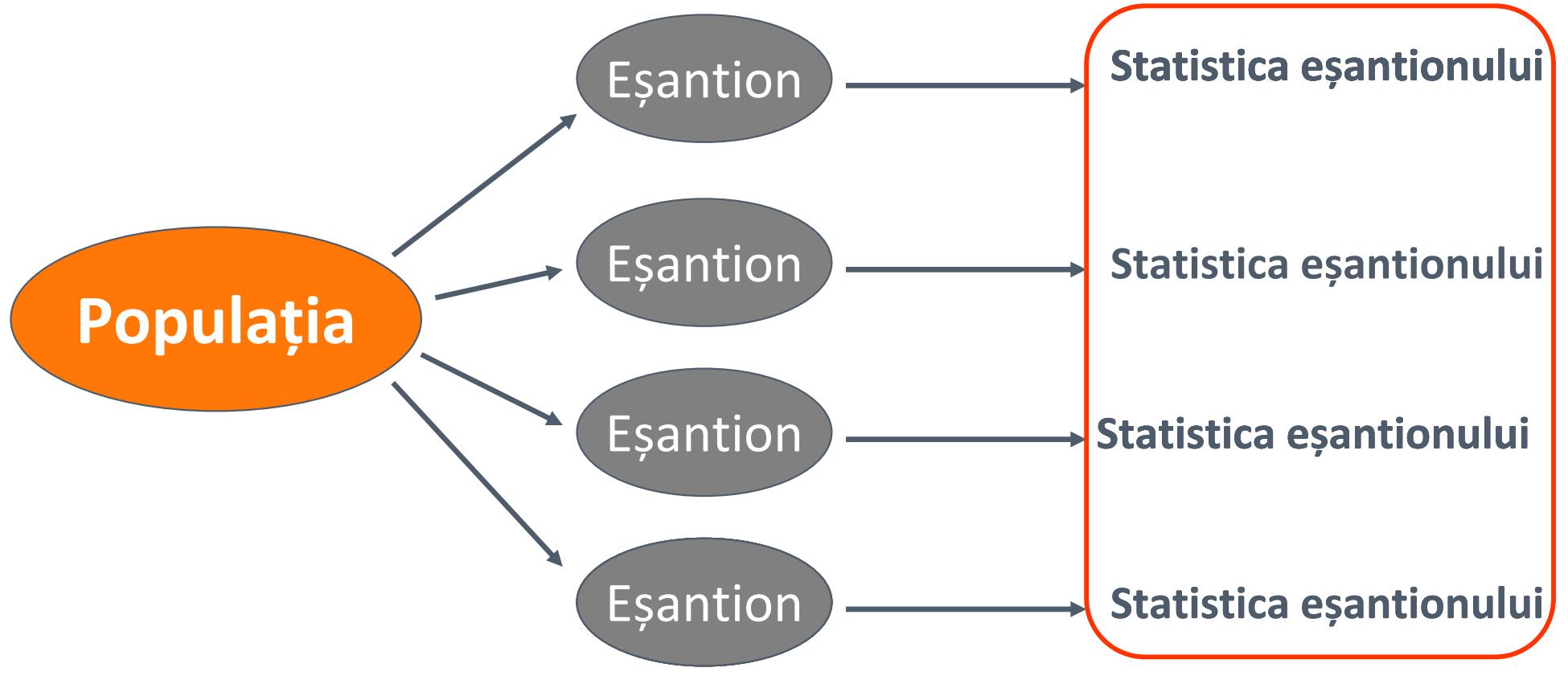


# Estimarea parametrilor statistici

- » Estimatorul punctual
- » Intervalul de încredere
  - Medie
  - Proporție

# Cuprins

>  
2



Distribuția  
eșantionului



Distribuția  
de eșantionare



Copii de 1 an  
din Ro

N = volumul  
populației

AB:  $x_{AB,1}, x_{AB,2}, \dots, x_{AB,1000}$

...

GJ:  $x_{GJ,1}, x_{GJ,2}, \dots, x_{GJ,1000}$

...

VN:  $x_{VN,1}, x_{VN,2}, \dots, x_{VN,1000}$

$\bar{x}_{AB}$

$\bar{x}_{GJ}$

$\bar{x}_{VN}$

$$\mu = \frac{x_{AB} + \dots + x_{GJ} + \dots + x_{VN}}{N}$$

$$\sigma = \sqrt{\frac{\sum_{i=1}^N (x_i - \bar{x})^2}{N}}$$

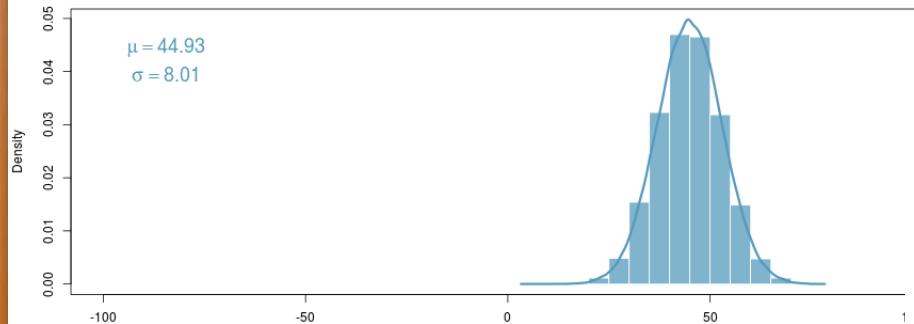
Distribuția  
de eșantionare

media  $(\bar{x}) \approx \mu$

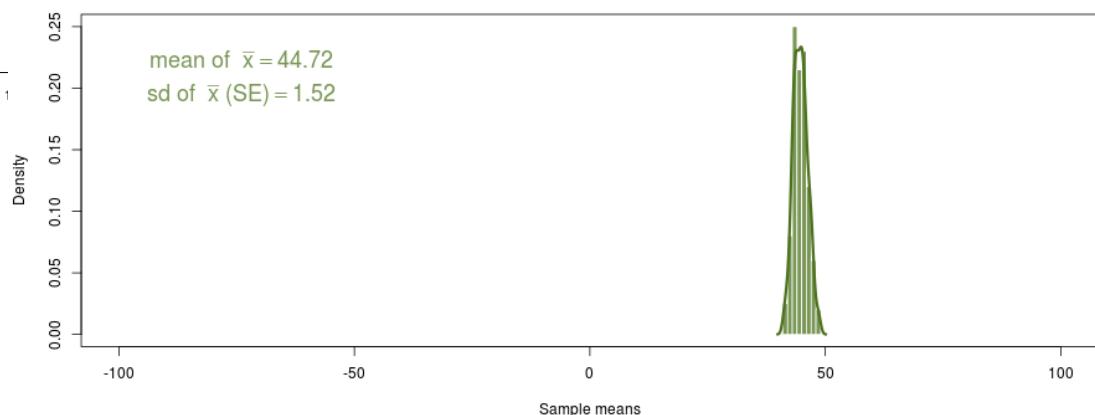
$n \uparrow \Rightarrow s < \sigma$

> 4

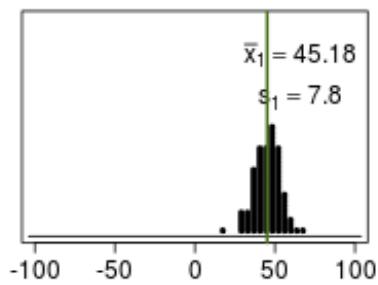
Population distribution: Normal

**n=30****200 eşantioane**

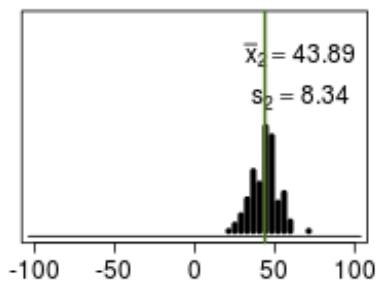
Sampling distribution



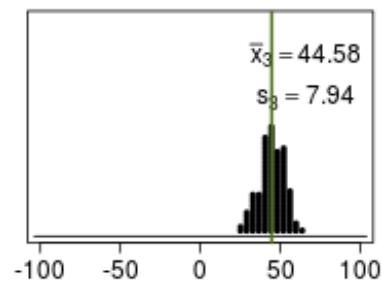
Sample 1



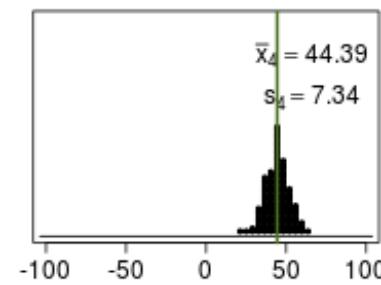
Sample 2



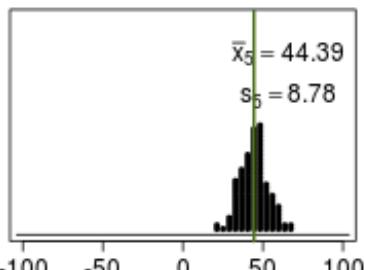
Sample 3



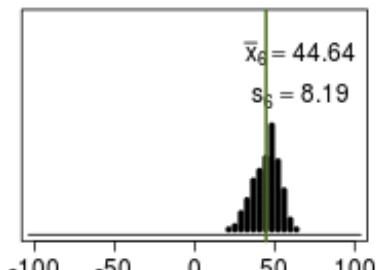
Sample 4



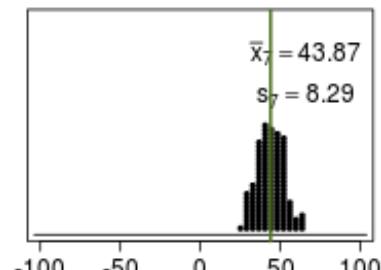
Sample 5



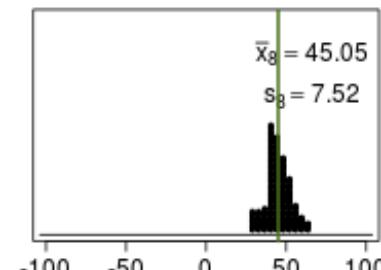
Sample 6



Sample 7



Sample 8



- » Distribuția statisticii eșantionului este aproape normală, cu media aproape egală cu cea a populației și cu deviația standard egală cu eroarea standard (deviația standard a populației împărțită la radical din volumul eșantionului).

$$\bar{x} \sim N \left( \text{mean} = \mu, SE = \frac{\sigma}{\sqrt{n}} \right)$$

↓ **forma**      ↓ **centralitatea**      ↓ **dispersia**

- » Condiții:
- » Independența: eșantioanele trebuie să fie independente (eșantion randomizat / asignare randomizată). – în caz de eșantionare fără înlocuire,  $n < 10\% \times N$ .
- » Volumul eșantionului/asimetrie: populația e normal distribuită / dacă e distribuția e asimetrică volumul eșantionului e mare ( $n > 30$ )

## Teorema limită centrală

6

## » Distribuția normală: De ce o folosim?

- Multe variabile biologice urmează o distribuție normală
- Distribuția normală este bine înțeleasă din punct de vedere matematic

## » Estimarea punctuală

- O valoare a parametrului teoretic estimat
  - +  $m$  (media eșantionului) este un estimator punctual al mediei populației ( $\mu$ )
- Este influențată de fluctuațiile de eșantionare
- Poate să fie foarte departe de valoarea reală a parametrului estimat

» Se recomandă ca estimarea unui parametru teoretic să se realizeze prin intermediul unui interval nu a unei singure valori

- Acum intervalul se numește interval de confidență/încredere
- Parametrul estimat aparține cu o probabilitate mare intervalului de confidență

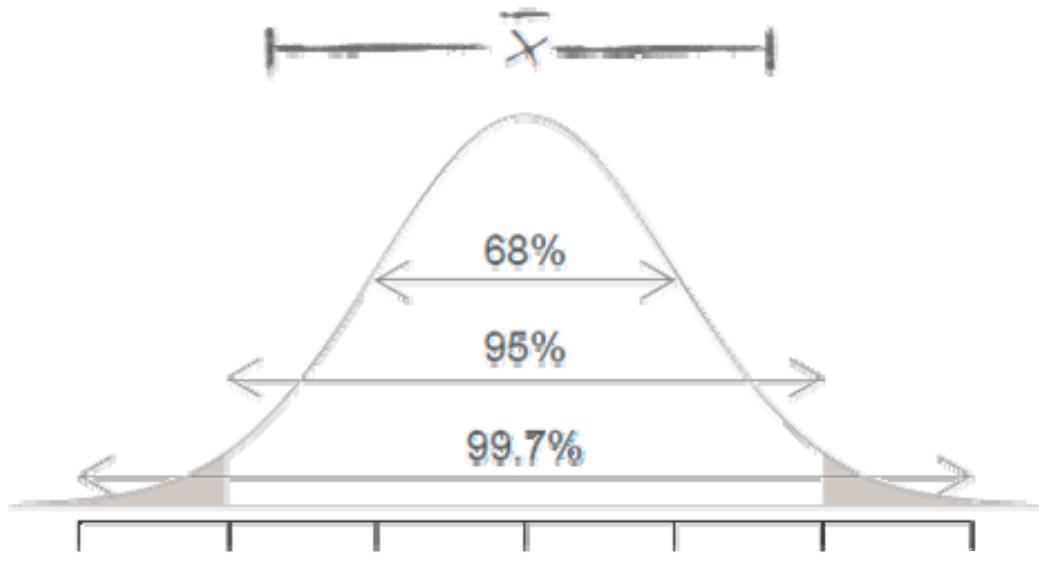
» Intervalul de confidență/încredere = un interval plauzibil de valori asociat unui parametru al populației

- Dacă raportăm un estimator punctual, cel mai probabil acesta nu va fi egal cu parametrul populației
- Dacă raportăm un interval avem o sansă ca acest interval să cuprindă valoarea parametrului populației

## Intervalul de încredere - De ce? > 8

## Teorema limită centrală

$$\bar{x} \sim N\left(\text{mean} = \mu, \text{SE} = \frac{\sigma}{\sqrt{n}}\right)$$



Marja de eroare (ME)

$$\text{~IC95\% (medie): } \bar{x} \pm 1,96 \times \text{SE}$$

## Intervalul de încredere

> 9

# Definiție

- » Un sir de valori al unui estimator de interes calculat astfel încât pentru o probabilitate de eroare aleasă să includă valorile adevărate ale variabilei.
- » **P[valoarea critică inferioară < estimatorul < valoarea critică superioară] = 1- $\alpha$** 
  - unde  $\alpha$  = nivelul de semnificație
- » Intervalul definit de valorile critice va cuprinde estimatorul populației cu o probabilitate de  $1-\alpha$
- » Se aplică în cazul variabilelor distribuite normal!

- » Unul din primele exemple de asimetrie comportamentală a omului este preferința de a întoarce capul spre dreapta nu spre stânga. Un studiu realizat pe un eșantion de 124 cupluri a pus în evidență că 64,5% din acestea întorc capul spre dreapta când se sărută. Eroarea standard asociată acestui estimator este egală cu aproximativ 4%. Care din următoarele sunt false?
1. Un volum de eșantion mai mare va determina o eroare standard mai mică.
  2. Marja de eroare pentru un IC de 95% asociată procentului de cupluri care întorc capul la dreapta când se sărută e aproximativ 8%.
  3. IC95% pentru procentul de cupluri care întorc capul spre dreapta când se sărută este  $64,5\% \pm 4\%$
  4. IC99,7% pentru procentul de cupluri care întorc capul spre dreapta când se sărută este  $64,5\% \pm 12\%$

# Intervalul de încredere pentru medie

- » Eroarea standard a mediei este egală cu deviația standard împărțită la radicalul volumului eşantionului
  - Dacă deviația standard este mare, șansa de eroare în estimator este mare
  - Dacă volumul eşantionului este mare, șansa erorii în estimator este mică.

$$\left[ \bar{x} - Z_{\alpha} \frac{s}{\sqrt{n}}, \bar{x} + Z_{\alpha} \frac{s}{\sqrt{n}} \right]$$

**Valoare critică**

$$\bar{x} \pm Z_{\alpha} \frac{s}{\sqrt{n}}$$

Z = valoare  
constantă pentru un  
prag de semnificație  
dat

Condiții:

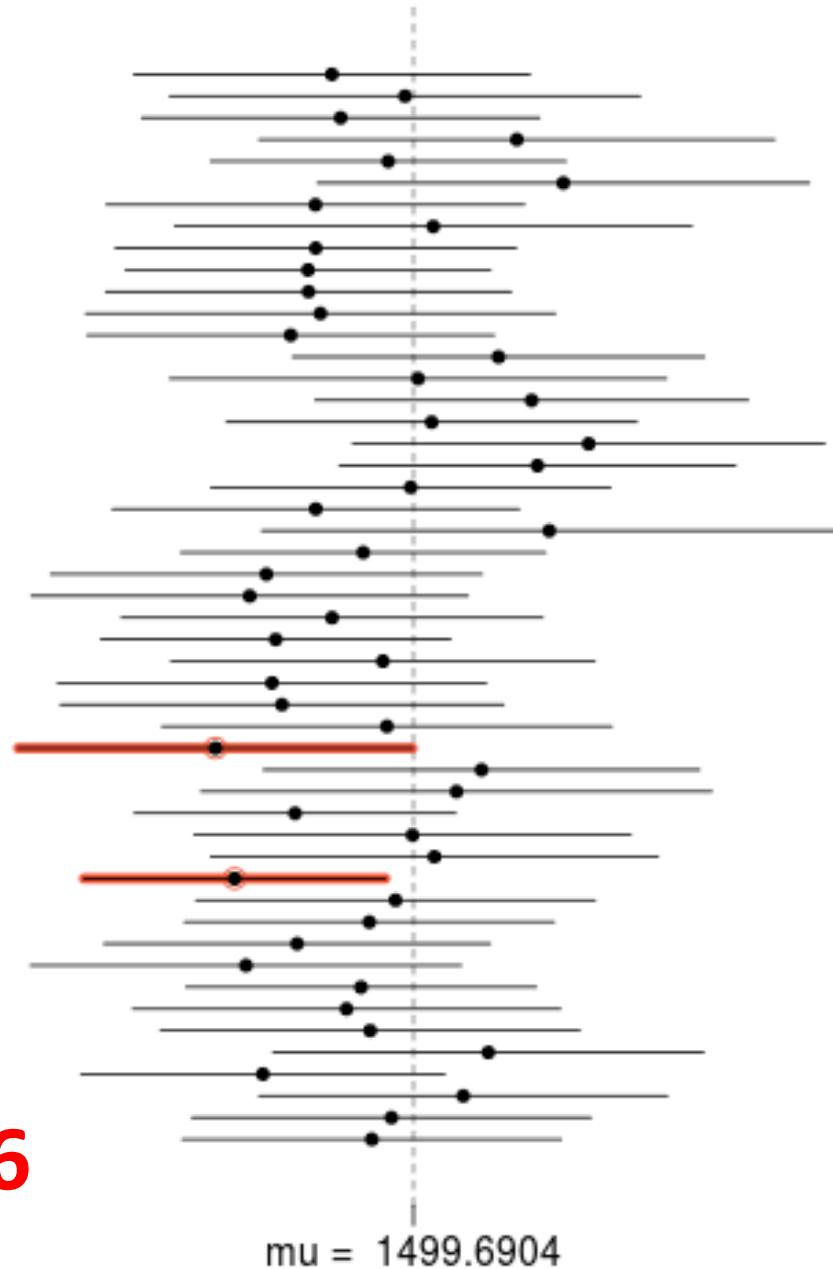
» Independența:

- Eșantioanele observate trebuie să fie independente
- Eșantion randomizat (fără înlocuire  $n < 10\% \times N$ ) sau asignare randomizată

» Asimetria distribuției:  $n \geq 30$ , sau mai mare pentru distribuții asimetrice

## Intervalul de încredere pentru medie > <sup>13</sup>

- » Luăm mai multe eșantioane și construim intervalele de confidență de 95%
- » ~ 95% din aceste intervale vor conține media adevărată a populației ( $\mu$ )
- » Intervale de confidență frecvent utilizate: 95%, 98% și 99%.

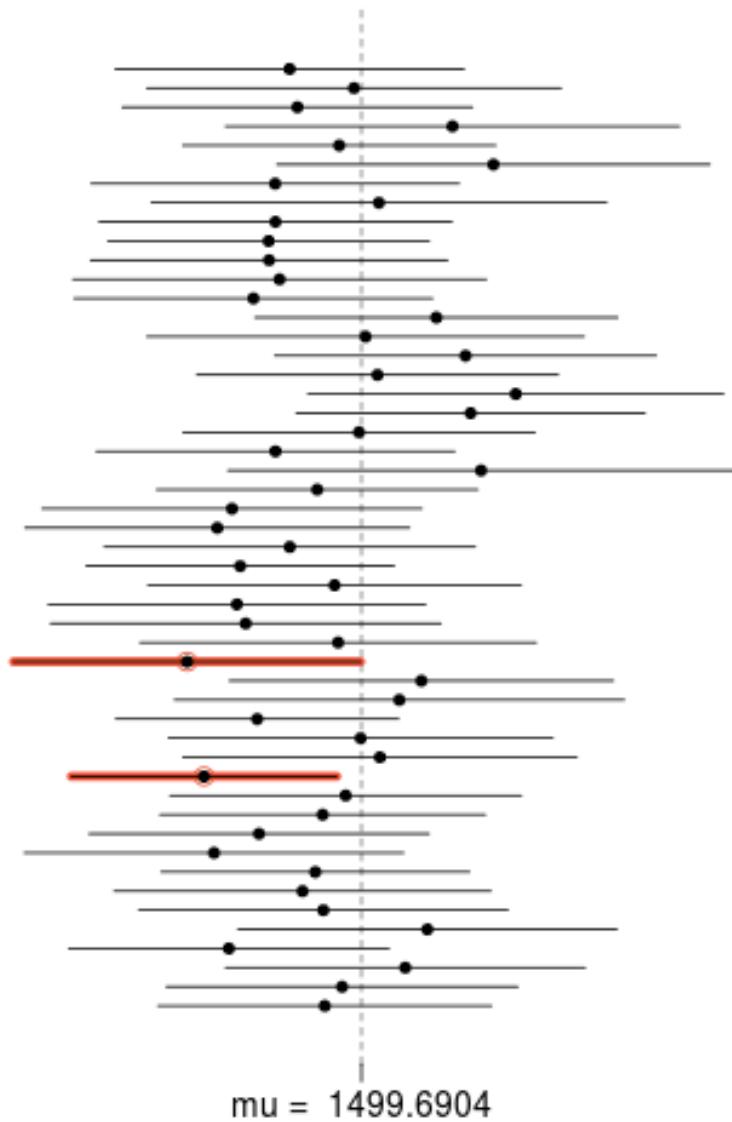
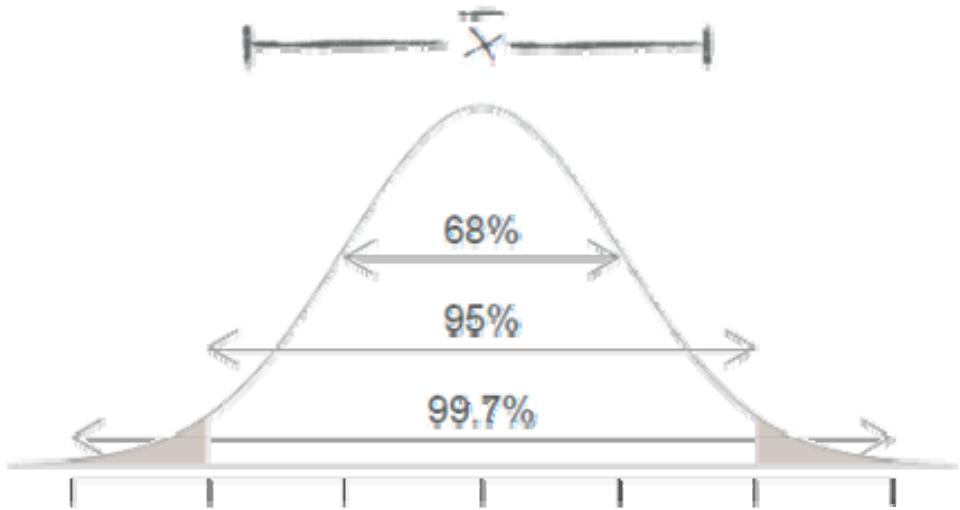


$$48/50 = 0,96$$

## Acuratețe vs precizie



- » Dacă dorim să fim siguri că media populației este cuprinsă în IC acesta trebuie să fie larg sau îngust?

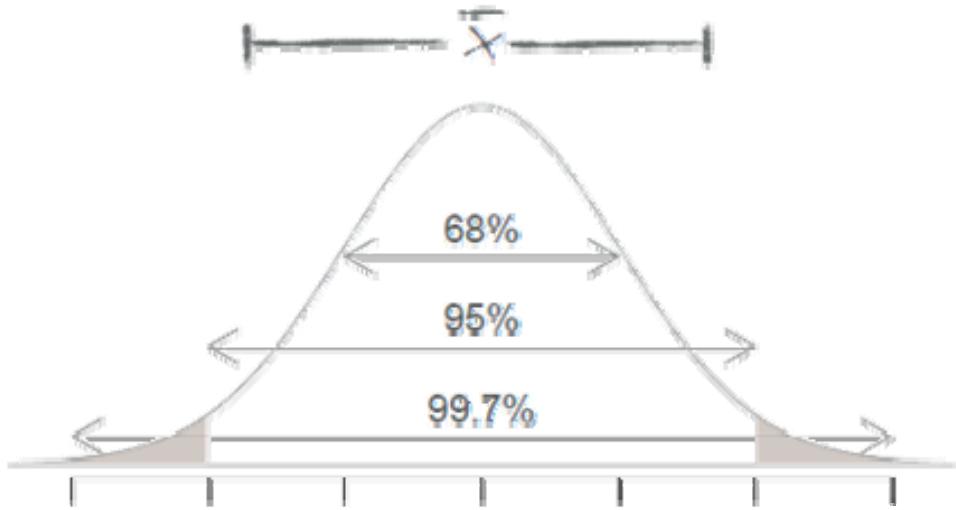


## Acuratețe vs precizie

15

- » Dacă dorim să fim siguri că media populației este cuprinsă în IC acesta trebuie să fie larg sau îngust?
  - 99% ( $\uparrow$ ): lărgimea ( $\uparrow$ ) + acuratețe ( $\uparrow$ ) + precizie ( $\downarrow$ )
- » Cum putem avea precizie mare și acuratețe mare?
  - $n \uparrow$

$$ME = Z_{\alpha} \frac{s}{\sqrt{n}} \rightarrow n = \left( \frac{Z_{\alpha} \times s}{ME} \right)^2$$



## Acuratețe vs precizie

16

- » Se dorește testarea efectului unui medicament folosit în tratamentul epilepsiei la mamă asupra dezvoltării cognitive a copilului. Dezvoltarea cognitivă se testează prin estimarea indicelui de inteligență a copilului de 3 ani născut de femei care au urmat în timpul sarcinii tratament cu medicamentul de interes.
- » Studii anterioare au arătat că deviația standard a indicelui de inteligență a copilului de 3 ani este egală cu 18 puncte.
- » Care este numărul de copii în vîrstă de 3 ani care trebuie incluși în studiu pentru a obține un interval de confidență de 90% cu o margine a erorii mai mică sau egală cu 4 puncte?

$$ME \leq 4$$

$$IC = 90\%$$

$$Z = 1,65$$

$$\Sigma = 18$$

$$ME = Z_{\alpha} \frac{\sigma}{\sqrt{n}}$$

$$n = \left( \frac{Z_{\alpha} \times \sigma}{ME} \right)^2 = \left( \frac{1,65 \times 18}{4} \right)^2$$

$$n = 55,13 \rightarrow n = 56$$

- » Se dorește testarea efectului unui medicament folosit în tratamentul epilepsiei la mamă asupra dezvoltării cognitive a copilului. Dezvoltarea cognitivă se testează prin estimarea indicelui de inteligență a copilului de 3 ani născut de femei care au urmat în timpul sarcinii tratament cu medicamentul de interes.
- » Studii anterioare au arătat că deviația standard a indicelui de inteligență a copilului de 3 ani este egală cu 18 puncte.
- » Care este numărul de copii în vîrstă de 3 ani care trebuie incluși în studiu pentru a obține un interval de confidență de 90% cu o margine a erorii mai mică sau egală cu 4 puncte?

**ME  $\leq$  4**

**IC = 90%**

**Z = 1,65**

**$\sigma = 18$**

**n=56**

**ME  $\leq$  2**

**IC = 90%**

**Z = 1,65**

**$\sigma = 18$**

**n=4×56=224**

**ME  $\leq$  6**

**IC = 90%**

**Z = 1,65**

**$\sigma = 18$**

**n=24,50**

**ME  $\leq$  4**

**IC = 95%**

**Z = 1,96**

**$\sigma = 18$**

**n=77,79=78**

**ME  $\leq$  2**

**IC = 95%**

**n=311,17=312**

**ME  $\leq$  5**

**IC = 95%**

**n=49,79=50**

- » Un eșantion de 49 studenți au fost întrebați în câte relații exclusive au fost implicați până la data studiului. Studenții din eșantion au avut în medie 3 relații exclusive, cu o deviație standard de 1,2. estimați media adevărată a numărului de relații exclusive bazată pe rezultatele acestui eșantion utilizând intervalul de confidență de 95%. Distribuția de eșantionare a fost aproximativ normală.
- » Pasul 1: verificarea condițiilor.
  - Numărul de relații exclusive ale unui student este independent de al altui student.  $n = 49 < 10\% \times N$  (numărul de studenți din universitate)
  - $n > 30 \rightarrow$  distribuția de eșantionare a numărului de relații exclusive dintr-un eșantion de volum egal cu 49 este aproximativ normală.

## Exemplu

19

- » Un eșantion de 49 studenți au fost întrebați în câte relații exclusive au fost implicați până la data studiului. Studenții din eșantion au avut în medie 3 relații exclusive, cu o deviație standard de 1,2. Estimați media adevărată a numărului de relații exclusive bazată pe rezultatele acestui eșantion utilizând intervalul de confidență de 95%. Distribuția de eșantionare a fost aproximativ normală.
- » Pasul 2: calculăm IC95%

$$n = 49 \quad SE = \frac{1.2}{\sqrt{60}} \approx 0.1714$$

$$\bar{x} = 3$$

$$s = 1.2$$

$$ME = 3 \pm 1,96 * 0.1714$$

$$IC95\% [2,66; 5,66]$$

Suntem 95% siguri că studenții au fost implicați în medie în 3,66 - 5,66 relații exclusive.

## Exemplu

# Intervale de încredere pentru diferență între doi estimatori: Interpretare

- » Dacă 0 este conținut în intervalul de încredere, diferența dintre cele două estimări (medii, proporții, rații, etc.) este zero
  - » Dacă zero nu este conținut în intervalul de încredere, diferența dintre cei 2 estimatori punctuali nu este egală cu zero.
- 
- <http://www.biomedcentral.com/1746-6148/8/68>
  - *BMC Veterinary Research* 2012, **8**:68 doi:10.1186/1746-6148-8-68

## Results

After optimising the cut-off values in order to avoid doubtful results without deteriorating the concordance between the results of the two tests, the I-ELISA appeared to be slightly more sensitive than CFT (Se<sub>I-ELISA</sub>=0.917 [0.822; 0.992], 95% Credibility Interval (CrI) compared to Se<sub>CFT</sub>=0.860 [0.740; 0.967], 95% CrI). However, CFT was slightly more specific than I-ELISA (Sp<sub>CFT</sub>=0.988 [0.947; 1.0], 95% CrI) compared to Sp<sub>I-ELISA</sub>=0.952 [0.901; 1.0], 95% CrI).

# Intervale de Încredere: Interpretare

» Când aceeași procedură se repetă pe mai multe eșantioane, intervalul de încredere (care va fi diferit pentru fiecare eșantion) va cuprinde în 95% din cazuri valoarea reală a estimatorului punctual.

# Intervalul de Încredere

» Se calculează în funcție de:

- Talia eşantionului sau a populației
- Tipul de variabilă (calitativă SAU cantitativă)

» Formula de calcul cuprinde 2 părți

- Un estimator al calității eşantionului pe baza căruia estimatorul populației s-a calculat (eroarea standard)

+ Eroarea standard:

- Cu cât  $n$  este mai mare cu atât eroarea standard este mai mică.
- Este întotdeauna mai mică decât deviația standard

- Gradul de încredere (confidență) al intervalului specificat (scorul  $Z_\alpha$ )

» Se poate calcula pentru orice estimator

# Intervalul de încredere pentru medie

» Media glicemiei la un eșantion de 121 pacienți este de 105 iar variația de 36. Care este intervalul de încredere al mediei glicemiei în populația din care s-a extras eșantionul cu un prag de semnificație  $\alpha=0,05$ , considerând că glicemia este normal distribuită și pentru acest prag  $Z = 1,96$ .

$$\gg n = 121$$

$$\bar{X} = 105$$

$$\gg s^2 = 36$$

$$\gg s = 6$$

$$\left[ 105 - 1,96 \frac{6}{\sqrt{121}}, 105 + 1,96 \frac{6}{\sqrt{121}} \right]$$

$$\gg [105 - 1.07, 105 + 1.07]$$

$$\gg [103.93 - 106.07]$$

$$\gg [104 - 106]$$

# Compararea mediilor cu ajutorul intervalului de încredere

Figure 1.

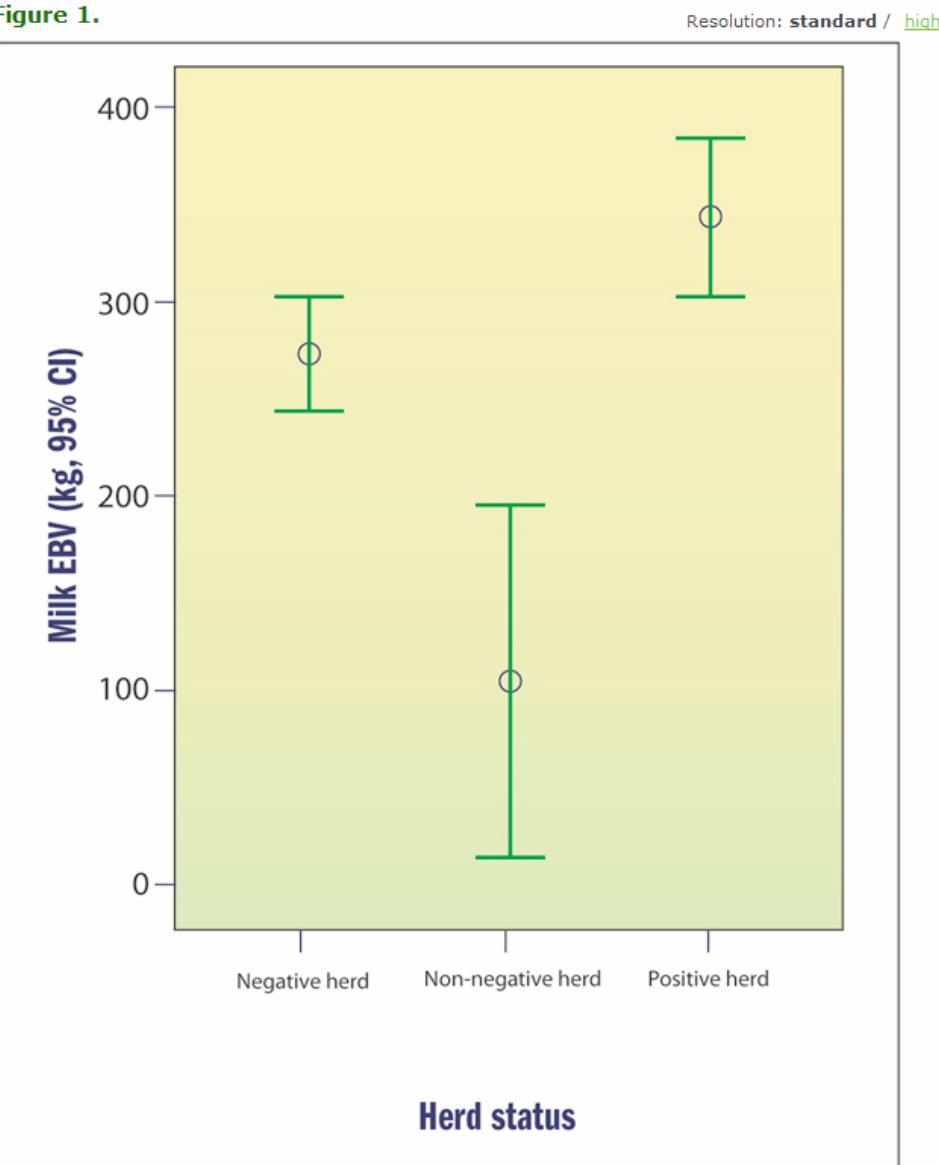
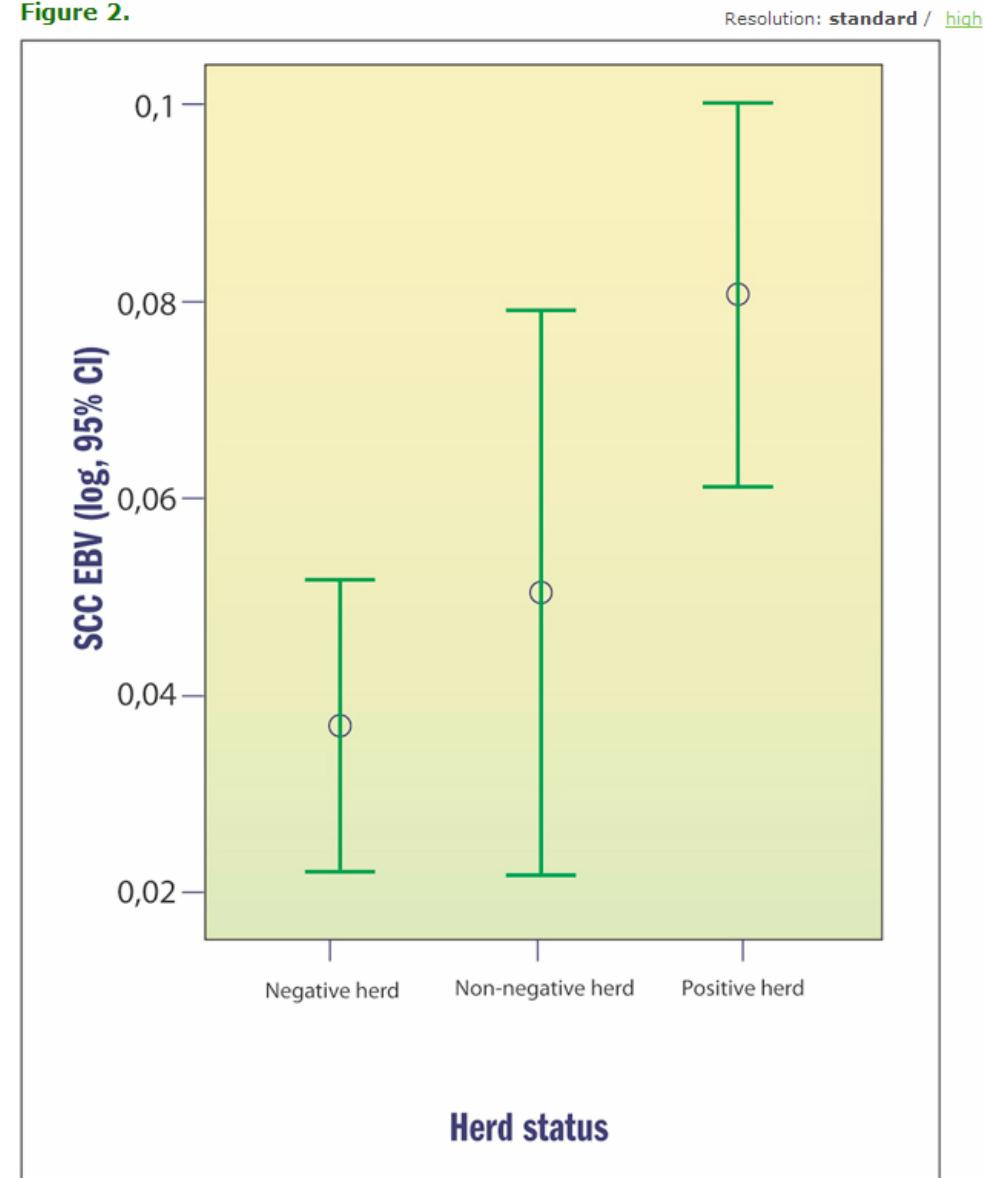


Figure 2.



# Intervalul de încredere pentru frecvențe

» Se calculează dacă:

$n \cdot f > 10$ , unde  $n$  = talia eșantionului,  $f$  = frecvența

$$\left[ f - Z_\alpha \sqrt{\frac{f(1-f)}{n}}; f + Z_\alpha \sqrt{\frac{f(1-f)}{n}} \right]$$

# Intervalul de încredere pentru frecvențe

- » Suntem interesați în estimarea frecvenței cancerului de sân la femeile între 50 și 54 de ani care au antecedente familiale pozitive. Într-un studiu randomizat la care au participat 10000 de femei, s-a constatat că 400 dintre acestea au fost diagnosticate cu cancer de sân.
- » Care este intervalul de încredere de 95% asociat frecvenței observate?

$$\» f = 400/10000 = 0.04$$

$$\left[ 0.04 - 1.96 \sqrt{\frac{0.04 \cdot 0.96}{10000}}; 0.04 + 1.96 \sqrt{\frac{0.04 \cdot 0.96}{10000}} \right]$$

- » [0,04-0,004; 0,04+0,004]
- » [0,036; 0,044]

$$\left[ f - Z_{\alpha} \sqrt{\frac{f(1-f)}{n}}; f + Z_{\alpha} \sqrt{\frac{f(1-f)}{n}} \right]$$

- » Estimarea corectă a unui parametru statistic se face cu ajutorul intervalului de încredere.
- » Intervalul de încredere depinde de volumul eşantionului și de eroarea standard.
- » Cu cât eroarea standard este mai mare cu atât intervalul de încredere este mai larg.
- » Cu cât volumul eşantionului este mai mic cu atât intervalul de încredere este mai larg.

# De reținut!

28