



UMF
UNIVERSITATEA DE
MEDICINĂ ȘI FARMACIE
IULIU HAȚIEGANU
CLUJ-NAPOCA

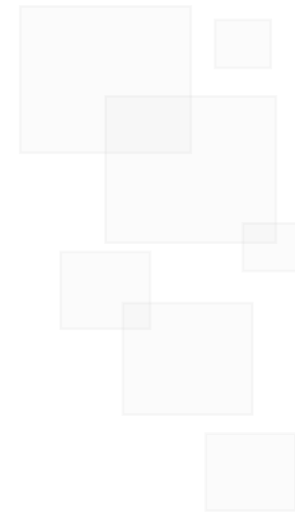
Testarea ipotezelor statistice

Analiza de corelație și regresie

Cuprins

- » Analiza corelației dintre **2 VARIABLE CANTITATIVE**:
 - Corelații liniare versus corelații monotone
 - Regresia liniară simplă

Corelații liniare



Formularea problemei:

Studiul **relației / legăturii / corelației** dintre 2 VARIABILE CANTITATIVE:

- Greutatea (kg) și Tensiunea arterială sistolică (mmHg)?
- Nivelul de inteligență (IQ) și mărimea creierului?

Greutate = X: X_1, X_2, \dots, X_n
TAS = Y: Y_1, Y_2, \dots, Y_n

0. Evaluarea grafică a relației:

- Nor de puncte (Scatter plot)

1. Intensitatea / forța/magnitudinea relației

- Coeficient de corelație Pearson sau Spearman

2. Testarea relației dintre variabilele X și Y pe populație

- Test statistic pentru coeficientul de corelație

3. Predictie : predicția valorilor unei variabile în funcție de valorile celeilalte

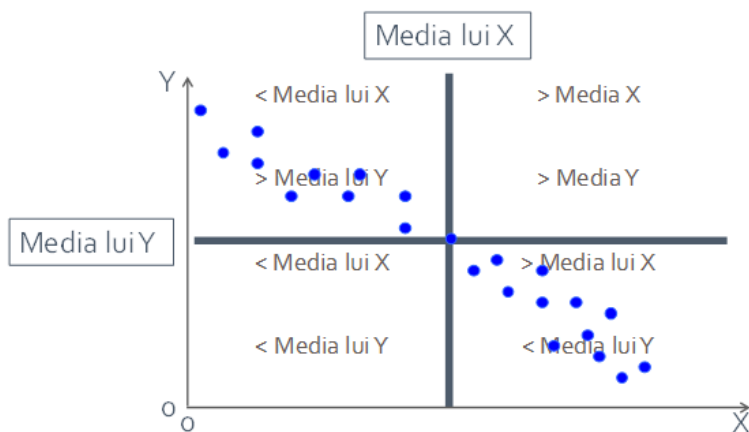
- Regresia liniară

Relația dintre variabile: coeficienți de corelație

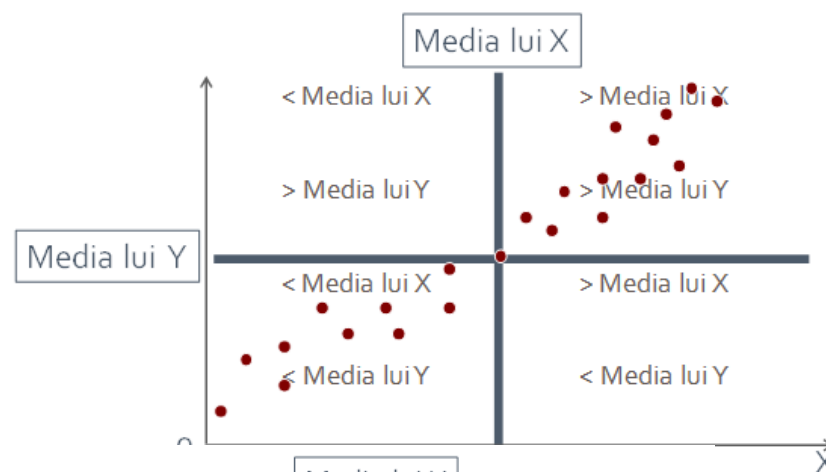
- 2 coeficienți de corelație:
- Pearson și Spearman

	Pearson	Spearman
Simbol	R / r	ρ (rho)
Distribuția	Fiecare dintre variabile are o distribuție normală de probabilitate	Orice distribuție
Relație	Liniară	Monotonă

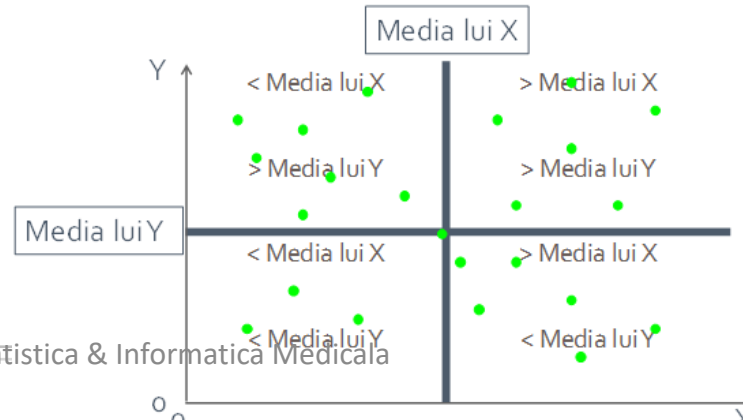
- Reprezentarea grafică: nor de puncte



corelație liniară negativă



corelație liniară pozitivă



lipsa corelației liniare

Calcul coeficient de corelație Pearson

Formula de calcul

Covarianța de eșantionare $COV(X,Y)$:
$$COV(X,Y) = \frac{1}{n-1} \sum_{i=1}^n (X_i - \bar{X})(Y_i - \bar{Y})$$

Coeficientul de corelație liniară Pearson:
$$r = \frac{COV(X,Y)}{S_X \cdot S_Y}$$

Coeficient de determinare $d = r^2$

X_i, Y_i = valorile celor două variabile X și Y , \bar{X} și \bar{Y} sunt mediile celor două variabile; n = talie esantion; S_x și S_y = deviațiile standard de eșantionare
 r = coefficient de corelatie Pearson, d = coefficient de determinare

Coeficient de corelație Pearson - interpretări

Coeficientul de corelație Pearson (notație r): direcția și intensitatea relației dintre două variabile cantitative;

Interpretarea direcției/sensului/tendinței:

- dacă $r > 0 \Rightarrow$ tendință crescătoare/ pantă ascendentă/ legătură de directă proporționalitate/ corelație pozitivă;
- dacă $r < 0 \Rightarrow$ tendință descrescătoare/ pantă descendentă/ legătură de inversă proporționalitate/ corelație negativă;
- dacă $r \approx 0 \Rightarrow$ nici o corelație liniară

Interpretarea intensității relației dintre variabile

- cu cât r este mai mare (în valoare absolută) cu atât relația dintre variabile este mai puternică;
- cu cât r este mai aproape de 0, cu atât relația este mai slabă/ puțin importantă,

Coeficient de corelație Pearson - interpretări

•Regulile lui Colton:

- dacă r aparține intervalului $[-0,25; 0,25]$ \Rightarrow Nu există nici o corelație liniară între cele două variabile sau există o corelație liniară neglijabilă
- dacă r aparține intervalului $[-0,50; -0,25)$ sau $(0,25; 0,50]$ \Rightarrow corelație liniară slabă spre acceptabilă
- dacă r aparține intervalului $[-0,75; -0,50)$ sau $(0,50; 0,75]$ \Rightarrow corelație liniară moderată spre bună
- dacă r aparține intervalului $[-1; -0,75)$ sau $(0,75; 1]$ \Rightarrow corelație liniară foarte bună

Exemplu de calcul

- Pe un eșantion de 10 pacienți adulți s-au măsurat variabilele Greutate (kg) și HDL colesterol (mg/dL)-vezi tabelul de mai jos, Se știe că valoarea covarianței dintre cele două variabile este $COV(Greutate, HDL) = 14,32$ iar deviațiile standard de eșantionare sunt $S_{HDL} = 1,07$ mg/dL; $S_{Greutate} = 14,6$ kg, Calculați coeficientul de corelație liniară.

HDL (mg/dL)	6,8	5,3	4,3	5,0	7,1	5,5	3,8	4,6	4,0	6,0
Greutatea (kg)	90	75	70	73	110	67	60	65	59	80

Soluție:

$$r = \frac{COV(Greutate, HDL)}{S_{Greutate} \cdot S_{HDL}} = \frac{14,32}{14,6 \cdot 1,07} = 0,91$$

Interpretare :

$r > 0 \Rightarrow$ tendință crescătoare / relație de directă proporționalitate,

$r > 0,75 \Rightarrow$ intensitatea corelației liniare este foarte bună (regulile lui Colton)

Test statistic de semnificație a coeficientului de corelație Pearson

- **Scopul testului:** testează dacă între cele două variabile cantitative există o relație de corelație (sau asociere) semnificativă statistic
- **Condiții de aplicabilitate:**
 - Observații independente în eșantion
 - Variabile cantitative
 - Cele două variabile sunt normal distribuite
 - Fără valori aberante (foarte îndepărtate de norul de puncte)

Test statistic de semnificație a coeficientului de corelație Pearson

- Ipoteza nulă (H_0)
 - Nu există o corelație liniară semnificativă statistic între variabilele X și Y pe populația de interes
- Ipoteza alternativă (H_1):
 - Există o corelație liniară semnificativă statistic între variabilele X și Y pe populația de interes
- Decizia testului cu ajutorul valorii p :
 - ✓ dacă $p < \alpha (0,05) \Rightarrow$ se respinge $H_0 \rightarrow$ suntem în favoarea lui $H_1 \rightarrow$ cu un risc de eroare de 5%, există o corelație liniară între variabilele X și Y pe populația de interes
 - ✓ dacă $p \geq \alpha (=0,05) \Rightarrow$ nu se respinge H_0

Test statistic de semnificație a coeficientului de corelație Pearson

Exemplul1:

Coeficientul de corelație pentru relația dintre Trigliceride (mg/dL) și Greutate(kg) pentru 50 de subiecți adulți cu diabet DZ2 este 0,72, iar valoarea p asociată este 0,001. Perechile de observații sunt independente, datele au avut distribuție normală, relația este liniară.

Interpretarea valorii p asociate coeficientului de corelație:

- $p < 0,05 \rightarrow$ există o corelație semnificativă statistic între variabilele trigliceride și greutate pe populația de adulți cu DZ2 (Corelația dintre trigliceride și greutate este semnificativă statistic)

Interpretarea direcției și intensității corelației:

- $r = 0,72 > 0 \rightarrow$ relația este de proporționalitate directă
- $r = 0,72$ este în $[0,50$ și $0,75)$ \rightarrow intensitatea corelației este moderată până la bună

Test statistic de semnificatie a coeficientului de corelatie Pearson

Exemplul2:

Coeficientul de corelație pentru relația dintre Trigliceride (mg/dL) și Glicemie (mg/dL) pentru 24 de subiecți cu vârsta între 20-30 ani este de 0,39, iar valoarea p asociată este de 0,35.

Perechile de observații sunt independente, datele sunt normal distribuite, relația este liniară.

Interpretarea valorii p asociate coeficientului de corelație:

- $p > 0,05 \rightarrow$ nu am găsit o **corelație semnificativă statistic** între Trigliceride și Glicemie pe populația de adulți cu vârsta între 20-30 ani (formulare echivalentă: nu avem suficiente dovezi statistice pentru a demonstra existența corelației liniare semnificative între trigliceride și glicemie pe populația de interes)

Interpretarea direcției și intensității corelației:

- interpretarea este inutilă deoarece relația nu a atins semnificația statistică - deci rezultatul poate fi influențat foarte mult de hazard (șansă). Pe eșantion, relația pare a fi de directă proporționalitate ($r = 0,39 > 0$), iar intensitatea corelației este scăzută / acceptabilă ($r = 0,39$ - este în $[0,25$ și $0,50)$).

Coeficientul de corelație Spearman

Scop: evaluarea **asocierii** dintre 2 variabile din punct de vedere al **direcției** și **magnitudinii** relației dintre acestea

Condiții de aplicare:

- Observații independente în eșantion;
- Cele 2 variabile sunt variabile ordinale/cantitative (normal distribuite sau fără distribuție normală)

Utilitate: evaluarea relației dintre

- două variabile ordinale
- o variabilă ordinală și una cantitativă
- două variabile cantitative care nu sunt normal distribuite

Coeficient de corelatie Spearman - interpretari

Interpretarea direcției/sensului/tendinței:

- dacă $\rho > 0 \rightarrow$ corelație pozitivă
- dacă $\rho < 0 \rightarrow$ corelație negativă
- dacă $\rho \approx 0 \rightarrow$ nici o corelație

Interpretarea intensității relației

- cu cât ρ este mai mare (în valoare absolută) cu atât relația dintre variabile este mai puternică
- cu cât ρ este mai aproape de 0, cu atât relația este mai slabă/ puțin importantă

Test statistic de semnificație a coeficientului de corelație Spearman

- Ipoteza nulă (H_0)
 - Nu există o **corelație semnificativă statistic** între variabilele **X** și **Y** pe **populația de interes**
- Ipoteza alternativă (H_1):
 - Există o **corelație semnificativă statistic** între variabilele **X** și **Y** pe **populația de interes**
- Decizia testului cu ajutorul valorii p :
 - ✓ dacă $p < \alpha (=0,05) \Rightarrow$ se respinge $H_0 \Rightarrow$ suntem în favoarea lui $H_1 \Rightarrow$ cu un risc de eroare de 5%, există o **corelație semnificativă statistic** între variabilele **X** și **Y** pe **populația de interes**

WOOC LAP -QUIZZ



- accesați link-ul:
www.wooclap.com
- introduceți codul:
- **CODE: AMGC9**



Regresia liniară



Tipuri de Regresii

Clasificare în funcție de:

- **tipul variabilei dependente**
 - **variabila cantitativă → regresie liniară**
 - Variabila calitativa dihotomială – regresie logistică
- **Numărul de variabile dependente:**
 - **regresie univariată (o singură variabilă dependentă)**
 - Regresie multivariată (≥ 2 variabile dependente)
- **Numărul de variabile independente :**
 - regresie **simplă** (o variabilă independentă)
 - regresie **multiplă** (≥ 2 variabile independente)

Regresia liniară simplă

Scop: evaluarea relației liniare dintre o variabilă independentă și o variabilă dependentă cantitativă

Obiective:

- **predicția valorilor** unei variabile cantitative ca funcție a valorilor variabilei independente
- **evaluează**
 - **Existența** legăturii/**dependenței**/relației dintre 2 variabile
 - **Direcția** legăturii/relației dintre 2 variabile
 - **Importanța/magnitudinea** legăturii/relației dintre 2 variabile

Regresia liniară simplă

- **Ecuatia dreptei de regresie:** $Y = b_0 + b_1X$

b_0 = ordonata în origine: valoarea lui Y când $X=0$

b_1 = panta dreptei de regresie,

Interpretarea lui b_1 = coeficientul variabilei X

- **Semnul lui b_1**

- dacă $b_1 > 0$ tendința crescătoare / pantă ascendentă / pantă pozitivă / legătură de directă proporționalitate
- dacă $b_1 < 0$ tendința descrescătoare / pantă descendentă / pantă negativă / legătură de inversă proporționalitate
- dacă $b_1 \approx 0$ nici o tendință liniară

- **Magnitudinea lui b_1**

- **La o creștere cu 1 unitate a variabilei independente X , valoarea variabilei dependente Y crește (sau scade) în medie cu b_1 unități**

Exemplu: corelație liniară

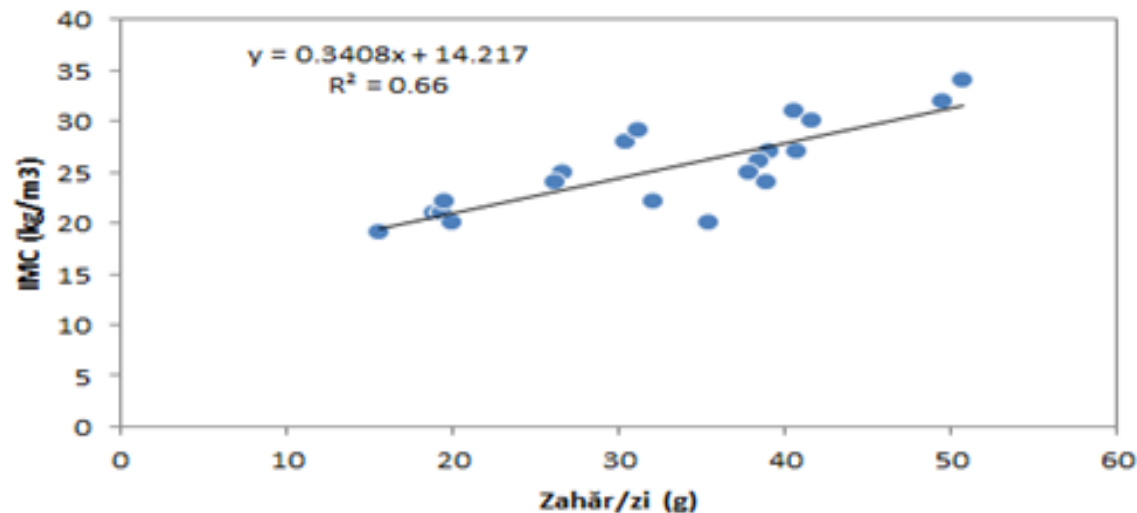
- **Un eșantion aleatoriu de 20 studenți AMG 1, UMF Cluj, a fost investigat pentru a identifica corelația liniară dintre indicele de masă corporală și consumul de zahăr pe zi (grame/zi).**
- Variabilele IMC(kg/m²) și Consumul zilnic de zahar (gr.) au avut distribuție normală
- Analiza corelației: $r = 0,812$ ($p < 0,001$)

Exemplu: regresie liniară

- Un eșantion aleatoriu de 20 studenți AMG 1, UMF Cluj, a fost investigat pentru a identifica dependența dintre indicele de masă corporală (IMC, kg/m²- variabilă dependentă) și consumul de zahăr pe zi ca variabilă independentă.
- Rezultate:
- Analiza de regresie:
- $\text{IMC}(\text{kg/m}^2) = 14,22 + 0,34 \times \text{Zahăr/zi}(\text{g})$

Exemplu: regresie liniară

- $\text{IMC}(\text{kg}/\text{m}^3) = 0,34 \times \text{Zahăr}/\text{zi}(\text{g}) + 14,22$
- $R^2 = 0,66$
- Coeficientul regresiei: la o creștere a consumului zilnic de zahăr cu 1 gram, IMC crește în medie cu $0,34 \text{ kg}/\text{m}^2$. În plus, 66% ($R^2 = 0,66$) din variația IMC-ului este explicată de relația liniară a acestuia cu consumul de zahar,



De reținut!

- Testarea **corelației** dintre două variabile se realizează prin:
 - Coeficientul de corelație Pearson: 2 variabile cantitative normal distribuite
 - Corelația Spearman: 2 variabile de orice tip (cantitative care nu urmează o distribuție normală sau 2 variabile calitative ordinale)
- Testarea **dependentei liniare** dintre 2 variabile se realizează prin regresia liniară simplă

Exemple de probleme

Q1. S-a realizat un studiu pe 2000 subiecți cu vârstă de peste 55 ani pentru a identifica relația liniară varsta (ani) și tensiunea arterială sistolică (TAS, mmHg), S-a obținut un coeficient de corelație liniară (r) între varsta și TAS de 0,55.

Care din următoarele afirmații sunt corecte?

- A. Există o relație de dependență negativă între vârsta și TAS
- B. Corelația este pozitivă; TAS tinde să fie mai mare pentru subiecții cu vârsta mai mare, intensitatea corelației fiind moderată spre bună
- C. Nu există relație între vârsta și TAS
- D. Corelația este negativă; TAS tinde să fie mai mare pentru subiecții cu vârsta mai mică, intensitatea corelației fiind moderată spre bună
- E. Coeficientul de corelație folosit este coeficientul de corelație al lui Pearson

R1: B, E

Exemple de probleme

Q2. S-a realizat un studiu pe 2000 subiecți cu vârstă de peste 55 ani pentru a identifica relația liniară vârsta (ani) și tensiunea arterială sistolică (TAS, mmHg). Covarianța dintre vârsta și TAS a fost egală cu 10,24 iar deviațiile standard de esantionare au fost egale cu 5 ani pentru vârsta respectiv 4 mmHg pentru TAS. Presupunem ca cele două variabile au o distribuție normală.

Care din următoarele afirmații sunt corecte?

- A. Coeficientul de corelație Pearson a fost egal cu $r=0,512$
- B. Nu este posibilă determinarea nici unui coeficient de corelație
- C. Nu există relație între vârsta și TAS
- D. În esanționul de studiu, ambele variabile variază în același sens
- E. În esanționul de studiu, ambele variabile variază în sens contrar

R2: A, D

Exemple de probleme

Q3. S-a realizat un studiu pe 2000 subiecți cu vârstă de peste 55 ani pentru a identifica relația liniară vârsta (ani) și tensiunea arterială sistolică (TAS, mmHg), S-a obținut un coeficient de corelație liniară între vârsta și TAS de 0,55 și o valoare p asociată $p < 0,001$.

Care din următoarele afirmații sunt corecte?

- A. Există o corelație liniară semnificativă statistic între vârsta și TAS pentru că $p < 0,05$
- B. În eșantionul de studiu există o corelație liniară pozitivă semnificativă statistic între vârsta și TAS pentru că $p < 0,05$
- C. Ipoteza nulă a testului statistic a cărui valoare p a fost $< 0,001$ afirmă că: Nu există corelație liniară semnificativă între vârsta și TAS
- D. Ipoteza alternativă a testului statistic a cărui valoare p a fost $< 0,001$ afirmă că: Există corelație liniară semnificativă între vârsta și TAS
- E. În eșantionul de studiu, corelația liniară dintre variabile a fost pozitivă iar intensitatea acesteia a fost moderată spre bună

R3: A, C, D, E

Exemple de probleme

Q4. S-a realizat un studiu pe 2000 subiecți cu vârstă de peste 55 ani pentru a identifica relația liniară vârstă (ani) și tensiunea arterială sistolică (TAS, mmHg). S-a obținut un model de regresie liniară între vârsta (ani) și TAS de forma: $TAS = 100,13 + 0,365 \times \text{varsta}$. Care din următoarele afirmații sunt corecte?

- A. Există o relație de dependență pozitivă între vârsta și TAS
- B. Există o relație de proporționalitate directă între vârsta și TAS pentru că panta regresiei este pozitivă
- C. În eșantionul de studiu, la o creștere a vârstei pacientului cu 1 an, TAS crește în medie cu 0,365 mmHg,
- D. Formula $TAS = 100,13 + 0,365 \times \text{varsta}$ reprezintă rezultatul unei regresii univariante multiple
- E. Pe eșantionul de studiu, nu există relație între vârsta și TAS

R4: A, B, C

MULȚUMESC PENTRU ATENȚIE!

