

Culegerea și stocarea datelor Analiza datelor I

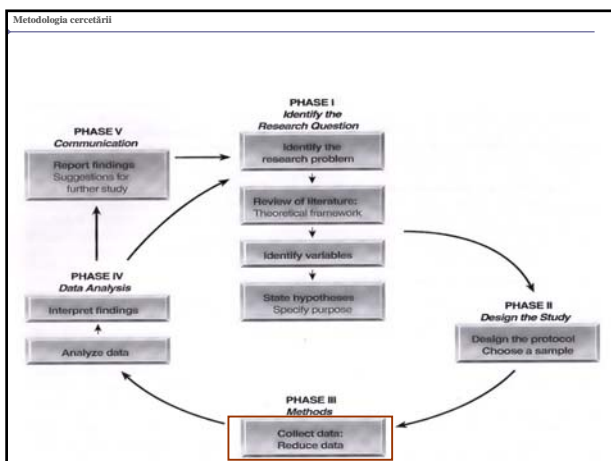
"To consult the statistician after an experiment is finished is often merely to ask him to conduct a *post mortem* examination. He can perhaps say what the experiment died of."

Presidential Address to the First Indian Statistical Congress, 1938



Cuprins

- Culegerea și stocarea datelor:
 - Baze de date Excel
 - Tipuri de date și formate
- Analiza datelor I



Culegerea și stocarea datelor

Scop:

- Organizarea datelor în formatul care să permită sumarizarea și prelucrarea acestora

Aplicabilitate pentru teză:

- Stocarea în format electronic va permite sumarizarea și prelucrarea statistică



Culegerea datelor

- Constituie un element esențial al cercetării
- De realizează pe fișe de culegere a datelor:
 - pe suport de hârtie
 - electronic
- Pentru prelucrarea datelor e necesar suportul electronic
- Conținutul fișei de culegere a datelor:
 - Pentru fiecare pacient toate caracteristicile urmărite
 - exemplu: numele, vârsta, greutatea, înălțimea, tensiunea arterială sistolică, ...



Gestiunea datelor medicale cu Microsoft Excel

- Pachet de programe destinat:
 - tratării datelor în formă tabelară
 - prelucrării statistice și reprezentării grafice a informației conținută în aceste tabele
 - este dotat cu o funcție pentru tratarea datelor tabelii ca baza de date



Baza de date Excel

- Regiune compactă de date care ocupă ca suprafață cel puțin două rânduri (denumite articole) și două coloane (denumite câmpuri) adiacente, coloane consecutive și rânduri ne-consecutive sau rânduri consecutive și coloane ne-consecutive
- Prima linie dintr-o astfel de regiune poartă numele de antet și conține denumiri de câmpuri



Baza de date Excel

- Dimensiunea maximă a unei baze de date Excel este dată de:
 - dimensiunea unei foi de calcul: 256 coloane × 65536 rânduri
 - numărul maxim de foi de calcul: 256



CATEGORIE	TOTAL	Iul	AUG	SEPT	OCT	NOV	DEC
1. Incidențe	100	10	10	10	10	10	10
2. Patru	100	10	10	10	10	10	10
3. Patru	100	10	10	10	10	10	10
4. Patru	100	10	10	10	10	10	10
5. Patru	100	10	10	10	10	10	10
6. Patru	100	10	10	10	10	10	10
7. Patru	100	10	10	10	10	10	10
8. Patru	100	10	10	10	10	10	10
9. Patru	100	10	10	10	10	10	10
10. Patru	100	10	10	10	10	10	10
11. Patru	100	10	10	10	10	10	10
12. Patru	100	10	10	10	10	10	10
13. Patru	100	10	10	10	10	10	10
14. Patru	100	10	10	10	10	10	10
15. Patru	100	10	10	10	10	10	10
16. Patru	100	10	10	10	10	10	10
17. Patru	100	10	10	10	10	10	10
18. Patru	100	10	10	10	10	10	10
19. Patru	100	10	10	10	10	10	10
20. Patru	100	10	10	10	10	10	10
21. Patru	100	10	10	10	10	10	10
22. Patru	100	10	10	10	10	10	10
23. Patru	100	10	10	10	10	10	10
24. Patru	100	10	10	10	10	10	10
25. Patru	100	10	10	10	10	10	10



Manipularea foilor de calcul

- Clic dreapta de mouse pe foaia de calcul și:

[Insert]	adăugarea unei noi foi de calcul
[Delete]	ștergerea foii de calcul selectată
[Rename]	schimbarea denumirii foii de calcul
[Move or Copy]	schimbarea ordinii foii de calcul selectate sau copierea acestora
[Select All Sheets]	selectarea tuturor foilor de calcul
[Format - Sheet - Background]	definirea background-ului pe pagina de lucru
[Format - Sheet - Hide]	ascunderea foii de calcul selectate



Editarea foilor de calcul

- Inserare unui rând sau a unei coloane:
 - clic dreapta de mouse pe eticheta de rând sau coloană și activarea opțiunii [Insert]
- Ștergere:
 - clic dreapta de mouse pe eticheta de rând sau coloană și activarea opțiunii [Delete]
- Definirea antetelor și a notelor de subsol: [View - Header and Footer] sau [View - PageSetup] – vizibile la Print Preview și la printare



Editarea foilor de calcul

- Selectarea:
 - unei celule: clic stânga pe celulă
 - unei coloane: clic stânga pe eticheta coloanei (idem și pentru selectarea unui rând, de această dată clic stânga pe eticheta rândului)
 - unui domeniu aleator de celule: clic stânga pe o celulă, simultan activarea tastei CTRL în combinație cu clic stânga celulele dorite a fi selectate



Editarea foilor de calcul

- Căutarea și înlocuirea datelor din celule:
 - [Edit - Find] pentru căutare
 - [Edit - Find - Replace] pentru înlocuire
- Copiere: [Edit - Copy - Paste]
- Mutarea: [Edit - Cut - Paste]
- Ștergerea: [Edit - Cut]

Atenție! Este necesară selectarea prealabilă a datelor!



[Edit - Paste Speciale...]

[Edit - Paste Speciale...]	[Edit - Paste Speciale...]
Conținutul și toate proprietățile acestuia	All
Formulele de calcul	Formulas
Valorile selecției fără formulele care au generat aceste valori	Values
Proprietățile de formatare	Formats
Comentariile anexate celulelor selectate	Comments
Regulile de validare a datelor din selecție	Validation
Total cu excepția formataților de chenare	All except borders
Doar formatațiile de lățime a coloanelor	Column widths
Formulele din celulele selectate și formatarea datelor	Formulas and number formats
Valorile din celulele selectate și formatarea datelor	Values and number formats



Formatare celule: [Format - Cells ...]

Denumire	Tip	Observații
Number	numeric	aliniat implicit la dreapta
Date	calendaristic	aliniat implicit la dreapta
Time	timp	aliniat implicit la dreapta
Text	non numeric	aliniat implicit la stânga
Scientific	1.0E-01	aliniat implicit la

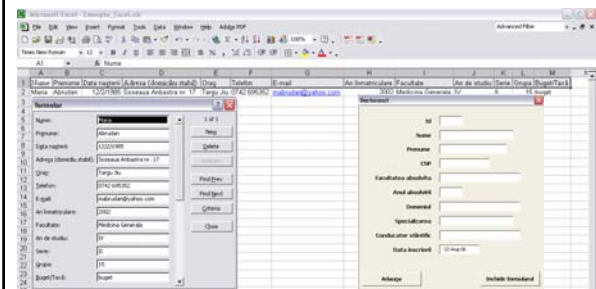
- Transformarea unui grup de celule într-o singură celulă:
[Format - Cells ... - Alignment - Merge cells]



Formatare celule: [Format - Conditional Formatting...]



Crearea unei baze de date Excel - Formulare



Validarea datelor dintr-o bază de date Excel

- [Data - Validation...]:
 - crearea unei liste de opțiuni
 - limitarea datelor de intrare la un anumit tip sau de o anumită mărime



Metodologia cercetării

Validarea datelor dintr-o bază de date Excel

- Criteriile de validare:

Criteriul	Semnificație
Whole Number	În funcție de operatorul ales: incluziune (minim - maxim), excluziune (minim - maxim), egal cu, diferit de, mai mare decât, mai mic decât, mai mare sau egal cu, mai mic sau egal cu
Decimal	Includerea sau excluziunea unui șir de numere care îndeplinesc condiția operatorului (vezi operatorii anteriori)
List	Includerea unei liste de opțiuni sau legătura cu o listă de opțiuni creată anterior
Date	Includerea sau excluziunea de date tip calendaristic (dd/mm/yyyy) în funcție de operatorul ales
Time	Includerea sau excluziunea de date tip timp (hh:mm:ss AM/PM) în funcție de operatorul ales
Text Length	Specificarea lungimii textului prin folosirea operatorilor
Custom	Folosirea formulelor de validare create de către utilizator

Sorana D. BOLBOACĂ Curs 4: Culegerea și stocarea datelor & Analiza datelor I.

Metodologia cercetării

Operațiunea de filtrare

- Afișarea conținutului care satisface criteriile de filtrare impuse de utilizator, articolele care nu satisfac criteriile fiind ascunse
- Tipuri:
 - auto-filtrarea [Data - Filter - AutoFilter]
 - filtrarea avansată [Data - Filter - Advanced Filter...]

Sorana D. BOLBOACĂ Curs 4: Culegerea și stocarea datelor & Analiza datelor I.

Metodologia cercetării

Operațiunea de filtrare

Sorana D. BOLBOACĂ Curs 4: Culegerea și stocarea datelor & Analiza datelor I.

Metodologia cercetării

Operatori în auto-filtrare

Operator	Denumire variabilă	Criteriu	Ce localizează în baza de date? înregistrări în care
egal	SEX	M	sexul este masculin (M)
ne-egal	ALCOOL	FALSE	variabila Alcool are valoarea TRUE
mai mare	VARSTA	40	persoanele au vârsta mai mare de 40 de ani
mai mare sau egal	DATA ANGAJĂRII	2/11/1987	data angajării este 2 octombrie 1987 (inclusiv)
mai mic	IMC	21	indicele de masă corporală are valoare mai mică de 21
mai mic sau egal	SALAR	1000	valorile din câmpul salar sunt mai mici de 1000
începe cu	TIPUL LAPTELUI	mi	valorile corespunzătoare câmpului TIPUL LAPTELUI încep cu mi
nu începe cu	TIPUL LAPTELUI	mi	valorile corespunzătoare câmpului TIPUL LAPTELUI nu încep cu mi
se termină cu	ORAȘ	na	valorile corespunzătoare câmpului ORAȘ care se termină cu na
nu se termină cu	ORAȘ	na	valorile corespunzătoare câmpului ORAȘ care nu se termină cu na
conține	ETNIE	ag	valorile corespunzătoare câmpului ETNIE care conțin literele ag
nu conține	ETNIE	ag	valorile corespunzătoare câmpului ETNIE care nu conțin literele ag

Sorana D. BOLBOACĂ Curs 4: Culegerea și stocarea datelor & Analiza datelor I.

Metodologia cercetării

Aplicații biomedicale

- Management medical
- Gestiunea datelor medicale
- Gestiunea în domeniul medical
 - Salarizare
 - Contabilitate
 - Facturarea serviciilor medicale

Sorana D. BOLBOACĂ Curs 4: Culegerea și stocarea datelor & Analiza datelor I.

Metodologia cercetării

De reținut!

Gestiunea datelor cu Microsoft Excel

- Ieftin și ușor de implementat!
- Pentru utilizare e nevoie de abilități minime de lucru cu calculatorul.
- Permite organizarea datelor după placul utilizatorului.
- Permite importul de date în alte programe:
 - De gestiune a informațiilor: Microsoft Access
 - De prelucrare statistică: EpiInfo

Sorana D. BOLBOACĂ Curs 4: Culegerea și stocarea datelor & Analiza datelor I.

Culegerea datelor

- Fișa de culegere a datelor în format electronic:
 - Fișier Microsoft Excel
 - Primul rând conține obligatoriu denumirea variabilei și unitatea de măsură (dacă este cazul)
 - Abrevierile folosite se introduc ca și comentariu la denumirea variabilei
 - În fiecare rând se introduc datele corespunzătoare unui singur pacient
 - Dacă nu avem toate datele specificăm în celulele corespunzătoare "lipsă"



Culegerea datelor

- Fișa de culegere a datelor în format electronic:
 - Se culeg **doar** date primare (măsurate sau observate)
 - Nu se culeg date derivate din date primare
 - Exemplu:
 - Dată primară = înălțime, greutate
 - Dată secundară = indicele de masă corporală
 - Exemplu:
 - Dată primară = tensiunea arterială sistolică, tensiunea arterială diastolică
 - Dată secundară = status hipertensiv – hipotensiv - normotensiv



De reținut! Culegerea datelor

- Culegerea datelor trebuie să se realizeze în conformitate cu protocolul de cercetare.
- Pentru sumarizarea și analiza statistică a datelor este necesară stocarea acestora în format electronic.
- Se culeg doar date primare! Datele secundare se obțin în urma prelucrării datelor primare.



Analiza datelor

Scop:

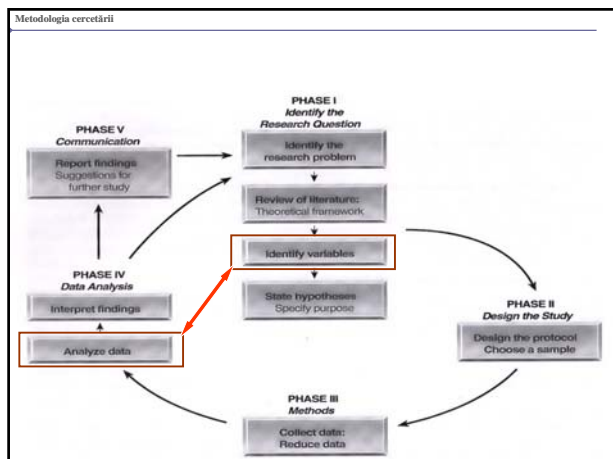
- Sumarizarea datelor
- Verificarea ipotezelor de cercetare

Finalitate:

- Generarea de noi ipoteze de cercetare

Aplicabilitate pentru teză:

- Capitolul Rezultate



De ce analiza statistică?

- 2 scopuri:
 - Descriptiv (statistica descriptivă):
 - Modalități de sumarizare a caracteristicilor importante ale unui set de date medicale
 - Inferențial (statistica inferențială):
 - Cum (și când) generalizăm rezultatele obținute pe un eșantion la populația generală



Metodologia cercetării

"There are three kinds of lies: lies, damned lies, and statistics."
Benjamin Disraeli

- Popularizată în SUA de Mark Twain
 - "...statement refers to the persuasive power of numbers, the use of statistics to bolster weak arguments, and the tendency of people to disparage statistics that do not support their positions..."

Sorana D. BOLBOACĂ Curs 4: Culegerea și stocarea datelor & Analiza datelor I.

Metodologia cercetării

SUMARIZAREA TABELARĂ ȘI REPREZENTAREA GRAFICĂ A DATELOR

Statistica descriptivă

Sorana D. BOLBOACĂ Curs 4: Culegerea și stocarea datelor & Analiza datelor I.

Metodologia cercetării

Conținut

- Principii de sumarizare tabelară
- Principii de reprezentare grafică
- Sumarizarea tabelară și/sau reprezentarea grafică a datelor:
 - Atribut (calitative): o variabilă
 - Atribut (calitative): două variabile
 - Numerice (cantitative): o variabilă
 - Numerice (cantitative): două variabile

Sorana D. BOLBOACĂ Curs 4: Culegerea și stocarea datelor & Analiza datelor I.

Metodologia cercetării

Principii de sumarizare tabelară

1. Simple: de preferat 2/3 tabele mai mici în loc de unul încărcat
2. Informative prin ele însele
 - Abrevieri sau simboluri explicate la subsolul tabelului
 - Etichete de rând și coloană
 - Unități de măsură
 - Titlul: ce? când? Unde?
 - Linii și/sau coloane de sinteză (total)
3. Dacă datele nu sunt originale trebuie să se menționeze sursa lor într-o notă de subsol

Sorana D. BOLBOACĂ Curs 4: Culegerea și stocarea datelor & Analiza datelor I.

Metodologia cercetării

Principii de reprezentare grafică

- Orice reprezentare grafică trebuie să aibă:
 - Titlul
 - Definierea axelor
 - Unități de măsură pentru fiecare axă (dacă este cazul)
 - Legendă (dacă este cazul)
- O reprezentare grafică trebuie să se "înțeleagă" singură!
 - Fără a se citi textul!!!

Sorana D. BOLBOACĂ Curs 4: Culegerea și stocarea datelor & Analiza datelor I.

Metodologia cercetării

Principii de reprezentare grafică

- Scopul unei reprezentări grafice este de a transmite o informație
- Când construim o reprezentare grafică trebuie să răspundem la întrebarea: Care este scopul acestei reprezentări?
- Datele trebuie reprezentate grafic în așa fel încât să fie utile în înțelegerea fenomenului clinic
- Atenție la compoziția culorilor (nu puneți fundaluri colorate) și la dimensiunea caracterelor!

Sorana D. BOLBOACĂ Curs 4: Culegerea și stocarea datelor & Analiza datelor I.

Variabile calitative: 1 variabilă Tabelul de frecvență

- Se ordonează datele crescător
- Se determine frecvența fiecărei valori
- Se includ valorile distincte și frecvențele într-un tabel pe două coloane:
 - Frecvența absolută (numărul de cazuri care îndeplinesc criteriul)
 - Frecvența relativă = raportul dintre frecvența absolută și volumul eșantionului/populației (simbol = n). Valorile se pot prezenta și procentual.



Variabile calitative: 1 variabilă Tabelul de frecvență

- Se pot alcătui tabele de frecvențe cu mai multe coloane care să cuprindă:
 - frecvențe absolute
 - frecvențe absolute cumulate crescător / descrescător
 - frecvențe relative
 - frecvențe relative cumulate crescător / descrescător
- Microsoft Excel:
 - funcția COUNTIF
 - Tabele Pivot [Data - Pivot Table and Pivot Chart Report ...]



Variabile calitative: 1 variabilă Tabelul de frecvență

Frecvența absolută

Frecvența relativă

Diagnostic	Nr. persoane	Procent (%)
Asfixia la naștere	527	26,1
Traumatisme obstetricale	92	4,6
Stare septică	7	0,3
Pneumonie	181	9,0
Diaree	8	0,4
Malformații congenitale	598	29,6
Alte cauze	606	30,0
Total	2019	100



Variabile calitative: 1 variabilă Tabelul de frecvență

Suma frecvențelor relative ale tuturor valorilor seriei care sunt mai mici sau egale decât x

Suma frecvențelor absolute ale tuturor valorilor seriei care sunt mai mici sau egale decât x

Diagnostic	f_a	f_r	f_a cumulat ↑	f_r cumulat ↑
Asfixia la naștere	527	26.10	527	26.10
Traumatisme obstetricale	92	4.56	619	30.66
Stare septică	7	0.35	626	31.01
Pneumonie	181	8.96	807	39.97
Diaree	8	0.40	815	40.37
Malformații congenitale	598	29.62	1413	69.99
Alte cauze	606	30.01	2019	100
Total	2019	100		



Variabile calitative: 1 variabilă Tabelul de frecvență

- Pentru seria statistică 5, 6, 7, 7, 8, 8, 5, 7, 8, 7 cărei din valorile de mai jos îi corespunde frecvența relativă cumulată crescător de 0.7:
 - A. 8
 - B. 6
 - C. 5
 - D. 7

Nici un răspuns nu este corect



Variabile calitative: 1 variabilă Tabelul de frecvență

- Pentru seria statistică 5, 6, 7, 7, 8, 8, 5, 7, 8, 7 cărei din valorile de mai jos îi corespunde frecvența relativă cumulată crescător de 0.7?

Valoare	f_a	f_r	f_a cc	f_r cc
5	2	0,20	2	0,20
6	1	0,10	3	0,30
7	4	0,40	7	0,70
8	3	0,30	10	1
Total	10	1		



Metodologia cercetării

Variabile calitative: 2 variabile Tabelul de contingență

	TBC=da	TBC=nu	Total
sex=F	2	10	12
sex=M	24	54	78
Total	26	64	90

Sorana D. BOLBOACĂ Curs 4: Culegerea și stocarea datelor & Analiza datelor I.

Metodologia cercetării

Variabile calitative: n variabile Tabel de frecvență

Tabelul 1. Distribuția patologiilor pulmonare asociate silicozei

	BrC	BPOC	Emfizem	CPC	TBC	Total
silicoza grad I	12	20	0	0	14	46
silicoza grad I/II	1	5	1	1	1	9
silicoza grad II	3	7	1	1	7	19
silicoza grad II/III	0	1	0	0	0	1
silicoza grad III	0	3	0	0	4	7
Total	16	36	2	2	26	82

BrC = bronșită cronică; BPOC = bronho-pneumopatie cronică obstructivă;
CPC = cord pulmonar cronic; TBC = tuberculoză pulmonară

Sorana D. BOLBOACĂ Curs 4: Culegerea și stocarea datelor & Analiza datelor I.

Metodologia cercetării

Variabile cantitative: 1 variabilă Tabele pe clase de frecvență

Greutate (g)	f_a	f_r	f_r cumulată ↑
(2800 – 3200]	151	18,60	18,60
(3200 – 3400]	299	36,82	55,42
(3400 – 3600]	300	36,95	92,37
(3600 – 3800]	0	0,00	92,37
(3800 – 4000]	62	7,64	100
Total	812	100	

Sorana D. BOLBOACĂ Curs 4: Culegerea și stocarea datelor & Analiza datelor I.

Metodologia cercetării

Reprezentarea grafică: 1 variabilă Plăcinta (PIE)

- Variabile calitative sau cantitative. Dacă este cantitativă trebuie să fie clase de frecvențe.
- Se folosește pentru a reprezenta frecvențe absolute sau relative:
 - Vizualizarea prevalenței relative a unui fenomen de sănătate
- Datele se culeg ca frecvențe absolute

Sorana D. BOLBOACĂ Curs 4: Culegerea și stocarea datelor & Analiza datelor I.

Metodologia cercetării

Reprezentarea grafică: 1 variabilă Plăcinta (PIE)

Distribuția patologiei cardiovasculare

Categorie	Pondere (%)
boala cardiovasculara +	55%
boala cardiovasculara -	45%

Sorana D. BOLBOACĂ Curs 4: Culegerea și stocarea datelor & Analiza datelor I.

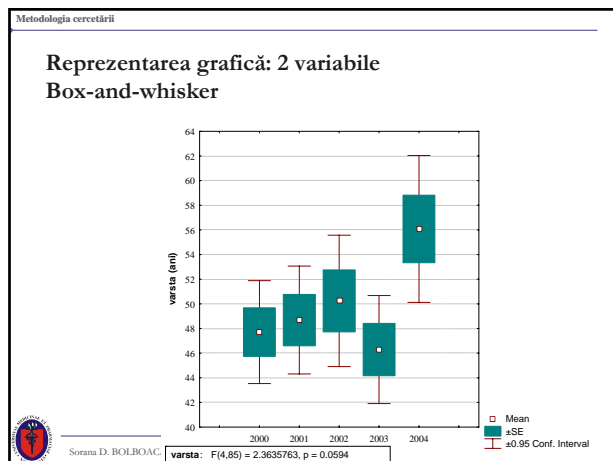
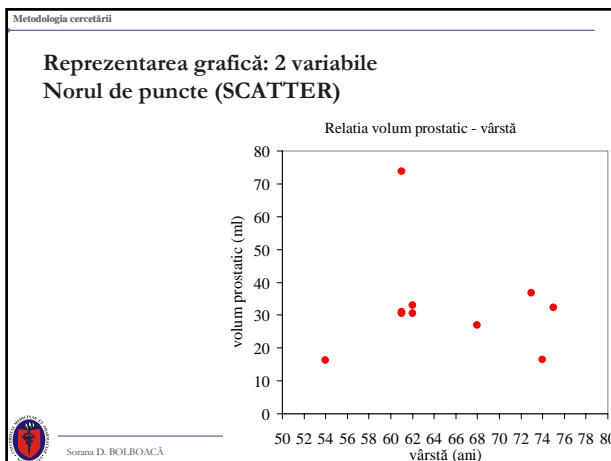
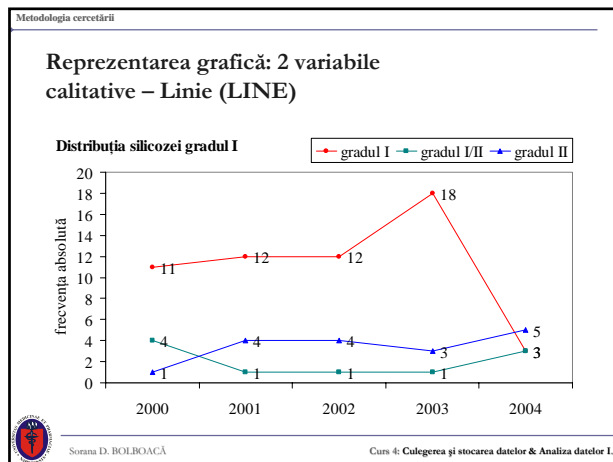
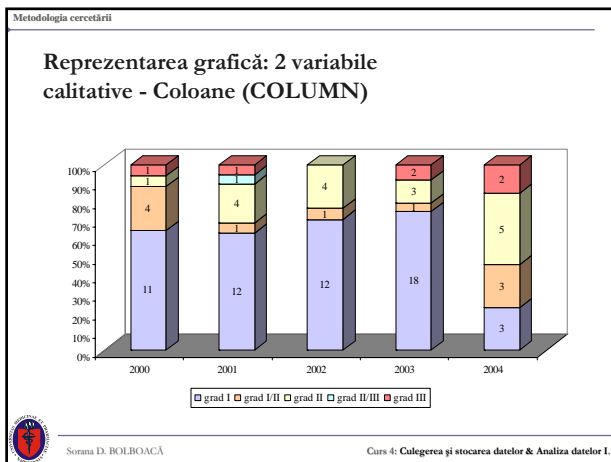
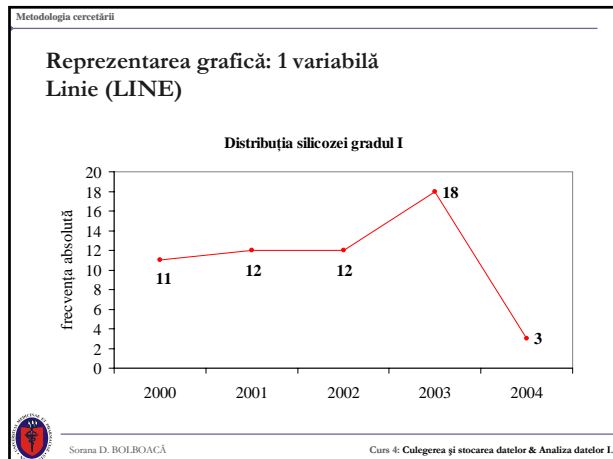
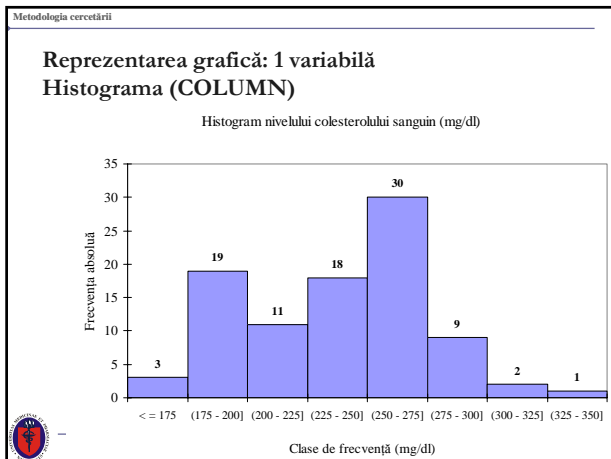
Metodologia cercetării

Reprezentarea grafică: 1 variabilă Coloane (COLUMN)

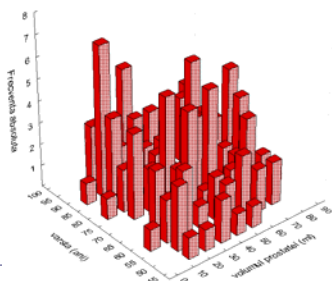
Modalitatea de implantare a cristinelor artificiale

tip implant	frecvența absolută
per primam	27
per secundam	15

Sorana D. BOLBOACĂ Curs 4: Culegerea și stocarea datelor & Analiza datelor I.



Reprezentarea grafică: 2 variabile Coloane (Histograme) bi-dimensionale



De reținut!

- Sumarizarea tabelară și reprezentarea grafică se realizează cu scopul transmiterii de informații.
- În realizarea lor trebuie să ținem cont de scop (putem distra atenția privitorului de la ceea ce dorim să transmitem).
- Asigurați-vă că aveți titluri informative, denumiri de rânduri și coloane; totaluri pe rânduri și/sau coloane.



De reținut!

- Asigurați-vă că axele au denumiri și unități de măsură.
- Minimizați numărul de culori.
- Evitați reprezentările grafice 2D și 3D:
 - Ceea ce se reprezintă ocupă o dimensiune mai mică din grafic.
 - Poate distorsiona imaginea în compararea a două distribuții.



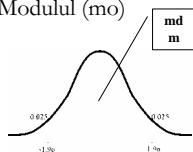
Parametrii statistici descriptivi



Date numerice, o variabilă

Măsuri ale tendinței centrale

- Cvartile
- Media aritmetică (m)
- Mediana (md)
- Modulul (mo)



Curba lui Gauss (distribuția normală)

Măsuri ale dispersiei sau variabilității

- Amplitudinea
- Varianța
- Eroarea standard a mediei

Caracteristici ale unei serii de date în jurul valorilor de mijloc



Date numerice, o variabilă: cvartile

- Divizăm setul de valori în 100 părți (maxim - minim):
Percentile

- Percentila de 90% = valoarea sub care regăsim 90% din valorile seriei [=PERCENTILE(argument)]
- [Insert - Function...]

- Divizăm scala în 10 părți: decile

- Divizăm scala în 4 părți: cvartile
[=QUARTILE(argument)]

- Permit aprecierea distribuției datelor analizate

- Exemplu: fie Q_1 (1/4), Q_2 (1/2) și Q_3 (3/4) primele 3 cvartile. Dacă $Q_2 - Q_1 \approx Q_3 - Q_2$ distribuția datelor este aproximativ simetrică. Dacă nu distribuția este asimetrică (spre dreapta sau spre stânga)



Metodologia cercetării

Date numerice, o variabilă: cuantile

- Cvartila 1: 3150 (Q_1)
- Cvartial 2: 3300 (Q_2)
- Cvartila 3: 3450 (Q_3)

Greutate la naștere (g)
2800
3100
3200
3300
3400
3500
3800

- Metoda de calcul: funcția QUARTILE
- Microsoft Excel
- Argumentele funcției:
 - Array: selectăm coloane unde avem date
 - Quart: 1 (cvartial 1), 2 (cvartial a doua), 3 (cvartila a treia)

$$Q_2 - Q_1 \approx Q_3 - Q_2$$

$$150 = 150$$

Sorana D. BOLBOACĂ Curs 4: Culegerea și stocarea datelor & Analiza datelor I.

Metodologia cercetării

Date numerice, o variabilă: Media aritmetică

- **Simboluri standard (universal recunoscute):**
 - μ = media aritmetică a unei populații
[=AVERAGE(argument)]
 - m = media aritmetică a unui eșantion (\bar{x})
 - Au aceeași formulă de calcul
 - Diferența constă în semnificația lui n (volumul populației sau volumul eșantionului)
 - Σ = sumă [=SUM(argument)]
 - Exemplu: fie variabila zile de supraviețuire (notată cu 'x') cu trei valori 1, 2, 4. $\Sigma x = 1+2+4 = 8$
 - $m(x) = \Sigma x / n = 8/3$

Sorana D. BOLBOACĂ Curs 4: Culegerea și stocarea datelor & Analiza datelor I.

Metodologia cercetării

Date numerice, o variabilă: Mediana

- Nu are simbol standard (uneori se notează cu 'md')
- Calculare (= MEDIAN(argument))
 - Așezăm valorile datelor în ordine crescătoare
 - Mediana este egală cu valoarea datei din mijlocul seriei dacă volumul eșantionului este impar
 - Valoarea medianei este egală cu media aritmetică a celor două valori din mijlocul seriei dacă volumul eșantionului este par

Sorana D. BOLBOACĂ Curs 4: Culegerea și stocarea datelor & Analiza datelor I.

Metodologia cercetării

Date numerice, o variabilă: Mediana

- Exemplu: calcularea medianei
 - Antigenul prostatic specific (simbol: PSA, unitate de măsură: ng/ml) pentru un eșantion de 10 pacienți (volumul eșantionului $n = 10$) cu manifestări prostatice: 7,6; 4,1; 5,9; 9,0; 6,8; 8,0; 7,7; 4,4; 6,1; 7,9
 - aranjăm datele în ordine crescătoare:
 - 4,1; 4,4; 5,9; 6,1; 6,8; 7,6; 7,7; 7,9; 8,0; 9,0
 - $n = 10 \rightarrow md = (6,8 + 7,6)/2 = 7,2$

Sorana D. BOLBOACĂ Curs 4: Culegerea și stocarea datelor & Analiza datelor I.

Metodologia cercetării

Date numerice, o variabilă: Media Aritmetică & Mediana

- Nu toate datele urmează o distribuție normală:
 - Deviere negativă: coadă la stânga
 - $md > m$ ($m < md < mo$)
 - Deviere pozitivă: coadă la dreapta
 - $m > md$ ($mo < md < m$)

Sorana D. BOLBOACĂ Curs 4: Culegerea și stocarea datelor & Analiza datelor I.

Metodologia cercetării

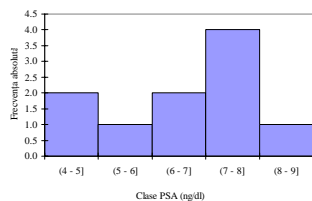
Măsuri de simetrie

- Într-o distribuție simetrică
media aritmetică = mediana = valoarea modală

Sorana D. BOLBOACĂ Curs 4: Culegerea și stocarea datelor & Analiza datelor I.

Date numerice, o variabilă: Modulul

- Nu are simbol standard (uneori se folosește 'mo')
- Valoarea cea mai frecvență a seriei
- Pe graficul de tip bare: modulul este valoarea barei cu frecvența cea mai mare



De reținut! Măsurile de centralitate

	+++	---
MODUL	<ul style="list-style-type: none"> ■ ușor de calculat ■ utilă pentru datele nominale 	<ul style="list-style-type: none"> ■ slabă stabilitatea de eșantionare
MEDIANANA	<ul style="list-style-type: none"> ■ nu e afectată de valorile extreme 	<ul style="list-style-type: none"> ■ Într-o oarecare măsură slabă stabilitate de eșantionare
MEDIA	<ul style="list-style-type: none"> ■ stabilitate de eșantionare ■ în legătură cu varianța 	<ul style="list-style-type: none"> ■ Nu este utilă pentru datele discrete ■ E afectată de distribuția asimetrică a datelor



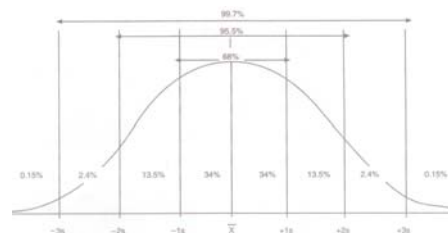
Date numerice, o variabilă: Variația & Deviația Standard

- Simbol standard variația populației: σ^2
 - Formula de calcul: $\sigma^2 = \sum(x - \mu)^2/n$
 - Alternativ: $\sigma^2 = [\sum(x^2 - n \times \mu^2)]/n$
- Simbol standard variația populației: s^2
 - Formula de calcul: $s^2 = [\sum(x^2 - n \times m^2)]/(n - 1)$
 - Exemplu:
 - $m_{PSA} = 6.75$
 - $s^2 = (74.1^2 + 4.4^2 + \dots + 9.0^2 - 10 \times 6.75^2)/9 = 2.585$
- Deviația standard: simbol s
 - Formula $s = \sqrt{s^2}$ [= SQRT(argument)]
 - Exemplu: $s = \sqrt{2.585} = 1.6078$



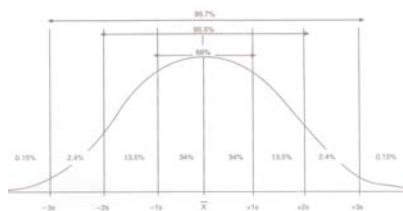
Date numerice, o variabilă: Media & Deviația standard

- Procentul de cazuri cu 1, 2 și 3 deviații standard ale mediei în distribuția normală



Date numerice, o variabilă: Media & Deviația standard

- Pe o curbă de distribuție normală care este procentul cazurilor care au valori mai mici decât $m+2 \times s^2$?
- Pe o curbă de distribuție normală care este procentul cazurilor care au valori mai mici decât $m+1 \times s^2$?



Date numerice, o variabilă: Eroarea standard a mediei

- Simbol standard:
 - Pentru media populației: σ_m
 - Pentru media eșantionului: s_m
 - Formule de calcul:
 - Populație: $\sigma_m = \sigma/\sqrt{n}$
 - Eșantion: $s_m = s/\sqrt{n}$
 - Exemplu PSA:
 - $s_m = s/\sqrt{n} = 1,6078/\sqrt{10} = 0,5084$



Date numerice, două variabile: Covarianța

- Analiza dependenței sau independenței dintre două variabile măsurabile
 - Covarianța:
 - Când urmărim două variabile măsurabile care variază simultan, una în relație cu cealaltă
 - Coeficientul de corelație:
 - Cuantificarea legăturii dintre cele două variabile



Date numerice, două variabile: Covariația

- Necesită înregistrarea unor date perechi
- Formula de calcul:
 - Populație: $\sigma_{xy} = (xy - n \times \mu_x \times \mu_y) / n$
 - Eșantion: $s_{xy} = (\sum xy - n \times m_x \times m_y) / (n - 1)$
 - Exemplu:
 - $s_{xy} = (75 \times 32.3 + 68 \times 27.0 + \dots + 74 \times 16.4 - 10 \times 65.1 \times 32.73) / (10 - 1) = -10.6478$
 - Excel (doar pentru populație):
 - [=COVAR(variabila1, variabila2)]



Date numerice, două variabile: Covariația

- Interpretare:
 - Dacă o variabilă crește odată cu cealaltă variabilă (așa cum este tensiunea arterială sistolică și diastolică) covarianța este pozitivă și valoarea acesteia este mare.
 - Dacă o variabilă crește în timp ce cea de-a doua variabilă scade (cum este de exemplu vârsta și volumul prostatei) covarianța este negativă și valoarea acesteia este mare.
 - Dacă creșterea sau descreșterea unei variabile nu este în legătură cu cealaltă variabilă de interes, valoarea covarianței este mică.



Măsuri de împrăștiere: Coeficientul de variație

- Interpretarea omogenității:

Coeficient de variație (CV)	Interpretare:
	populația poate fi considerată
CV < 10%	omogenă
10% ≤ CV < 20%	relativ omogenă
20% ≤ CV < 30%	relativ eterogenă / relativ heterogenă
> 30%	eterogenă / heterogenă

Date numerice, două variabile:
Coeficientul de corelație

- Formula de calcul:
 - Populație: $\rho_{xy} = \sigma_{xy} / (\sigma_x \times \sigma_y)$
 - Eșantion: $r_{xy} = s_{xy} / (s_x \times s_y)$
- Interpretare:
 - Dacă o variabilă în relație perfect lineară directă cu cea de-a doua variabilă coeficientul de corelație ia valoarea +1 (variabila 1 ↑ + variabila 2 ↑) / -1 (variabila 1 ↑ + variabila 2 ↓ sau variabila 1 ↓ + variabila 2 ↑)
 - Regulile lui Colton

Date numerice, două variabile:
Coeficientul de corelație

- Interpretare: Regulile lui COLTON:
 - un coeficient de corelație de la -0.25 la 0,25 înseamnă o corelație slabă sau nulă,
 - un coeficient de corelație de la 0.25 la 0.50 (sau de la -0.25 la -0.50) înseamnă un grad de asociere acceptabil
 - un coeficient de corelație de la 0.5 la 0.75 (sau de la -0.5 la -0.75) înseamnă o corelație moderată spre bună
 - un coeficient de corelație mai mare decât 0.75 (sau mai mic decât -0.75) înseamnă o foarte bună asociere sau corelație



Date numerice, două variabile: Coeficientul de corelație

- Exemplu:
 - $s_{x(\text{vârsta})} = 6,9992$
 - $s_{y(\text{volum})} = 15,9351$
 - $s_{xy} = -10,6478$
 - $r_{xy} = -10,6478 / (6,9992 \times 15,9351) = -0,0956$
 - Acest rezultat ne spune că pe eșantionul studiat format din 10 pacienți, volumul prostatei tinde să descrească odată cu creșterea vârstei, dar relația dintre cele două variabile este foarte slabă



Sumarizarea și analiza datelor



Sumarizarea și analiza datelor

- Generalizarea rezultatelor obținute pe eșantion asupra populației
1. Intervalul de confidență (pentru medie, pentru frecvențe)
 2. Statistică inferențială



Sumarizarea și analiza datelor: intervalul de confidență

- Definiție. Scop
- Interpretare
- Intervalul de încredere pentru medie
- Intervalul de încredere pentru frecvență



Sumarizarea și analiza datelor: intervalul de confidență De ce intervalul de încredere?

- Estimarea punctuală
 - = o valoare pentru parametrul teoretic estimat
 - Influențată de fluctuațiilor de eșantionare
 - poate fi la o mare distanță de valoarea reală a parametrului estimat
- Este recomandabil să se estimeze un parametru teoretic nu printr-o singură valoare ci printr-un interval, numit interval de încredere (în care să se poată afirma că parametrul estimat se găsește cu o probabilitate ridicată).



Sumarizarea și analiza datelor: intervalul 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 $1 - \alpha$
- Se aplică în cazul variabilelor distribuite normal!



Sumarizarea și analiza datelor: intervalul de confidență Interpretare

- Dacă intervalul de încredere pentru diferența dintre o medie observată și una teoretică cuprinde valoarea 0, datele sunt compatibile cu o diferență a mediei populației egală cu 0.
- Dacă intervalul de încredere pentru diferența dintre o medie observată și una teoretică nu cuprinde valoarea 0, datele nu sunt compatibile cu egalitatea mediilor populației.



Sumarizarea și analiza datelor: intervalul de confidență

- Se calculează în funcție de:
 - Talia eșantionului sau a populației
 - Variabila de studiat (calitativă, cantitativă)
- Formula de calcul cuprinde 2 părți:
 - Un estimator al calității eșantionului pe baza căruia estimatorul populației s-a calculat (eroarea standard)
 - Gradul de încredere (confidență) al intervalului specificat (scorul Z_α)
- Cel mai frecvent utilizat este intervalul de încredere pentru medie



Sumarizarea și analiza datelor: intervalul de confidență 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.

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



Sumarizarea și analiza datelor: intervalul de confidență Intervalul de încredere pentru medie

- Scorul Z este scorul distribuției normale de medie 0 și deviație standard de 1. Orice distribuție poate fi transformată în scorul Z utilizând formula:

$$Z = \frac{(X - \bar{X})}{s}$$

- Scorul pozitiv este mai mare decât media
- Scorul negativ este mai mic decât media
- Pentru intervalul de confidență de 95%: $Z_{5\%} = 1,96$
- Pentru intervalul de confidență de 99%: $Z_{1\%} = 2,58$

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



Sumarizarea și analiza datelor: intervalul de confidență 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 $Z = 1,96$.

$$\bar{X} = 105$$

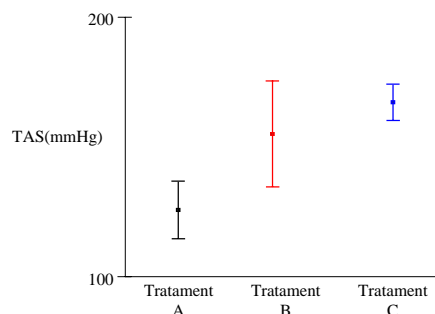
- $n = 121$
- $s^2 = 36$
- $s = 6$

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

- [105-1,07, 105+1,07]
- [103.93 - 106.07]
- [104-106]



Sumarizarea și analiza datelor: intervalul de confidență Compararea mediilor cu ajutorul intervalului de încredere



Sumarizarea și analiza datelor: intervalul de confidență Intervalul de încredere pentru frecvențe

- Dacă $n \cdot p > 10$

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



Sumarizarea și analiza datelor: intervalul de confidență 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]

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



Sumarizarea și analiza datelor: intervalul de confidență

- Rata șansei, respectiv riscul relativ

- Dacă intervalul de confidență al ratei șansei sau al riscului relativ conține valoarea 1 → expunerea nu este factor de risc pentru patologia de interes



Sumarizarea și analiza datelor: intervalul de confidență 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.

