

## Seminar 6      Corelatii si regresii

Deschideți fișierul "P6.xls" si salvati-l sub denumirea: "Nume\_P6.xls" in directorul dedicat cursului. Fișierul conține mai multe foi (worksheets)!

**0.** În worksheet-ul 'Corelatie 0' determinati coeficientul de corelatie Pearson intre lungimea femurului si al humerusului introducand manual formule in celulele indicate:

femur	humerus	$(x-x_{med})$	$(y-y_{med})$	$(x-x_{med})^2$	$(y-y_{med})^2$	$(x-x_{med}) \cdot (y-y_{med})$
38	41					
56	63					
59	70					
64	72					
74	84					
$x_{med} =$	$y_{med} =$					
$n =$				$s_x =$	$s_y =$	COV =

$$COV_{xy} = \frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y}) \quad r_{xy} = \frac{\frac{1}{n} \cdot \sum_{i=1}^n (x_i - \bar{x}) \cdot (y_i - \bar{y})}{s_x \cdot s_y} \quad s = \sqrt{\frac{\sum_{i=1}^n (x_i - \bar{x})^2}{n}}$$

**I.** În worksheet-ul 'Corelatie 1' aveți date pentru nivelul hormonilor FSH (hormon de stimulare foliculara) si TSH (hormone de stimulare tiroidiana) pentru un lot de 6 pacienti.

1. Folosiți opțiunea "chart" pentru a crea un grafic tip "scatter" a nivelului TSH în funcție de nivelul FSH. Dați nume adecvate graficului și axelor. Răspundeți la urmatoarele întrebări privind natura relației dintre cele două variabile:

- există o corelație între cele două variabile?
- corelația este liniară?

2. Folosind programul EXCEL, calculați coeficientul Pearson  $r_{xy}$  pentru datele indicate.  
*Există două funcții: CORREL și PEARSON, care dau exact același rezultat!*

3. Calculați coeficientul Spearman pentru aceleasi date. Pentru aceasta trebuie să urmați următorii pași:

- Calculați ordinul "rank" pentru datele din fiecare set: într-o parte goală a worksheet-ului folosiți funcția RANK pentru a calcula ordinal corespunzător fiecărei valori FSH. Efectuați aceiași pași pentru TSH.

Atentie: programul EXCEL atribuie greșit ordinul datelor prea apropiate!. (Ex: *pentru datele 3; 4; 5; 5; 6 - ordinul perechii de date "5" trebuie să fie 3,5!*)

Pentru a calcula corect coeficientul Spearman  $r_s$ , ordinele datelor apropiate trebuie să fie ajustate manual!

- Calculați pătratul diferențelor dintre fiecare pereche de ordine.
- Folosiți funcția SUM pentru a însuma pătratele diferențelor calculate. (Acest pas și cel anterior se pot face într-unul singur dacă se folosește funcția SUMXMY2).

- Calculați  $r_s$  folosind formula 
$$r_s = 1 - \frac{6 \cdot \sum d^2}{n^3 - n}$$

*Indicație:* puteți folosi funcția COUNT pentru a determina mărimea eșantionului automat.

4. Testați semnificația fiecărei statistici calculate pentru întrebările 1 și 2 considerând 5% nivel de încredere. Enunțați cele două ipoteze clar, precum tipul de test-t folosit ("one-tailed" sau "two-tailed").

- Pentru coeficientul Pearson  $r_{xy}$ , calculați  $t_{calc}$  folosind formula 
$$t_{calc} = r \cdot \sqrt{\frac{n-2}{1-r^2}}.$$

- Calculați probabilitatea asociată valorii  $t_{calc}$ , pentru gradul de libertate  $n-2$ , folosind funcția TDIST.

- Deoarece mărimea eșantionului este mai mică decât 10, pentru testarea semnificației statistice a coeficientul Spearman  $r_s$ , trebuie folosit comparata valoarea  $r_{S_{calc}}$  cu  $r_{S_{crit}}$ , folosind tabelul adecvat (se găsește în curs).

**II. Determinarea valorilor FSH și TSH s-a făcut pentru încă un lot de pacienți. Datele obținute se găsesc în worksheet-ul 'Corelație 2'.**

- Trasați noul grafic și repetați calculele de la punctele 1 - 4 pentru setul de date mai larg.

*Avertizare:* în noile date sunt ordine apropiate pe care EXCEL nu le calculează adecvat, deci trebuie ajustate ordinele calculate de EXCEL. Acest lucru se poate face în modul următor:

- Să presupunem că ordinul pentru FSH se găsește în domeniul: D3:D20 (când urmați instrucțiunile de mai jos folosiți referința adecvată datelor proprii).

- Într-o zonă goală a worksheet-ului tastați următoarea formulă:

$$=D3+(COUNTIF(D$3:D$20,D3)-1)/2$$

Cu această formulă se adaugă 0,5 ordinelor corespunzătoare la două datelor apropiate, 1,0 pentru ordinele corespunzătoare la trei date apropiate, etc.

- Copiați celula în jos la următoarele 19 rânduri pentru a calcula ordinul *real* pentru fiecare valoare a debitului.

- Repetați pașii de mai sus pentru datele TSH.

- Continuați ca de obicei pentru a calcula coeficientul de corelație Spearman.

Comparați diferențele dintre rezultatele celor două teste înainte și după eșantionul adițional, ținând cont de puterea lor relativă și de presupunerile ce stau la baza fiecărei statistici.

### III.

A. Să se determine ecuația drepte de regresie (prin metoda celor mai mici pătrate) pentru seturile de date ce reprezintă valoarea calciului seric (y) și valoarea parathormonului (x) din worksheet-ul “regresie 1”.

B. Determinați valoarea coeficientului de regresie și apoi folosind formulele date la curs determinați ecuația drepte de regresie.

IV. Pentru 9 pacienți cu anemie aplastică se evaluează procentul de reticulocite și numărul limfocitelor (worksheet-ul “regresie 2”).

a) Trasati dreapta de regresie dintre procentul de reticulocite și numărul limfocitelor.

b) Calculați coeficientul de corelație și coeficientul de determinare, și interpretați nivelul de corelație.

c) Calculați variația reziduală.

Aplicați analiza "Regression" din "Data analysis" și încercați să înțelegeți cât mai multe din datele afișate.

V. Folosind tabelul din worksheet-ul “regresie 3” calculați:

a) coeficientul de corelație dintre durata spitalizării și vârstă

b) cea mai bună relație liniară dintre durata spitalizării și vârstă.

c) verificați semnificația statistică a acestei relații

Pentru a răspunde la aceste întrebări puteți aplica analiza "Regression" din "Data analysis"