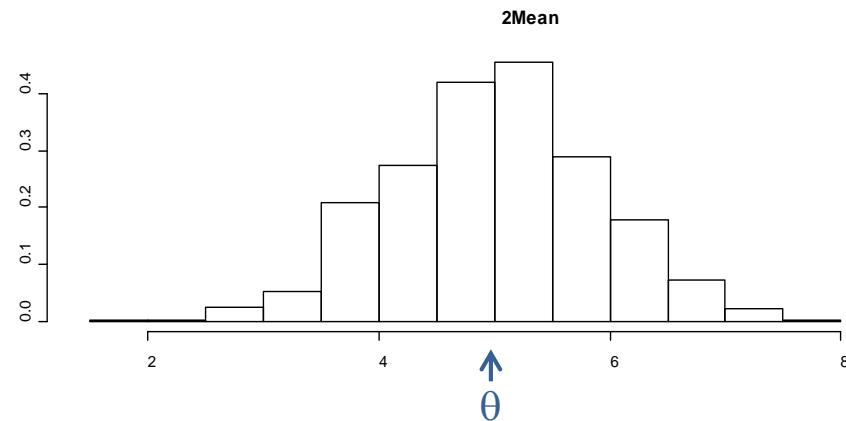
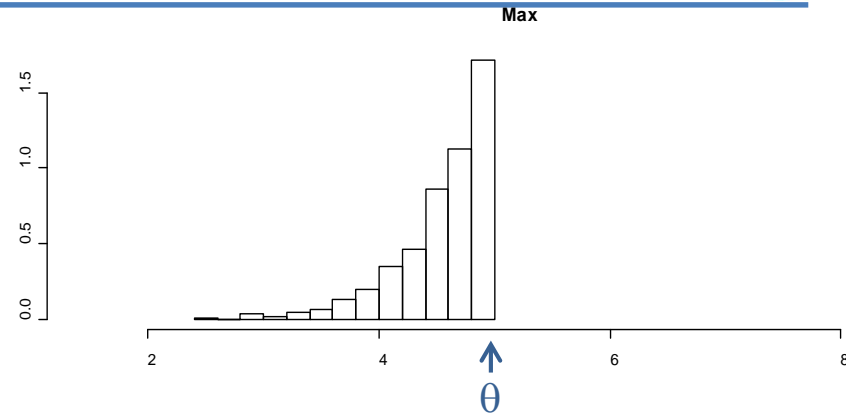
# Estimació puntual. Exemple

Estimar el límit superior  $\theta$  d'una llei uniforme, amb mínim zero:  $U[0, \theta]$



Tres estimadors possibles:

- Màxim de la mostra (esbiaixat)
- Doble de la mitjana mostral (no esbiaixat)
- Doble de la mediana mostral (no esbiaixat)



# Teorema del Límit Central (repàs)

- El T.L.C. estableix que, si s'agafen mostres de grandària  $n$  d'una població de mitjana  $\mu$  i desviació típica  $\sigma$ , a mesura que creix  $n$ , la distribució de la mitjana mostral  $\bar{X}$  s'aproxima a la d'una normal de mitjana  $\mu$  i desviació típica  $\sigma/\sqrt{n}$ :

$$\bar{X}_n \xrightarrow{n \text{ gran}} N\left(\mu, \frac{\sigma}{\sqrt{n}}\right)$$

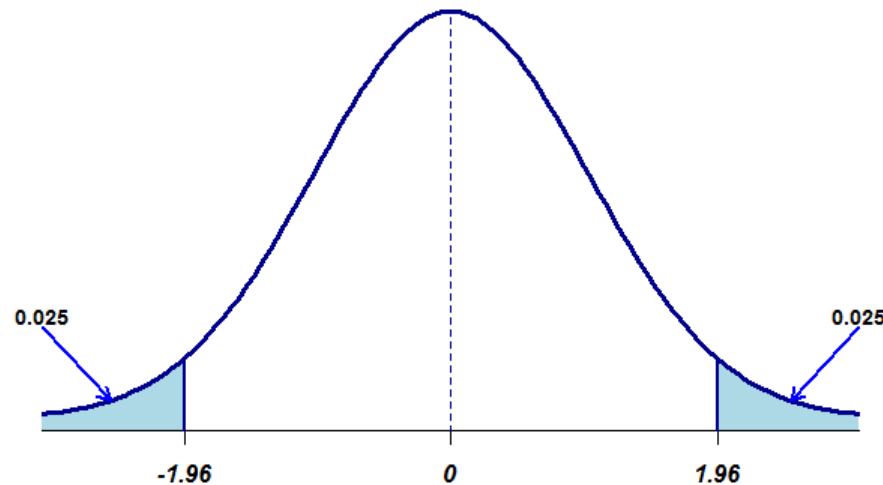
- ¿Quina grandària ha de tenir  $n$  per a que es compleixi el TLC?*
  - Si  $X \sim N \rightarrow \bar{X}_n \sim N \forall n$  (Les combinacions lineals de Normals i.i.d. són Normals)
  - Si  $X$  és quantitativa  $\rightarrow \bar{X}_n \sim N$  si  $n > 30$  (Com més s'assembli  $X$  a la Normal abans convergirà a una Normal)
  - Si  $X$  té una distribució **discreta** i/o **asimètrica**, la convergència requereix una grandària mostral ( $n$ ) més gran



# Quantils en la $Z \sim N(0,1)$ (repàs)

**Definició:** El quantil  $\alpha$  és el valor fins al qual s'acumula una probabilitat  $\alpha$

**Notació:** Aquell valor  $a$  tal que  $F_Z(a) = P(Z \leq a) = \alpha$  l'indicarem per  $z_\alpha$



**Exemples:**

$$z_{0.025} = -1.960$$

$$z_{0.975} = 1.960$$

$$z_{0.95} = 1.645$$

La Normal, al ser simètrica, compleix que:

$$z_\alpha = -z_{1-\alpha}$$

$$z_{\alpha/2} = -z_{1-\alpha/2}$$

# Interval per a la mitjana mostral. Exemple

- Situació:

$X$ : “temps d’execució d’un algoritme”  $\sim N(\mu = 100\text{ms}, \sigma = 10\text{ms})$ .

- Plantegem les distribucions de les VA:  $X$ ,  $\bar{X}_9$  i  $\bar{X}_{100}$

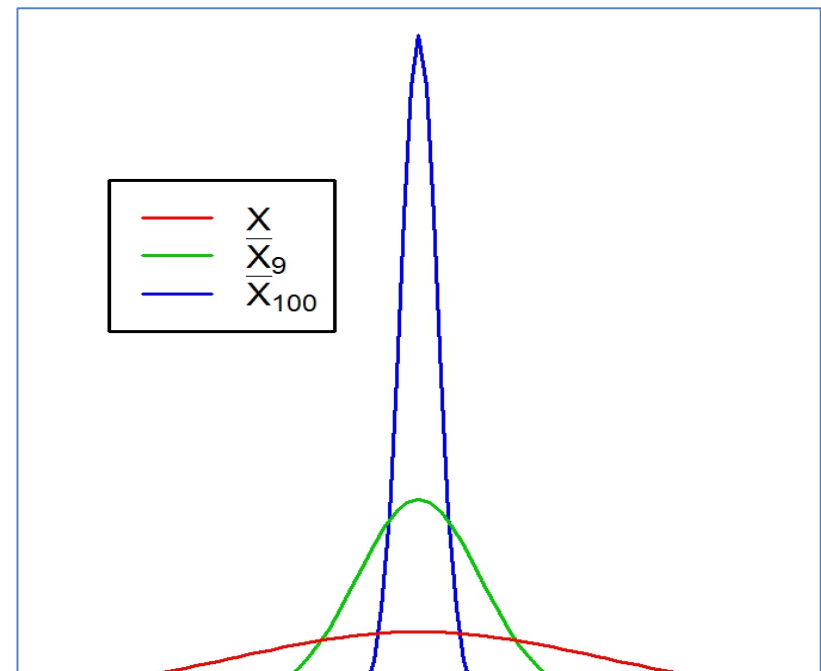
- Calcularem, per  $n=1$ ,  $n=9$  i per  $n=100$ , els intervals amb probabilitats *grans* (95%) d’assegurar que les possibles  $\bar{x}$  hi pertanyeran [deixem fora només una petita proporció ( $\alpha$ : 0.05 o 5%)]

- Distribucions:

$$X = \bar{X}_1 \sim N(\mu = 100 \text{ ms}, \sigma = 10 \text{ ms})$$

$$\bar{X}_9 \sim N\left(\mu = 100 \text{ ms}, \sigma = \frac{10}{\sqrt{9}} = 3.33 \text{ ms}\right)$$

$$\bar{X}_{100} \sim N\left(\mu = 100 \text{ ms}, \sigma = \frac{10}{\sqrt{100}} = 1 \text{ ms}\right)$$



# Interval per a la mitjana mostral. Exemple

Els límits  $v, w$  dels intervals els podem calcular utilitzant les taules de la  $N(0,1)$ :

$$z_{0.975} = 1.96 \quad z_{0.025} = -1.96$$

- Rang que conté el 95% de les infinites execucions de l'algoritme ( $X$ )

$$v, w = \mu \pm z_{0.975} \cdot \sigma = 100 \pm 1.96 \cdot 10 = 100 \pm 19.60 = [80.40, 119.60]$$

- Rang que conté el 95% de les mitjanes de les infinites mostres de  $n = 9$  execucions ( $\bar{X}_9$ )

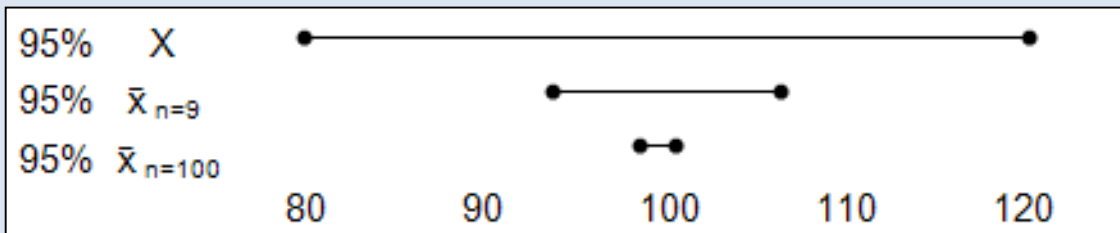
$$v, w = \mu \pm z_{0.975} \sigma / \sqrt{n} = 100 \pm 1.96 \cdot 10 / 3 = 100 \pm 6.53 = [93.47, 106.53]$$

- Rang que conté el 95% de les mitjanes de les infinites mostres de  $n = 100$  execucions ( $\bar{X}_{100}$ )

$$v, w = \mu \pm z_{0.975} \sigma / \sqrt{n} = 100 \pm 1.96 \cdot 10 / 10 = 100 \pm 1.96 = [98.04, 101.96]$$

Representació esquemàtica:

Com l'amplada de l'interval depèn inversament de l'arrel de  $n$ , passar de  $n=1$  a  $n=9$ , fa l'interval 3 vegades més estret



# Estimació per interval de $\mu$

- Hem vist que sabem calcular un “interval” que contingui  $\bar{x}$  a partir de  $\mu$ . Però el problema real és **calcular interval per  $\mu$ , coneixent  $\bar{x}$**  (és a dir, passar d’un interval per a la mitjana mostral  $\bar{x}$  a un per a la mitjana poblacional  $\mu$ )

- A partir d’una probabilitat  $1 - \alpha$  entre dos valors  $a$  i  $b$  (simètrics): *(amb  $\sigma$  coneguda)*

$$P(a \leq \bar{X}_n \leq b) = 1 - \alpha \rightarrow P\left(\frac{a - \mu}{\sigma/\sqrt{n}} \leq \frac{\bar{X}_n - \mu}{\sigma/\sqrt{n}} \leq \frac{b - \mu}{\sigma/\sqrt{n}}\right) = 1 - \alpha \rightarrow P\left(z_{\frac{\alpha}{2}} \leq \frac{\bar{X}_n - \mu}{\sigma/\sqrt{n}} \leq z_{1 - \frac{\alpha}{2}}\right) = 1 - \alpha$$

- Obtenim l’interval de la v. a.  $\bar{X}_n$  amb **probabilitat  $1 - \alpha$**

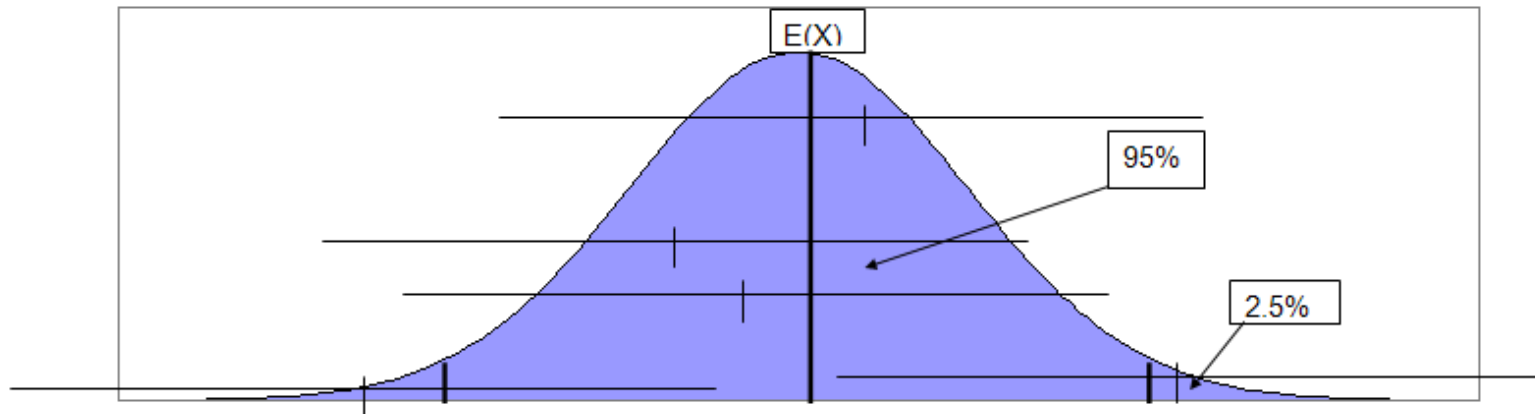
$$P\left(\mu + z_{\frac{\alpha}{2}} \cdot \frac{\sigma}{\sqrt{n}} \leq \bar{X}_n \leq \mu + z_{1 - \frac{\alpha}{2}} \cdot \frac{\sigma}{\sqrt{n}}\right) = 1 - \alpha$$

- I reordenant obtenim **l’interval de confiança  $1 - \alpha$  del paràmetre  $\mu$**

$$P\left(\bar{X}_n + z_{\frac{\alpha}{2}} \cdot \frac{\sigma}{\sqrt{n}} \leq \mu \leq \bar{X}_n + z_{1 - \frac{\alpha}{2}} \cdot \frac{\sigma}{\sqrt{n}}\right) = 1 - \alpha$$

# Estimació per interval de $\mu$

- $P\left(\bar{X}_n + z_{\frac{\alpha}{2}} \cdot \frac{\sigma}{\sqrt{n}} \leq \mu \leq \bar{X}_n + z_{1-\frac{\alpha}{2}} \cdot \frac{\sigma}{\sqrt{n}}\right) = 1 - \alpha$  significa que amb una confiança de  $1 - \alpha$  podem assegurar que  $E(X) = \mu$  estarà en el rang calculat.
- Si  $1 - \alpha$  és 95% ( $\alpha = 5\%$ ): **el 95% dels intervals (IC) contindran  $\mu$**



- *Aquest procediment encerta el  $100 \cdot (1 - \alpha)\%$  de les vegades!*
- Denotem **IC( $\mu$ ,  $1 - \alpha$ )** a l'**INTERVAL DE CONFIANÇA**  $1 - \alpha$  de  $\mu$ , i l'expresssem:

$$IC(\mu, 1 - \alpha) = \bar{x} \mp z_{1-\frac{\alpha}{2}} \cdot \frac{\sigma}{\sqrt{n}}$$

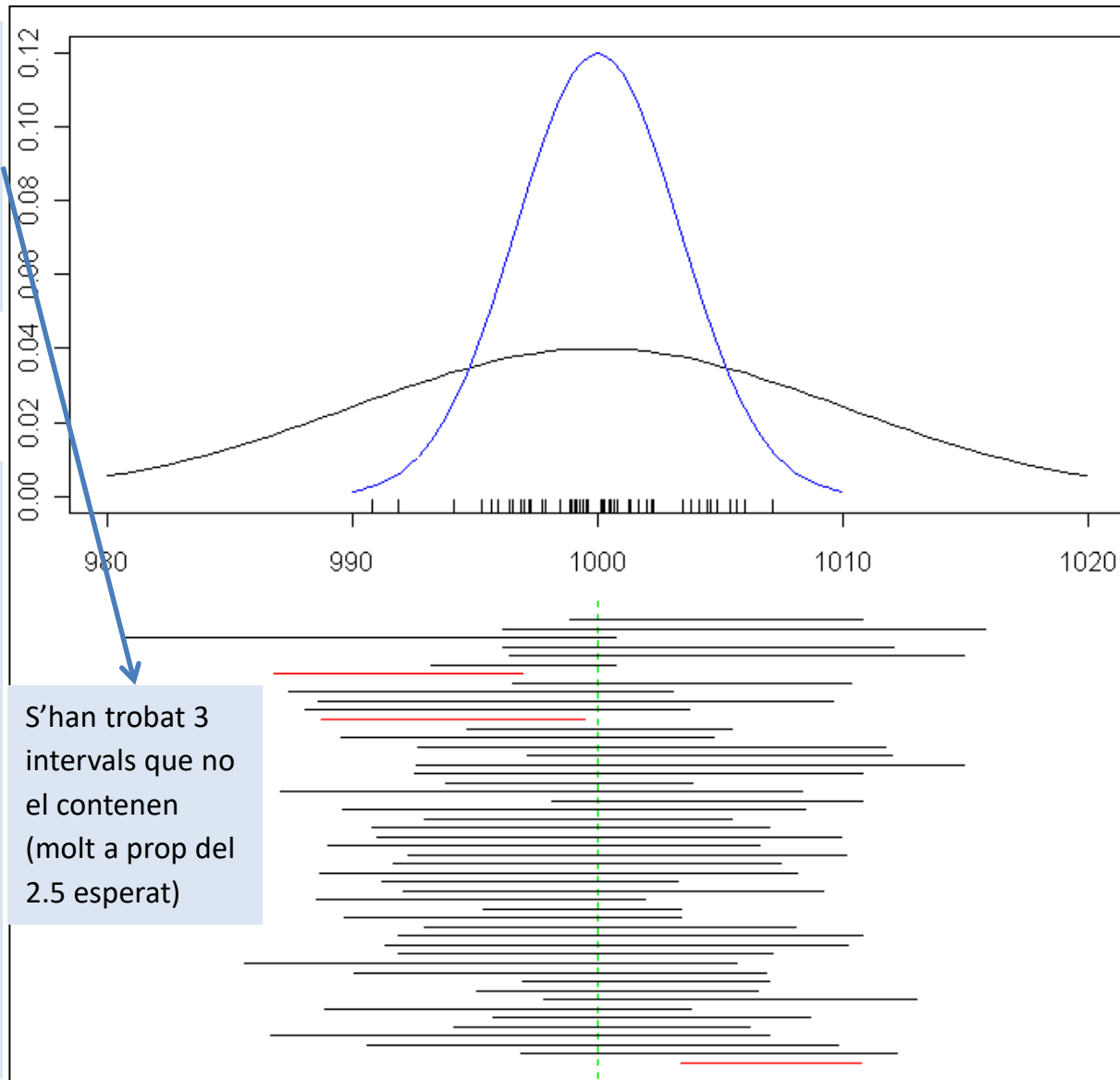
Com que la Normal és simètrica, llavors  $z_{\alpha/2} = -z_{1-\alpha/2}$

**Atenció:** nosaltres només observarem una mostra, i no sabrem si l'IC trobat conté o no  $\mu$ , però sí sabem que aquest procediment a la llarga dóna un  $100 \cdot (1 - \alpha)\%$  d'encerts.

# Estimació per interval de $\mu$ . Simulació

S'han simulat 50 mostres amb  $n=9$  provinents d'una  $N(\mu=1000, \sigma)$ . Calculant el IC95%, esperem que aproximadament el 5% (2.5) d'ells no continguin el valor real de  $\mu$

**Nota tècnica:** Amb un IC determinat (p.ex., [985, 1004]), s'ha de dir **“tenim un alt grau de confiança (p.ex., 95%) de que el paràmetre es trobi entre aquest dos valors”**, però no és correcte parlar d'una probabilitat de 0.95 que el paràmetre estigui entre els dos valors trobats, perquè el paràmetre no es considera un element aleatori. Serà desconegut, però no és incert!



# Interval de confiança per $\mu$ ( $\sigma$ coneguda). Exemple

- Així doncs, l'interval de confiança  $1-\alpha$  de  $\mu$  (amb  $\sigma$  coneguda) es calcula com:

$$IC(\mu, 1 - \alpha) = \bar{x} \mp z_{1-\frac{\alpha}{2}} \cdot \frac{\sigma}{\sqrt{n}}$$

- Recordeu que ens basem en el TCL i perquè es complís calia que la variable X inicial fos Normal o que n fos més gran de 30. Per tant, els requisits per realitzar aquest càlcul són:  **$n > 30$  o  $X \sim N$**

## EXEMPLE:

Una embotelladora d'ampolles de litre té una dispersió de  $\sigma = 10\text{cc}$ . En una mostra a l'atzar de  $n = 100$  ampolles d'aquesta màquina, la mitjana observada ha sigut  $\bar{x} = 995\text{cc}$ . Calculeu un interval de confiança del 95% de  $\mu$ .

$$IC(\mu, 0.95) = \bar{x} \mp z_{1-\frac{\alpha}{2}} \cdot \frac{\sigma}{\sqrt{n}} = 995 \mp 1.96 \cdot \frac{10}{\sqrt{100}} = 995 \mp 1.96 = [993.04, 996.96]$$

**Resultat:** amb una confiança del 95%,  $\mu$  es troba entre 993.04 i 996.96

# Interval de confiança per $\mu$ ( $\sigma$ coneguda). Exercici

1. La glicèmia en mmol/L té una desviació típica de  $\sigma = 1$  en una mostra de  $n = 9$  pacients, la mitjana  $\bar{x}$  val 5. Calculeu el IC( $\mu, 0.95$ ).

Amb una "força" del 95%, creiem que l'autèntic valor poblacional està entre aquests límits

2. Sense canviar la confiança, com podríem reduir l'interval a la meitat?
3. Calculeu l' IC amb una confiança del 99%

**ATENCIÓ:** quan  $n$  augmenta la precisió dels IC augmenta (interval més estret). Si augmenta la confiança (disminuint el risc  $\alpha$  d'error), la precisió dels IC disminueix (interval més ample)

**ATENCIÓ:** En aquest cas, per estimar  $\mu$  necessitem conèixer  $\sigma \rightarrow$  *situació molt particular i infreqüent*



# Interval de confiança. Mecànica

Passos	Esquema de Resolució
<b>1</b>	Definir l'estadístic a ser utilitzat
	Especificar la seva distribució
	Indicar les premisses necessàries per dir que segueix la distribució
	Delimitar el nivell de confiança (usualment $1-\alpha=95\%$ )
<b>2</b>	Calcular l'interval
<b>3</b>	Interpretar el resultat

# Interval de confiança per $\mu$ ( $\sigma$ desconeguda)

- Si **desconeixem**  $\sigma$ , la podem substituir per  $S$ , i llavors l'estadístic  $Z = \frac{\bar{x} - \mu}{\sigma/\sqrt{n}}$  passa a ser  $t = \frac{\bar{x} - \mu}{s/\sqrt{n}}$  que és el quocient de 2 VAs i ja no es pot suposar que segueix una distribució  $N(0,1)$ .
- Tal com diu el pas 2 de la mecànica de construcció de l'IC, cal conèixer la llei de probabilitats que segueix l'estadístic. En aquest cas, es tracta d'una nova distribució anomenada ***t-Student***.
- Per aquest estadístic, la distribució de probabilitat concreta és  $t_{n-1}$  ( $n-1$  graus de llibertat). Els percentils es poden trobar a taules específiques o amb R.
- Així doncs, l'interval de confiança  $1-\alpha$  de  $\mu$  (amb  $\sigma$  desconeguda) és:

$$IC(\mu, 1 - \alpha) = \bar{x} \mp t_{n-1, 1-\frac{\alpha}{2}} \cdot \frac{s}{\sqrt{n}}$$

**NOTA:** la situació de no conèixer la  $\sigma$  de la població és més freqüent

# Interval de confiança per $\mu$ ( $\sigma$ desconeguda)

- Podem utilitzar-la per conèixer el quocient informació/soroll utilitzant  $S$  en lloc de  $\sigma$

$$\hat{t} = \frac{\bar{X} - \mu}{\sqrt{\frac{s^2}{n}}} \sim t_{n-1}$$

- Demostració:

$$\begin{aligned} \hat{t} &= \frac{\bar{X} - \mu}{\sqrt{\frac{s^2}{n}}} = \frac{(\bar{X} - \mu)/\sqrt{\sigma^2/n}}{\sqrt{\frac{s^2/n}{\sigma^2/n}}} = \frac{Z}{\sqrt{\frac{s^2}{\sigma^2}}} = \frac{Z}{\sqrt{\frac{\sum_{i=1}^n (x_i - \bar{x})^2 / (n-1)}{\sigma^2}}} = \frac{Z}{\sqrt{\frac{\sum_{i=1}^n (x_i - \bar{x})^2 / \sigma^2}{n-1}}} \\ &= \frac{Z}{\sqrt{\frac{Y_n}{n-1}}} \sim t_{n-1} \end{aligned}$$

- Coneixem que el quocient informació/soroll segueix una 't' de Student.
- Això ens ajudarà a conèixer el IC95% per la  $\mu$  amb  $\sigma$  desconeguda

**Nota:** Observeu que es requereix la Normalitat de les  $X$ 's

# Interval de confiança per $\mu$ ( $\sigma$ desconeguda). Ex.

En 9 dies consecutius s'ha observat el nombre de terminals en una Universitat connectats a internet: 587, 470, 676, 451, 436, 672, 584, 697 i 408.

**R:** `nterm <- c(587, 470, 676, 451, 436, 672, 584, 697, 408)`

Recordem que havíem calculat les estimacions puntuals: Una estimació del IC al 95% de confiança de la mitjana poblacional, assumint que la desviació poblacional ( $\sigma$ ) val 100:

$1-\alpha$	$\sigma$	IC( $\mu, 1-\alpha$ )	Resolució amb R
95%	Coneguda ( $\sigma=100$ )	[488.11; 618.78]	<pre>n &lt;- 9 ; sigma &lt;- 100 mean(nterm) - qnorm(0.975)*sigma/sqrt(n) mean(nterm) + qnorm(0.975)*sigma/sqrt(n)</pre>
99%	Coneguda ( $\sigma=100$ )	[467.58 ; 639.31]	<pre>n &lt;- 9 ; sigma &lt;- 100 mean(nterm) - qnorm(0.995)*sigma/sqrt(n) mean(nterm) + qnorm(0.995)*sigma/sqrt(n)</pre>
95%	Desconeguda	[465.74; 641.15]	<pre>n &lt;- 9 mean(nterm) - qt(0.975, 8)*sd(nterm)/sqrt(n) mean(nterm) + qt(0.975, 8)*sd(nterm)/sqrt(n)</pre>
99%	Desconeguda	[425.83 ; 681.06]	<pre>n &lt;- 9 mean(nterm) - qt(0.995, 8)*sd(nterm)/sqrt(n) mean(nterm) + qt(0.995, 8)*sd(nterm)/sqrt(n)</pre>

# Interval de confiança per $\mu$ . Premisses

Per garantir el nivell de confiança de l'IC, s'ha de complir certes premisses

- Si sigma és coneguda, exigirem una de les condicions:
  - **$X \sim N$**   $\rightarrow$  la combinació lineal de Normals és Normal ( $\bar{X} \sim N$ )
  - **Tenir una mostra gran ( $n \geq 30$ )**  $\rightarrow$  Pel TCL,  $\bar{X} \sim N$
- Si sigma no és coneguda, exigirem una de les condicions:
  - **$X \sim N$**   $\rightarrow (\bar{x} - \mu) / \sqrt{s^2/n} \sim t_{n-1}$
  - **Tenir una mostra gran ( $n \geq 100$ )**  $\rightarrow$  Pel TCL,  $\bar{X} \sim N$

La premissa que sempre s'ha de complir és que **l'origen de la mostra ha de ser aleatori** (v.a.i.i.d)

Dist. de referència	$\sigma$ coneguda	$\sigma$ desconeguda
<b>X Normal</b>	Normal	t de <i>Student</i>
<b>X no Normal</b>	Normal <i>si <math>n \geq 30</math></i>	Normal <i>si <math>n \geq 100</math></i>

Amb grans mostres la variació de “s” serà limitada (s estima molt bé  $\sigma$ ), i podem considerar que

$$(\bar{x} - \mu) / \sqrt{s^2/n} \approx (\bar{x} - \mu) / \sqrt{\sigma^2/n} \sim N(0,1)$$

[Aplicable al cas  $X \sim N$ ,  $\sigma$  desconeguda i n gran]

# Interval de confiança per $\pi$ en una Binomial ( $n, \pi$ )

- Sigui  $X \sim B(n, \pi) \rightarrow$ 

$$E(X) = \pi \cdot n$$

$$V(X) = \pi \cdot (1 - \pi) \cdot n$$
- Sigui  $P = X/n \rightarrow$ 

$$E(P) = E(X/n) = E(X)/n = \pi \cdot n / n = \pi$$

$$V(P) = V(X/n) = V(X)/n^2 = \pi \cdot (1 - \pi) \cdot n / n^2 = \pi \cdot (1 - \pi) / n$$
- Per construir l'IC es pot recorre a la convergència de la Binomial a la Normal [amb la premissa de  $n$  gran i  $\pi$  no extrema  $\rightarrow \pi \cdot n \geq 5$  i  $(1 - \pi) \cdot n \geq 5$ ]:

$$P \rightarrow N\left(\mu_P = \pi, \sigma_P = \sqrt{\frac{\pi(1-\pi)}{n}}\right)$$

- Així, l'interval de confiança s'assemblaria al de  $\mu$ :

$$IC(\pi, 1 - \alpha) = P \mp z_{1-\frac{\alpha}{2}} \sigma_P = P \mp z_{1-\frac{\alpha}{2}} \sqrt{\frac{\pi(1-\pi)}{n}}$$

La **paradoxa** de que necessitem conèixer  $\pi$  per estimar el IC de  $\pi$  es pot solucionar de 2 maneres:

a) Substituint  $\pi$  per  $P \rightarrow IC(\pi, 1 - \alpha) = P \mp z_{1-\alpha/2} \cdot \sqrt{(P(1 - P))/n}$

b) Aplicant el màxim de  $\pi \cdot (1 - \pi)$  que correspon a  $\pi = 0.5 \rightarrow IC(\pi, 1 - \alpha) = P \mp z_{1-\alpha/2} \cdot \sqrt{(0.5(1 - 0.5))/n}$

## IC per $\pi$ en una Binomial( $n, \pi$ ). Exemple

Llencem 100 vegades una moneda a l'aire i observem 56 cares ( $P = 56/100 = 0.56$ ).

Les dues solucions per l'IC segons com estimem  $\pi$ :

$$IC(\pi, 0.95) = P \mp z_{1-\frac{\alpha}{2}} \sqrt{\frac{P \cdot (1-P)}{n}} = 0.56 \mp 1.96 \sqrt{\frac{0.56 \cdot 0.44}{100}} \approx 0.56 \mp 0.10 = [0.46, 0.66]$$

$$IC(\pi, 0.95) = P \mp z_{1-\frac{\alpha}{2}} \sqrt{\frac{\pi_m \cdot (1-\pi_m)}{n}} = 0.56 \mp 1.96 \sqrt{\frac{0.50 \cdot 0.50}{100}} \approx 0.56 \mp 0.10 = [0.46, 0.66]$$

Donen el mateix IC fins al 2n decimal. El motiu és que la probabilitat estimada (0.56) és molt similar a la probabilitat de màxima indeterminació (0.50)

# Interval de confiança per $\sigma^2$ . Distribució $\chi^2$

- Gràcies a aquesta distribució, coneixem la distribució de  $S^2$  (estimador de  $\sigma^2$ ):

$$s^2 = \hat{\sigma}^2 = \frac{\sum_{i=1}^n (x_i - \mu)^2}{n-1} \rightarrow (n-1) \cdot \frac{s^2}{\sigma^2} \sim \chi_{n-1}^2$$

- Demostració: si  $X_i \rightarrow N$ , llavors

Les  $x_i$  han de provenir d'una Normal

$$n \cdot \frac{\hat{\sigma}^2}{\sigma^2} = n \cdot \frac{(\sum_{i=1}^n (x_i - \mu)^2)/n}{\sigma^2} = \frac{\sum_{i=1}^n (x_i - \mu)^2}{\sigma^2} = \sum_{i=1}^n \left( \frac{x_i - \mu}{\sigma} \right)^2 \sim \chi_n^2$$

- També, si  $X_i \rightarrow N$

$$(n-1) \cdot \frac{\hat{\sigma}^2}{\sigma^2} = (n-1) \cdot \frac{(\sum_{i=1}^n (x_i - \bar{x})^2)/(n-1)}{\sigma^2} = \frac{\sum_{i=1}^n (x_i - \bar{x})^2}{\sigma^2} = \sum_{i=1}^n \left( \frac{x_i - \bar{x}}{\sigma} \right)^2 \sim \chi_{n-1}^2$$

- Per tant, **coneixem la distribució de  $s^2$ !!** i podrem relacionar la distribució  $\chi_{n-1}^2$  amb  $s^2$  per definir IC de  $\sigma^2$ , tal com fem amb les distribucions  $Z$  i  $t_{n-1}$  amb  $\bar{x}$  per definir IC de  $\mu$ .



# Interval de confiança per $\sigma^2$

- Hem vist que:  $(n - 1) \cdot \frac{S^2}{\sigma^2} \sim \chi_{n-1}^2$  sempre que  $\mathbf{X}_i \sim \mathbf{N}$  (premissa)
- Per tant:

$$P\left(\chi_{n-1, \frac{\alpha}{2}}^2 \leq \frac{S^2 \cdot (n-1)}{\sigma^2} \leq \chi_{n-1, 1-\frac{\alpha}{2}}^2\right) = 1 - \alpha$$

$$P\left(\frac{1}{\chi_{n-1, 1-\frac{\alpha}{2}}^2} \leq \frac{\sigma^2}{S^2 \cdot (n-1)} \leq \frac{1}{\chi_{n-1, \frac{\alpha}{2}}^2}\right) = 1 - \alpha$$

$$P\left(\frac{S^2 \cdot (n-1)}{\chi_{n-1, 1-\frac{\alpha}{2}}^2} \leq \sigma^2 \leq \frac{S^2 \cdot (n-1)}{\chi_{n-1, \frac{\alpha}{2}}^2}\right) = 1 - \alpha$$

**Atenció:** és un IC per  $\sigma^2$ , no per  $\sigma$ !!

$$IC(\sigma^2, 1 - \alpha) = \left[ \frac{s^2(n-1)}{\chi_{n-1, 1-\frac{\alpha}{2}}^2}, \frac{s^2(n-1)}{\chi_{n-1, \frac{\alpha}{2}}^2} \right]$$

**Nota:** no és un interval simètric, ja que la distribució no ho és. Això implica calcular dos valors en la distribució en lloc de fer  $\pm$ .

# IC per $\sigma^2$ en una Normal ( $\mu, \sigma$ ). Exemple

En les 25 execucions d'un mateix programa s'ha observat una variabilitat  $s^2 = 8^2$

$$\begin{aligned}
 IC(\sigma^2, 0.95) &= \left[ \frac{s^2(n-1)}{\chi_{n-1, 1-\frac{\alpha}{2}}^2}, \frac{s^2(n-1)}{\chi_{n-1, \frac{\alpha}{2}}^2} \right] = \\
 &= \left[ \frac{8^2(25-1)}{12.401}, \frac{8^2(25-1)}{39.364} \right] = [123.86, 39.02] \rightarrow \text{Oops! M'he equivocat} \\
 &= \left[ \frac{8^2(25-1)}{39.364}, \frac{8^2(25-1)}{12.401} \right] = [39.02, 123.86] \rightarrow \text{Ara sí!}
 \end{aligned}$$

Resultat:

$$IC(\sigma^2, 0.95) = [39.02, 123.86]$$

$$IC(\sigma, 0.95) = [6.25, 11.13]$$

← Fent l'arrel quadrada,  
obtenim un interval per  $\sigma$

# Formulari per IC

TCL:  $X_1, \dots, X_n$  i.i.d. ( $n \rightarrow \infty$ ), amb  $E(X_i) = \mu$  i  $V(X_i) = \sigma^2$ ,

llavors  $\frac{\sum_{i=1}^n X_i}{n} = \bar{X}_n \approx N(\mu, \sigma^2/n)$  i també  $\sum_{i=1}^n X_i \approx N(n\mu, \sigma^2 n)$

Estadístic mitjana mostral ( $\bar{x}$ ):  $\frac{(\bar{x} - \mu)}{\sqrt{\sigma^2/n}} \approx N(0,1)$   $\frac{(\bar{x} - \mu)}{\sqrt{s^2/n}} \approx t_{n-1}$  on  $\bar{x} = \sum_{i=1}^n x_i / n$

Estadístic variància mostral ( $s^2$ ):  $s^2 \frac{n-1}{\sigma^2} \approx \chi_{n-1}^2$  on  $s^2 = \frac{\sum_{i=1}^n (x_i - \bar{x})^2}{n-1} = \frac{\sum_{i=1}^n x_i^2 - n(\bar{x})^2}{n-1} = \frac{\sum_{i=1}^n x_i^2 - \frac{(\sum_{i=1}^n x_i)^2}{n}}{n-1}$

Paràmetre	Estadístic	Premisses	Distribució	Interval de Confiança $1-\alpha$ (Risc $\alpha$ )
$\mu$	$\hat{z} = \frac{(\bar{x} - \mu)}{\sqrt{\sigma^2/n}}$	[ $X \rightarrow N$ ò $n \geq 30$ ] i $\sigma$ coneguda	$\hat{Z} \rightarrow N(0,1)$	$\mu \in (\bar{x} \pm z_{1-\alpha/2} \sqrt{\frac{\sigma^2}{n}})$
$\mu$	$\hat{t} = \frac{(\bar{x} - \mu)}{\sqrt{s^2/n}}$	$X \rightarrow N$	$\hat{t} \rightarrow t_{n-1}$	$\mu \in (\bar{x} \pm t_{n-1, 1-\alpha/2} \sqrt{\frac{s^2}{n}})$
$\mu$	$\hat{z} = \frac{(\bar{x} - \mu)}{\sqrt{s^2/n}}$	$n \geq 100$	$\hat{Z} \rightarrow N(0,1)$	$\mu \in (\bar{x} \pm z_{1-\alpha/2} \sqrt{\frac{s^2}{n}})$
$\sigma$ (normal)	$\hat{X}^2 = \frac{s^2(n-1)}{\sigma^2}$	$X \rightarrow N$	$\hat{X}^2 \rightarrow \chi_{n-1}^2$	$\sigma^2 \in \left( \frac{S^2(n-1)}{\chi_{n-1, 1-\alpha/2}^2}, \frac{S^2(n-1)}{\chi_{n-1, \alpha/2}^2} \right)$
$\pi$ (Binomial)	$\hat{z} = \frac{(p - \pi)}{\sqrt{\pi(1-\pi)/n}}$	$(1-\pi)n \geq 5$ $\pi n \geq 5$	$\hat{Z} \rightarrow N(0,1)$	$\pi \in (P \pm z_{1-\alpha/2} \sqrt{\frac{\hat{\pi}(1-\hat{\pi})}{n}})$ $\hat{\pi} = P$ o $\hat{\pi} = 0.5$
$\lambda$ (Poisson)	$\hat{z} = \frac{(L - \lambda)}{\sqrt{\lambda}}$	$\lambda \geq 5$	$\hat{Z} \rightarrow N(0,1)$	$\lambda \in (L \pm z_{1-\alpha/2} \sqrt{L})$

El color gris indica un IC secundari

# Proves d'hipòtesi. Exemple

Afirmo que encerto el 80% dels meus tirs lliures a basket, i un amic em posa a prova. Dels 20 tirs, solament faig 8. *“Fuà! Algú que encerta un 80% gairebé mai faria 8 de 20; així que no em crec la teva afirmació”*.

El raonament de l'amic es basa en demanar-se què passaria si l'afirmació fos certa i es repetís la mostra de 20 tirs moltes vegades. Segurament poques vegades encertaria un nombre tan baix com 8. Un resultat com aquest és tan poc freqüent que aporta certa *evidència* per rebutjar la meva afirmació.

De fet, **aquesta evidència es pot quantificar**:

$$\text{Sigui } M = n^{\circ} \text{ encerts: } M \sim B(20, 0.8) \rightarrow P(M \leq 8) = 0.000102$$

Es a dir, si jo fos tan bo, solament 1 de cada 10000 vegades obtindria una puntuació tan nefasta (o pitjor).

Dues opcions: o he tingut mala sort o l'afirmació era falsa!

Source: The Basic Practice of Statistics. David S. Moore. 4th Ed.

# Proves d'hipòtesi. Raonament

- Al igual que amb els IC, hem de pensar com si l'experiència es pogués repetir un gran nombre de vegades.
- Però ara no volem l'IC que contengui el paràmetre, ara partim d'una afirmació (una **hipòtesi** de partida, o **nul·la**), i volem estudiar si les dades proporcionen proves en contra seu.
- Una repetició intensa (una mostra infinita = la població) seria definitiva.
- Però, amb una mostra finita, quina informació aporten les dades?

## Nota tècnica:

- Formalment, es distingeix entre les proves de Fisher (per aportar coneixement o evidència o inferència) i els contrastos de Neyman-Pearson (per minimitzar els errors al prendre decisions). Els primers son rellevants per la Ciència (p.e., la Física) i els segons per la Tècnica (p.e. la Arquitectura). Però en aquest curs no distingirem i ho englobarem tot sota Proves de Hipòtesi (PH).
- Read more in: <http://onlinestatbook.com/chapter9/significance.html>

# Proves d'hipòtesi. Raonament

La hipòtesi nul·la ( $H_0$ ) es planteja formalment amb un paràmetre (o varis). El paràmetre en qüestió pren un valor que representem:

$$H_0: \pi = 0.80$$

$\pi$  representa la probabilitat poblacional d'encertar un tir lliure, i volem comprovar si aquest valor és coherent amb les observacions.

Al igual que amb els IC, la mostra es *concentra* en un estadístic, que segueix una distribució de probabilitat coneguda si s'assumeix certa la  $H_0$ .

Addicionalment a  $H_0$  afegim la hipòtesi **alternativa**  $H_1$ , que pot ser totalment complementària a la nul·la (enfoc bilateral), o parcialment (unilateral):

$$H_1: \pi \neq 0.80$$

$$H_1: \pi < 0.80$$

$H_1$  determina el(s) sentit(s) més oposat(s) a  $H_0$ : per exemple, el nombre de encerts pot ser l'estadístic, i si  $H_1$  fos " $\neq$ " serien *sospitosos* tant els nombres d'encerts que van cap a 0 com els que van cap a 20. Si  $H_1$  fos " $<$ " serien *sospitosos* només els que van cap a 0 (que és el que hem pres, donat que el meu amic no confia molt en les meves habilitats).

# Proves d'hipòtesi. P-valor

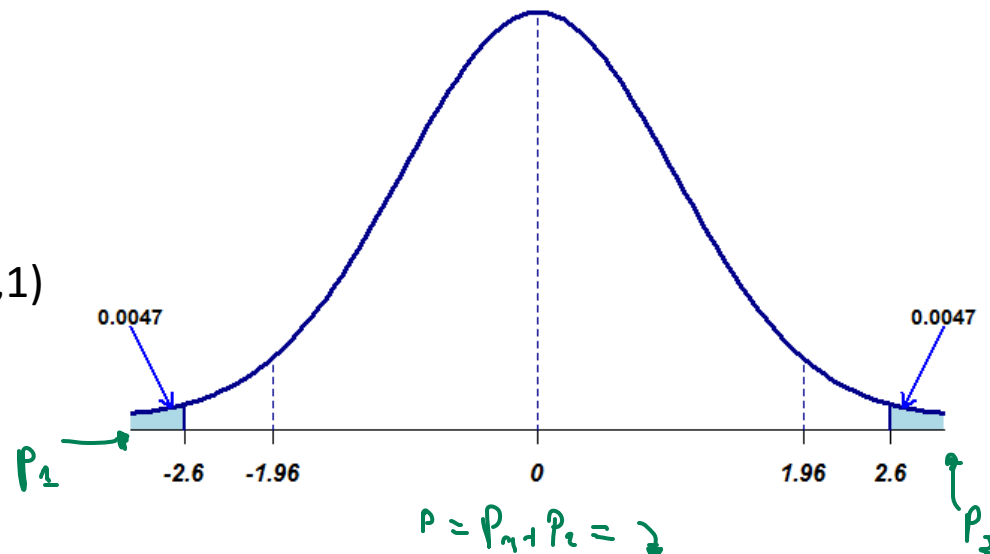
**P**, o **p-valor**, és la probabilitat de, sota  $H_0$ , obtenir resultats igual o més *extremes* que el observat.

**Exemple:**

**Estadístic:**  $Z = 2.6$

**Distribució de l'estadístic (sota  $H_0$ ):**  $Z \sim N(0,1)$

**Contrast:** bilateral



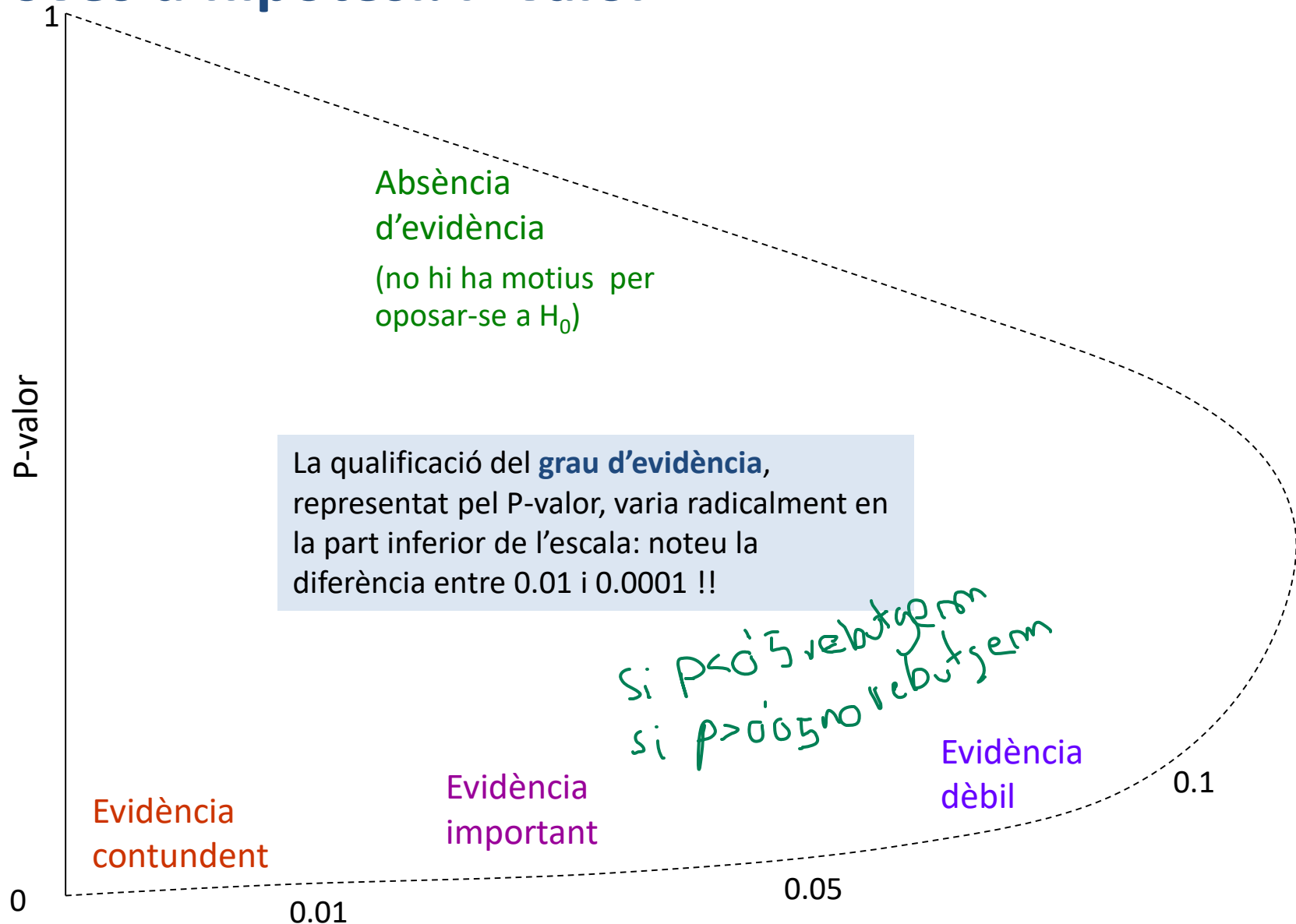
$$\left. \begin{aligned} P(Z < -2.6) &= 0.0047 \\ P(Z > +2.6) &= 0.0047 \end{aligned} \right\} \rightarrow p\text{-valor} = P(|Z| > 2.6) = 2 \cdot 0.0047 = 0.0094$$

[**Taules:**  $2 \cdot (1 - 0.9953)$  ; **R:** `pnorm(-2.6) + (1 - pnorm(2.6))` )

$P = p\text{-valor} = P(\text{VA } Z \text{ "més lluny" de } H_0 \text{ que el valor observat } z)$

RECORDEU: "sota  $H_0$ " = suposem ('temptativament') que és cert que els paràmetres poblacionals valen el que diu  $H_0$

# Proves d'hipòtesi. P-valor





# Malentessos amb el p-valor

- El  $P$ -valor diu amb quina freqüència poden passar events com el de la mostra (o més extrems) quan la hipòtesi  $H_0$  és correcta:
  - Si el  $P$ -valor és petit  $\rightarrow$  tenim evidència en contra de  $H_0$
  - Si el  $P$ -valor no és petit, **NO** demostra la “veritat” de  $H_0$
- *Incorrecte:  $P = 0.000$   $\rightarrow$  Correcte:  $P < 0.001$ .*
- $P$ -valor **NO** és cap de les següents probabilitats:
  - la probabilitat d’haver-se equivocat
  - la probabilitat que la hipòtesi nul·la sigui certa
  - la probabilitat d’haver rebutjat erròniament la hipòtesi nul·la
- $1 - P$ -valor **NO** és la probabilitat que la hipòtesi alternativa sigui certa
- Trobareu més a “[Frequent misunderstandings](#)” a [Wikipedia](#)

# Proves d'hipòtesi. Resolució

1. Escollir una **variable** segons els objectius de l'estudi
2. Escollir un disseny i un **estadístic**
3. Definir una **hipòtesi**  $H_0$  que es vol posar a prova, enfront una hipòtesi alternativa  $H_1$
4. Especificar la **distribució** de l'estadístic si  $H_0$  fos certa (i les premisses adients)
5. Amb les dades, calcular el **valor de l'estadístic** ( $z$ )
6. **Contrastar  $H_0$** . Dues alternatives per fer-ho:
  - a. Si  $|z| > z_{1-\alpha}$  (unilateral) o  $|z| > z_{1-\alpha/2}$  (bilateral) llavors rebutjar  $H_0$  ( $H_0$  és poc versemblant)
  - b. Calcular el valor de  $P \rightarrow$  Si  $P < \alpha$ , llavors rebutjar  $H_0$  ( $H_0$  és poc versemblant)
7. Afegir l'estimació per interval **IC(1- $\alpha$ )**

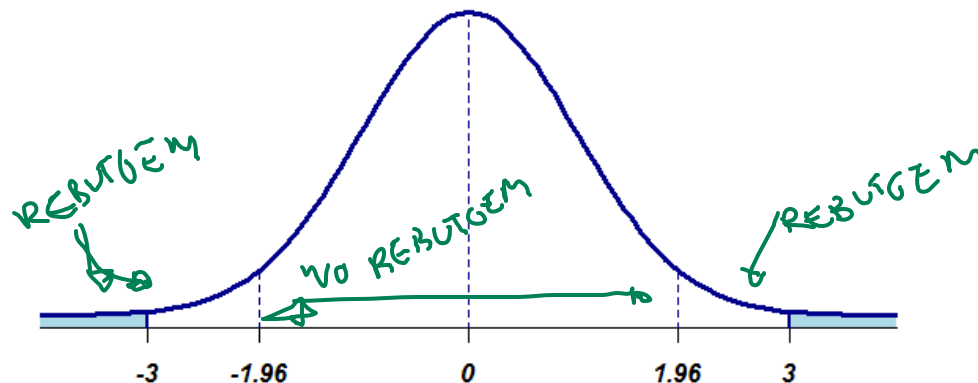
**Problema:** Com fixem  $\alpha$ ?

- Valorar el cost que representa una conclusió equivocada
- Definir un llindar arbitrari per fixar què es considera rellevant

# Proves d'hipòtesi. Llimars d'acceptació

Es poden trobar els límits d'una regió crítica

- **Bilateral:** a l'esquerra de  $z_{\alpha/2}$  (-1.96), i a la dreta de  $z_{1-\alpha/2}$  (1.96)
- **Unilateral:**
  - per l'esquerra: per sota de  $z_{\alpha}$  (-1.645)
  - per la dreta: per sobre de  $z_{1-\alpha}$  (1.645)
- Si l'estadístic cau a la regió crítica (en blau al dibuix), llavors la hipòtesi nul·la ( $H_0$ ) és dubtosa.
- Si l'estadístic cau a la regió d'acceptació (en blanc al dibuix), llavors no podem rebutjar  $H_0$



# Proves d'hipòtesi sobre la $\mu$ . Estadístic

- Si  $\sigma$  és coneguda, l'estadístic de referència és:

$$z = \frac{\bar{x} - \mu}{\frac{\sigma}{\sqrt{n}}} \sim N(0,1)$$

*mitj. mostra* →  $\bar{x}$        $\mu$  ← *mitjana poblacional*

- Si  $\sigma$  és desconeguda, l'estadístic de referència és

$$t = \frac{\bar{x} - \mu}{\frac{s}{\sqrt{n}}} \sim t_{n-1}$$

*desv. tipus mostra* →  $s$        $t_{n-1}$  ← *t. estudent*

- Sota  $H_0$ , tindrà:

- distribució Normal estàndard, si la grandària mostral és suficient:

$$t = \frac{\bar{x} - \mu}{\frac{s}{\sqrt{n}}} \sim N(0,1) \quad \text{si } n \geq 100$$

- distribució t Student amb  $(n-1)$  g.l si la variable estudiada  $X$  és Normal

$$t = \frac{\bar{x} - \mu}{\frac{s}{\sqrt{n}}} \sim t_{n-1} \quad \text{si } X \sim N$$

# Proves d'hipòtesi sobre $\mu$ . Exemple

desviació hip

Si en el cas d'una embotelladora de 1 litre tenim:  $\bar{x} = 997$ ,  $s = 10$  i  $n = 100$ . Podem pensar que la mitjana poblacional és  $\mu = 1000$  cc?

1. **Variable:** contingut en envasos de 1000cc
2. **Estadístic:**  $\hat{t} = (\bar{x} - \mu) / (s / \sqrt{n})$
3.  **$H_0$ :**  $\mu = 1000$ cc      vs.       **$H_1$ :**  $\mu \neq 1000$ cc
4. **Distribució de l'estadístic** sota  $H_0$ :  $(\bar{x} - \mu) / (s / \sqrt{n}) \sim N(0,1)$  ja que  $n = 100$ ;
5. **Càlculs:**  $\hat{t} = (\bar{x} - \mu) / (s / \sqrt{n}) = \frac{997 - 1000}{10 / \sqrt{100}} = -3$
6. **P-valor** =  $P(|z| > |-3|) = 0.0027$       [**Taules:**  $2 * (1 - 0.9987)$  ; **R:** `pnorm(-3) + (1 - pnorm(3))`]
7. **Conclusió:** com que  $P$  és menor que  $\alpha$ , es rebutja  $H_0$ :  $\mu = 1000$ cc

**Conclusió pràctica:** ens estan estafant!

$$8. \text{IC}(\mu, 0.95) = \bar{x} \mp z_{1-\frac{\alpha}{2}} \cdot \frac{\sigma}{\sqrt{n}} = 997 \mp 1.96 \cdot \frac{10}{\sqrt{100}} = [995.04, 998.96]$$

**Nota:** per graus llibertat propers a 100, D. t Student és molt propera a DN

**R:** funció **t.test**

# Proves d'hipòtesi sobre $\mu$ . Exercici

En 9 fitxers, la diferencia  $D$  entre els temps d'execució de dos programes de compressió de fitxers ha estat de mitjana 6.71 i desviació 6.00. Acceptant que  $D \sim N$ , ¿es pot acceptar que  $E(D) = \mu = 0$ ? (és a dir, acceptar que els dos compressors tarden el mateix en mitjana?)

1. **Variable:**  $D$  (diferència en temps)
2. **Estadístic:**  $\hat{t} = (\bar{d} - \mu) / (s / \sqrt{n})$
3.  $H_0 : \mu = 0$  vs.  $H_1 : \mu \neq 0$
4. **Distribució de l'estadístic sota  $H_0$ :**  $(\bar{d} - \mu) / (s / \sqrt{n}) \sim t_{n-1}$  ja que  $D \sim N$
5. **Càlculs:**  $\hat{t} = (\bar{d} - \mu) / (s / \sqrt{n}) = \frac{6.71 - 0}{6 / \sqrt{9}} = 3.355$
6. **P-valor:**  $P(|\hat{t}| > |3.355|) = 0.01$  [R: `pt(-3.355) + (1-pt(3.355))`]
7. **Conclusió:** es rebutja  $H_0$ :  $E(D) = \mu = 0$  ja que el p-valor (0.01) és  $< \alpha$

**Conclusió pràctica:** no triguen el mateix

8. **IC( $\mu$ , 0.95):**  $\bar{x} \mp t_{n-1, 1-\frac{\alpha}{2}} \cdot \frac{s}{\sqrt{n}} = 6.71 \mp 2.306 \cdot \frac{6}{\sqrt{9}} = [2.1, 11.3]$

**R:** funció ***t.test***

# Proves d'hipòtesi sobre $\pi$ . Exemple

En el **exemple** anterior del basket, amb  $n=20$  i 8 encerts.

1. **Variable**: resultat de cada tir (canasta o no)
2. **Estadístic**: nombre d'encerts  $X$
3.  $H_0 : \pi = 0.80$  (sóc un magnífic tirador) vs.  $H_1 : \pi < 0.80$  (no sóc tan bo)
4. Si  $H_0$  es certa:  $X \sim B(n, \pi) = B(20, 0.80)$

**Premisses**: mostra de tirs independents, amb la mateixa probabilitat d'encert

[Si  $n$  és gran (i  $\pi$  lluny de 0 i de 1), pot ser més simple utilitzar l'aproximació de la Binomial a la Normal (en aquest cas, treballem amb  $P=X/n$  enlloc de  $X$ ): Si  $H_0$  es certa:  $P \sim N(\pi, \pi(1-\pi)/n)$ . Requereix una premissa addicional:  $n \cdot \pi > 5$  i  $n \cdot (1-\pi) > 5$  (No s'aplicarà en l'exemple)]

5. **Càlcul de l'estadístic** :  $x = 8$
6. **P-valor** =  $P(X \text{ més lluny de } H_0 \text{ que } x) = 0.0001$  (*unilateral*)
7. **Conclusió**: **SÍ**, hi ha (forta) evidència en contra de  $H_0: \pi=0.80$
8. El càlcul de l'**IC** quan la  $n$  és petita no és senzill, però es pot trobar amb la instrucció de R **binom.test**. En aquest exemple val  $[0, 0.61]$

Si la mostra és petita obtindreu intervals molt amples... Com ha de ser!

# Proves d'hipòtesi sobre $\pi$ . Exercici

Llencem una moneda 100 vegades i obtenim 63 cares. Està “trucada” la moneda?

1. **Variable:** resultat de cada llançament(cara o creu)
2. **Estadístic:**  $Z = (P - \pi) / \sqrt{\pi \cdot (1 - \pi) / n}$
3.  **$H_0$  :**  $\pi = 0.50$  (moneda equilibrada) vs.  **$H_1$  :**  $\pi \neq 0.50$  (moneda trucada)
4. **Distribució si  $H_0$  es certa:**  $Z \sim N\left(\pi, \sqrt{\pi \cdot (1 - \pi) / n}\right) = N(0.5, 0.05)$

**Premisses:** 1) llançaments independents, amb la mateixa prob. de cara i (2)  $n\pi > 5$ ;  $n(1-\pi) > 5$

5. **Càlcul de l'estadístic:**  $= \frac{0.63 - 0.5}{\sqrt{0.5 \cdot (1 - 0.5) / 100}} = 2.6$
6. **P-valor:**  $P(X \text{ més lluny de } H_0 \text{ que } x) = 0.0094$  (bilateral)
7. **Conclusió:** **SÍ**, hi ha (certa) evidència en contra de  $H_0$ :  $\pi = 0.50$
8. **IC95%:**

$$IC(\pi, 0.95) = p \mp z_{0.975} \cdot \sqrt{(\pi \cdot (1 - \pi)) / n} = 0.63 \mp 1.96 \cdot \sqrt{(0.63 \cdot 0.37) / 100} = [0.53, 0.73]$$



# Proves d'hipòtesi vs IC

- Per exemple, sobre el valor del paràmetre  $\pi$  en la població:
  - **PH** fa una pregunta “tancada”: **¿és  $\pi = 0.5$ ?**
  - **IC** fa una pregunta “oberta”: **¿quin es el valor de  $\pi$ ?**
- Donar els resultats sempre amb IC implica:
  - Si es rebutja  $H_0$ , dir on es troba el paràmetre
  - Si no es rebutja  $H_0$ , quantificar la informació de que es disposa
  - L'IC proporciona informació més fàcil d'interpretar que el P-valor.

# Formulari. Proves d'hipòtesi

Paràmetre	Hipòtesi	Estadístic	Premisses	Distribució sota $H_0$	Criteri Decisió (Risc $\alpha$ )
$\mu$	$H : \mu = \mu_0$	$\hat{z} = \frac{\bar{y} - \mu_0}{\sqrt{\sigma^2/n}}$	$Y \rightarrow N$ ò $n \geq 30$ i $\sigma$ coneguda	$\hat{z} \rightarrow N(0,1)$	Rebutjar $H$ si $ \hat{z}  > z_{1-\alpha/2}$ ( $ \hat{z}  > 1.96$ amb $\alpha=5\%$ )
$\mu$	$H : \mu = \mu_0$	$\hat{t} = \frac{\bar{y} - \mu_0}{\sqrt{s^2/n}}$	$Y \rightarrow N$	$\hat{t} \rightarrow t_{n-1}$	Rebutjar $H$ si $ \hat{t}  > t_{n-1, 1-\alpha/2}$ ( $ \hat{t}  > t_{n-1, 0.975}$ amb $\alpha=5\%$ )
$\mu$	$H : \mu = \mu_0$	$\hat{z} = \frac{\bar{y} - \mu_0}{\sqrt{s^2/n}}$	$n \geq 100$	$\hat{z} \rightarrow N(0,1)$	Rebutjar $H$ si $ \hat{z}  > z_{1-\alpha/2}$ ( $ \hat{z}  > 1.96$ amb $\alpha=5\%$ )
$\pi$ (Binomial)	$H : \pi = \pi_0$	$\hat{z} = \frac{p - \pi_0}{\sqrt{\frac{\pi_0(1 - \pi_0)}{n}}}$	$(1-\pi_0)n \geq 5$ $\pi_0 n \geq 5$	$\hat{z} \rightarrow N(0,1)$	Rebutjar $H$ si $ \hat{z}  > z_{1-\alpha/2}$ ( $ \hat{z}  > 1.96$ amb $\alpha=5\%$ )
Anexe: $\lambda$ (Poisson)	$H : \lambda = \lambda_0$	$\hat{z} = \frac{f - \lambda_0}{\sqrt{\lambda_0}}$	$\lambda_0 \geq 5$	$\hat{z} \rightarrow N(0,1)$	Rebutjar $H$ si $ \hat{z}  > z_{1-\alpha/2}$ ( $ \hat{z}  > 1.96$ amb $\alpha=5\%$ )
$\sigma$ (normal)	$H : \sigma = \sigma_0$	$\hat{\chi}^2 = \frac{s^2 \cdot (n-1)}{\sigma^2}$	$Y \rightarrow N$	$\hat{\chi}^2 \rightarrow \chi^2_{n-1}$	Rebutjar $H$ si $\hat{\chi}^2 < \chi^2_{n-1, \alpha/2}$ o $\hat{\chi}^2 > \chi^2_{n-1, 1-\alpha/2}$
En les proves unilaterals s'acumula el valor de P-valor a un sol costat:				$H : \mu \leq \mu_0 \rightarrow$ Rebutjar $H$ si $\hat{z} > z_{1-\alpha}$ $H : \mu \geq \mu_0 \rightarrow$ Rebutjar $H$ si $\hat{z} < -z_{1-\alpha}$	

# Annexe: Premissa de Normalitat

- Hi ha dues mesures que ajuden a valorar el grau d'ajustament, afinitat o similitud a una certa distribució de referència

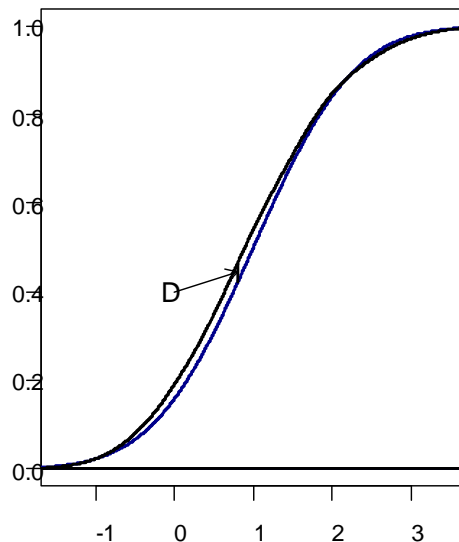
Kolmogorov-Smirnov (Estadístic D)	Shapiro-Wilk (Estadístic W)
<p>Distància màxima entre la funció de distribució empírica i la teòrica.</p> <p>Valors elevats indiquen No Normalitat</p> <p>Entre 0 i 1 (usualment, prop de 0).</p> <p><b>Valors alts indiquen des-ajustament</b></p>	<p>Mesura la correlació entre els quantils observats i els teòrics.</p> <p>Valors elevats indiquen Normalitat</p> <p>Entre 0 i 1 (usualment, prop de 1).</p> <p><b>Valors alts indiquen bon ajustament</b></p>

- Tots dos estadístics fluctuen a les mostres i han de interpretar-se amb prudència. Proporcionen **P-valors**, però farem només un anàlisi descriptiu i visual (qqplot)
- A continuació mostrem alguns exemples generats a partir de distribucions conegudes i amb diferents mides mostrals.

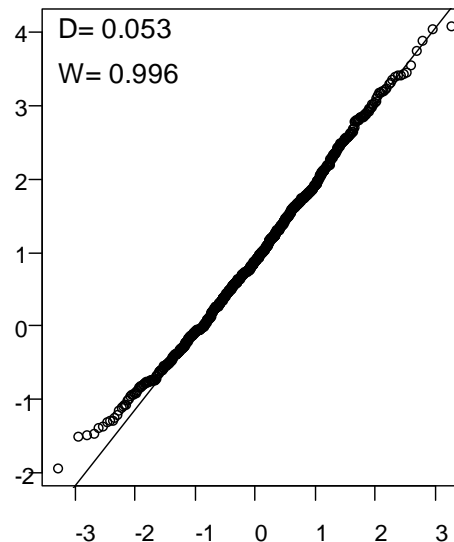
# Annexe: Premissa de Normalitat

- Els paquets solen proporcionar P-valors de la  $H_0: X \sim N$
- No ens hem de fiar massa dels P-valors obtinguts d'aquestes proves
- Es millor emprar les mesures i les eines gràfiques (QQ-Norm, PP-Norm)

Distribution (n = 965)



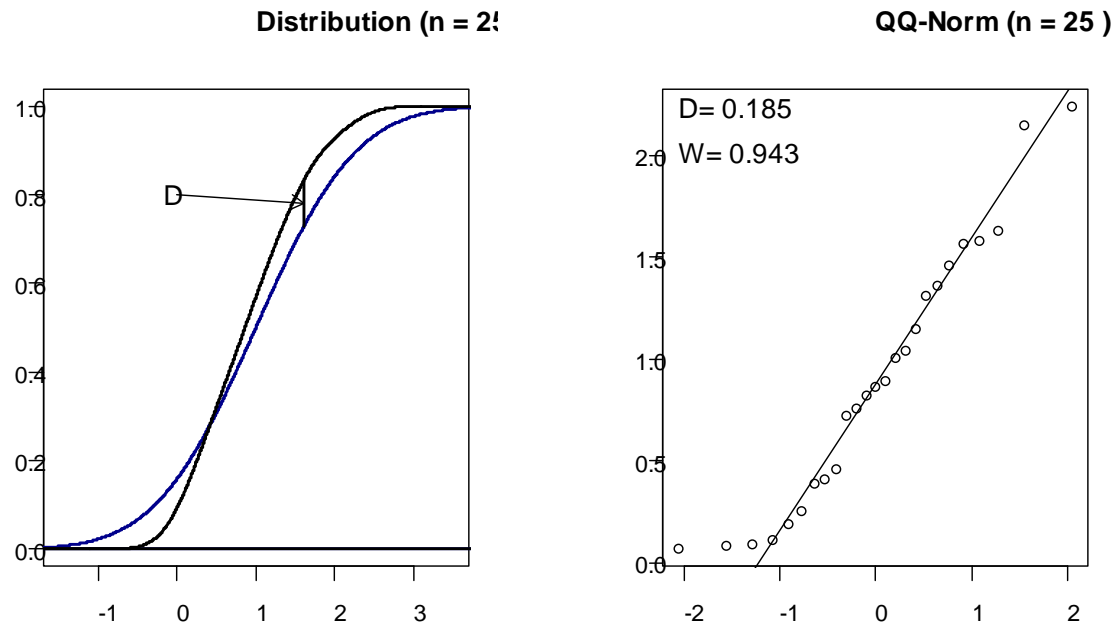
QQ-Norm (n = 965)



- Aquestes 965 observacions van estar generades seguint una distribució Normal
- D, W i el QQ-Norm mostren que les dades s'ajusten prou bé a una Normal
- No obstant, els p-valors de les proves (0.008 i 0.013) ens farien rebutjar la hipòtesi de Normalitat

# Annexe: Premissa de Normalitat

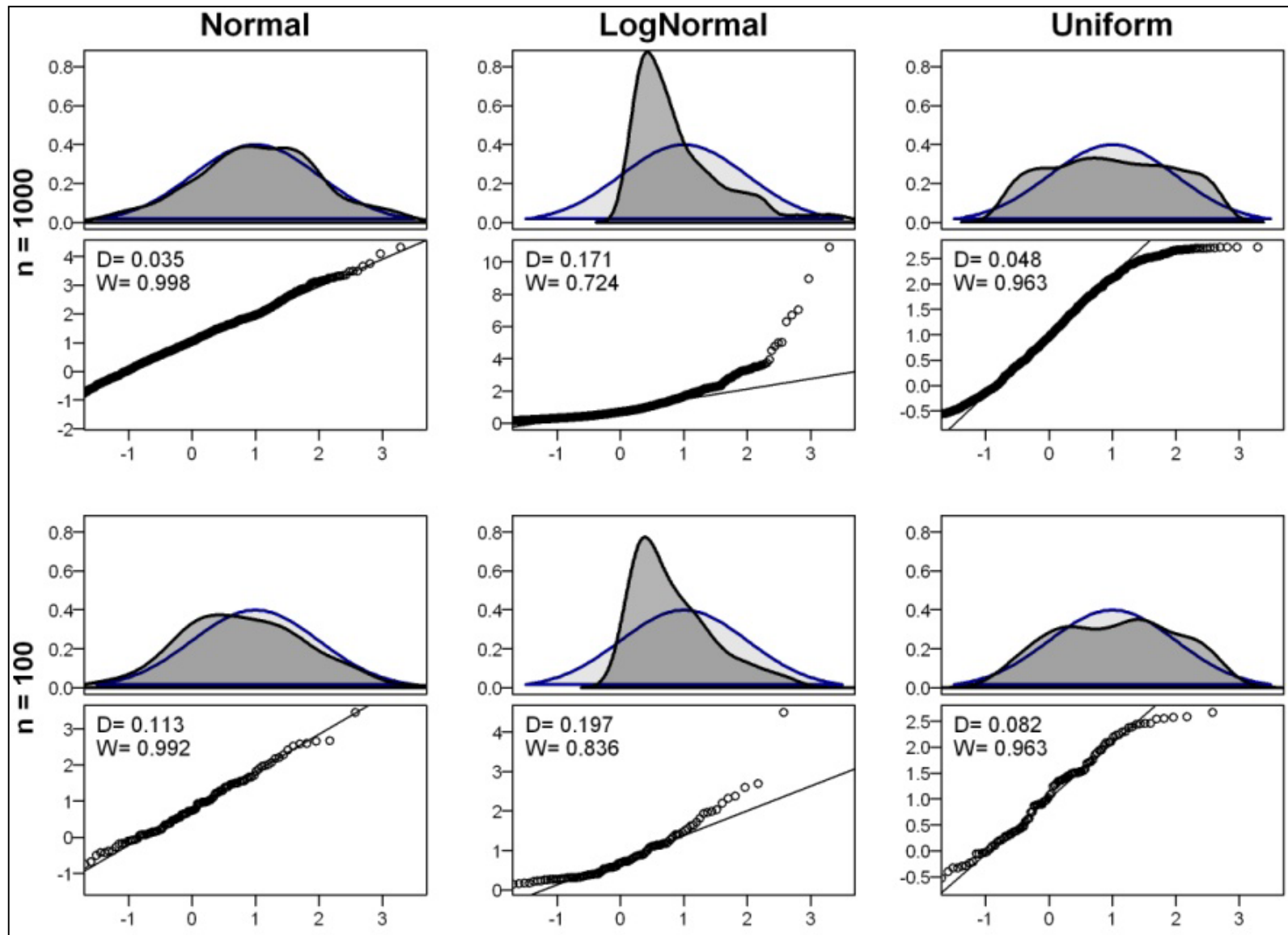
- Aquí tenim un exemple del cas contrari



- Aquestes 25 observacions van estar generades seguint una distribució Exponencial
- $F_X$  i D mostren la distància amb la teòrica, encara que W té una bona correlació.
- No obstant, els P-valors de les proves (0.32 i 0.17) no ens fan rebutjar la hipòtesi de Normalitat

Noteu la paradoxa: quan la mostra és petita, és important detectar la No Normalitat, però els P valors fallen. En canvi, quan la mostra és gran, és poc important detectar-la (pel bon comportament asimptòtic dels estimadors) però el P-valor detecta desviacions irrelevantes.

# Annexe: Premissa de Normalitat (n gran)



# Annexe: Premissa de Normalitat (n petita)

