

Autor: Bondor Cosmina-loana

Asociere și predicție

- A** ALWAYS
- S** SEEK
- K** KNOWLEDGE

Objective

- Corelație
- Regresie liniară
- Exerciții

Scenariu

- 36 de pacienți cu artroplastie bilaterală totală a genunchiului
- A fost măsurată:
 - Greutatea
 - Densitatea minerală osoasă
- Scop: Există asociere între greutate și densitatea minerală osoasă

Corelație

- Relația dintre două caracteristici
- Cum este relația?
- Putem prezice un eveniment?
- Ce eroare putem să acceptăm?

Asociere

- Termeni similari: relație, dependență.
- Termenul **corelație** se utilizează numai în cazul a două variabile numerice sau ordinale
- Dacă Y poate fi calculat din X = asociere
- Ex. Y este de 3 ori mai mare decât X . Dacă îl știm pe X , putem să-l calculăm pe Y .

Dependență

- Dacă X este cauza lui Y, atunci Y este **dependent** de X

Greutatea asociată cu densitatea minerală osoasă

- X = variabila independentă (greutatea)
- Y = variabila dependentă (densitatea minerală osoasă)
- Mai întâi X, apoi Y – demonstrația cauzalității
- ! Ca să demonstrăm cauzalitatea avem nevoie de protocol de studiu specific, de mai multe studii, un studiu de corelație nu este suficient

Relația dintre două caracteristici

- Două variabile cantitative sau calitative ordinale
 - relație liniară: coeficient de corelație Pearson
 - relație neliniară: coeficient de corelație Spearman
 - Coeficient de corelație interclase
- Două variabile calitative:
 - Risc relativ
 - Rata șansei
 - Corelația interobservator - Cohen Kappa coeficient de concordanță
- Dacă este implicată variabila de timp ca variabilă dependentă
 - Rata hazardului



Coeficient de corelație Pearson

Dacă X și Y sunt două variabile cantitative sau ordinale. Coeficientul de corelație Pearson:

$$r_{XY} = \frac{\sum_{i=1}^n (X_i - \bar{X})(Y_i - \bar{Y})}{\sqrt{\sum_{i=1}^n (X_i - \bar{X})^2} \sqrt{\sum_{i=1}^n (Y_i - \bar{Y})^2}}$$

Unde n - numărul de pacienți (dimensiunea eșantionului), \bar{X}, \bar{Y} - media aritmetică a variabilelor X respectiv Y



Coeficient de corelație Pearson

- Indică asocierea dintre X și Y
- $r=1$ sau $r=-1$ corelație perfectă

- Este între -1 și 1

$$r \in [-1,1]$$

Cu cât $|r|$ se apropie de 1 cu atât asocierea este mai puternică

Cu cât $|r|$ se apropie de 0 cu atât asocierea este mai slabă

! Vorbim despre relație liniară



Graficul de corelație – XY Scatter

- Două variabile – pot fi cu două unități de măsură diferite
- Variabila dependentă pe axa OY și variabila independentă pe axa OX

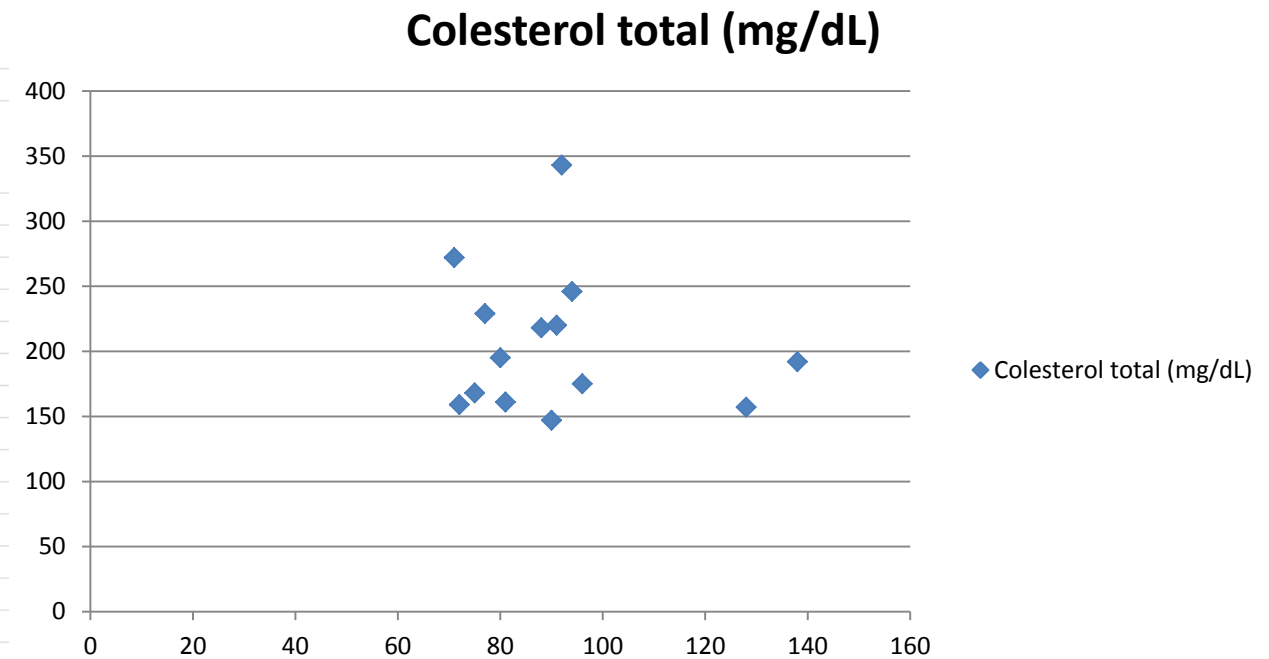
5		
6		
7		
8		
9		
10		
11		
12		
13		
14		
15		
16		
17		
18		
19		
20		
21		

	Glicemie (mg/dL)	Colesterol total (mg/dL)
7	75	168
8	92	343
9	77	229
10	128	157
11	81	161
12	138	192
13	88	218
14	72	159
15	71	272
16	80	195
17	91	220
18	94	246
19	90	147
20	96	175

5		
6		
7		
8		
9		
10		
11		
12		
13		
14		
15		
16		
17		
18		
19		
20		
21		

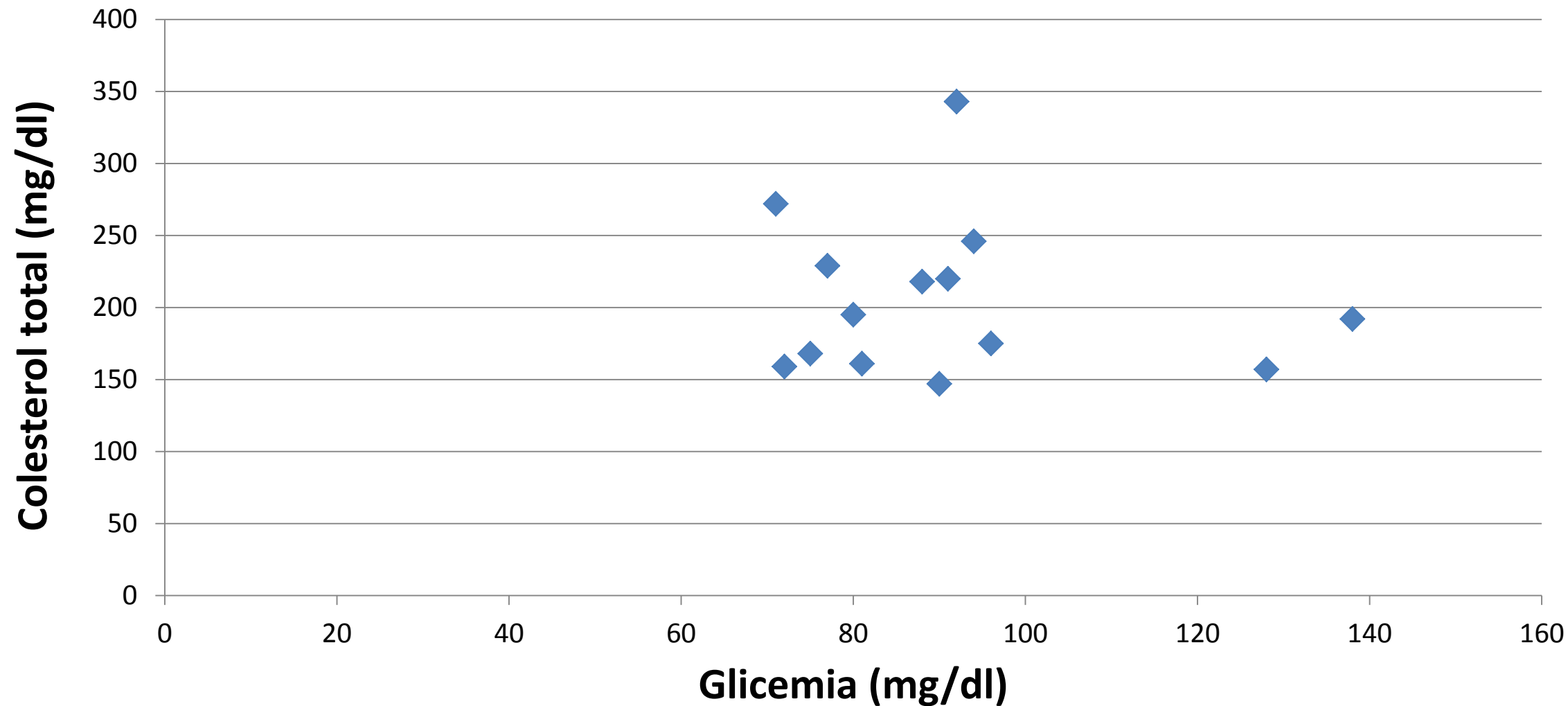
	Glicemie (mg/dL)	Colesterol total (mg/dL)
7	75	168
8	92	343
9	77	229
10	128	157
11	81	161
12	138	192
13	88	218
14	72	159
15	71	272
16	80	195
17	91	220
18	94	246
19	90	147
20	96	175

Selectăm tot tabelul

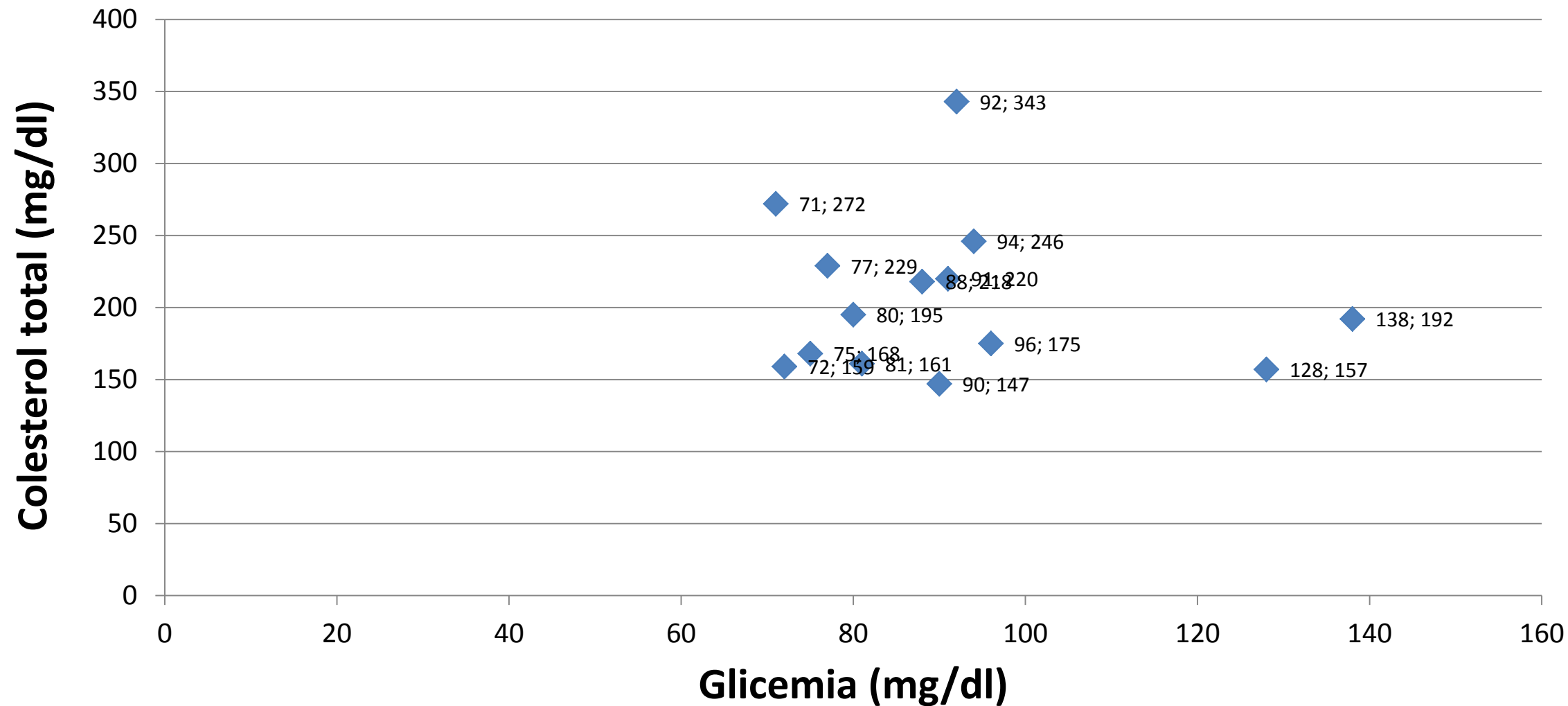


Inserăm grafic Scatter

Corelatia dintre colesterol si glicemie

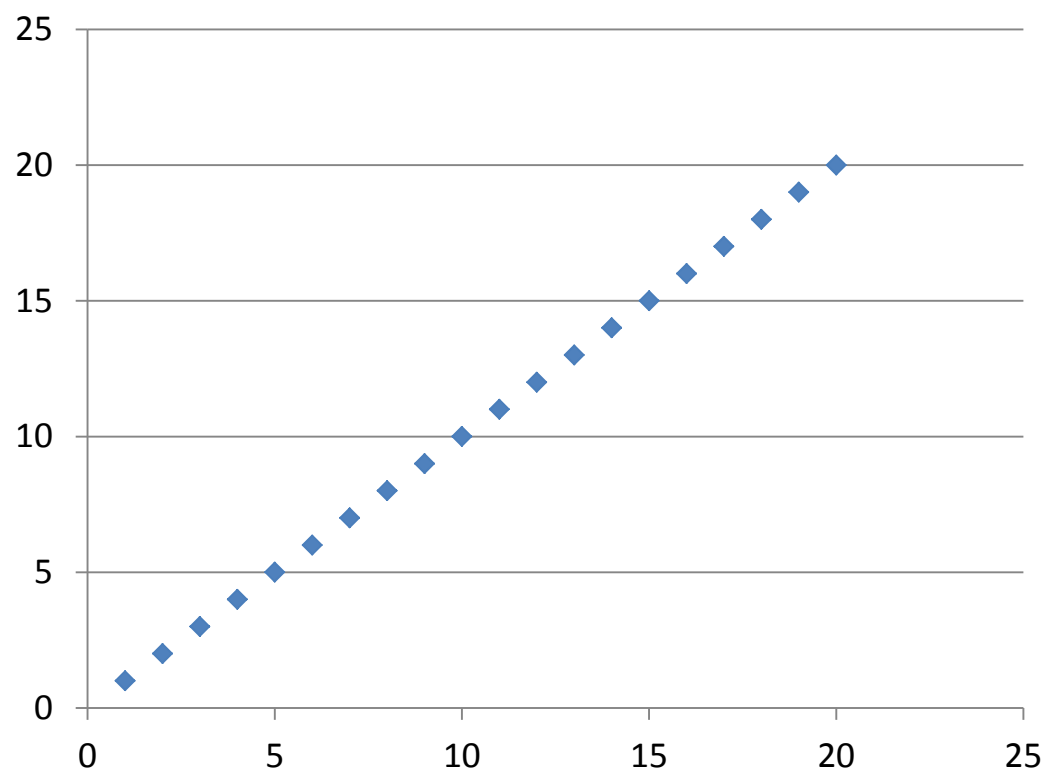


Corelatia dintre colesterol si glicemie

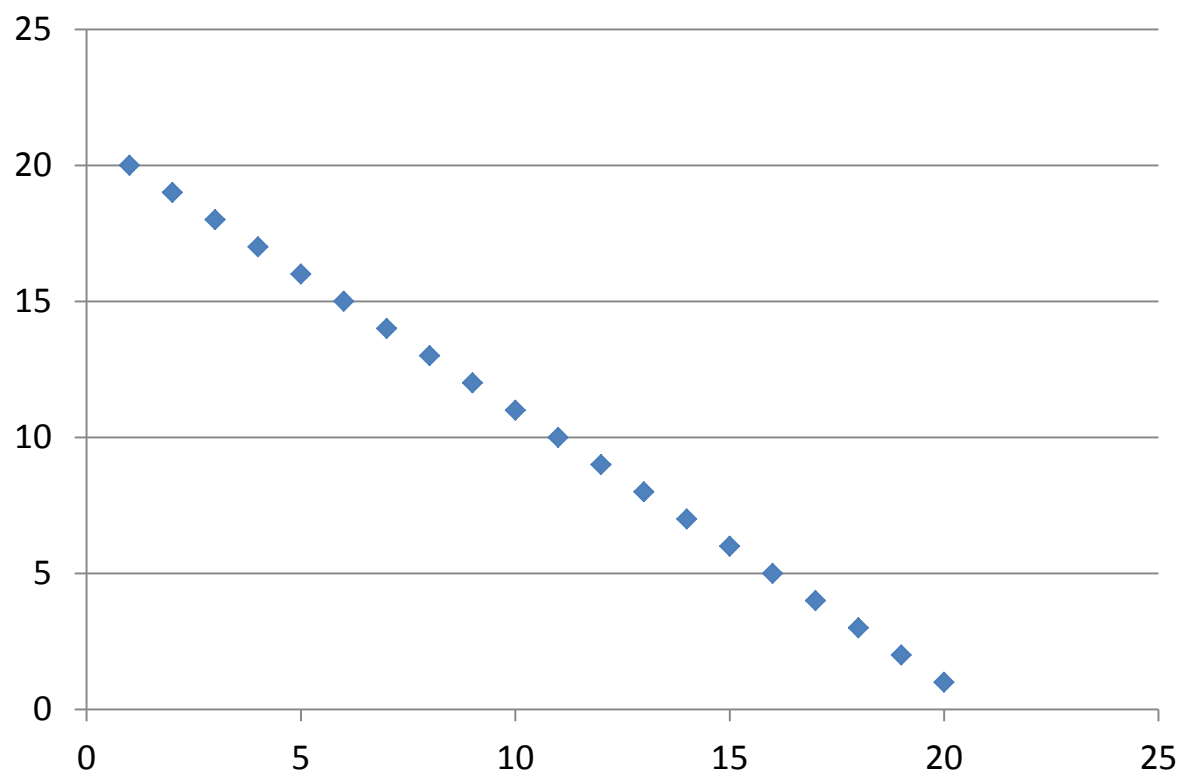


Corelația perfectă

- $r=1$



$r=-1$



Coeficient de corelație Pearson

- Dacă $r > 0$ atunci relația dintre X și Y este **pozitivă** (directă)
(valorilor mici ale lui X – le corespund valori mici ale lui Y ,
valorilor mari ale lui X – le corespund valori mari ale lui Y)
- Dacă $r < 0$ atunci relația dintre X și Y este **negativă** (inversă)
(valorilor mici ale lui X – le corespund valori mari ale lui Y ,
valorilor mari ale lui X – le corespund valori mici ale lui Y)

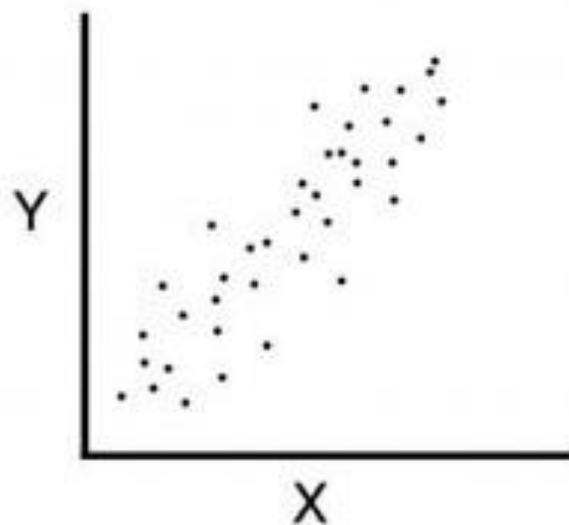


Graficul de corelație – XY Scatter

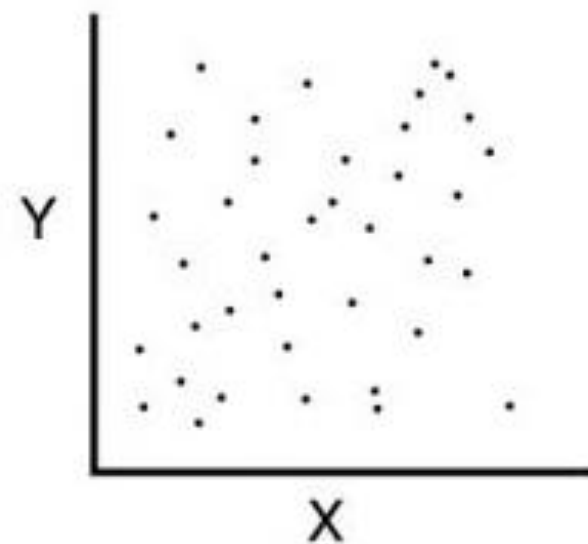
- Relație imperfectă, directă, pozitivă



Corelație puternică
 $r \approx 1$



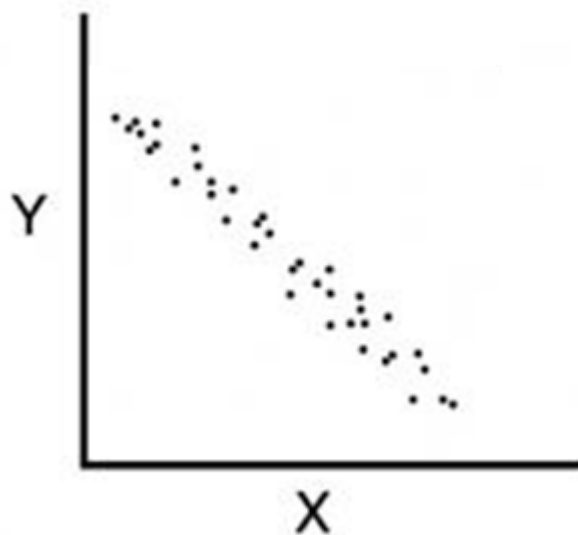
Corelație moderată
 $r \approx 0.5$



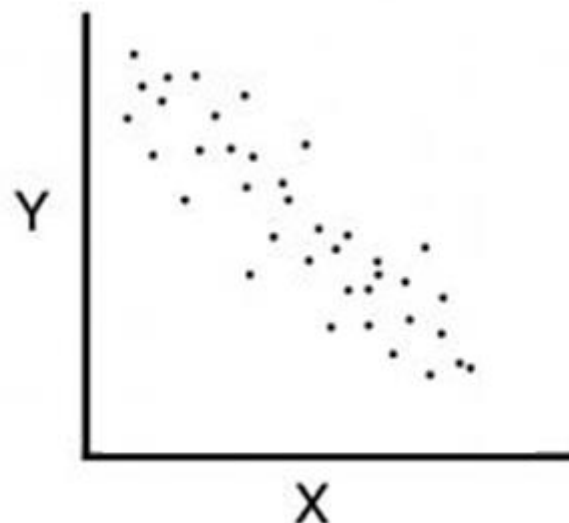
Corelație nulă
 $r \approx 0$

Graficul de corelație – XY Scatter

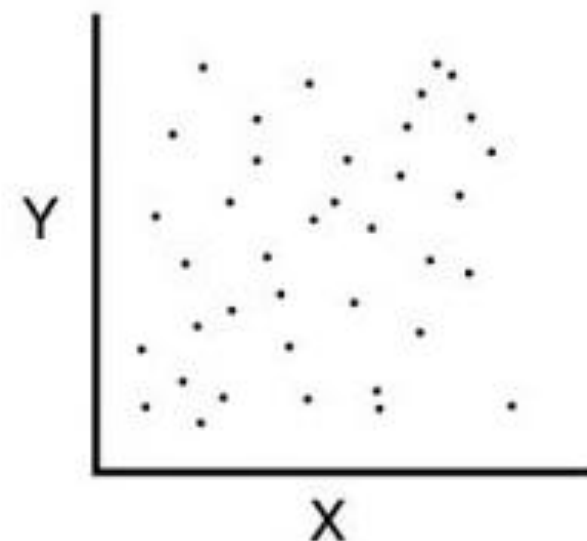
- Relație imperfectă, inversă, negativă



Corelație puternică
 $r \approx -1$



Corelație moderată
 $r \approx -0.5$



Corelație nulă
 $r \approx 0$

Regulile lui Colton

- [Colton T. Statistics in Medicine. Little Brown and Company, New York, NY 1974]:
 - $R \in [-0,25 ; +0,25] \rightarrow$ Nu există asociere spre asociere slabă
 - $R \in (0,25 ; +0,50] \cup (-0,25 ; -0,50] \rightarrow$ Asociere acceptabilă
 - $R \in (0,50 ; +0,75] \cup (-0,50 ; -0,75] \rightarrow$ Asociere moderată sau bună
 - $R \in (0,75 ; +1] \cup (-0,75 ; -1] \rightarrow$ Asociere foarte bună



Vom reveni la
această noțiune în
cursurile următoare

Exemplu

- Ex. Corelația dintre Greutate și densitatea minerală osoasă ($r=0.43$, $p<0.001$) la pacienții cu artroplastie bilaterală totală a genunchiului
- **Pozitivă** – Greutatea mare – corespunde cu valori mari de densitate osoasă (DMO crește cu greutatea subiectului)

Regulile lui Colton → **relație acceptabilă**

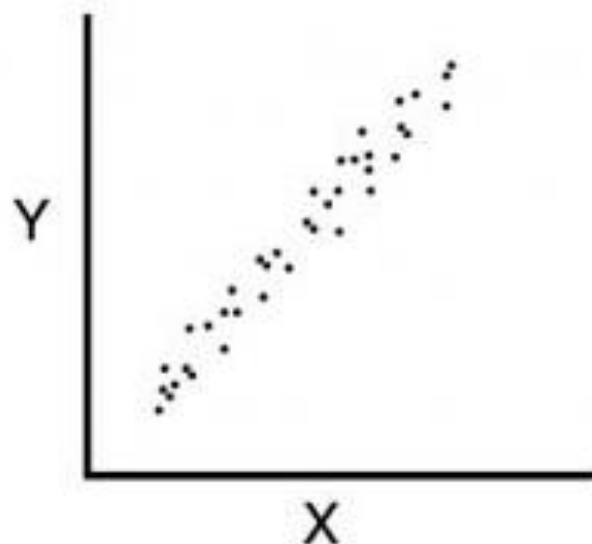
Vom reveni la
această noțiune în
cursurile următoare

Exemplu

- Ex. Corelația dintre Vârsta pacientului și densitatea minerală osoasă ($r = -0.32$, $p = 0.007$) la pacienții cu artroplastie bilaterală totală a genunchiului
- **Pozitivă** – Vârstă mare – corespunde cu valori mici de densitate osoasă (DMO scade cu vârsta)

Regulile lui Colton → **relație acceptabilă**

Graficul de corelație – XY Scatter



Interpretare



Coeficient de corelație $r > 0$

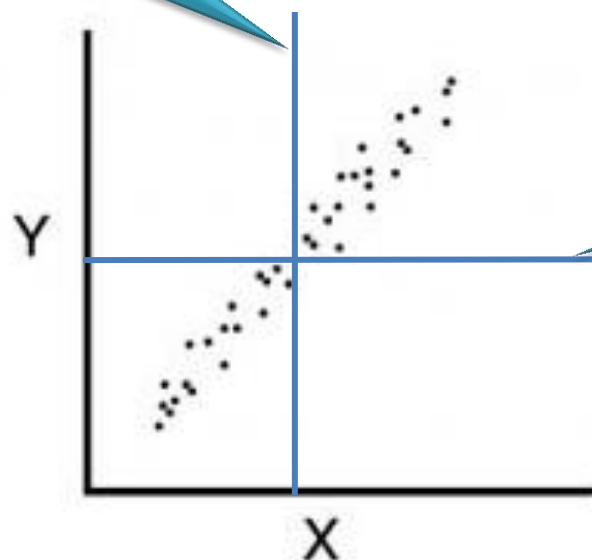
X crește, crește și Y, dacă r este mare (aproape de 1)

Valori mici ale X corelează cu valori mici ale Y

Valori crescute ale X corelează cu valori crescute ale Y

Media lui X

Graficul de corelație – XY Scatter



Media lui Y

Interpretare

Coeficient de corelație $r > 0$

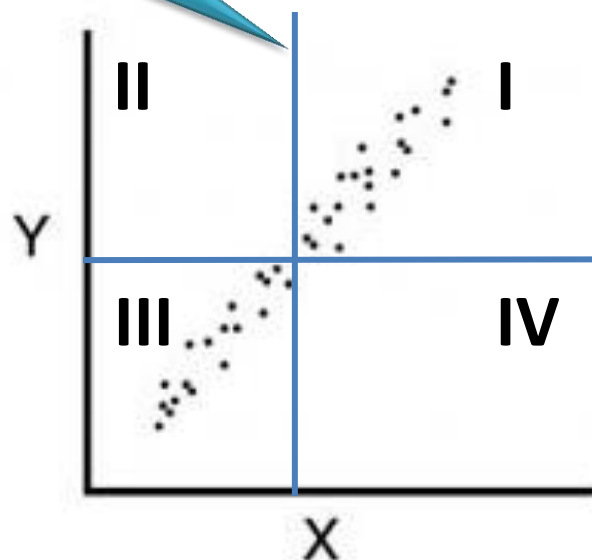
X crește, crește și Y, dacă r este mare (aproape de 1)

Valori mici ale X corelează cu valori mici ale Y

Valori crescute ale X corelează cu valori crescute ale Y

4 cadrane

Graficul de corelație – XY Scatter



Interpretare

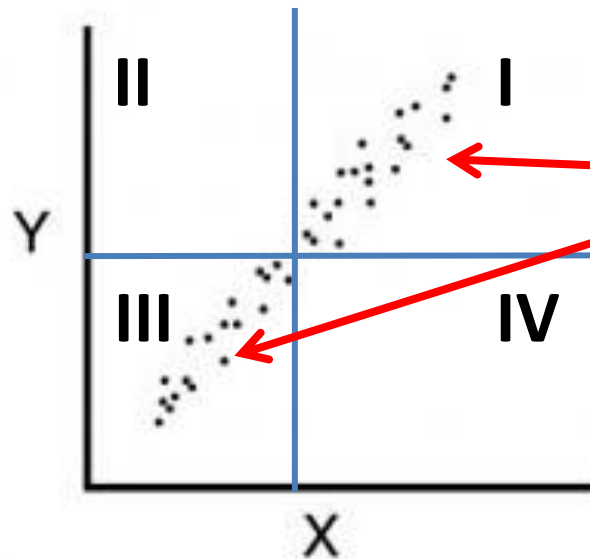
Coeficient de corelație $r > 0$

X crește, crește și Y, dacă r este mare (aproape de 1)

Valori mici ale X corelează cu valori mici ale Y

Valori crescute ale X corelează cu valori crescute ale Y

Graficul de corelație – XY Scatter



Când punctele sunt în
cadrantul I și III atunci
corelația este puternică

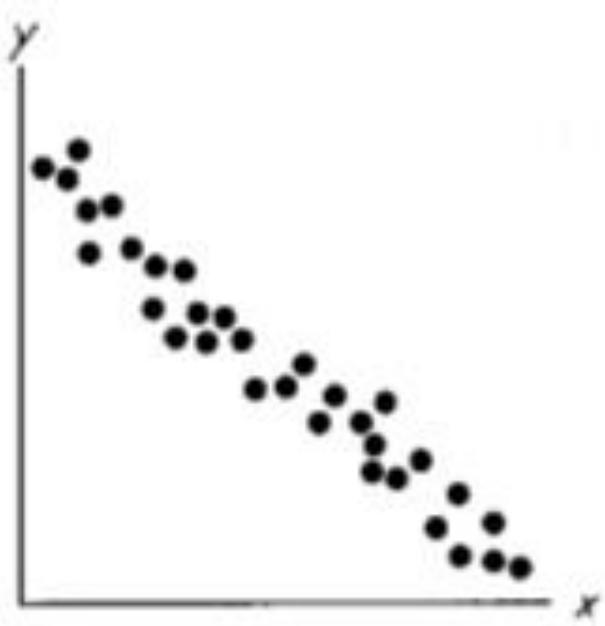
Coeficient de corelație $r > 0$

X crește, crește și Y, dacă r este mare (aproape de 1)

Valori mici ale X corelează cu valori mici ale Y

Valori crescute ale X corelează cu valori crescute ale Y

Graficul de corelație – XY Scatter



Interpretare

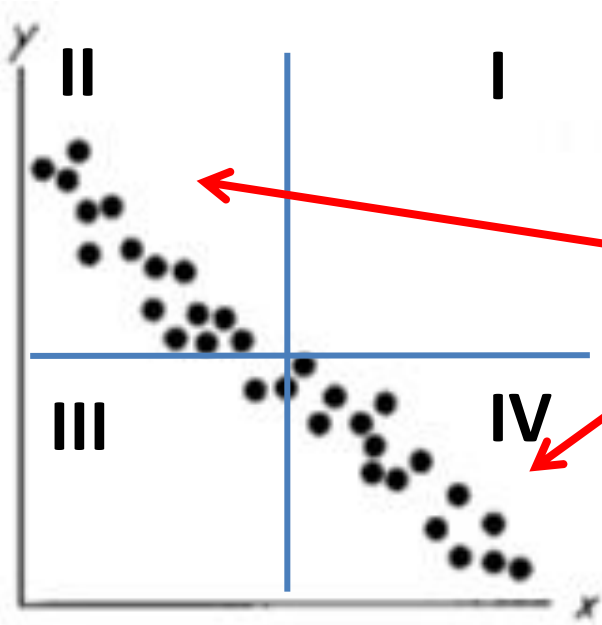
Coeficient de corelație $r < 0$

X crește și Y descrește, dacă r este mic (aproape de -1)

Valori mici ale X corelează cu valori crescute ale Y

Valori crescute ale X corelează cu valori mici ale Y

Graficul de corelație – XY Scatter



Când punctele sunt în
cadrantul II și IV atunci
corelația este puternică

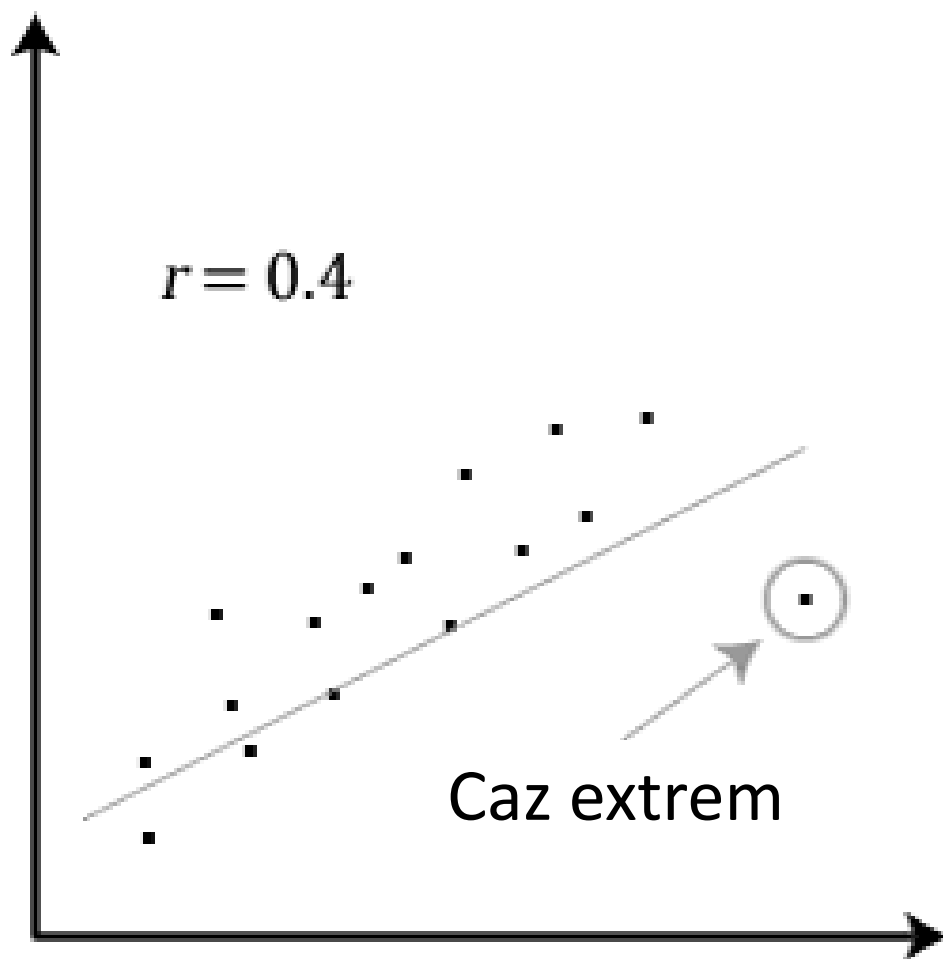
Coeficient de corelație $r < 0$

X crește și Y descrește, dacă r este mic (aproape de -1)

Valori mici ale X corelează cu valori crescute ale Y

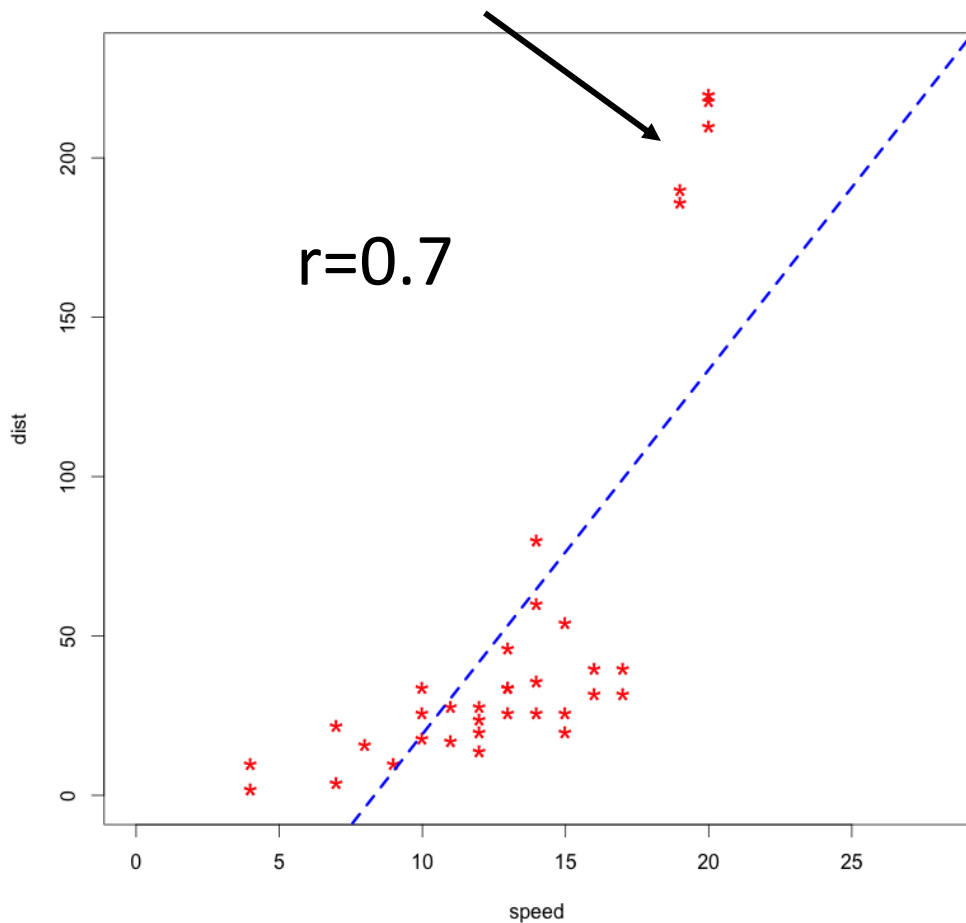
Valori crescute ale X corelează cu valori mici ale Y

Efectul cazurilor extreme asupra r – coeficientul de corelație Pearson

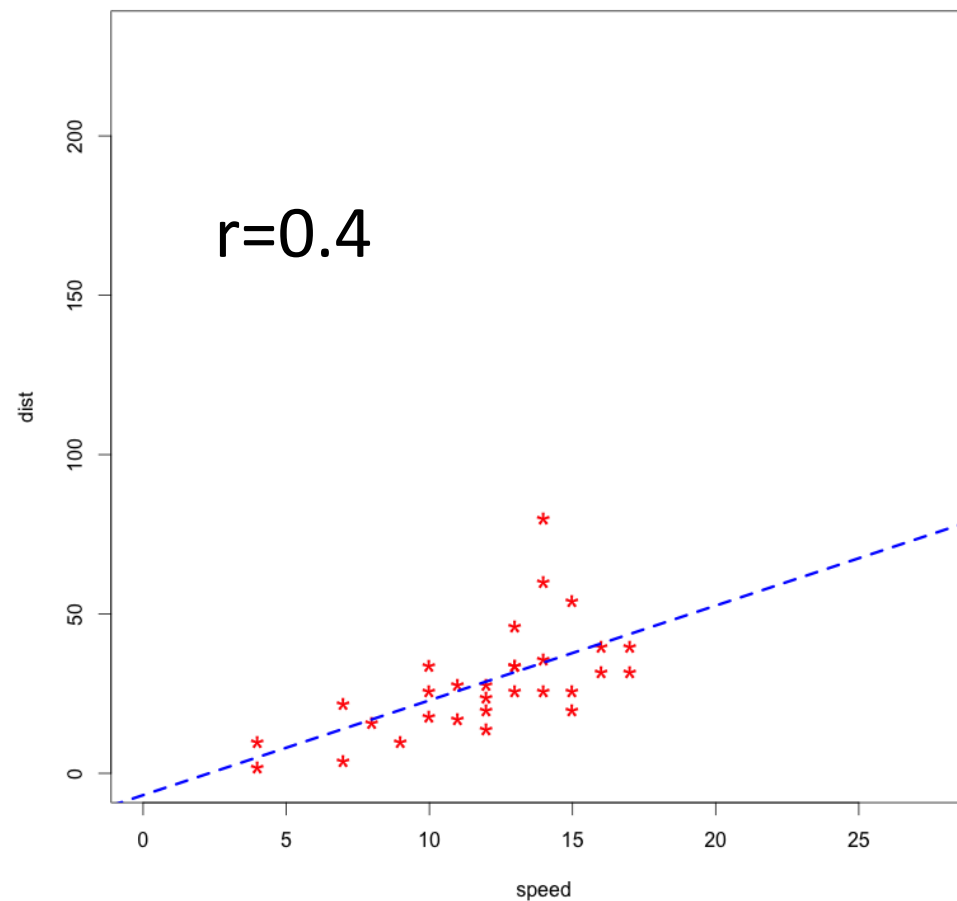


Efectul cazurilor extreme asupra r – coeficientul de corelație Pearson

Cazuri extreme



Fără cazurile extreme



Mărimea efectului – coeficientul de determinare d

$$d=r^2$$



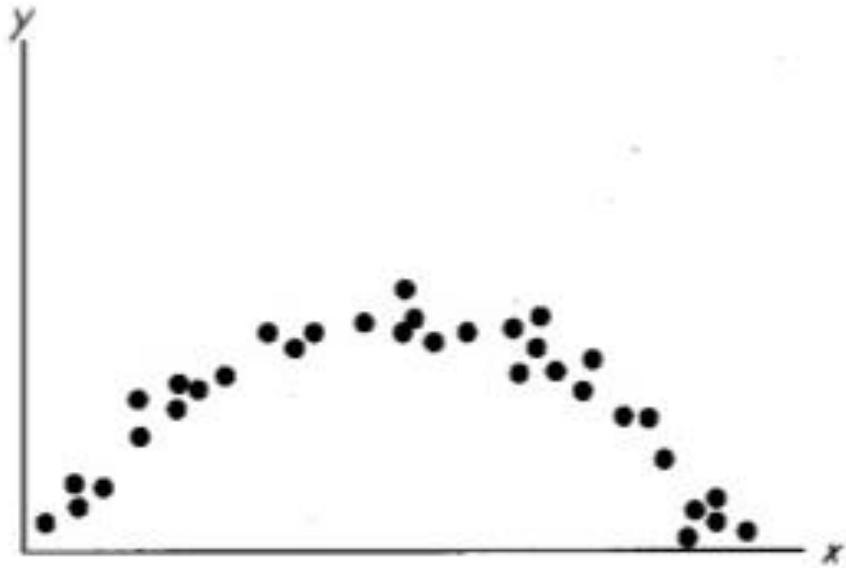
Ex. $r=0.5$, atunci $d=0.25$

25% din variația lui Y poate fi explicat prin relația liniară cu X

Cerințe pentru calcularea coeficientului de corelație Pearson

- X și Y să nu aibă cazuri extreme
- X și Y să fie variabile cantitative
- Relația să fie liniară
- Soluții posibile:
 - Transformări: logaritmare, normalizare
 - Scoaterea cazurilor extreme din analiză
 - calcularea coeficientului de corelație Spearman

Graficul de corelație – XY Scatter



Relație neliniară

- Nu există o relație liniară între X și Y – asocierea respectivă nu este adecvată pentru calcularea coeficientului de corelație Pearson, ar trebui să calculăm coeficientul de corelație Spearman

Coeficientul de corelație Spearman

- Fie R_X și R_Y rangurile a două variabile X și Y :

$$r_{XY} = \frac{\sum_{i=1}^n (R_{X,i} - \overline{R_X})(R_{Y,i} - \overline{R_Y})}{\sqrt{\sum_{i=1}^n (R_{X,i} - \overline{R_X})^2} \sqrt{\sum_{i=1}^n (R_{Y,i} - \overline{R_Y})^2}}$$

- Poate fi calculat și pentru două variabile ordinale

Coeficientul de corelație Spearman

- Indică asocierea dintre X și Y
- Întotdeauna între -1 și 1

$$r \in [-1,1]$$

când $|r|$ se apropie de 1 asocierea este mai mare

când $|r|$ se apropie de 0 asocierea este mai mică

! Vorbim despre relație neliniară



Coeficientul de corelație Spearman

- Dacă $r > 0$ atunci asocierea dintre X și Y este **pozitivă** (directă, neliniară)
- Dacă $r < 0$ atunci asocierea dintre X și Y este **negativă** (inversă, neliniară)
- Se interpretează după regulile lui Colton



Regresie

- Predicție = regresie – modelarea matematică a relației dintre Y și X
- univariată - modelul cu o singură variabilă independentă X
- multivariată - modelul cu mai multe variabilă independentă X_i , $i=2,n$

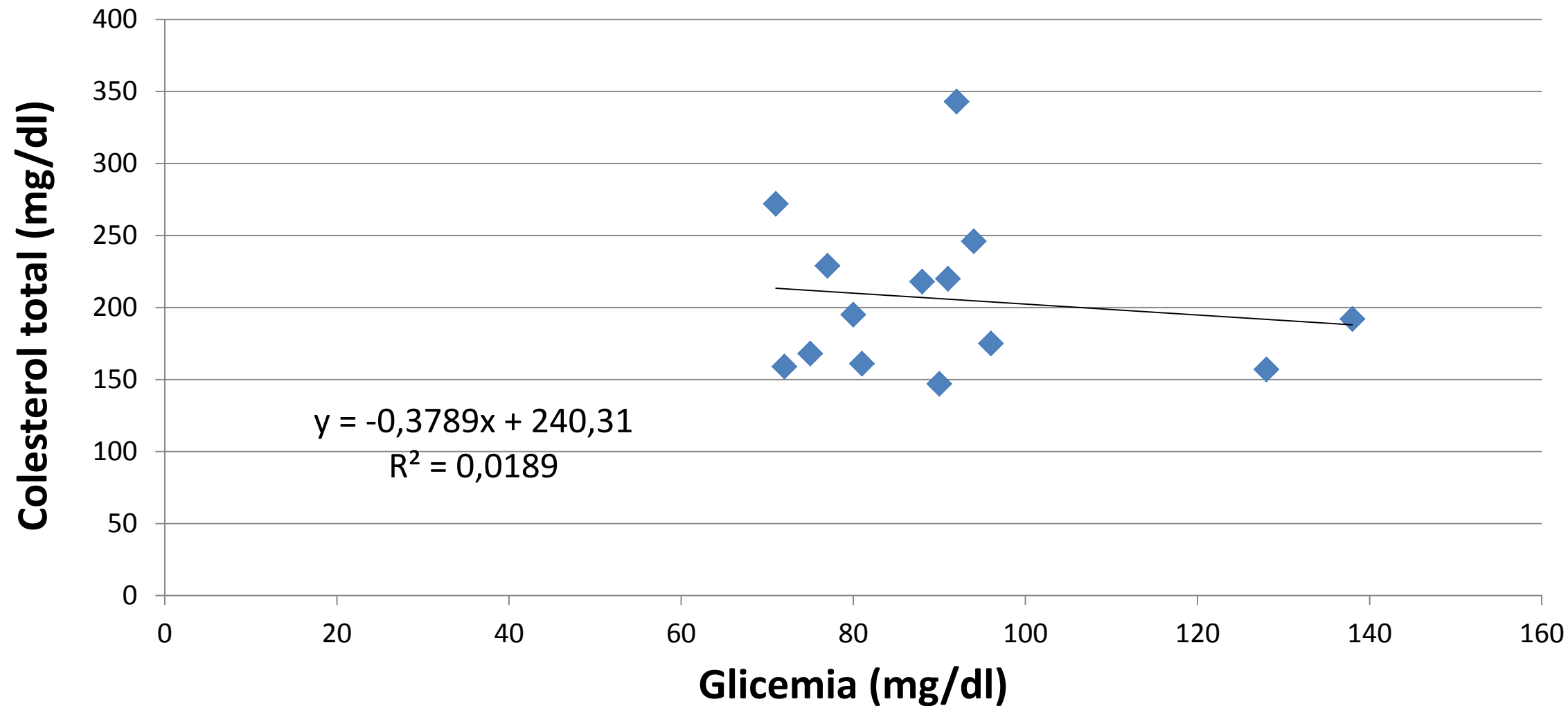
Regresie liniară

- Regresie liniară - funcția de predicție liniară: $Y=aX+b$
- **Metoda celor mai mici pătrate** - pentru a obține ecuația matematică - modalitatea de a determina ecuația celei mai potrivite linii pentru norul de date
- Punctul în care linia traversează axa Y este notat cu a și panta liniei cu b, ecuația drepte:

$$Y=aX+b$$



Corelatia dintre colesterol si glicemie



Regresie liniară ($Y=aX+b$)

- Y crește cu a de câte ori X crește cu 1

$$Y = 3X + 5$$

- $X=1$, atunci $Y = 3 \cdot 1 + 5 = 8$
- $X=2$, atunci $Y = 3 \cdot 2 + 5 = 11$
- $X=3$, atunci $Y = 3 \cdot 3 + 5 = 14$
-
- Dacă a este pozitiv, atunci când X crește, Y crește

Regresie liniară ($Y=aX+b$)

- Y crește cu $-a$ de câte ori X crește cu 1

$$Y = -3X + 5$$

- $X=1$, atunci $Y = -3 \cdot 1 + 5 = 2$
- $X=2$, atunci $Y = -3 \cdot 2 + 5 = -1$
- $X=3$, atunci $Y = -3 \cdot 3 + 5 = -4$
-
- Dacă a este negativ, când X crește Y descrește

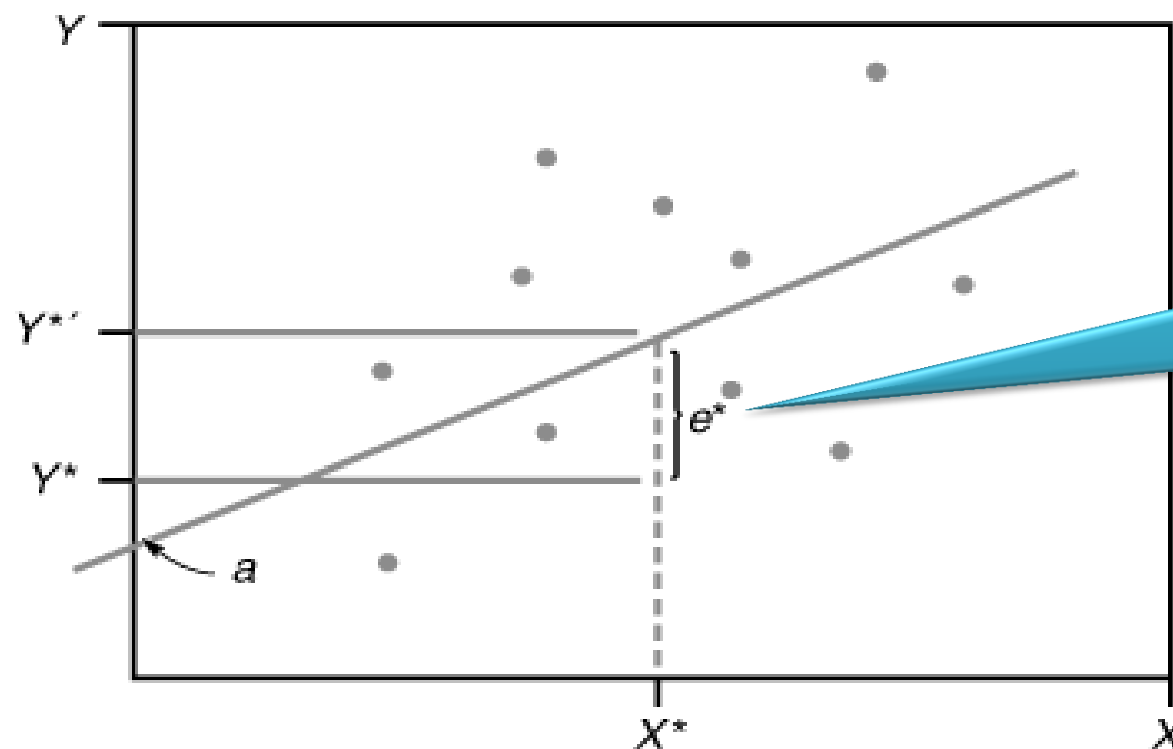
- Formula pentru coeficienții liniei de regresie:

$$b = \frac{\Sigma(X - \bar{X})(Y - \bar{Y})}{\Sigma(X - \bar{X})^2}$$

$$a = \bar{Y} - b\bar{X}$$

Regresie liniară ($Y=aX+b$)

- Când modelul liniar nu se potrivește perfect (nu toate punctele sunt pe linia de regresie) --> avem o eroare



Eroarea pentru
punctul
respectiv

Cum realizăm predicții?

- Avem nevoie de ecuația de regresie

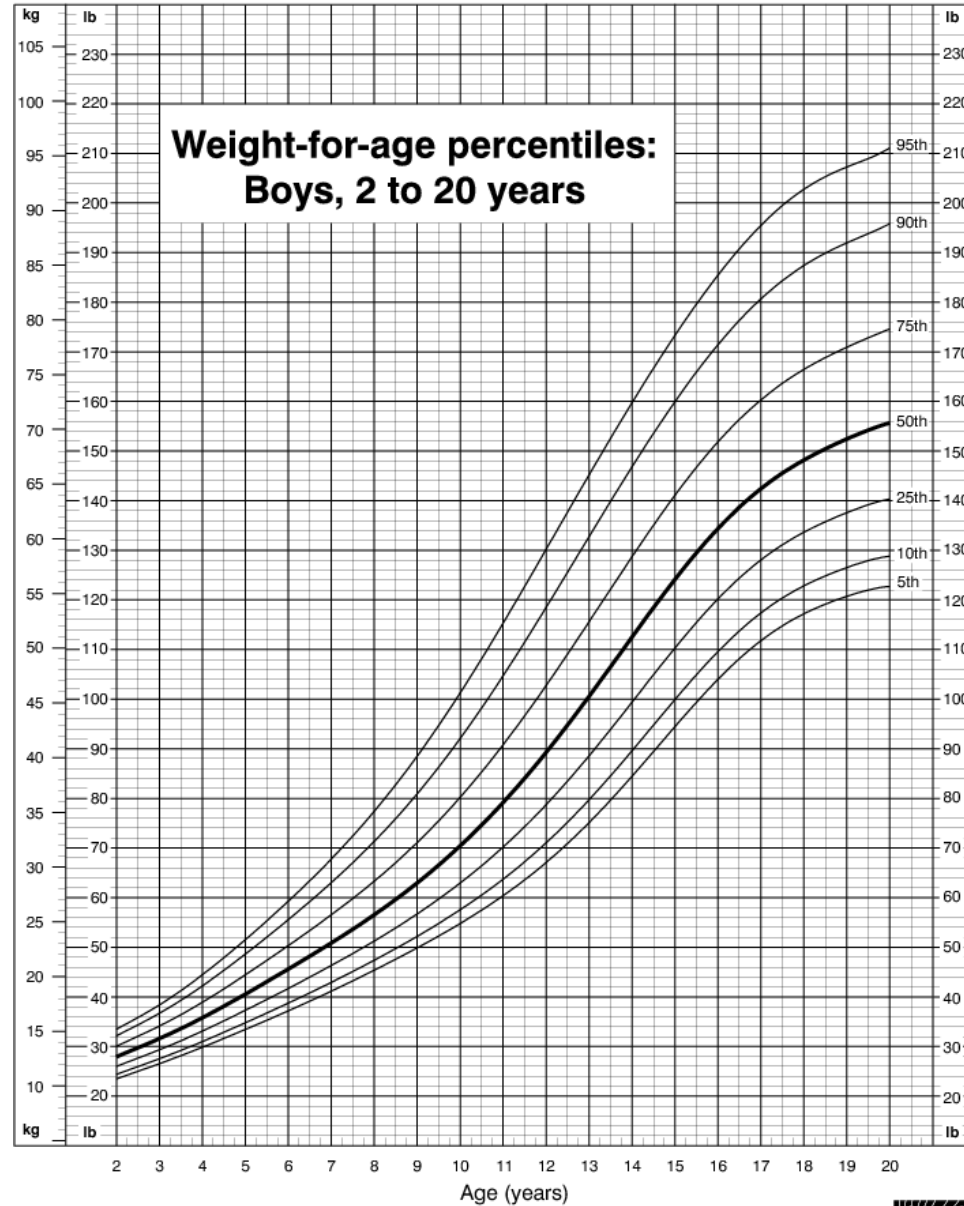
Cum putem obține ecuația de regresie?

- Ex. Greutatea crește cu vârsta
- Greutatea variază de la 1 la 200
- Vârsta variază de la 0 la 110
- Considerăm vârsta. Pentru fiecare valoare a vârstei luăm minim 3 indivizi. Ca să ne asigurăm că avem minim 3 indivizi pe fiecare valoare a vârstei selectăm la întâmplare în loc de 330 de indivizi: 3000. Obținem greutatea fiecărui individ.

Ce constatăm?

- Femeile și bărbații au greutate diferite
- Copiii și adolescenții au dependențe diferite față de adulți
- Este nevoie să renunțăm la unii indivizi
- Relația nu este liniară
- Ne focusăm numai pe adolescenți și copii, băieți

CDC Growth Charts: United States



Published May 30, 2000.

SOURCE: Developed by the National Center for Health Statistics in collaboration with the National Center for Chronic Disease Prevention and Health Promotion (2000).



SAFER • HEALTHIER • PEOPLE™

Acum să realizăm predicții

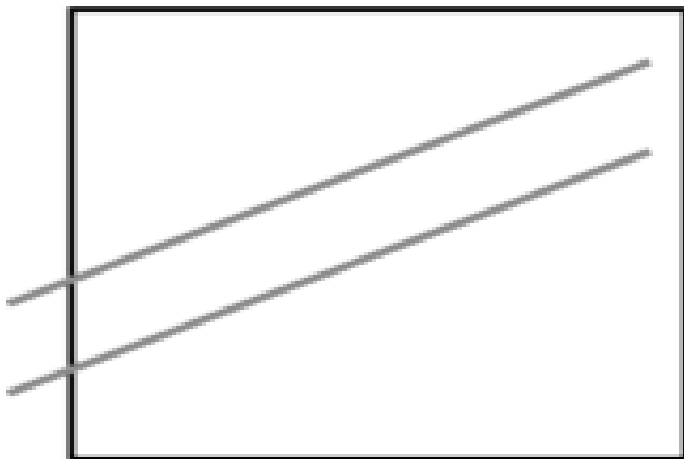
- Câte kg va avea un băiat la 10 ani?
 - Între 29 kg și 36 kg cu 50% eroare
 - Între 24 kg și 46 kg cu 5% eroare
- Câte kg va avea un băiat la 5 ani?
 - Între 17 kg și 19 kg cu 50% eroare
 - Între 15 kg și 24 kg cu 5% eroare

Regresie liniară

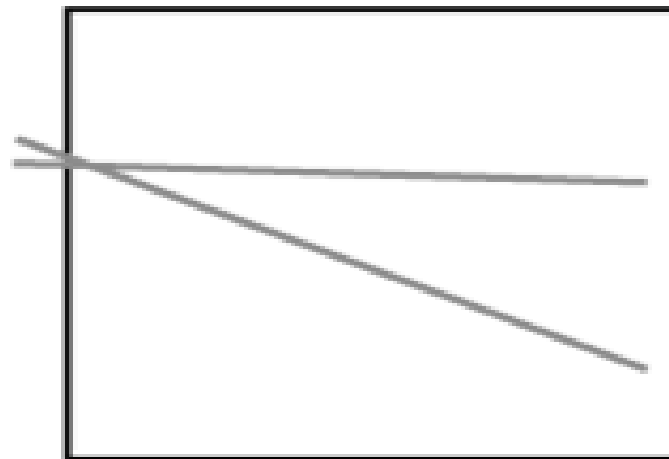
- Eroarea standard – media erorilor de predicție pentru fiecare punct în parte

$$ES_{X,Y} = \sqrt{\frac{\sum_{i=1}^n (Y_{i-predicted} - Y_{i-observed})^2}{n-2}}$$

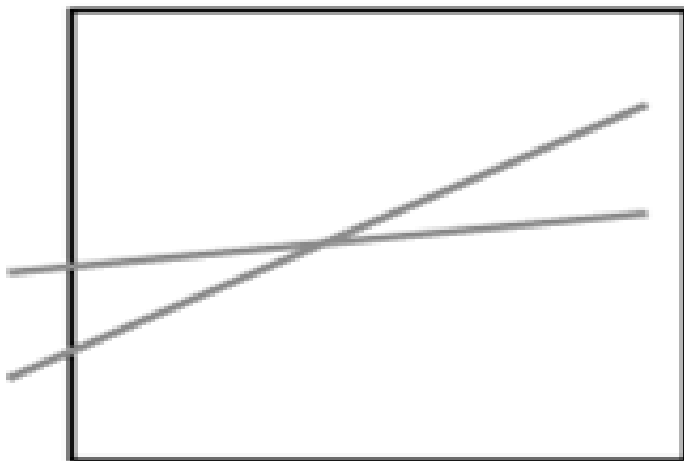
Dacă repetăm studiul, linia de regresie este diferită



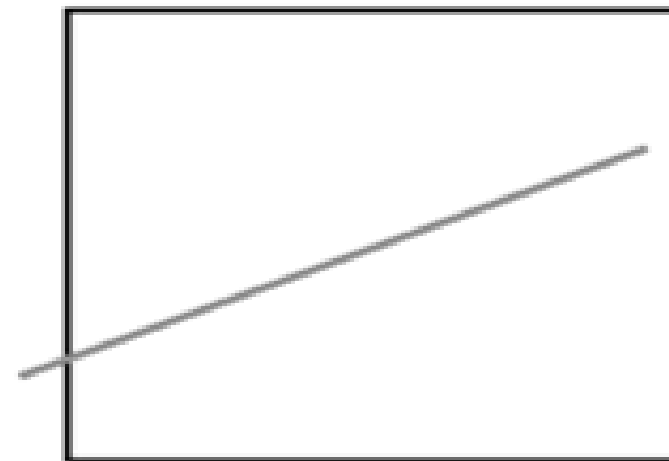
A



B



C



D



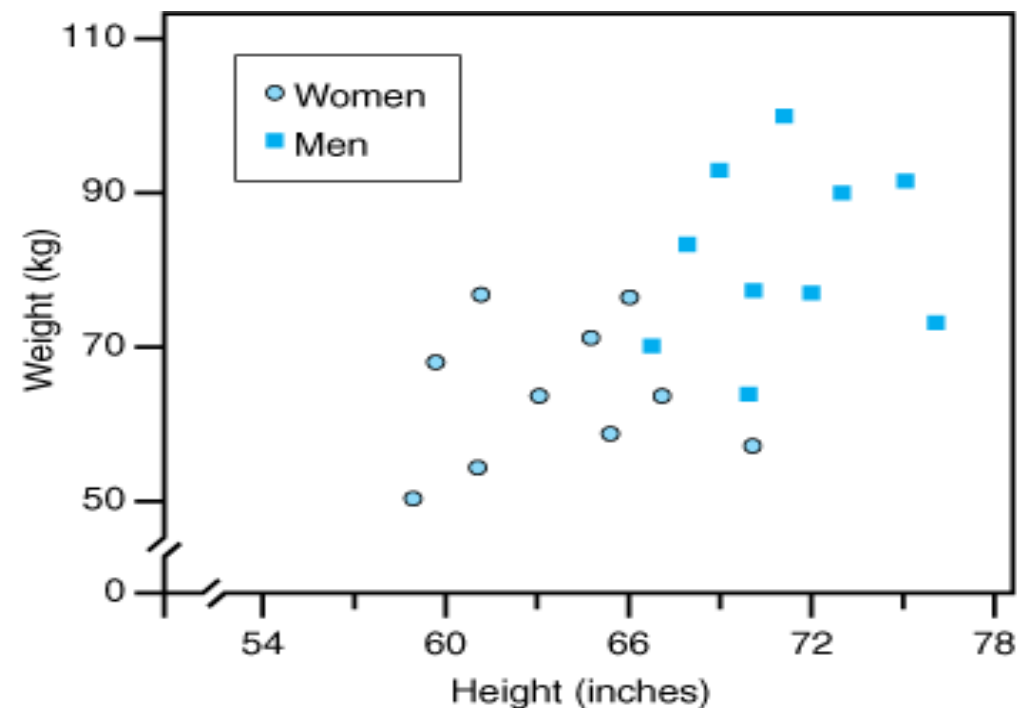
- A. Pantă egală, constantă diferită; B. Pantă diferită, constantă egală;
- C. Pantă diferită și constantă diferită; D. Aceeași constantă și pantă

Relație neliniară - ce putem face?

- Transformarea datelor
- Regresie polinomială
- Regresie pe mai multe secțiuni (intervale)

Erori în utilizarea regresiei

- Amestecarea a două populații diferite
- Amestecarea datelor rezultate din măsurători repetate pe aceeași axă (Y nu are toate datele independente)



Coeficientul de corelație sau regresia liniară ?

- Coeficientul de corelație este independent de unitatea de măsură, regresia liniară nu
- Atunci când greutatea este măsurată în kg sau grame, linia de regresie are un coeficient diferit, în timp ce coeficientul de corelație este același
- Coeficientul de corelație și coeficientul liniei de regresie au același semn (+ sau -)

- Nu are sens calcularea coeficientul de corelație între X și Y dacă Y este calculat din X
- Ex. IMC este corelat cu înălțimea și greutatea (a fost calculat pe baza lor - avem deja ecuația nu mai este necesar să o calculăm)
- $X \cdot Y \cdot Z$ se corelează cu Z
- X / Z se corelează cu Y / Z

Corelația nu este cauzalitate

- În Cluj-Napoca s-a observat o corelare între creșterea numărului de fracturi ale picioarelor și creșterea numărului de filme vizionate
- Primarul urmează să interzică firmele furnizoare de televiziune prin cablu și internet

Corelația nu este cauzalitate

- De fapt există o altă variabilă care influențează cele două creșteri concomitente:
- frigul care se instalează odată cu venirea toamnei și a iernii, acesta duce
 - la apariția gheții care înmulțește numărul de fracturi
 - oamenii stau mai mult în casă și se uită la mai multe filme

Corelația nu este cauzalitate

- De la Ecuator la Pol scade temperatura, iar această scădere se corelează cu creșterea înălțimii persoanelor care locuiesc în zonele respective, astfel că s-a ajuns la concluzia că înălțimea persoanelor depinde de temperatura mediului

Corelația nu este cauzalitate

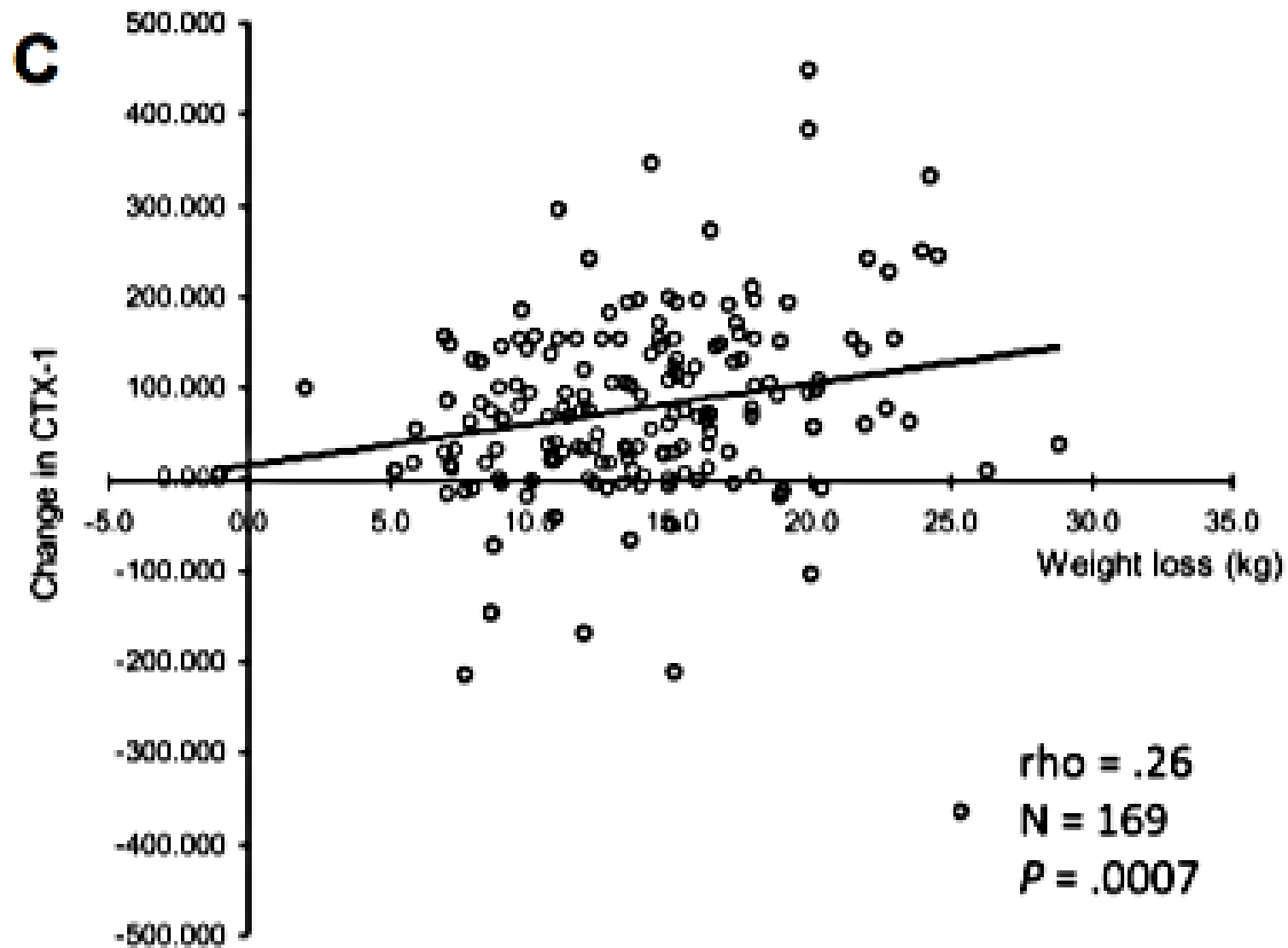
- La Ecuator locuiesc pigmeii - un trib de oameni de înălțime mică
- Europa este locuită de europeni - o populație de înălțime medie
- În nordul Europei trăiesc persoane mai înalte (urmașii vikingilor)
- O coincidență: înălțimea persoanelor depinde de temperatura mediului

- Corelația trebuie să aibă și sens.



Scenariu

- 192 obezi cu osteoartrita genunchiului (OA)
- OA evaluat cu un biomarker uCTX-I (Urine C-terminal telopeptide of collagen type I)
- Program intensiv de scădere în greutate
- Scopul studiului: scăderea în greutate se asociază cu creșterea biomarkerului uCTX-I
- uCTX-I a fost direct corelată cu scăderea în greutate ($r = 0.22$, $P = 0.0007$)



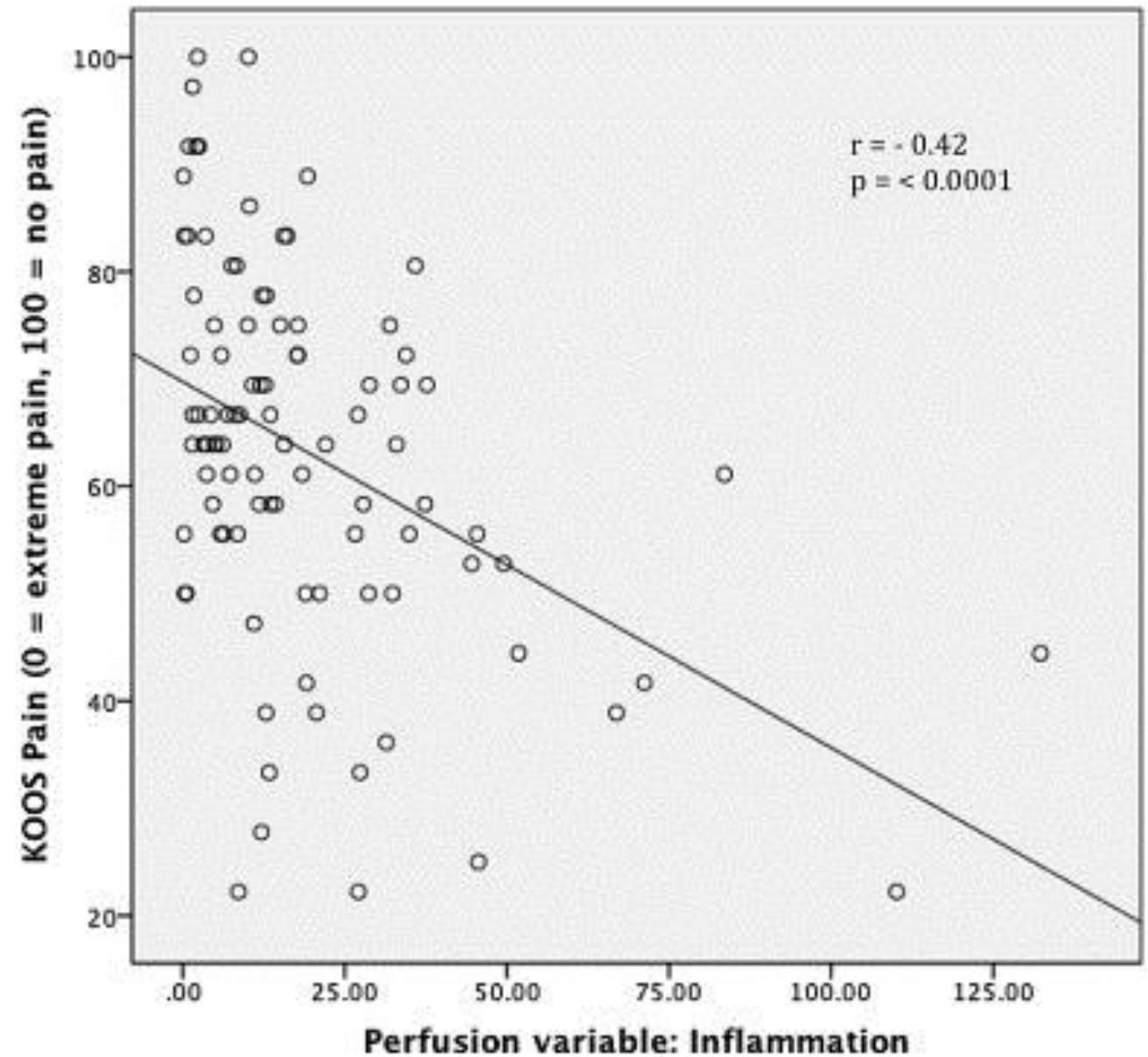
Bartels EM, Christensen R, Christensen P, Henriksen M, Bennett A, Gudbergson H, Boesen M, Bliddal H. Effect of a 16 weeks weight loss program on osteoarthritis biomarkers in obese patients with knee osteoarthritis: a prospective cohort study. *Osteoarthritis Cartilage*. 2014 Nov;22(11):1817-25.

Scenariu

- Scop - asocierea dintre durerea resimțită la nivelul genunchiului și semnele de inflamație în padul de grăsime infrapatelar (IPFP) la **pacienții obezi cu osteoartrita genunchiului**
- Inflamația în padul de grăsime infrapatelar evaluată cu RMN (CE-MRI și DCE-MRI)
- KOOS - scorul final al osteoartritei (durerea și alte simptome) 100 = fără durere, 0 = durere extremă
- 95 de pacienți

[Ballegaard C, Riis RG, Bliddal H, Christensen R, Henriksen M, Bartels EM, Lohmander LS, Hunter DJ, Bouert R, Boesen M. Knee pain and inflammation in the infrapatellar fat pad estimated by conventional and dynamic contrast-enhanced magnetic resonance imaging in obese patients with osteoarthritis: a cross-sectional study. Osteoarthritis Cartilage. 2014 Jul;22(7):933-40.]

- Coeficient de corelație Spearman
- $r = -0,42$



[Ballegaard C, Riis RG, Bliddal H, Christensen R, Henriksen M, Bartels EM, Lohmander LS, Hunter DJ, Bouert R, Boesen M. Knee pain and inflammation in the infrapatellar fat pad estimated by conventional and dynamic contrast-enhanced magnetic resonance imaging in obese patients with osteoarthritis: a cross-sectional study. *Osteoarthritis Cartilage*. 2014 Jul;22(7):933-40.]

- Mulțumesc!