

PRELUCRAREA DATELOR BIOMEDICALE PRIN METODE BIOSTATISTICE. PROBLEME

1. Introducere. Obiectivele ședințelor de exerciții practice de prelucrări de date

Biostatistica oferă metodele și instrumentele de descriere și analiză a datelor medicale, a rezultatelor experimentelor și observațiilor biomedicale. Totodată, abordarea problemelor medicale din punct de vedere statistic permite atât planificarea studiilor experimentale și observaționale, cât și descrierea fenomenelor biologice prin metode matematice.

Domeniul este vast și necesită un studiu aprofundat, dincolo de o introducere succintă pe care o putem acoperi în spațiul rezervat din cadrul cursurilor pentru studenții de la Universitatea de Medicină. În selecția de probleme prezentată în continuare (ca o completare a secțiunii de biostatistică din cursul Mihalaș & Lungeanu 2009) ne-am propus familiarizarea cu noțiunile de bază privind prelucrările statistice, acoperind în același timp aspectele fundamentale:

- descrierea datelor calitative și cantitative;
- estimarea statistică și testarea statistică a ipotezelor de cercetare biomedicală;
- analiza și corelarea;
- studii epidemiologice de risc.

Materialul începe cu prezentarea unor probleme rezolvate și apoi propune o serie de probleme la care nu sunt date soluțiile, lăsând cititorilor plăcerea de a le explora pe cont propriu. Problemele propuse sunt preluate și adaptate din: (a) Knapp & Clinton Miller 1992; (b) Rosner 2000.

Toate problemele propuse se pot rezolva cu *EpiInfo* și/sau *Microsoft Excel*. Soluțiile sunt prezentate în varianta rezolvării lor cu unul dintre cele două pachete, dar aceste abordări nu sunt în nici un fel restrictive. La fel de bine, studenții pot încerca utilizarea altor pachete statistice: *R*, *Statistica*, *SPSS*, etc.

2. Descrierea datelor. Estimarea statistică



Problema 2.1

Într-un studiu privind asocierea dintre fumat și boala coronariană, o echipă de cercetători a obținut următoarea distribuție a valorilor presiunii arteriale sistolice (în mm Hg) pentru un grup de 37 de fumători.

Presiunea arteriala	Frecvente	Frecvente relative	Frecvente
89.5-109.5	5	13.5	13.5
109.5-129.5	15	40.5	54.0
129.5-149.5	A	27.0	C
149.5-169.5	3	8.1	89.1
169.5-189.5	2	B	94.5
189.5-209.5	D	5.4	99.9
TOTAL	37		100.0

1. Valoarea **A** din tabel este:
2. Valoarea **B** din tabel este:
3. Valoarea **C** din tabel este:
4. Probabilitatea ca un individ selectat la întâmplare din acest eșantion de 37 de fumători să aibă presiunea sistolică între 89.5 mm Hg și 129.5 mm Hg este:



Soluții 2.1

1. Valoarea **A** din tabel este: $37 \cdot 0.27 = 9.99 \Rightarrow \mathbf{A=10}$
2. Valoarea **B** din tabel este: $\frac{2}{37} \cdot 100 = 5.4 \Rightarrow \mathbf{B=5.4}$
3. Valoarea **C** din tabel este: $54.0 + 27.0 = 81.0 \Rightarrow \mathbf{C=81.0}$
4. Probabilitatea ca un individ selectat la întâmplare din acest eșantion de 37 de fumători să aibă presiunea sistolică între 89.5 mm Hg și 129.5 mm Hg este: **54.0%** (coloana de frecvențe cumulate, rândul 2).



Problema 2.2

În procesul de pregătire pentru un test național, un lot de 200 de studenți au dat un test simulat la care au avut de răspuns la un număr de 100 de întrebări. Fiecare student a răspuns corect la un număr de întrebări între 35 și 59. Media

aritmetică a numărului de răspunsuri corecte a fost 47, iar deviația standard de 4. Presupunem că numărul de răspunsuri corecte urmează o distribuție normală.

1. Procentul de studenți care a răspuns corect la un număr de întrebări cuprins între 43 și 51 este de aproximativ:

2. Procentul de studenți care a răspuns corect la cel puțin 55 de întrebări este de aproximativ:



Soluții 2.2

1. Dacă media este $m=47$, iar deviația standard $s=4$, atunci în intervalul $(47 \pm 4) = (43; 51)$ se vor găsi aproximativ 68% dintre valori (în situația în care presupunem că valorile urmează o distribuție normală).

Concluzie: 68% dintre studenți au răspuns corect la un număr de întrebări cuprins între 43 și 51.

2. Dacă media este $m=47$, iar deviația standard $s=4$, atunci în intervalul $(47 \pm 2 * 4) = (39; 55)$ se vor găsi aproximativ 95% dintre valori (în situația în care presupunem că valorile urmează o distribuție normală). Înseamnă că doar aproximativ $(1 - 0.95)/2 = 0.025$ dintre valori sunt mai mari sau egale cu 55.

Concluzia: aproximativ 2.5% dintre studenți au răspuns corect la cel puțin 55 de întrebări.



Problema 2.3

Centrul de diabet și boli de nutriție din Timișoara are în observație un număr mare de pacienți din care în atenția noastră intră doar un eșantion de 32 de pacienți. Glicemia a fost determinată cu un aparat portabil (glucotest, cu bandetele). Datele sunt prezentate în tabelul de mai jos.

	SEX	GLICEMIE mg%
1	M	127
2	F	130
3	F	126
4	M	111
5	F	109
6	M	99
7	F	116
8	F	116
9	F	150
10	M	95
11	F	120
12	M	100
13	M	116
14	M	116
15	M	114
16	M	108
17	M	112
18	M	98
19	M	115
20	F	109
21	F	106
22	F	103
23	M	97
23	F	95
25	M	119
26	M	117
27	M	105
28	M	137
29	F	106
30	F	101
31	F	129
32	M	132

1. Pentru un nivel de încredere $1-\alpha=0.68$, determinați intervalul de încredere pentru valorile de glicemie. Interpretați rezultatul obținut (formulați concluzia în cuvinte).
2. Estimați media pentru valorile de glicemie ale acestor pacienți, cu $1-\alpha=0.95$. Interpretați rezultatul obținut (formulați concluzia în cuvinte).



Soluții 2.3

Pentru simplificare, mai jos vom presupune că valorile de glicemie urmează o distribuție normală. Sunt considerate corecte atât versiunile de

rezolvare riguroasă (utilizând distribuția t), cât și cele simplificate (utilizând distribuția normală și scorurile z).

Media valorilor de glicemie este $m=113.56$, iar deviația standard $s=13.09$ (se pot determina cu *Epi*, cu *Microsoft-Excel*, sau chiar prin calcul manual).

1. Pentru $1-\alpha=0.68$, $z \approx 1$, deci intervalul de încredere pentru valorile individuale de glicemie va fi aproximativ

$$(113.56 \pm 1 * 13.09) = (100.47; 126.65)$$

Concluzia: putem spune că aproximativ 68% dintre valorile de glicemie se vor încadra între 100.47 și 126.65

2. Pentru $1-\alpha=0.95$, $z \approx 2$, deci intervalul de încredere pentru media valorilor de glicemie va fi aproximativ

$$\left(113.56 \pm 2 * \frac{13.09}{\sqrt{32}} \right) = (108.94; 118.19)$$

Concluzia: putem spune că media valorilor de glicemie (pentru populația din care a fost extras acest eșantion) se încadrează între 108.94 și 118.19, cu o probabilitate de 95%.

3. Teste statistice



Problema 3.1

Cu datele de la Problema 2.3, ce test statistic se aplică pentru a vedea dacă există diferențe între glicemia celor două sexe? Formulați ipotezele statistice.

Aplicați testul și trageți concluzia. ($\alpha = .05$). Formulați în termeni statistici și medicali.



Soluții 3.1

Glicemia este o variabilă numerică despre care presupunem că urmează o distribuție normală - rezultă că tendința centrală va fi descrisă de medie, iar ipotezele statistice se vor referi la acest parametru.

Ipotezele statistice sunt:

$$H_0: \mu_M = \mu_F$$

media glicemiei este aceeași pentru cele două sexe;

$$H_a: \mu_M \neq \mu_F$$

media glicemiei este diferită pentru cele două sexe;

Se aplică testul **t nepereche** (avem două eșantioane independente - unul de femei și unul de bărbați) în varianta bi-direcțională (2-tailed).

Se aplică testul **t-nepereche 2-tailed**: se obține statistica $t=0.706$, căreia îi corespunde valoarea $p=0.486 > \alpha=0.05$.

Concluzie: Se acceptă H_0 și vom spune că NU există diferențe statistice semnificative între mediile glicemiei pentru cele două sexe.

Concluzia medicală: nivelul de glicemie este același pentru pacienții din acest centru, indiferent de sex – diferențele observate se încadrează în limitele variabilității biologice.



Problema 3.2

S-a făcut un studiu menit să investigheze efectele contraceptivelor orale (CO) asupra bolilor de inimă la femeile cu vârsta între 40 și 44 de ani. S-a găsit că între cele 5000 de femei care inițial erau utilizatoare curente de CO, 13 femei au dezvoltat infarct miocardic (IM) în decursul unei perioade de 3 ani. În același timp, dintre cele 10000 de femei care nu utilizau CO, doar 7 au dezvoltat IM în perioada de 3 ani.

Cum putem determina dacă aceste diferențe sunt semnificative sau ele se datorează doar șansei?

1. Sintetizați rezultatele studiului într-un tabel de contingență.
2. Ce test statistic se utilizează pentru a determina semnificația presupusei asocieri dintre utilizarea CO și bolile de inimă?
3. Formulați ipotezele statistice.
4. Aplicați testul ales. Trageți concluzia statistică.
5. Formulați concluzia medicală.



Soluții 3.2

Tabelul de contingență va fi:

	IM +	IM -	Totaluri
utiliz CO	13	4987	5000
NU utiliz CO	7	9993	10000
Totaluri	20	14980	15000

Se aplică **testul χ^2** – se compară proporțiile celor care dezvoltă infarct miocardic în rândul utilizatoarelor CO, respectiv a celor care nu utilizează CO.

Ipoteza de zero: proporția celor care dezvoltă IM este aceeași în cele 2 grupe de femei (utilizatoare, respectiv ne-utilizatoare CO).

Ipoteza alternativă: există proporții diferite de cazuri de IM pentru cele două grupe de femei.

Se aplică testul χ^2 (Epi aplică implicit varianta bi-direcțională).

Se obține:

$$\chi^2 (\text{uncorrected}) = 9.04, p = 0.0026$$

$$\chi^2 (\text{Yates corrected}) = 7.67, p = 0.0056$$

(corecția de continuitate – în această situație concluziile statistice sunt aceleași).

Concluzie: H_0 este respinsă (riscul erorii de tipul I este sub pragul de semnificație α), iar H_a va fi acceptată.

Vom spune că **sunt diferențe statistic semnificative** între proporția celor care dezvoltă IM pentru cele două grupe de femei. Cum $p < 0.01$, putem spune că **diferențele sunt chiar foarte semnificative din punct de vedere statistic**. Concluzia medicală care se impune este ca aceste diferențe (implicit riscul respectiv) nu pot fi ignorate, ele nu se datorează doar șansei.

Observați însă ce loturi mari au fost implicate în acest studiu pentru a se atinge puterea statistică necesară punerii în evidență a acestor diferențe între proporții.



Problema 3.3

S-a făcut un studiu de comparare a efectului analgezic al unor produse farmaceutice utilizate de pacienți cu dureri lombare: ibuprofen (400 mg), codeină (60 mg), codeină (30 mg) și placebo. Cei 20 de participanți la studiu au fost distribuiți aleatoriu în patru grupe ($n=5$ pentru fiecare grup). La două ore după administrarea produsului, pacienților li s-a cerut să descrie senzația de îndepărtare a durerii pe o scară între 0 și 100. Rezultatele sunt prezentate în tabel:

Scorul atribuit			
Ibuprofen	Codeină	Codeină	Placebo
82	80	77	65
89	70	69	75
77	72	67	67
72	90	65	55
92	68	57	63

1. Care sunt **ipotezele statistice**? Formulați-le în termeni statistici și în cuvinte.
2. Ce **test statistic** se aplică pentru a vedea dacă există diferențe semnificative între efectele celor patru substanțe? Motivați-vă decizia.
3. Aplicați testul ales și formulați concluzia (în termeni statistici și medicali).
4. Comparați apoi separat doar **Ibuprofen (400 mg)** cu **Placebo**.

Care sunt **ipotezele statistice**? Formulați-le în termeni statistici și în cuvinte.

Aplicați testul statistic potrivit și trageți concluzia (în termeni statistici și medicali).



Soluții 3.3

Ipoteza de zero: toate cele 4 produse farmaceutice au același efect - media scorurilor va fi aceeași (aceeași tendință centrală)

$$H_0: \mu_A = \mu_B = \mu_C = \mu_D$$

Ipoteza alternativă: cel puțin unul dintre produse diferă de celelalte - cel puțin o medie diferă de celelalte

$$H_a: \mu_i \neq \mu_j \text{ pentru cel puțin una din cele } C_4^2 = 6 \text{ combinații posibile.}$$

Se aplică testul ANOVA (Analiza Variantei) – se compară mediile pentru 4 eșantioane.

După aplicarea testului ANOVA se obține statistica $F=5.16$, careia îi corespunde valoarea $p=0.011 < \alpha=0.05$.

Concluzie: Se respinge H_0 și vom spune că există diferențe statistice semnificative în efectul analgezic al celor 4 produse farmaceutice – cel puțin unul dintre cele 3 produse medicamentoase (ibuprofen, codeină-60 mg, respectiv codeină-30mg) are un efect semnificativ diferit de placebo.

Ipoteza de zero: cele 2 produse farmaceutice au același efect - media scorurilor va fi aceeași (aceeași tendință centrală).

$$H_0: \mu_{\text{ibuprofen}} = \mu_{\text{placebo}}$$

Ipoteza alternativă: cele două produse au efecte diferite - mediile vor fi diferite.

$$H_a: \mu_{\text{ibuprofen}} \neq \mu_{\text{placebo}}$$

Presupunem că întrebările 9-10 sunt separate de punctele anterioare (pentru simplificare).

Se aplică testul **t- nepereche 2-tailed**: se obține statistica $t=3.547$, careia îi corespunde valoarea $p=0.008 < \alpha=0.05$.

Concluzie: Se respinge H_0 și vom spune că există diferențe statistic semnificative în efectul analgezic al celor 2 produse farmaceutice –ibuprofenul are un efect semnificativ diferit de placebo.

Observație (suplimentar):

Dacă suntem riguroși și aplicăm un *test post-hoc* (Bonferroni de exemplu - un test t la care α a fost împărțit între cele 6 teste t posibile) nivelul de semnificație atins va fi $p=0.02 < \alpha=0.05$, deci concluzia va fi aceeași.



Problema 3.4

Tabelul de mai jos prezintă concentrația de digoxină în sânge la 4 ore și la 8 ore după injectarea acesteia intravenos pe un lot de 9 bărbați sănătoși.

Concentrația de digoxină (mg/l)			
Pers.	4 ore	8 ore	Diferențe (d_i)
1	1.0	1.0	0.0
2	1.3	1.3	0.0
3	0.9	0.7	-0.2
4	1.0	1.0	0.0
5	1.0	0.9	-0.1
6	0.9	0.8	-0.1
7	1.3	1.2	-0.1
8	1.1	1.0	-0.1
9	1.0	1.0	0.0
	$m_1=1.056$	$m_2=0.99$	$m_d=-0.067$
	$s_1=0.15$	$s_2=0.18$	$s_d=0.017$

1. Ce test se aplică pentru a vedea dacă digoxina se elimină – concentrația de digoxină măsurată la 4 ore de la injectare diferă semnificativ față de cea măsurată la 8 ore? Explicați de ce.
2. Formulați ipotezele statistice, atât în termeni matematici, cât și în cuvinte.
3. Aplicați testul statistic și trageți concluzia ($\alpha= .05$) ? Formulați în termeni statistici și medicali.



Soluții 3.4

Se aplică **testul t-pereche** – se verifică existența unui fenomen măsurat printr-o variabilă numerică + valorile din cele două eșantioane de valori (la 4 ore și la 8 ore) se pun în corespondență – ne interesează tendința diferențelor dintre valoarea concentrației la 4 ore și cea de la 8 ore.

Tendința pentru variabile numerice se exprimă prin medie – deci ipotezele statistice se vor referi la medie.

Ipoteza de zero: digoxina nu se elimină – concentrația în sânge rămâne constantă \Rightarrow pentru populație diferențele sunt nule.

$$H_0: \mu_d=0$$

Ipoteza alternativă: diferențele nu sunt nule – mai mult, dacă digoxina se elimină, ne așteptăm ca diferențele pentru populație să fie negative.

Valoarea $p=0.01$ (sau $p=0.02$) $< \alpha=0.05$ reprezintă probabilitatea ca H_0 să fie adevărat în condițiile în care pentru acest eșantion s-au obținut valorile precizate.

Concluzie: Se respinge H_0 și vom spune că diferențele de concentrație sunt statistic semnificativ diferite de zero.

Concluzia în termeni medicali: Eliminarea digoxinei în intervalul de 4 ore (măsurat la 4 ore, respectiv la 8 ore de la injectare) este statistic semnificativă, deci nu poate fi neglijată.

4. Analiza corelației și regresiei



Problema 4.1

Obstetricienii obișnuiesc să ceară sumarul de urină pentru a testa nivelul de estriol la femeile gravide care se apropie de termen deoarece s-a descoperit că nivelul de estriol este relaționat cu greutatea fătului. Testul poate furniza o dovadă indirectă pentru o greutate mică (sub nivelul normal) a fătului.

În 1963, doi obstetricieni (J. Greene și J. Touchstone) au făcut un studiu în care au înregistrat nivelul de estriol în apropierea termenului, respectiv greutatea copilului la naștere. Tabelul de mai jos prezintă datele publicate de ei, iar pe pagina următoare sunt prezentate rezultate la prelucrarea statistică.

i	Estriol (mg/24 hr) x_i	Greutatea la naștere (g/100) y_i	i	Estriol (mg/24 hr) x_i	Greutatea la naștere (g/100) y_i
1	7	25	17	17	32
2	9	25	18	25	32
3	9	25	19	27	34
4	12	27	20	15	34
5	14	27	21	15	34
6	16	27	22	15	35
7	16	24	23	16	35
8	14	30	24	19	34
9	16	30	25	18	35
10	16	31	26	17	36
11	17	30	27	18	37
12	19	31	28	20	38
13	21	30	29	22	40
14	24	28	30	25	39
15	15	32	31	24	43
16	16	32			

1. Ce tip de studiu este cel prezentat. Motivați-vă afirmația și explicați care sunt caracteristicile acestui studiu.
2. Formulați ipotezele statistice ale acestui studiu în termeni matematici și în cuvinte.
3. Reprezentați grafic datele. Alegeți cea mai potrivită reprezentare grafică, astfel încât ea să fie în concordanță cu afirmațiile făcute la punctele 1 și 2.
4. Prelucrați datele din punct de vedere statistic, potrivit afirmațiilor pe care le-ați făcut la punctele 1, 2 și 3.

Indicație. Puteți încerca rezolvarea problemei cu Microsoft-Excel, dar și cu alte pachete statistice.

5. Interpretați rezultatele obținute. Formulați în termeni statistici și medicali.



Soluții 4.1

Studiul și-a propus *analiza corelației* dintre nivelul de estriol la femeile gravide care se apropie de termen și greutatea nou-născutului. În studiu au fost implicate 31 de viitoare mame și pentru fiecare dintre ele s-a consemnat valoarea pentru cele două variabile (sub formă de perechi de valori). Interesează investigarea relației dintre cele două variabile numerice și exprimarea tăriei acesteia printr-un *coeficient de corelație*.

Dacă se dorește să se poată face o predicție a greutateii nou-născutului (respectiv a fătului) în funcție de nivelul de estriol, atunci se face *analiza regresiei* – nivelul de estriol va fi considerat ca fiind variabila independentă. De obicei, pentru două variabile numerice (cum este cazul aici), se investighează o posibilă relație liniară exprimabilă printr-o *dreaptă de regresie*. Panta acestei drepte se numește *coeficient de regresie* și va permite calculul greutateii fătului în funcție de nivelul de estriol, chiar pentru valori ale acestuia care nu au apărut în studiul inițial (dar între limitele min. și max. din studiul făcut).

Din punctul de vedere al tipurilor de studii epidemiologice, nu ni se dau suficiente informații ca să putem identifica tipul de studiu. Am putea presupune că este un *studiu transversal*.

Ipotezele statistice sunt:

H₀: Ipoteza de zero: cele două variabile sunt independente - nu există nici o relație între nivelul de estriol și greutatea nou-născutului

$$H_0: \rho=0$$

H_a: Ipoteza alternativă: cele două variabile sunt corelate – există o relație între nivelul de estriol și greutatea nou-născutului

$$H_a: \rho \neq 0$$

Grafic de tip „scatter” (nor de puncte) – variabila independentă pe abscisă (Figura 1).

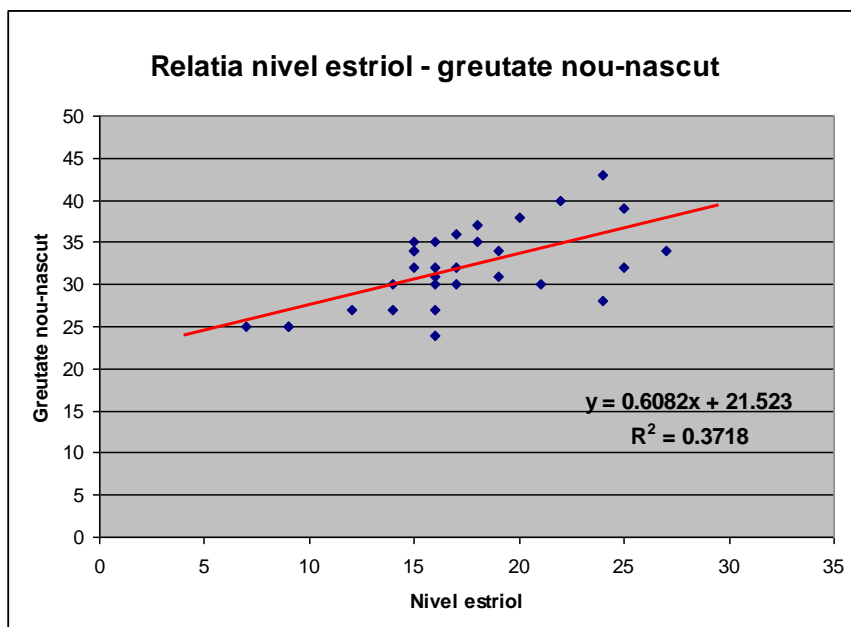


Figura 1. Cea mai potrivită reprezentare grafică este cea de tip „scatter”, care scoate în evidență relația dintre cele două variabile.

În Microsoft-Excel coeficientul de corelație se poate determina cu funcția CORREL sau cu *Tools* → *Data analysis* → *Regression*.

Coeficientul de corelație Pearson este $r = 0.6097$ cu $N=31-2=29$ grade de libertate.

Dacă se face calculul statisticii t se obține:

$$r = 0.609731$$

$$t = 4.142656$$

$$p = 0.000271$$

În Microsoft-Excel:

$$t = r \sqrt{\frac{n-2}{1-r^2}} \text{ adică } t = r * \text{SQRT}(7/(1-r^2))$$

$p < 0.001$ s-a obținut cu funcția TDIST.

Rezultat similar se obține și cu *Tools* → *Data analysis* → *Regression*.

O altă alternativă de rezolvare - se utilizează comanda *Regress* în programul *Analysis* din **EpiInfo**. În acest caz se obține un interval de încredere (cu $1-\alpha=0.95$) pentru coeficientul de corelație și se interpretează acest interval.

Judecând doar după r , putem spune că există o corelație directă (pozitivă) și relativ puternică (r peste 0.5) – corespunzătoare unei corelații puternice între nivelul de estriol și greutatea nou-născutului.

Pentru o discuție mai nuanțată și pentru a ne pronunța asupra semnificației statistice, trebuie luată în considerare **valoarea lui p** (indiferent prin ce metodă s-a obținut).

Concluzie: corelația între nivelul de estriol și greutatea nou-născutului este extrem de semnificativă din punct de vedere statistic ($p < 0.001$).

Din punct de vedere medical, există dovezile că nivelul de estriol al mamei la apropierea termenului este corelat direct și semnificativ cu greutatea fătului. Cunoscând nivelul de estriol, se poate estima greutatea fătului și pot fi luate din timp măsurile medicale necesare.



Problema 4.2

Un grup de cercetători și-au propus să investigheze relația dintre nivelul plasmatic de amfetamină și psihoza indusă de această substanță. În studiul făcut au fost cuprinși 10 consumatori cronici, cărora li s-a determinat nivelul plasmatic de amfetamină, determinare urmată imediat de o evaluare psihiatrică (în urma

căreia s-a obținut un scor al intensității psihozei). Datele rezultate sunt sintetizate în tabelul următor:

Pacient	Intensitatea psihozei (Y)	Nivelul de amfetamină (mg/ml) (X)
1	10	150
2	30	300
3	20	250
4	15	150
5	45	450
6	35	400
7	50	425
8	15	200
9	40	350
10	55	475

1. Reprezentați grafic datele din acest studiu – alegeți tipul de grafic care ilustrează cel mai bine relația dintre cele două variabile din studiu.
2. Determinați coeficientul de corelație Pearson.
3. Interpretați valoarea obținută - trageți concluzia pentru studiul în discuție.
4. Este această corelație statistic semnificativă? Argumentați.



Soluții 4.2

Se face un grafic de tip „scatter” (nor de puncte) – variabilă independentă pe abscisă (Figura 2).

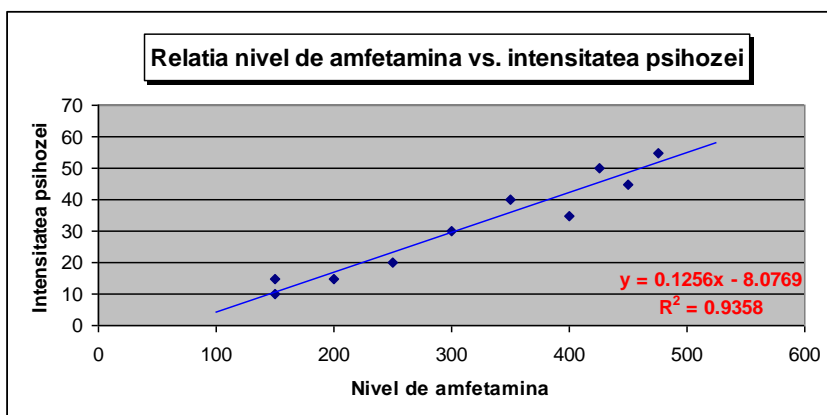


Figura 2. Cea mai potrivită reprezentare grafică este cea de tip „scatter”, care scoate în evidență relația dintre cele două variabile.

În Microsoft-Excel coeficientul de corelație se poate determina cu funcția CORREL sau cu *Tools* → *Data analysis* → *Regression*.

Coeficientul de corelație Pearson este $r = 0.97$ cu $N=10-2=8$ grade de libertate.

Judecând doar după r , putem spune că există o corelație directă (pozitivă) și puternică (r aproape de valoarea 1 – corespunzătoare unei corelații “perfecte”) între nivelul plasmatic de amfetamină și intensitatea psihozei.

Pentru o discuție mai nuanțată – punctul 2.4.

Dacă se face calculul statisticii t se obține:

$$r = 0.967381949$$

$$t = 10.80114078$$

$$p = 4.76491 \cdot 10^{-6}$$

În Microsoft-Excel:

$$t = r \sqrt{\frac{n-2}{1-r^2}} \text{ adică } t = r \cdot \text{SQRT}(8/(1-r^2))$$

$$p = 4.76491 \text{E-}06 \text{ obținut cu funcția TDIST}$$

Concluzie: corelația între nivelul plasmatic de amfetamină și intensitatea psihozei este extrem de semnificativă din punct de vedere statistic ($p < 0.0001$).

Rezultat similar se obține și cu *Tools* → *Data analysis* → *Regression*.

O altă alternativă de rezolvare - se utilizează comanda *Regress* în secțiunea *Analysis* din **EpiInfo**. În acest caz se obține un interval de încredere (cu $1-\alpha=0.95$) pentru coeficientul de corelație și se interpretează acest interval.

5. Studii observaționale. Analiza riscului



Problema 5.1

Într-un studiu privind frecvența de apariție a efectelor secundare la un medicament anti-artritic, s-au observat două grupuri de pacienți: unii care urmau tratament cu noul medicament și alții care nu urmau acest tratament. Au fost observați 400 de pacienți selectați dintr-un spital. Rezultatele sunt prezentate în tabelul ce urmează.

TRATAMENT	EFECTE SECUNDARE	FARA EFECTE SECUNDARE	TOTALURI
Nu au urmat tratament	39	271	310
Au urmat tratament	7	83	90
TOTALURI	46	354	400

1. Ce tip de studiu a fost utilizat? Explicați care sunt caracteristicile acestui studiu.
2. Calculați **RR** de apariție a efectelor secundare, asociate cu medicamentul anti-artritic aflat în studiu.
3. Interpretați **RR** calculat la punctul anterior – luați în considerare intervalul de încredere pentru **RR** și exprimați concluzia din punct de vedere medical.



Soluții 5.1.

Studiul este de tip cohort prospectiv (longitudinal; “*follow-up*”) – în studiu au fost implicate 2 grupuri de subiecți (unii “*expuși*” tratamentului în discuție, ceilalți nu) care au fost urmăriți un timp pentru a se depista eventualele efecte secundare. S-au numărat cazurile cu probleme adverse pentru ambele grupuri și s-a determinat riscul de apariție a efectelor asociate cu boala artritică la cei care au urmat tratamentul în discuție în comparație cu pacienții care nu au urmat acel tratament.

Dacă se face calculul utilizând programul Epi-STATCALC – tabelul de contingență va arăta astfel:

Tratament Exposure	-	Efecte secundare - Disease		Totaluri
		(+)	(-)	
	(+)	7	83	90
	(-)	39	271	310
	Totaluri	46	354	400

Se obține:

$$RR = 0.62$$

Intervalul de estimare pentru riscul relativ este:

$$0.29 < RR < 1.33$$

$$1-\alpha = 0.95$$

Intervalul de estimare îl conține pe 1 \Rightarrow nu putem spune că tratamentul în discuție ar implica un risc mai crescut sau mai scăzut în apariția efectelor secundare la pacienții ce suferă de artrită și urmează tratamentul anti-artritic în discuție.

(RR se definește ca fiind raportul între incidența efectelor secundare pentru cei ce urmează tratamentul în discuție și incidența efectelor pentru cei care nu urmează tratamentul respectiv, urmând probabil alte tratamente).

6. Probleme propuse



Până în anii 1970, emisia de compuși de plumb în aer nu era considerată o problemă de sănătate publică. O contribuție importantă în evidențierea efectelor poluării cu plumb a adus-o un studiu de la începutul anilor 1970 făcut în Texas pentru a investiga efectele expunerii la compuși de plumb asupra funcțiilor neurologice și psihologice [1, 2]. Datele au fost furnizate de către profesorul Bernard Rosner [3], cărui i-au fost date de către dr. Philip Landrigan [4].

Pe scurt, a fost studiat un grup de copii care au locuit în apropierea topitoriei de plumb din El Paso, Texas, cărora li s-a determinat nivelul de plumb din sânge. A fost determinat un grup de 46 copii “expuși” care aveau nivelul de plumb în sânge $\geq 40 \mu\text{g/ml}$ în 1972 (pentru câțiva dintre ei determinarea s-a făcut în 1973). În același timp, cei 78 de copii cu nivel de plumb în sânge $< 40 \mu\text{g/ml}$ (atât în 1972 cât și în 1973) au format grupul de control.

- [1] Landrigan, PJ, Whitworth, RH, Baloh, RW, Staehling, NW, Barthel, WF, Rosenblum, BH (1975). Neuropsychological dysfunction în children with chronic low-level lead absorption. *Lancet* 1: 708-715.
- [2] Morse, DL, Landrigan, PJ, Rosenblum, BH, Hubert JS, Housworth J (1979). El Paso revisited. Epidemiologic follow-up of an environmental lead problem. *JAMA* 249 (8). <http://jama.ama-assn.org/cgi/content/abstract/242/8/739> (Last access 13th April 2008).
- [3] Bernard, R (2000). *Fundamentals of biostatistics* (5th Ed.). Duxbury: Pacific Grove.
- [4] CDC. epidemic Intelligence Service. <http://www.cdc.gov/EIS/about/landrigan.htm> (Last access 13th April 2008).

O parte din datele obținute în acest studiu sunt furnizate în fișierul **lead_hw2345.xls** ce se găsește pe CD-ul ce însoțește această carte. Semnificația câmpurilor este:

Id – codul de identificare;

Area – distanța față de topitorie

1=0-1 mile;

2=1-2.5 mile;

3=2.5-4.1 mile.

Sex – sexul (1=masculin, 2=feminin);

Lead_type – descrie cele trei grupuri de copii în funcție de nivelul de plumb în sânge în 1972 și 1973:

1 – grupul de CONTROL, cu nivel de plumb în sânge < 40 $\mu\text{g/ml}$ atât în 1972, cât și în 1973;

2 – grupul EXPUS CURENT, cu nivelul de plumb în sânge ≥ 40 $\mu\text{g/ml}$ în 1973;

3 – grupul EXPUS ÎN TRECUT, cu nivelul de plumb în sânge ≥ 40 $\mu\text{g/ml}$ în 1972, dar cu nivel normal în 1973.

Pica și Colic – exprimă prezența simptomelor de pica și colici (1=da, 2=nu);

FWT_R și FWT_L - performanța copiilor la un test de tipul “*finger-wrist tapping*”;

MAXFWT – scorul obținut pentru mîna dominantă (maximul dintre FWT_R și FWT_L);

Age_y – vârsta exprimată în ani;

Group – identifică două grupuri de intoxicație:

0=neintoxicați;

1=intoxicați (cei pentru care Lead_type este 2 sau 3);

Acolo unde nu apar valori, înseamnă că datele lipsesc.

Întrebările 1-5

Distribuția valorilor pentru vîrstă este prezentată în tabelul de mai jos, unde intervalele sunt considerate închise la dreapta și deschise la stînga:

Vârsta ani	Frecvențe absolute	Frecvențe relative	Frecvențe cumulate%
0 - 2	0	0.000	0.00%
2 - 4	8	0.065	6.45%
4 - 6	A	0.153	B
6 - 8	31	C	46.77%
8 - 10	21	0.169	63.71%
10 - 12	14	0.113	75.00%
12 - 14	19	0.153	90.32%
14 - 16	12	0.097	D
TOTAL	124	1	

Graficul din Figura 3 sintetizează distribuția valorilor.

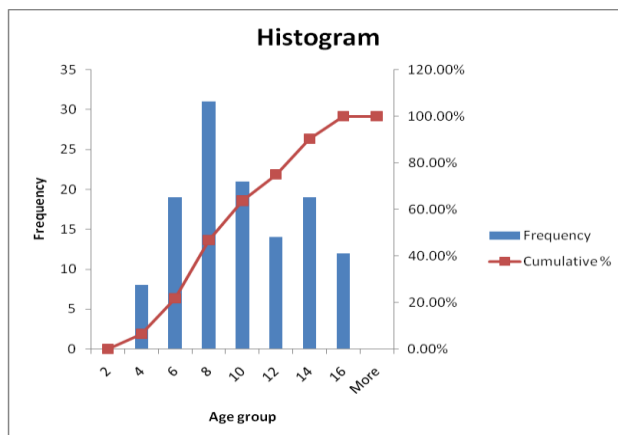


Figura 3. Distribuția valorilor de vârstă pentru datele din fișierul **lead_hw2345.xls**

1. Valoarea **A** din tabel este:
2. Valoarea **B** din tabel este:
3. Valoarea **C** din tabel este:
4. Valoarea **D** din tabel este:
5. Probabilitatea ca un individ selectat la întâmplare din acest eșantion de 124 de copii să aibă vârsta între 6 și 12 ani este:

Întrebările 6 - 8

Pentru următoarele trei întrebări, considerați eșantionul reprezentativ pentru zona respectivă.

6. Pentru un nivel de încredere $1-\alpha=0.68$, determinați intervalul de încredere pentru vârsta copiilor din zona respectivă. Interpretați rezultatul obținut (formulați concluzia în cuvinte).
7. Estimați media pentru vârsta copiilor din zona respectivă, cu $1-\alpha=0.95$. Interpretați rezultatul obținut (formulați concluzia în cuvinte).
8. Estimați procentul de copii care au vârsta sub 5 ani (cu aproximație).

Întrebările 9-13

9. Pentru un nivel de încredere $1-\alpha=0.95$, estimați media pentru vârsta copiilor de sex feminin. Interpretați rezultatul obținut (formulați concluzia în cuvinte).
10. Estimați media vârstei pentru pacienții de sex masculin, cu $1-\alpha=0.95$. Interpretați rezultatul obținut (formulați concluzia în cuvinte).

11. Ce test statistic se aplică pentru a vedea dacă există diferențe între vârstele copiilor de cele două sexe? Formulați ipotezele statistice.
12. Aplicați testul și trageți concluzia statistică ($\alpha=0.05$).
13. Formulați concluzia în termeni medicali.

Întrebările 14 - 15

Investigați asocierea dintre intoxicația cu plumb și prezența colicilor la copiii din studiul descris. Estimați riscul pe care-l constituie intoxicația cu plumb (descrisă prin variabila GROUP) pentru apariția colicilor (descrisă prin variabila COLIC).

Informații ajutătoare: în grupul copiilor neintoxicați cu plumb (GROUP=0) au fost 78 de copii, din care 11 au prezentat colici; în grupul copiilor intoxicați (GROUP=1) au fost 46 de copii, din care 12 au prezentat colici.

14. Ce tip de studiu statistic este recomandat pentru această investigație?
15. Investigați asocierea dintre intoxicația cu plumb și prezența colicilor.
 - 15.1. Formulați ipotezele statistice (matematic și în cuvinte);
 - 15.2. Precizați ce test statistic aplicați și de ce;
 - 15.3. Aplicați testul ales;
 - 15.4. Formulați decizia statistică;
 - 15.5. Formulați concluzia medicală în cuvinte.

Întrebările 16-17

Investigați performanța copiilor la testul “*finger-wrist tapping*” (FWT), exprimată de variabila numită MAXFWT, pentru cele trei grupuri de copii identificați prin variabila LEAD_TYPE.

16. Determinați intervalele de încredere pentru estimarea mediei performanței pentru fiecare grup în parte.
17. Comparați cele trei grupuri din punctul de vedere al performanței la testul FWT (variabila MAXFWT). Analizați atât din punct de vedere statistic, cât și al interpretării medicale a rezultatelor statistice:
 - 17.1. Formulați ipotezele statistice (matematic și în cuvinte);
 - 17.2. Precizați ce test statistic aplicați și de ce;
 - 17.3. Aplicați testul ales;
 - 17.4. Formulați decizia statistică;
 - 17.5. Formulați concluzia medicală în cuvinte.

Întrebările 18-19

Reluați investigarea asocierii dintre intoxicația cu plumb și prezența colicilor la copiii din studiul descris. Estimați riscul pe care-l constituie intoxicația cu plumb (descrisă prin variabila GROUP) pentru apariția colicilor (descrisă prin variabila COLIC).

[Informații ajutătoare: în grupul copiilor neintoxicați cu plumb (GROUP=0) au fost 78 de copii, din care 11 au prezentat colici; în grupul copiilor intoxicați (GROUP=1) au fost 46 de copii, din care 12 au prezentat colici.]

18. Rediscutați tipul de studiu statistic recomandat pentru această investigație, într-un context mai complex (luând în considerare și studiile epidemiologice).

19. Estimați riscul pe care-l constituie intoxicația cu plumb pentru apariția colicilor:

- 19.1.** Formulați ipotezele statistice pentru această abordare;
- 19.2.** Aplicați metoda pe care ați ales-o și precizați rezultatul obținut;
- 19.3.** Trageți concluzia statistică;
- 19.4.** Formulați concluzia medicală în cuvinte.

Întrebările 20-22

Similar cu modul în care ați investigat asocierea dintre intoxicația cu plumb și prezența colicilor la copiii din studiul descris (atât în tema anterioară, cât și în cea prezentă), investigați și asocierea intoxicării (descrisă prin variabila GROUP) cu prezența simptomelor de pica (descrisă prin variabila PICA).

20. Ce tip de studiu statistic este recomandat pentru aceasta investigație?

21. Investigați asocierea dintre intoxicația cu plumb și prezența simptomelor de pica.

22. Estimați riscul pe care-l constituie intoxicația cu plumb pentru apariția simptomelor de pica.

7. Concluzii

În această ședință ați învățat noțiuni de bază privind prelucrările statistice ale datelor biomedicale:

- ✚ descrierea datelor obținute în urma unor studii experimentale sau observaționale;
- ✚ estimarea statistică;
- ✚ testarea statistică a ipotezelor de cercetare;
- ✚ analiza regresiei și corelației;
- ✚ studii epidemiologice, estimarea riscului.

Referințe

- Vernic CV, Apostol SA, Frandes M, Mada L, Lungeanu D. APLICATII PRACTICE DE INFORMATICA SI BIOSTATISTICA MEDICALA IN NURSING. Editura Eurostampa, ISBN 978-606-32-0487-6, Colectia Derzelas, 2017:1-216
- Vernic CV, Mada L, Lungeanu D, Muntean C, Apostol SA, Catu CO, Ursoniu S. Aplicații practice de Informatică Medicală și Biostatistică. Editura Victor Babes, ISBN 978-606-8054-09-4, 2010: 1-229.
- Mihalaș GI, Lungeanu D. Informatică medicală și biostatistică. Timișoara: EVB, 2009.
- Sheskin DJ. Handbook of parametric and nonparametric statistical procedures (3rd Edition.). Boca Raton: Chapman & Hall/CRC, 2004.
- Rosner B. Fundamentals of biostatistics (5th Edition). Pacific Grove: Duxbury. Thomson Learning, 2000.
- Knapp RG, Clinton Miller M: Clinical epidemiology and biostatistics. Baltimore: Williams & Wilkins, 1992.