

Statistici Descriptive (2)

Statistici descriptive

Măsuri de centralitate <ul style="list-style-type: none">♦ Media♦ Mediana♦ Modulul♦ Valoarea centrală	Măsuri de dispersie / împrăștiere <ul style="list-style-type: none">♦ Amplitudine♦ Intervalul intercvartilic [Q1, Q3]♦ Variația (sau varianța)♦ Deviația standard♦ Coeficientul de variație♦ Eroarea standard
Măsuri de simetriei <ul style="list-style-type: none">♦ Asimetria♦ Boltirea	Măsuri de localizare <ul style="list-style-type: none">♦ Cvartile (decile; percentile)

Statistici descriptive de dispersie

- Calcularea parametrilor de dispersie/împrăștiere:
 - Amplitudinea
 - Intervalul intercvartilic [Q1, Q3]: IQR
 - Distanța intercvartile: $IQR = Q3 - Q1$
 - Deviația standard de la medie
 - Eroarea standard
 - Coeficientul de variație

Amplitudinea

$$A = X_{\max} - X_{\min}$$

- Nu ne spune nimic despre modalitatea în care datele variază în jurul valori centrale
- Valorile extreme afectează semnificativ valoarea amplitudinii
- Excel: RANGE (Descriptive Statistics)

	A	B	C	D	E	F	G	H	I	J
1	X_1	X_2	X_3	X_4	X_5	X_6	X_7	X_8	X_9	X_{10}
2	3	9	9	5	8	6	10	7	8	5



$$A = X_{\max} - X_{\min} = 10 - 3 = 7$$

Abaterea de la medie (de la mediană)

- Măsoară media distanțelor dintre fiecare valoare a variabilei și valoarea centrală (media aritmetică sau mediana).
- Dă o pondere egală fiecărei observații
- Sensibilitate mai mare decât amplitudinea sau intervalul intercuartilic dintre cuartila 3 și 1 (IQR)
- Se utilizează în calculul variației și respectiv a deviației standard

$$AD_m = \frac{|X_1 - m| + |X_2 - m| + |X_3 - m| + \dots + |X_n - m|}{n}$$

$$AD_{Me} = \frac{|X_1 - Me| + |X_2 - Me| + |X_3 - Me| + \dots + |X_n - Me|}{n}$$

- Unde \bar{X} este media aritmetică (Ma), Me=mediana datelor

Variația standard

- **VARIAȚIA (descriptivă) a EȘANTIONULUI** = media sumei pătratelor abaterilor de la medie

$$s^2 = \frac{(x_1 - m)^2 + (x_2 - m)^2 + \dots + (x_n - m)^2}{n} \text{ unde } n = \text{talie eșantionului}$$

Avantaje

- ✓ pozitivă

Dezavantaje

- ✗ unitățile varianței = pătratul unităților variabilei
- ✗ varianță mare => dispersie mare în jurul mediei => greu de inteles și de facut comparații
- ✗ este sensibilă la valori extreme
- ✗ are sens doar pentru o variabilă CANTITATIVĂ

Deviația standard

DEVIAȚIA STANDARD (abaterea standard) a eșantionului se notează cu **s**

- Are aceeași unitate de măsură ca și media și datele seriei

$$s = \sqrt{\frac{\sum_{i=1}^n (X_i - m)^2}{n}} = \sqrt{\frac{(x_1 - m)^2 + (x_2 - m)^2 + \dots + (x_n - m)^2}{n}} \quad \text{unde } n = \text{talie eșantionului}$$

$$s = \sqrt{s^2}$$

- **DEVIAȚIA STANDARD** (abaterea standard) a populației se notează cu **σ**

unde N = talia populației; μ = media pe populație

$$\sigma = \sqrt{\frac{\sum_{i=1}^N (X_i - \mu)^2}{N}} = \sqrt{\frac{(x_1 - \mu)^2 + (x_2 - \mu)^2 + \dots + (x_N - \mu)^2}{N}}$$

Deviația standard de eșantionare

- **Deviația standard de eșantionare** (deviația standard a eșantionului optimizată pentru a aproxima cel mai bine deviația standard a populației:

- $S = \sqrt{\frac{n}{n-1}} s$ unde n = talia eșantionului



$$S = \sqrt{\frac{(x_1 - m)^2 + (x_2 - m)^2 + \dots + (x_n - m)^2}{n-1}}$$

- Va fi utilizată în statistica inferențială
- Notăție S (majuscule) pentru deviația standard de eșantionare
- respectiv S^2 = variația de eșantionare

Proprietățile deviației standard

AVANTAJE

- ✓ Aceleași unități de măsură ca și variabila în sine.
- ✓ Pozitivă
- ✓ Dacă înmulțim valorile unei serii statistice cu o constantă, abaterea standard se înmulțește cu aceeași constantă
- ✓ Dacă adăugăm o valoare la fiecare valoare a unei serii statistice, abaterea standard nu se modifică

DEZAVANTAJE

- ✓ este sensibilă la valori extreme

Variația și deviația standard: exemplu de calcul

Id_student	Nota	Deviația de la medie ($X_i - m$)	$(X_i - m)^2$
X_1	9	$= 9 - 7,3 = 1,67$	$= 1,67^2 = 2,78$
X_2	6	-1,33	1,78
X_3	4	-3,33	11,11
X_4	9	1,67	2,78
X_5	4	-3,33	11,11
X_6	8	0,67	0,44
X_7	8	0,67	0,44
X_8	9	1,67	2,78
X_9	7	-0,33	0,11
X_{10}	4	-3,33	11,11
X_{11}	10	2,67	7,11
X_{12}	10	2,67	7,11

Variația (descriptivă) a eșantionului

$$s^2 = \frac{(9 - 7,33)^2 + (6 - 7,33)^2 + \dots + (10 - 7,33)^2}{12}$$

Media aritmetică (m): = 7,33

$$s^2 = 58,77 / 12 = 4,90$$

$$s = \sqrt{4,90}$$

Deviația standard a eșantionului

$$s = 2,21$$

Variația și deviația standard: exemplu de calcul

Id_student	Nota	Deviația de la medie ($X_i - m$)	$(X_i - m)^2$
X_1	9	$= 9 - 7,3 = 1,67$	$= 1,67^2 = 2,78$
X_2	6	-1,33	1,78
X_3	4	-3,33	11,11
X_4	9	1,67	2,78
X_5	4	-3,33	11,11
X_6	8	0,67	0,44
X_7	8	0,67	0,44
X_8	9	1,67	2,78
X_9	7	-0,33	0,11
X_{10}	4	-3,33	11,11
X_{11}	10	2,67	7,11
X_{12}	10	2,67	7,11

Variația de eșantionare

$$S^2 = \frac{(9 - 7,33)^2 + (6 - 7,33)^2 + \dots + (10 - 7,33)^2}{11}$$

Media aritmetică (m): $= 7,33$

$$S^2 = 58,77 / 11 = 5,34$$

$$s = \sqrt{5,34}$$

Deviația standard de eșantionare

$$S = 2,31$$

Variația și deviația standard: exemplu de calcul

Id_subiect	Greutate(kg)	Deviația de la medie ($X_i - m$)	$(X_i - m)^2$
1	71	-0.92	0.85
2	76	3.57	12.72
3	62	-9.87	97.45
4	72	-0.26	0.07
5	75	3.49	12.19
6	71	-1.36	1.85
7	70	-1.58	2.51
8	70	-1.80	3.24
9	70	-2.41	5.83
10	66	-5.84	34.06
11	60	-12.44	154.85
12	81	8.99	80.91
13	67	-4.80	23.02
14	65	-6.96	48.49
15	83	10.56	111.53
16	68	-4.35	18.89
17	77	4.70	22.12
18	74	1.74	3.01
19	65	-7.43	55.22
20	75	2.65	7.00
21	81	9.25	85.55
22	69	-2.74	7.50
23	73	1.21	1.47
...

Variația de eșantionare

$$S^2 = \frac{(71 - 72)^2 + (76 - 72)^2 + \dots (80 - 72)^2}{99}$$

Media aritmetică (m): = 72 kg

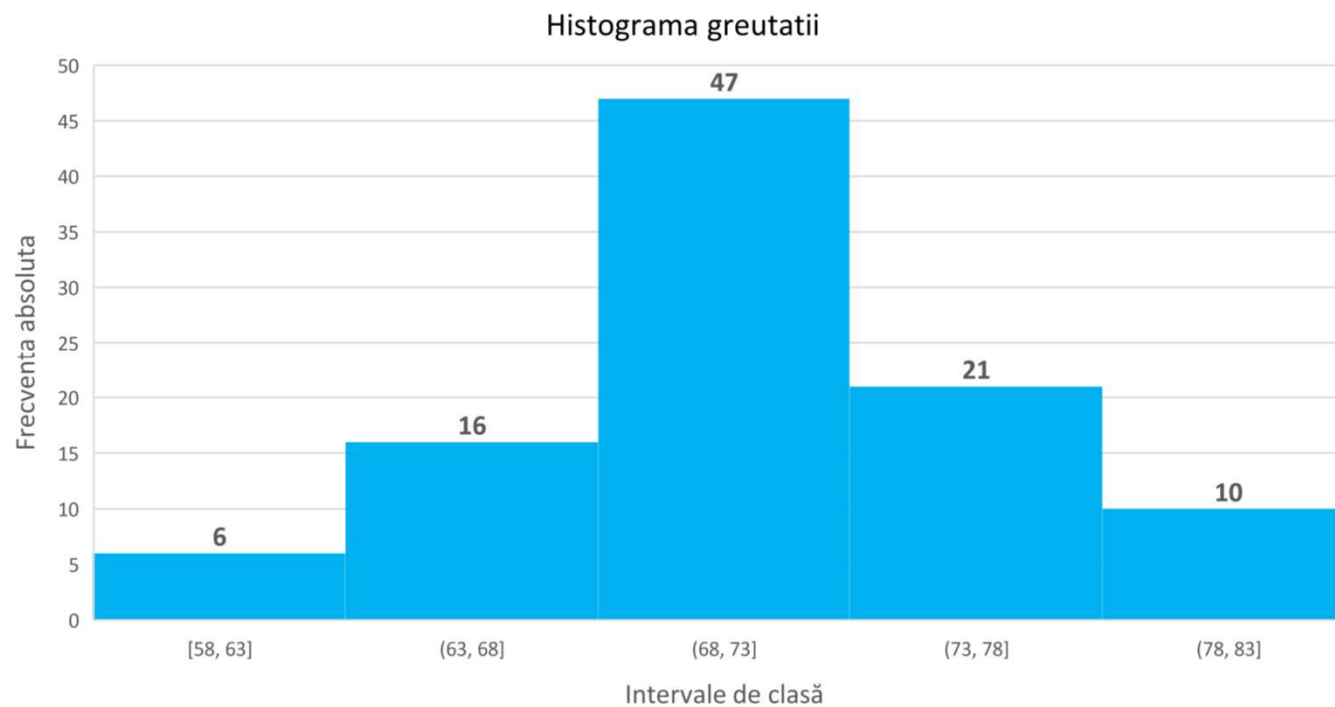
$$S^2 = 2414,10/99 = 24,38$$

$$s = \sqrt{24,38}$$

Deviația standard de eșantionare

$$S = 4,94 \text{ kg}$$

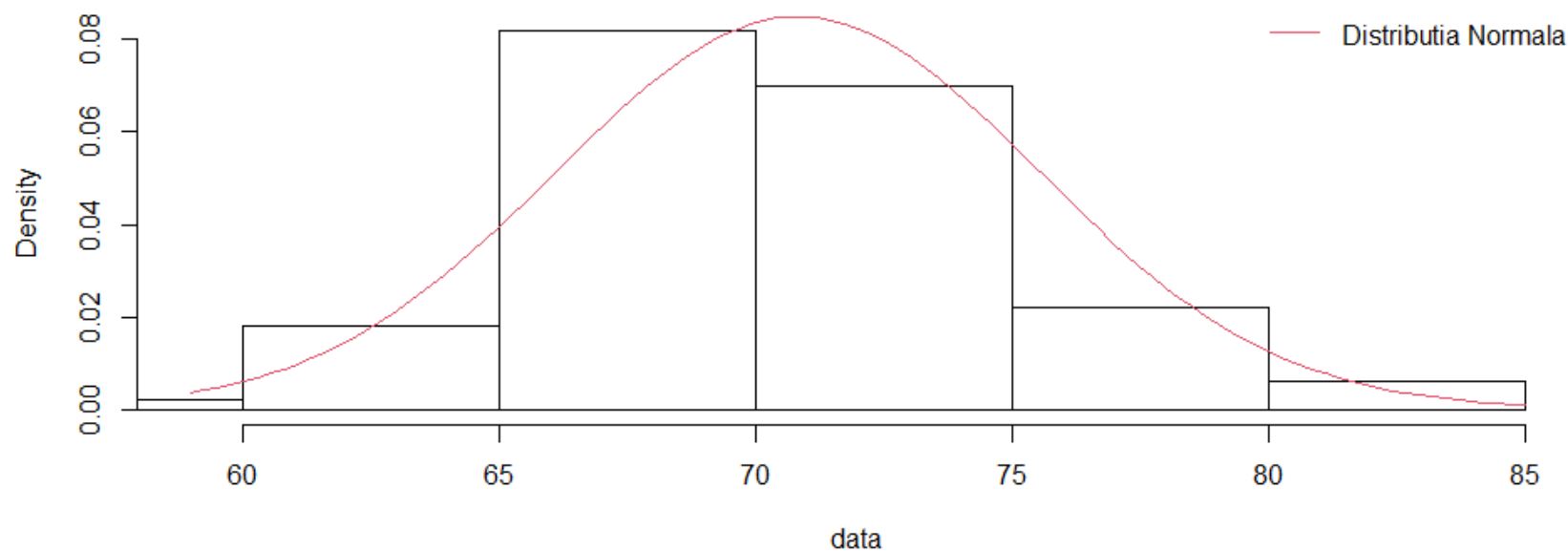
Histograma și curba normală (gaussiană)



Histograma și curba normală (gaussiană)

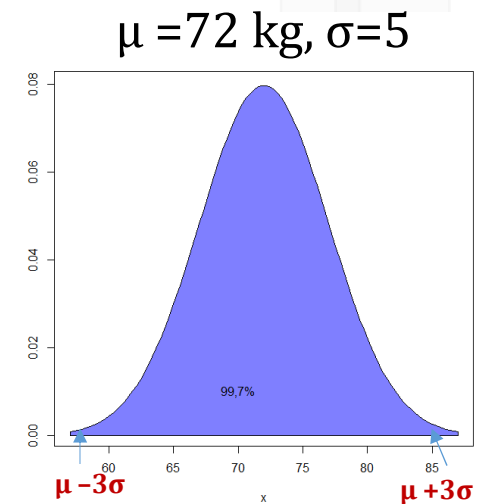
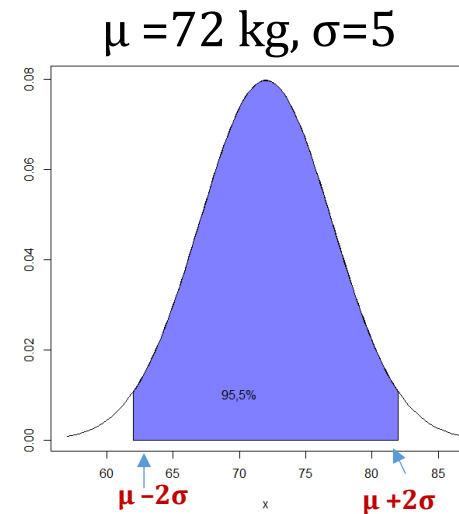
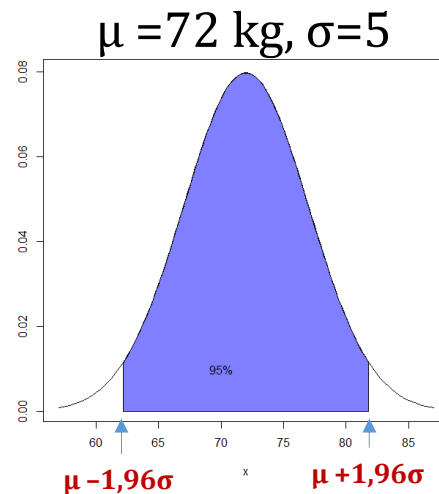
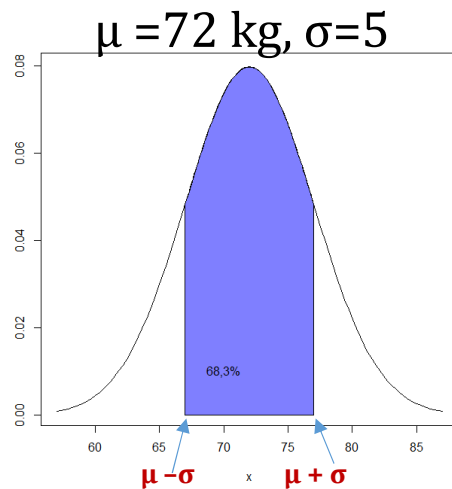
Histograma greutății

Histogram and theoretical densities



Curba Normală

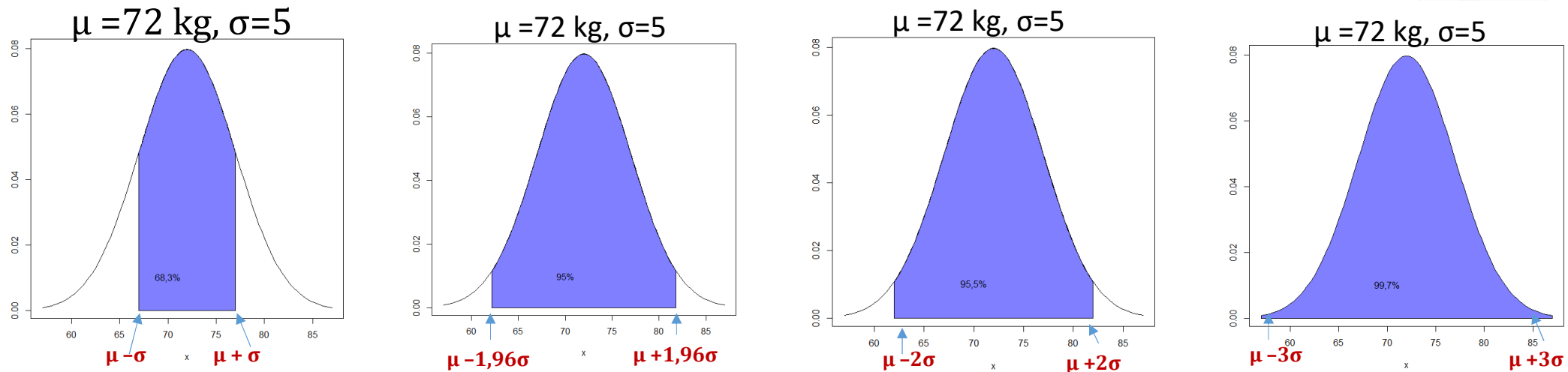
- Distribuție simetrică în jurul mediei



- Intervalul de valori $[72 - 5; 72 + 5] = [67, 77]$ conține aproximativ **~68,3 %** din valorile distribuției.
- Intervalul de valori $[72 - 2 \cdot 5; 72 + 2 \cdot 5] = [62, 82]$ conține aproximativ **~95,5 %** din valorile distribuției.
- Intervalul de valori $[72 - 1,96 \cdot 5; 72 + 1,96 \cdot 5] = [62,2, 81,8]$ conține **95 %** din valorile distribuției.
- Intervalul de valori $[72 - 3 \cdot 5; 72 + 3 \cdot 5] = [57, 87]$ conține aproximativ **~99,7 %** din valorile distribuției.

Curba Normală

- Distribuție simetrică în jurul mediei



- Intervalul de valori $[\mu - \sigma; \mu + \sigma]$ conține aproximativ **~68,3 %** din valorile distribuției.
- Intervalul de valori $[\mu - 2\sigma; \mu + 2\sigma]$ conține aproximativ **~95,5 %** din valorile distribuției.
- Intervalul de valori $[\mu - 1,96\sigma; \mu + 1,96\sigma]$ conține **95 %** din valorile distribuției.
- Intervalul de valori $[\mu - 3\sigma; \mu + 3\sigma]$ conține aproximativ **~99,7 %** din valorile distribuției.

Coeficientul de variație

- Măsură relativă a dispersiei datelor
- Formula de calcul:

$$CV(\%) = \frac{s}{m} \times 100$$

Reguli empirice de interpretare a coeficientului de variație [100]:

- $CV < 10\%$ → populația poate fi considerată omogenă
 - $10\% \leq CV < 20\%$ → populația poate fi considerată relativ omogenă
 - $20\% \leq CV < 30\%$ → populația poate fi considerată relativ eterogenă
 - $CV \geq 30\%$ → populația poate fi considerată eterogenă
- Evaluare a abaterii standard în raport cu valoarea medie
 - Are avantajul de a fi un indicator independent de unitățile de măsură
 - Se poate exprima și/sau procentual

Coeficientul de variație

- Măsură a variabilității relative utilizată pentru:
 - Măsurarea modificărilor care au apărut în populație în timp
 - Compararea variabilității a două populații când unitățile de măsură sunt diferite (mg/dL vs mmol/L – colesterol)
 - Frecvent exprimat procental

	Greutate (kg)	Înălțime (cm)
Media aritmetică	72,6	168
Deviația standard	13,6	10,2

- » Care din variabilele de mai sus are împrăștierea mai mare?
 - > Nu se poate răspunde la întrebare
- » Care din variabilele de mai sus are împrăștierea **relativă la medie** mai mare ?

Greutate: $CV = 13,6/72,6 \cdot 100 = 19\%$

Înălțime: $CV = 10,2/168 \cdot 100 = 6,1\%$

Dispersia vs. dispersia relativă

- Aplicabilitate: dispersia în două seturi de date

- *A*: 12, 13, 16, 18, 18, 20
- *B*: 120, 130, 160, 180, 180, 200

	Grup A	Grup B
Media	16	162
Deviația standard	3	29
CV (%)	18	18

- Aplicabilitate: dispersia în două seturi de date

- *A*: 12, 13, 16, 18, 18, 20
- *B*: 2, 3, 160, 18, 200, 300

	Grup A	Grup B
Media	16	114
Deviația standard	3	114
CV (%)	18	100

Eroarea standard

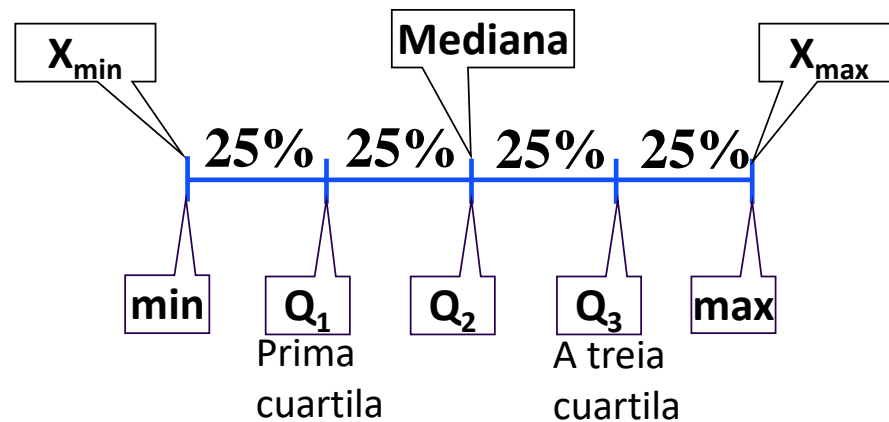
- este utilizată pentru a măsura precizia statistică a unei estimări
- este utilizată în statistica inferențială (teste de ipoteze și estimarea intervalului de încredere)

$$ES = \frac{s}{\sqrt{n}}$$

unde S= deviația standard de eșantionare-vezi cursul viitor-intervale de încredere

n=talie eșantion

Măsuri de localizare: Cuartilele și Intervalul intercuartilic (IQR)



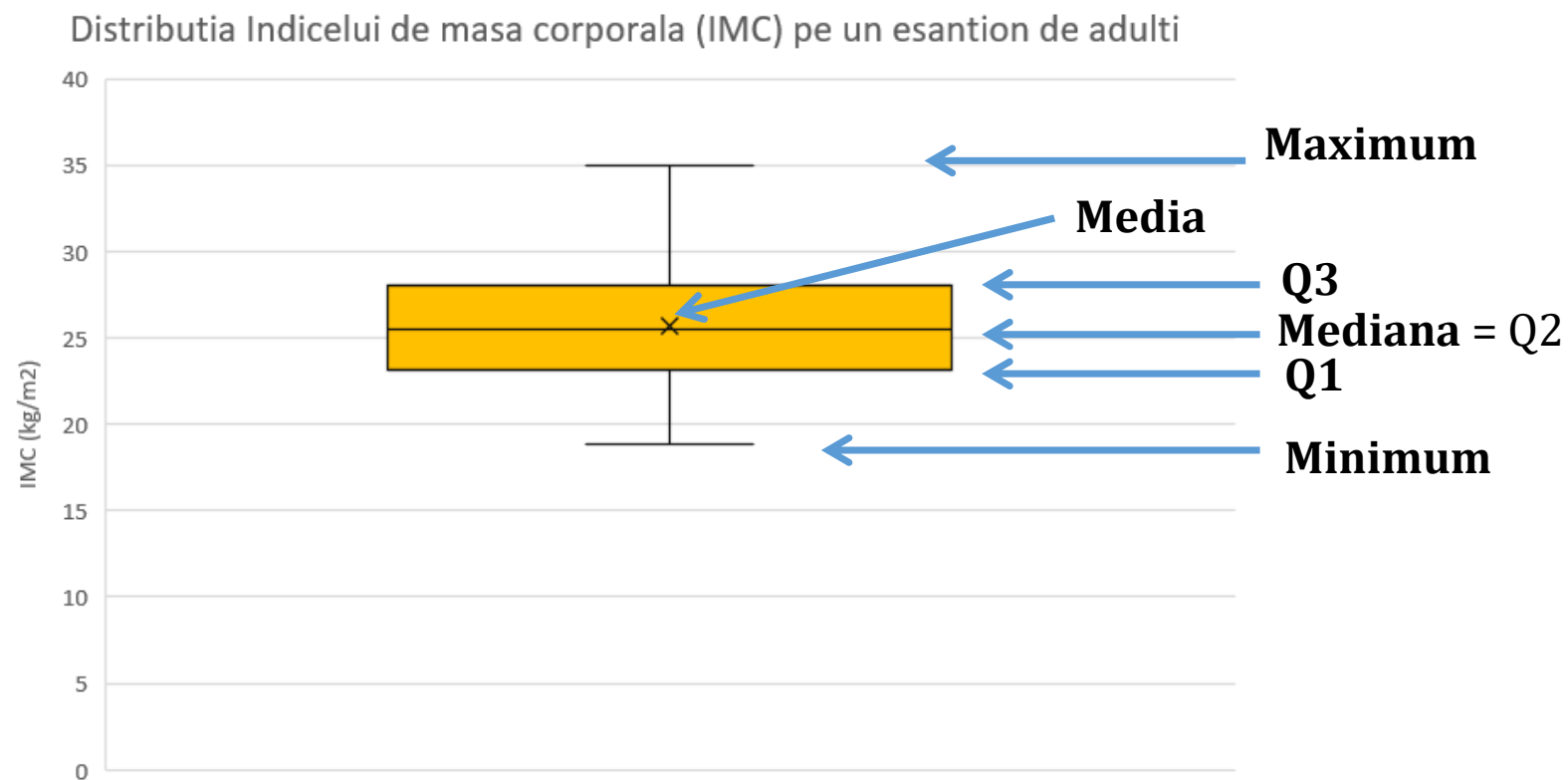
Măsură a dispersiei pentru 50%
din datele de mijloc.

$IQR = Q_3 - Q_1$, unde Q_3 = cvartila 3 (percentila 75%), Q_1 = cvartila 1 (percentila de 25%)

Sau

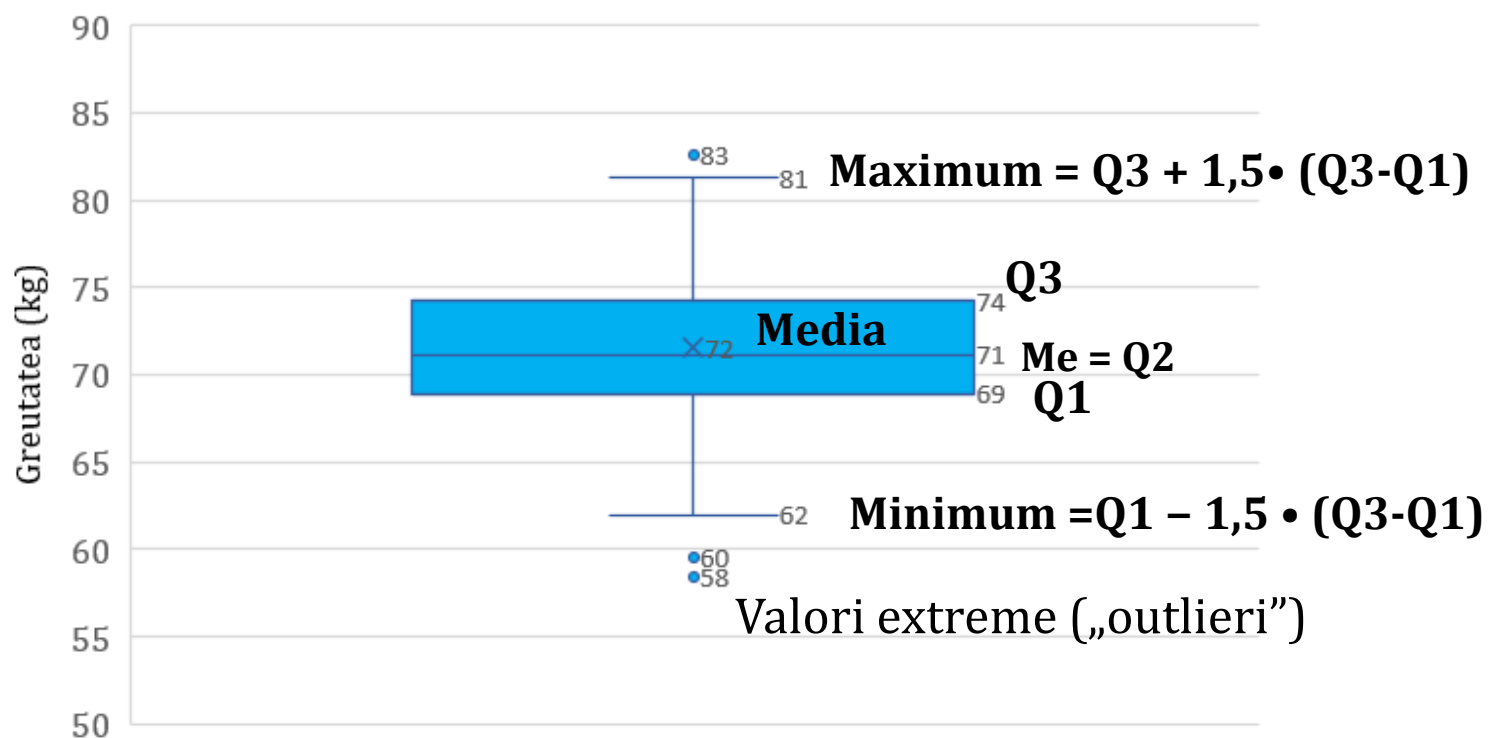
$IQR: [Q_1; Q_3]$

Reprezentare grafică: box-plot



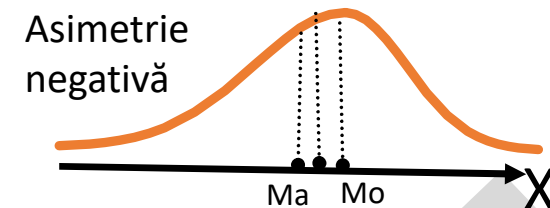
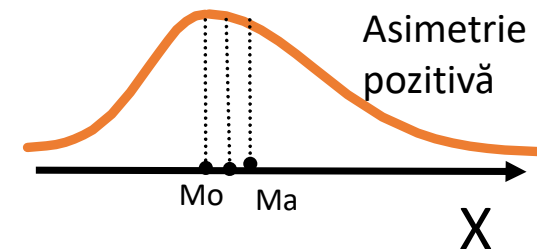
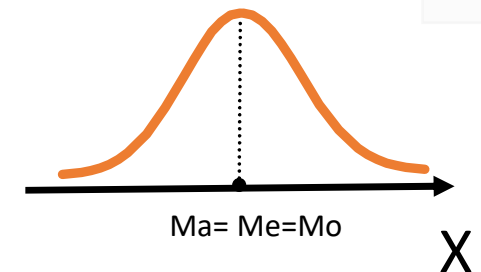
Reprezentare grafică: box-plot

Distribuția Greutății pe un eșantion de adulți



Măsuri ale simetriei: Coeficientul de asimetrie

- Coeficientul de asimetrie (notație: α_3): Măsoară gradul de asimetrie al unei distribuții de date
- **Semnul acestuia indică direcția asimetriei: pozitivă sau negativă**
 - ✓ $\alpha_3 \approx 0 \Rightarrow$ distribuție simetrică.
 - ✓ $\alpha_3 > 0 \Rightarrow$ distribuție alungită la dreapta– asimetrie pozitivă
 - ✓ $\alpha_3 < 0 \Rightarrow$ distribuție alungită la stanga– asimetrie negativă
- **Valoarea acestuia indică gradul asimetrie:**
 - ✓ $\alpha_3 \in [-0,5; 0,5]$ distribuție aproximativ simetrică
 - ✓ $\alpha_3 \in [-1; -0,5]$ sau $[0,5; 1]$ distribuție cu asimetrie moderată
 - ✓ $\alpha_3 < -1$ sau > 1 distribuție cu asimetrie importantă



Măsuri ale simetriei și boltirii: Coeficientul de boltire

- **Coeficientul de boltire** (notație: α_4): Măsoară gradul de aplatizare al unei distribuții de date în raport cu distribuția normală (Gaussiană): proprietatea distribuției de a fi mai ascuțită sau mai aplatizată decât curba normală.

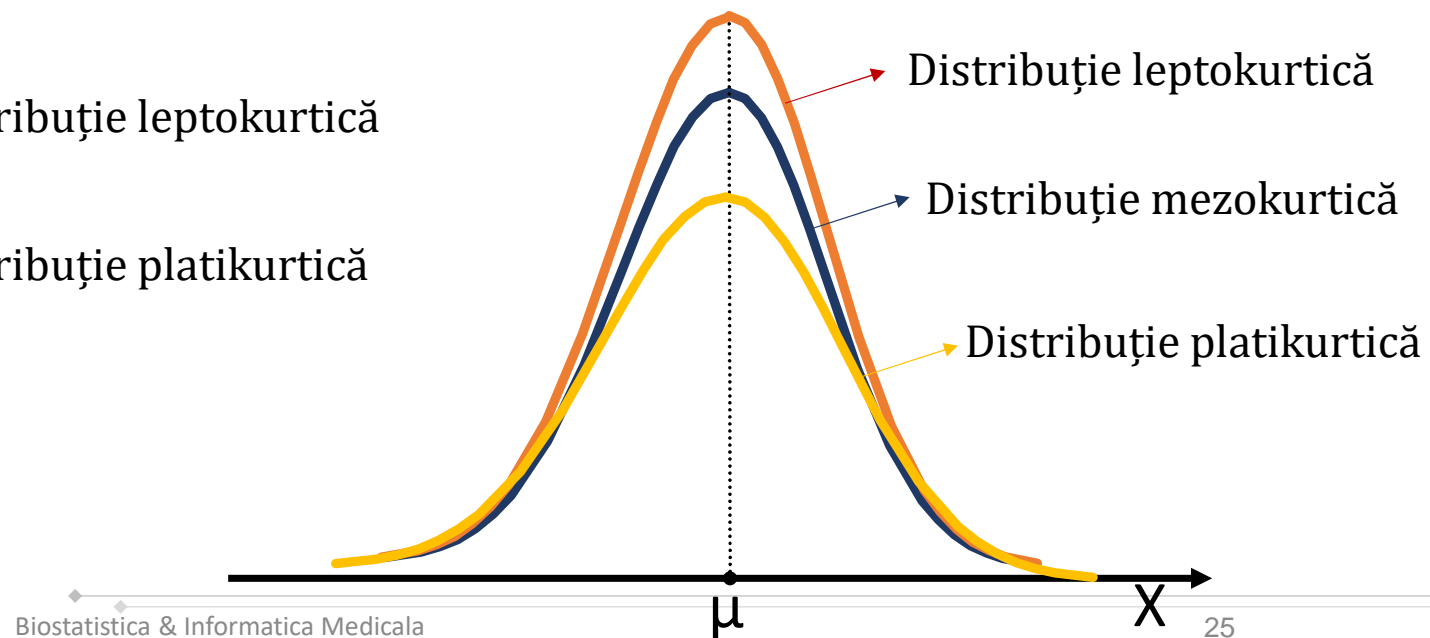
- Dacă:

- ✓ $\alpha_4 \approx 0 \Rightarrow$ distribuție mezokurtică (normală)

- ✓ $\alpha_4 > 0 \Rightarrow$ distribuție leptokurtică

- ✓ $\alpha_4 < 0 \Rightarrow$ distribuție platikurtică

-

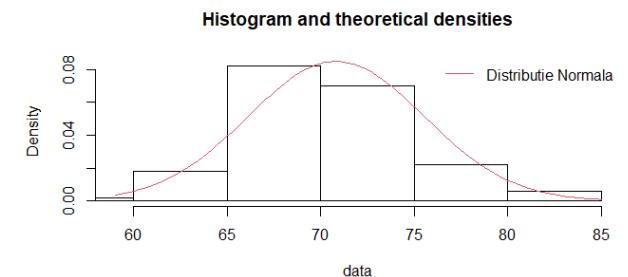
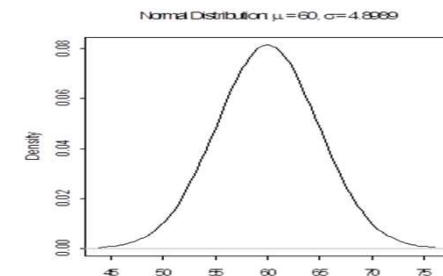


Verificarea condiției de normalitate a datelor (existența distribuției Normale)

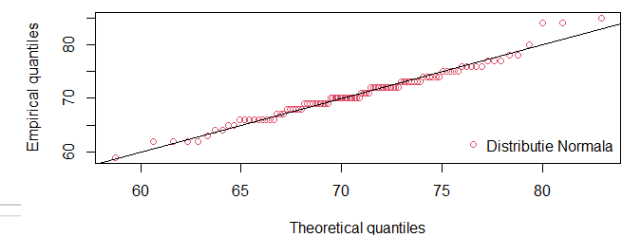
- **Datele unei variabile cantitative sunt considerate a avea o distribuție Normală (de probabilitate) dacă frecvențele valorilor acestei variabile sunt similare cu frecvențele teoretice generate de o funcție matematică (legea Normală).**

- **Utilitate:**

- ✓ Important pentru alegerea statisticilor descriptive
- ✓ De exemplu: pentru datele normale distribuite folosim media și abaterea standard, pentru datele distribuite nenormale folosim mediana și quartilele
- ✓ Important pentru alegerea tehnicilor statistice analitice/inferențiale

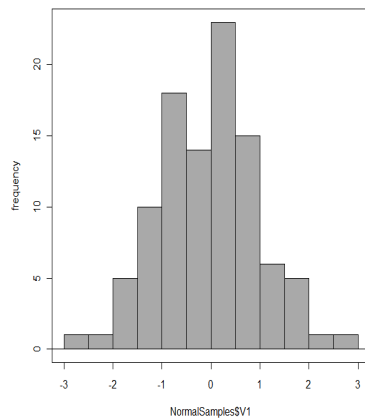


Grafic de cuantile Q-Q plot

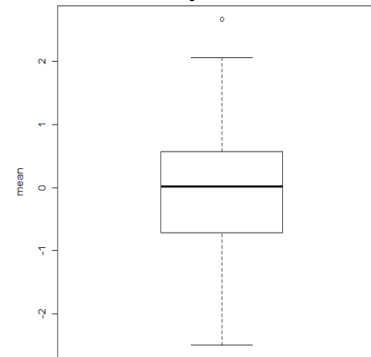


Date normal distribuite versus Date care au deviații de la normalitate

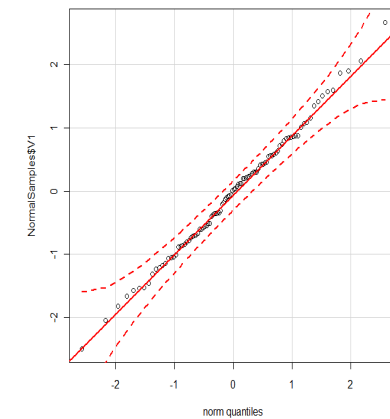
Histograma



Box plot

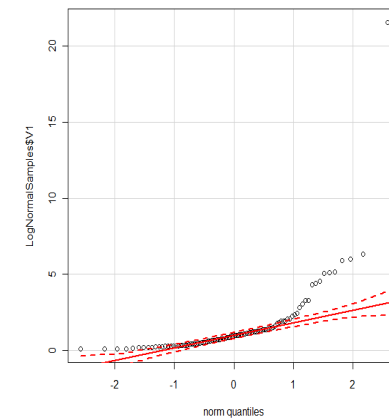
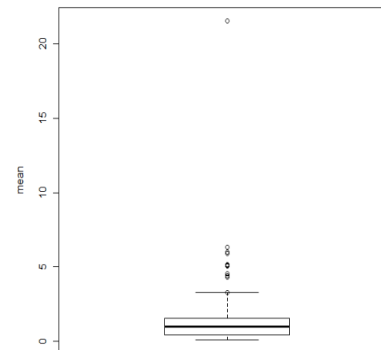
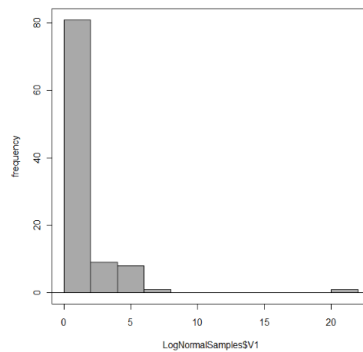


Grafic de cuantile



Date normal distribuite

Date care NU sunt normal distribuite



Verificarea condiției de normalitate a datelor (existența distribuției Normale)

Metode de verificare:

1. modalități grafice (cele mai bune)

- ≈ Histograma (simetrică, ca o pălărie)
- ≈ Boxplot (simetric în jurul mediei)
- ≈ **Density plot**
- ≈ **Graficul cuantilelor Q-Q plot**

2. statistici descriptive (nu foarte fiabile)

- ✓ Dacă media este \sim = mediană
- ✓ Dacă coeficientul de aplatizare/boltire \sim = 0 / aparține intervalului $[-1, 1]$ (kurtoză)
- ✓ Dacă coeficientul de simetrie \sim = 0 / aparține intervalului $[-1, 1]$ (asimetrie)

3. Teste de normalitate: testul Shapiro-Wilk

De reținut....

Amplitudinea	<ul style="list-style-type: none"> • se calculează pentru Variabile cantitative discrete sau continue • Slab informativă deoarece nu ne spune nimic cu privire la distribuția datelor
Deviația standard	<ul style="list-style-type: none"> • se calculează pentru Variabile cantitative • Cel mai frecvent raportată relativ la medie • Informații cu privire la cât de departe sunt datele față valoare centrală (cel mai frecvent media aritmetică – date normal distribuite)
Cvartilele	<ul style="list-style-type: none"> • se calculează pentru Variabile cantitative si calitative ordinale • Utile ca mărimi ale dispersiei dacă datele nu urmează distribuția normală
Coeficientul de variație	<ul style="list-style-type: none"> • se calculează pentru Variabile cantitative • Mărime a variabilității relative • Permite compararea variabilității pe două eșantioane diferite

STATISTICI DESCRIPTIVE: DATE CALITATIVE

RAPORTUL / PROPORTIA / RATA



Raportul

- Numere raționale pozitive a și b , $b \neq 0$
- **Raportul:** a/b
- Simbolică
 - $a:b$
 - a/b
- Numitorul nu include în mod obligatoriu subiecții numărătorului
- Într-o grupă de 12 studenți avem 4 fumători.
 - Raportul nefumători/fumători = $8/4 = 2/1 = 2 \Rightarrow$ adică la 2 nefumători există un fumător
 - Raportul fumători/nefumători = $4/8 = 0,5$

Proporția

» O **proporție** este un raport în care numărătorul face parte din numitor. Astfel o proporție are forma generală:

$$a/(a+b)$$

» Ia valori între:

> 0 și 1

> 0 și 100 dacă se exprimă procentual

» Toți indivizii de la numărător sunt incluși la numitor

» Frecvența relativă este o proporție (Prevalența unei boli este o proporție)

» Într-o grupă de 12 studenți avem 4 fumători.

» Proporția fumătorilor = $4/12=0,33$

» Proporția nefumătorilor = $8/12$

Proporția: exemplu

- La serviciul de urgențe ale unui spital județean s-au prezentat pentru consultație 1200 pacienți. Dintre aceștia 420 au fost internați (200 femei și 220 bărbați):
 - » Proporția pacienților internați = $420/1200 \cdot 100 = 35\%$
 - » Proporția pacienților de sex feminin internați = $200/420 \cdot 100 = 48\%$
 - » Proporția pacienților de sex masculin internați = $220/420 \cdot 100 = 52\%$

Prevalența

- = proporția de indivizi dintr-o populație care au boala la un moment dat
- estimează probabilitatea ca un individ să aibă boala la un moment dat
- Formula:

$$\text{Prevalența} = (\text{numărul de cazuri de boală}) / (\text{total populație})$$

Rata

- O rată calculată reflectă riscul de a surveni în timp un anumit eveniment.
- Ia valori de la 0 la infinit
- Număr de indivizi raportat la unitatea de timp (oră / zi / săptămână / lună / an etc.)
 - Rate de morbiditate
 - Rate de mortalitate
 - Rate de natalitate

Rata

- Într-un oraș cu populație de 100000 locuitori s-au înregistrat 200 născuți vii în anul 1999
- Rata de natalitate = $200/100000 * 1000 = 2$ nou născuți la mia de locuitori
- Rata de fertilitate = $(\text{nr. nașteri})/(\text{nr. femei cu vârsta între 15-45 ani}) * 1000$
- Rata de morbiditate = frecvența cazurilor de îmbolnăvire pentru o populație specificată într-o perioadă de timp
- Rata de mortalitate = $(\text{nr. decese})/(\text{populație}) * 1000$

Exemple de probleme pentru examen

E1. Valorile tensiunii arteriale sistolice (mmHg) măsurate pe un eșantion de 10 pacienti sunt următoarele: 120, 100, 110, 120, 130, 160, 130, 120, 140, 160.

Media aritmetica, mediana, modulul si valoarea centrala sunt următoarele:

- A. 129 – 125 – 120 – 130
- B. 130 – 125 – 130 – 125
- C. 120 – 130 – 120 – 125
- D. 129 – 130 – 120 – 130
- E. 125 – 125 – 120 – 130

R1: A

Rezolvare-> Media aritmetică(Ma) =? $Ma = (120+100+110+120+130+160+130+120+140+160) / 10 = 129 \text{ mmHg}$

-> Mediana (Me)=? Seria statistica ordonata: 100, 110, 120,120,120, 130,130,140,160,160

Talie esantion: numar par = 10=> $\frac{X_{\frac{10}{2}}+X_{\frac{10}{2}+1}}{2} = \frac{X_5+X_6}{2}=125$

-> Modulul (Mo)=? Valoarea din seria statistică a TAS cu cea mai mare frecvență (cea mai repetitiva valoare)=120

-> Valoare centrala=? $\frac{x_{max}+x_{min}}{2}=130$

Exemple de probleme pentru examen

E2. A fost evaluat un eșantion de studenți AMG1 cuprinzând 40 de băieți și 40 de fete. Pentru fiecare student s-au colectat înălțimea și greutatea și s-a calculat indicele de masă corporală (IMC). Valoarea medie și deviația standard ale IMC au fost de 20 kg/m^2 respectiv 4 kg/m^2 . Din următoarele afirmații, selectați cele pe care le considerați adevărate:

- A. Volumul eșantionului în acest studiu a fost de 40.
- B. Dacă valorile IMC au avut o distribuție simetrică atunci jumătate dintre studenți au avut valorile $\text{IMC} \leq 20 \text{ kg/m}^2$.
- C. Coeficientul de variație a fost egal cu 20%.
- D. IMC este relativ eterogen.
- E. Coeficientul de variație a fost independent de unitățile de măsură

R2. B, C, D, E

Exemple de probleme pentru examen

E4. Valorile tensiunii arteriale sistolice (mmHg) măsurate pe un eșantion de 10 pacienți sunt următoarele: 120, 100, 110, 120, 130, 160, 130, 120, 140, 160. Se cunoaște ca prima cvartilă (Q1) este egală cu 120, a doua cvartilă (Q2) este egală cu 125 și a treia cvartilă (Q3) este egală cu 137.5. Care dintre următoarele afirmații sunt corecte

- A. $Q2 - Q1 = 5$
- B. $Q3 - Q2 = 12$
- C. $Q3 - Q2 = 12,5$
- D. Distribuția este aproximativ simetrică
- E. Distribuția este simetrică

R4: A, C, E

Rezolvare: -> Într-o distribuție simetrică distanța dintre cvartilele (Q1 și Q3) și mediana (Q2) este aceeași => $Q2 - Q1 = Q3 - Q2$; în cazul nostru, cei doi termeni ($Q3 - Q2$ și $Q2 - Q1$) nu sunt egali deci seria statistică a valorilor TAS nu este simetrică.

Exemple de probleme pentru examen

E5. Pentru o variabila cu distribuție simetrică:

- A. Media (aritmetică) este egală cu prima cvartilă
- B. 50% din date sunt superioare mediei
- C. $Q3 - Q2 = Q2 - Q1$
- D. $Q3 - Q2 = 3 \cdot (Me - Q1)$
- E. Modul \approx Medie \approx Mediana

R5: B, C, E

MULȚUMESC PENTRU ATENȚIE!

