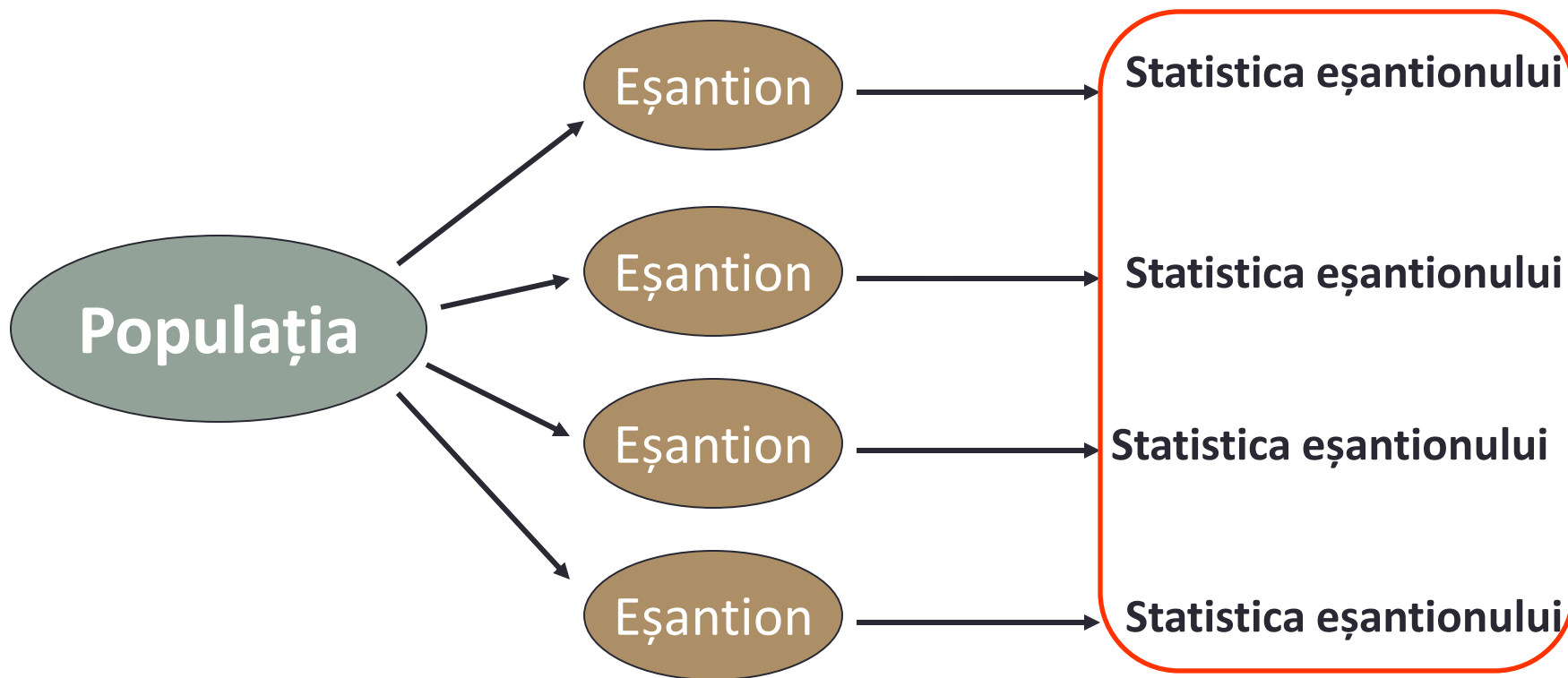


ESTIMATORI PUNCTUALI ȘI INTERVALE DE ÎNCREDERE

Despre ...

- Estimatorul punctual
- Intervalul de încredere:
 - Definiții
 - Pentru
 - medie
 - diferența dintre medii
 - frecvență
 - diferența dintre frecvențe
 - riscul relativ



**Distribuția
eșantionului**

≠

**Distribuția
de eșantionare**



N = volumul populației

Copii de 1 an din Ro

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}

$$\mu = \frac{\bar{x}_{AB} + \dots + \bar{x}_{GJ} + \dots + \bar{x}_{VN}}{N}$$

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

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

$n \uparrow \Rightarrow s < \sigma$

Estimarea punctuală

- O valoare a parametrului teoretic estimat
 - media eșantionului (\bar{X}) este un estimator punctual al mediei populației (μ)
- Este influențată de fluctuațiile de eșantionare
- Poate să fie foarte departe de valoarea reală a parametrului estimat

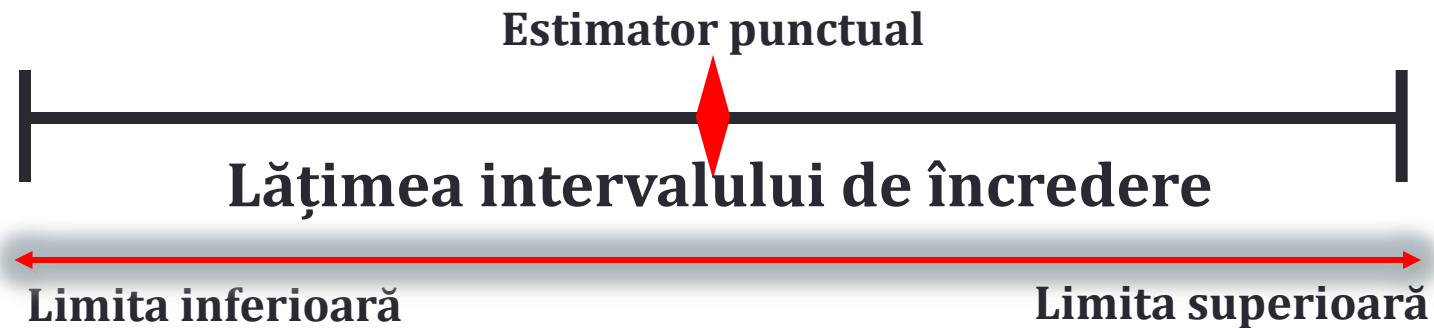
- Se recomandă ca estimarea unui parametru teoretic să se realizeze prin intermediul unui interval nu a unei singure valori
 - Acest interval se numește interval de încredere
 - Parametrul estimat aparține cu o probabilitate mare intervalului de confidență

Definiție

- Un șir de valori al unui estimator de interes calculat astfel încât pentru o probabilitate de eroare aleasă să includă valorile adevărate ale variabilei.
- **$P[\text{valoarea critică inferioară} < \text{estimatorul} < \text{valoarea critică superioară}] = 1 - \alpha$**
 - unde α = nivelul de semnificație
- Intervalul definit de valorile critice va cuprinde estimatorul populației cu o probabilitate de $1 - \alpha$

Estimatorul punctual vs. intervalul de încredere

- Estimatorul punctual = valoarea unei statistici obținută pe un eșantion
 - Cât de multă incertitudine este asociată estimatorului punctual?
- Un interval oferă mai multe informații despre o caracteristică a populației decât un estimator punctual



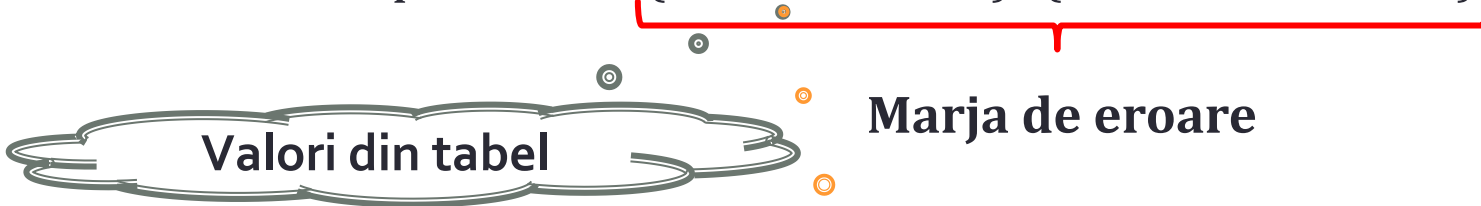
Intervalul de încredere

- **Intervalul de încredere:**

- Ia în considerare variabilitatea de eșantionare → are valori diferite pentru fiecare eșantion
- Se poate calcula pe baza observării unui singur eșantion
- Oferă informații despre parametrul necunoscut al populației

- **Formula generală:**

Estimator punctual \pm (Valoare critică) \times (Eroarea standard)

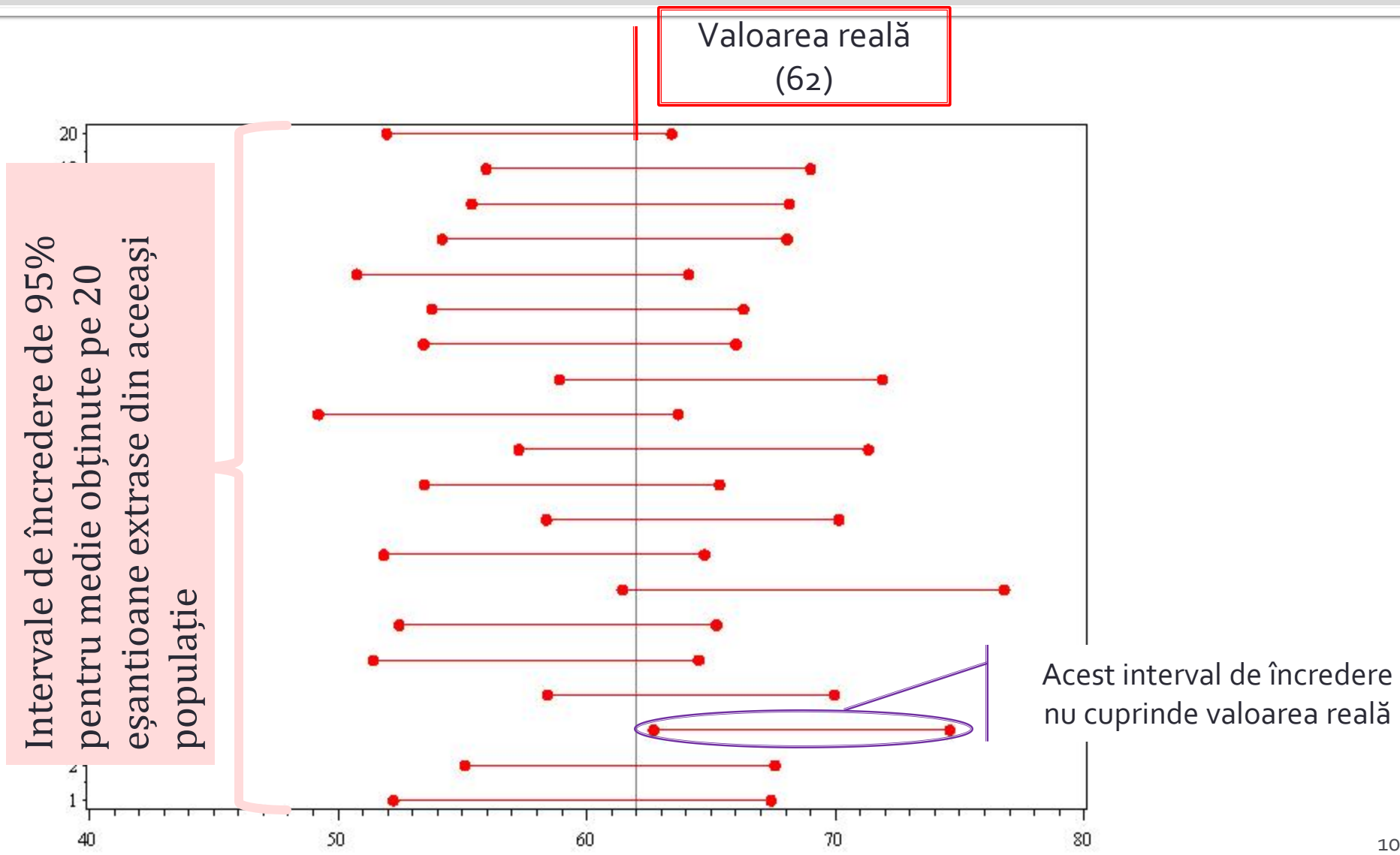


Depinde de α și df (grade de libertate)

Intervalul de încredere

- Marja de eroare și respectiv lățimea intervalului de încredere
 - este cu atât mai mică cu cât volumul eșantionului e mai mare.
 - variază cu valoarea nivelului de semnificație (α)
- Nivel de semnificație $\alpha = 5\%$ → interval de încredere este de 95% (IC95%) - $IC = (1 - \alpha) = 0,95$
- Interpretare:
 - Dacă toate eșantioanele posibile de volum n s-ar extrage din populație și mediile și intervalele de încredere asociate ar fi calculate, 95% din intervalele de încredere vor conține valoarea reală a parametrului populației
 - Un interval de încredere poate să conțină sau poate să nu conțină valoarea reală a parametrului (datorită riscului de 5%)

Intervalul de încredere



Intervalul de încredere

- Se calculează în funcție de:
 - Talia eșantionului
 - Tipul de variabilă (calitativă SAU cantitativă)
- Formula de calcul cuprinde 2 părți
 - Un estimator al calității eșantionului 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
- Se poate calcula pentru orice estimator
- Nivele de confidență utilizate curent: 90%, 95%, 98%, și 99%

Intervalul de încredere pentru medie

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

$$\bar{X} \pm Z_{\alpha} \frac{S}{\sqrt{n}}$$

- Condiții de aplicare a intervalului de încredere:
 - Independența: Observațiile trebuie să fie independente
 - Eșantion randomizat / Asignare randomizată
 - Pentru un eșantion fără înlocuire, $n < 10\% * N$
 - Volumul eșantionului: $n \geq 30$ sau mai mare dacă distribuția în populație nu este simetrică

Dacă dorim să fim siguri că intervalul conține parametrul populației, acesta va trebui să fie îngust sau larg?

IC larg \rightarrow acuratețea \uparrow precizie \downarrow
n \uparrow \uparrow \rightarrow Acuratețe \uparrow & precizie \uparrow

- Parametrul populației: nivelul colesterolului la femeile din România cu HTA și obezitate
- S-a extras din această populație un eșantion randomizat de volum 40 și s-a determinat o statistică a estimatorului punctual medie egală cu 220 mg/dl, deviația standard (S) egală cu 18 mg/dl
 - Media 220 mg/dl este estimatorul punctual al parametrului necunoscut al populației
 - Datorită variabilității de eșantionare, media va fi însoțită de intervalul de încredere asociat pentru estimarea corectă a parametrului populației

$$IC_{95\%} = \left[220 - 1,96 \frac{18}{\sqrt{49}}; 220 + 1,96 \frac{18}{\sqrt{49}} \right] = [215; 225]$$

Lățimea = 225 - 215 = 10

$$IC_{99\%} = \left[220 - 2,58 \frac{18}{\sqrt{49}}; 220 + 2,58 \frac{18}{\sqrt{49}} \right] = [213; 227]$$

Lățimea = 227 - 213 = 14

- 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 vârstă 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?

$$ME \leq 4$$

$$IC = 90\%$$

$$Z = 1,65$$

$$s = 18$$

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

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

Avem nevoie de cel puțin 56 subiecți pentru a obține o margine de eroare de până la 4 puncte

- 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 49 este aproximativ normală.

- Pasul 2: calculăm IC95%

$$n = 49$$

$$\bar{x} = 3$$

$$s = 1,2$$

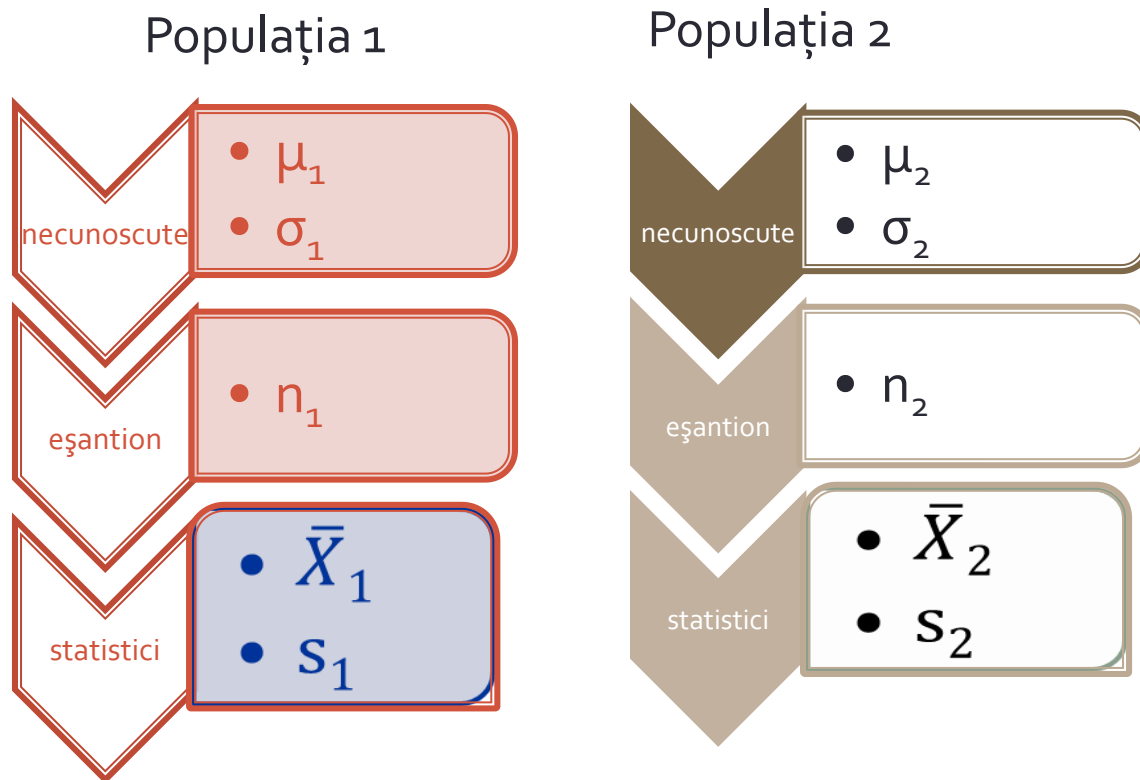
$$ES = \frac{1,2}{\sqrt{60}} \approx 0,1714$$

$$ME = 3 \pm 1,96 * 0,1714$$

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

- Suntem 95% siguri că studenții au fost implicați în medie în 3 până la 6 relații exclusive.

Intervale de încredere pentru diferența dintre medii



Estimăm $(\mu_1 - \mu_2)$ cu $\bar{X}_1 - \bar{X}_2$

Interpretare:

- Dacă valoarea zero este în intervalul de încredere → diferența dintre medii nu este semnificativ diferită de zero
- Dacă valoarea zero este în intervalul de încredere → diferența dintre medii este semnificativ diferită de zero

Intervale de încredere pentru diferența dintre medii

$$(\bar{X}_1 - \bar{X}_2) \pm t_{critic} \sqrt{\frac{(n_1 - 1)S_1^2 + (n_2 - 1)S_2^2}{n_1 + n_2 - 2}}$$

Grupa 10	7	7	8	8	8	6	9	6	5
Grupa 3	8	10	9	6	10	8	9	7	8

	Grupa 1	Grupa 2
Media	7,11	8,33
Variația	1,27	1,32
Deviația standard	1,61	1,75

df=15,97

pentru $\alpha = 0,05 \rightarrow t_{15,97} = 2,13$

$(7,11 - 8,33) \pm 2,13 \sqrt{(0,18 + 0,19)}$

$-1,22 \pm 2,13 * 0,61$

$-1,22 \pm 1,30 \rightarrow [-2,52; 0,08]$

Intervalul de încredere pentru frecvențe

- O frecvență:

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

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

- Diferența dintre două frecvențe

$$(f_1 - f_2) \pm Z_{\text{critic}} \times \text{sqrt}((f_1 \cdot (1 - f_1) / n_1) + (f_2 \cdot (1 - f_2) / n_2))$$

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?

- $f = 400/10000 = 0.04$

$$\left[0,04 - 1,96 \sqrt{\frac{0,04 \cdot 0,96}{10000}}; 0,04 + 1,96 \sqrt{\frac{0,04 \cdot 0,96}{10000}} \right]$$

- $[0,04 - 0,004; 0,04 + 0,004]$

- $[0,036; 0,044]$

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

Intervalul de încredere pentru RR

	Boală+	Boală-
Expunere +	a	b
Expunere -	c	d

$$RR = \frac{a/(a+b)}{c/(c+d)}$$

$$\ln(\widehat{RR}) \pm Z_{crit} \sqrt{\frac{b/a}{a+b} + \frac{d/c}{c+d}} \rightarrow$$

$$\left[\exp \left(\ln(\widehat{RR}) - Z_{crit} \sqrt{\frac{\frac{b}{a}}{a+b} + \frac{\frac{d}{c}}{c+d}} \right); \exp \left(\ln(\widehat{RR}) + Z_{crit} \sqrt{\frac{\frac{b}{a}}{a+b} + \frac{\frac{d}{c}}{c+d}} \right) \right]$$

	CCPulmonar+	CCPulmonar-
Fumat+	30	20
Fumat-	15	35

$$RR = \frac{30/(30+20)}{15/(15+35)} = 2$$

$$\ln(RR) = 0,69$$

$$Z_{critic} = 1,96$$

$$ES = \sqrt{0,06} = 0,24$$

$$[\exp(0,69 - 1,96 * 0,24); \exp(0,69 + 1,96 * 0,24)] \rightarrow [1,24; 3,23]$$

COMPARAREA MEDIILOR CU AJUTORUL INTERVALULUI DE ÎNCREDERE

<http://www.biomedcentral.com/content/pdf/1471-2458-12-1013.pdf>

Table 1 Living conditions of the MS-MV and the immigrant population (CASEN survey 2006)

	IMMIGRANT POPULATION 1% total sample, n = 154 431 weighted population (1877 real observations)		MS-MV GROUP 0.67% total sample, n = 108 599 weighted population (1477 real observations)	
	% or mean	95% CI	% or mean	95% CI
<i>DEMOGRAPHICS</i>				
Mean age**	X = 33.41	31.81–35.00	X = 26.13	23.41–28.26
Age categories:				
<16 years old**	13.60	11.29–16.28	45.25	39.53–51.10
16-65 years old**	79.08	75.92–81.93	47.26	41.64–52.94
>65 years old	7.32	5.33–9.97	7.49	5.31–10.46
Sex (female = 1)	45.21	41.74–48.72	51.27	47.99–55.41
Marital status:				
Single**	45.81	42.06–49.62	64.30	59.36–68.95
Married**	45.49	41.66–49.36	29.39	25.09–34.10

De reținut!

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

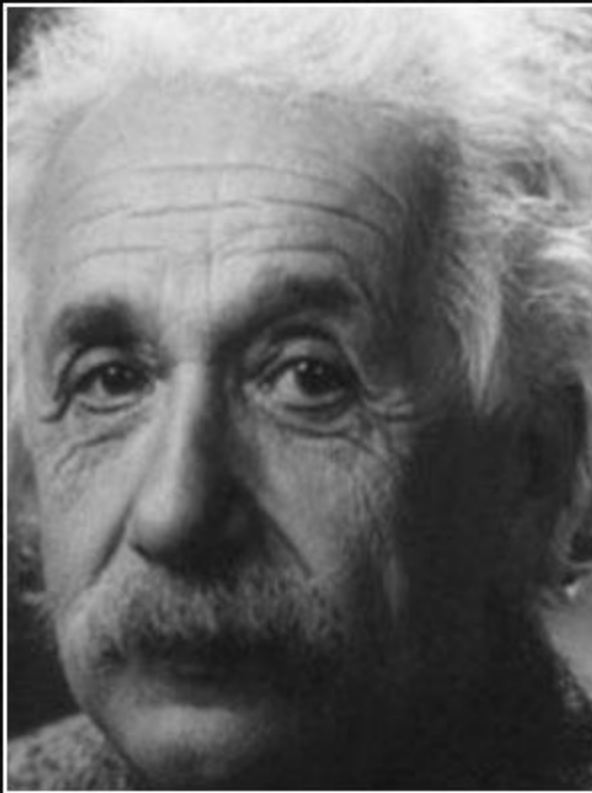
$$\left[\bar{X} - Z_{\alpha} \frac{S}{\sqrt{n}}; \bar{X} + Z_{\alpha} \frac{S}{\sqrt{n}} \right]$$

$$(\bar{X}_1 - \bar{X}_2) \pm t_{critic} \sqrt{\frac{(n_1 - 1)S_1^2 + (n_2 - 1)S_2^2}{n_1 + n_2 - 2}}$$

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

$$(f_1 - f_2) \pm Z_{critic} \times \text{sqrt}((f_1 * (1 - f_1) / n_1) + (f_2 * (1 - f_2) / n_2))$$

MULȚUMESC PENTRU ATENȚIE!



Intelligence is not the ability to store
information, but to know where to
find it.

— *Albert Einstein* —

AZ QUOTES