

CORELAȚII SI REGRESII

LEGĂTURA ÎNTRE MAI MULTE
VARIABLE

Tipuri de “relații”

- ⦿ două sau mai multe variabile cantitative
- ⦿ două variabile de ordine
- ⦿ două variabile calitative
- ⦿ o variabilă cantitativă cu o variabilă calitativă

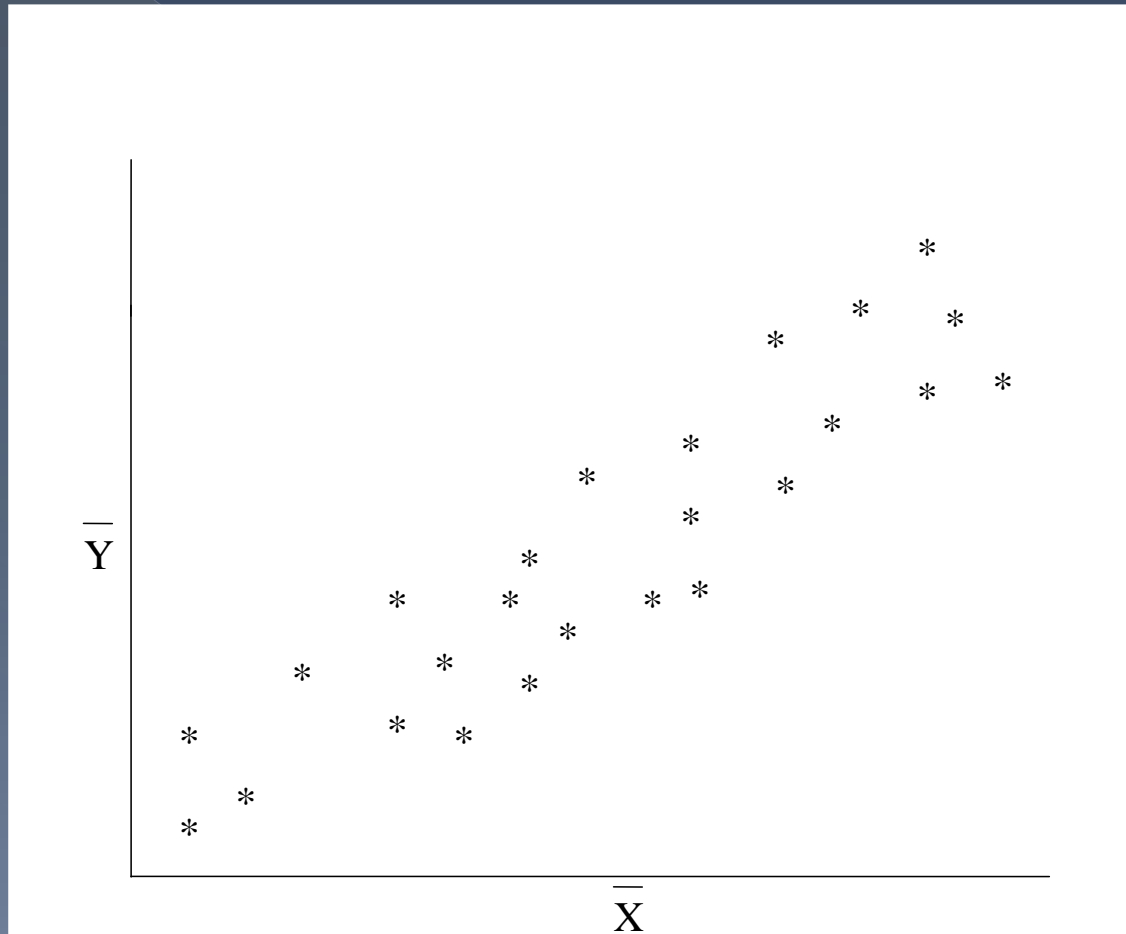
Variabile cantitative

Varsta $X: X_1, X_2, \dots, X_n$

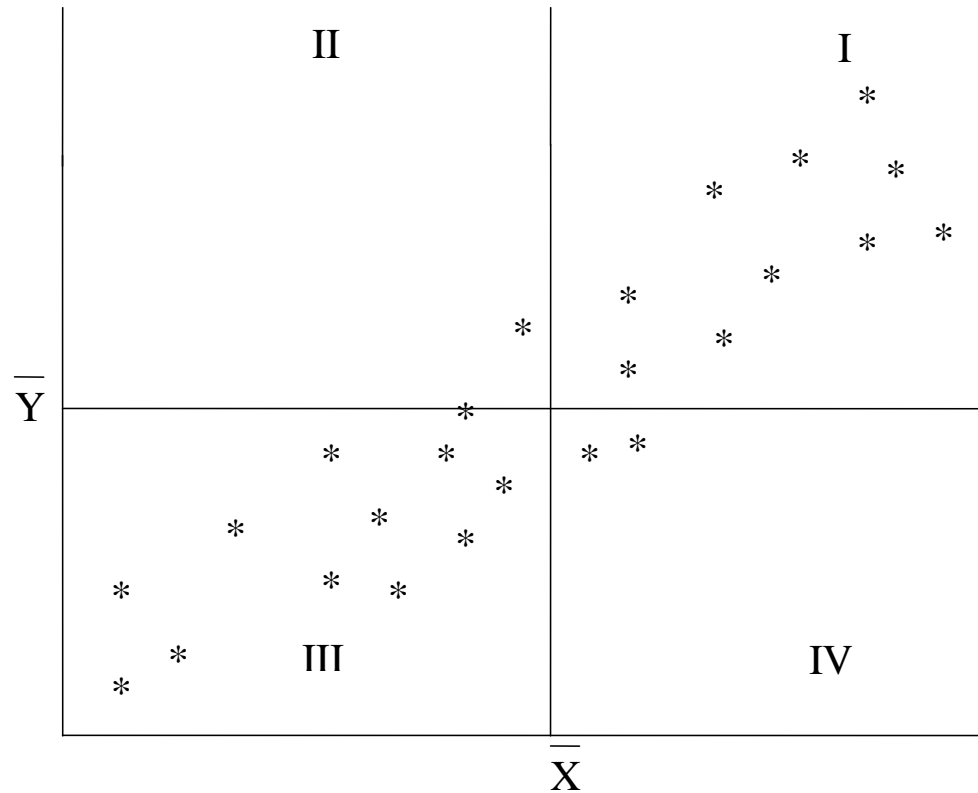
TAS $Y: Y_1, Y_2, \dots, Y_n$

- 1. Să se stabilească dacă există o legătură între variabilele X și Y (cantitative continue) și să se determine o modalitate de a măsura intensitatea acestei legături.
 - > Coeficientul de corelație
- 2. Să se stabilească dacă Y depinde de X și dacă da în ce formă se realizează această dependență.
 - > Funcția de regresie

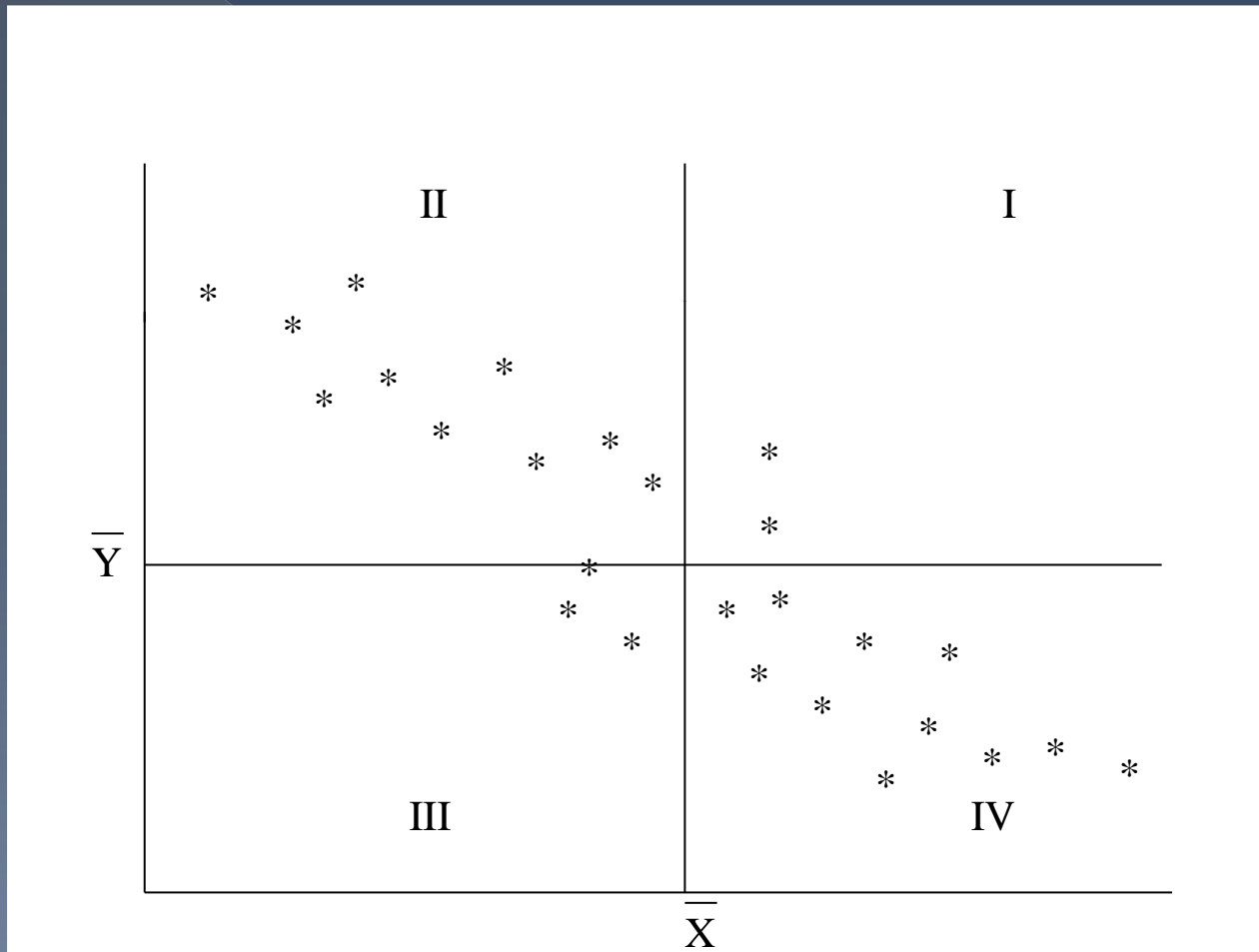
Statistici descriptive in două dimensiuni. Diagrama de dispersie



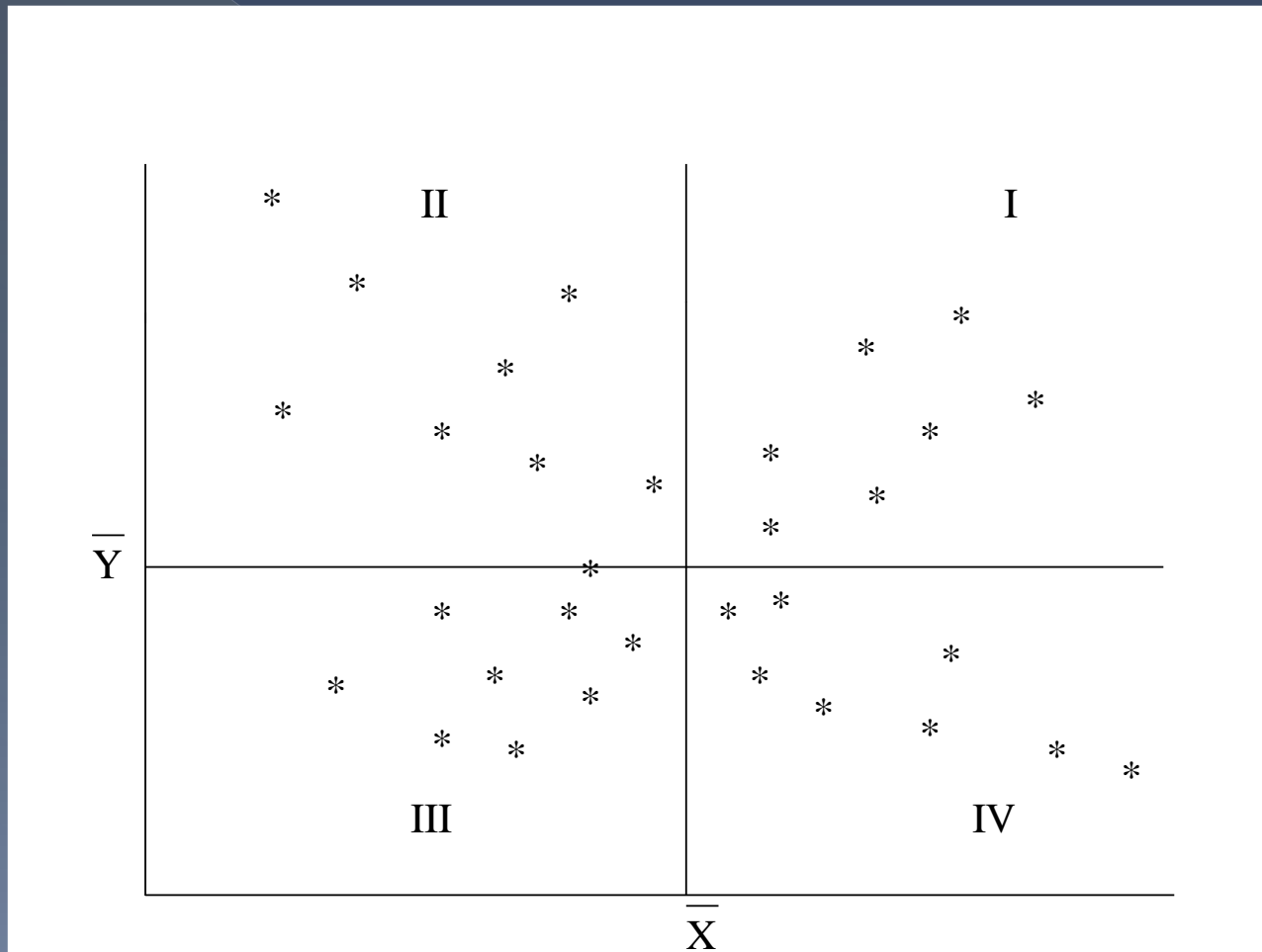
Statistici descriptive in două dimensiuni. Diagrama de dispersie



Statistici descriptive in două dimensiuni. Diagrama de dispersie



Statistici descriptive in două dimensiuni. Diagrama de dispersie



Indici de corelație. Suma produselor ecart

- Descrierea "intensității" relației dintre variabilele X și Y:

- > (X_i, Y_i) -în cadranele I sau III: $(X_i - \bar{X})(Y_i - \bar{Y}) \geq 0$
- > (X_i, Y_i) -în cadranele II sau IV: $(X_i - \bar{X})(Y_i - \bar{Y}) \leq 0$

$$SPE = \sum_{i=1}^n (X_i - \bar{X})(Y_i - \bar{Y})$$

Indici de corelație. Covarianța

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

Avantaje: mărime independentă față de volumul seriei statistice.

Indici de corelație.

Coeficientul de corelație

Pentru a obține un indicator independent și de unitățile de măsură ale celor două variabile se utilizează coeficientul de corelație sau coeficientul Bravais-Pearson:

$$r = \frac{COV(X, Y)}{S_x \cdot S_y}$$

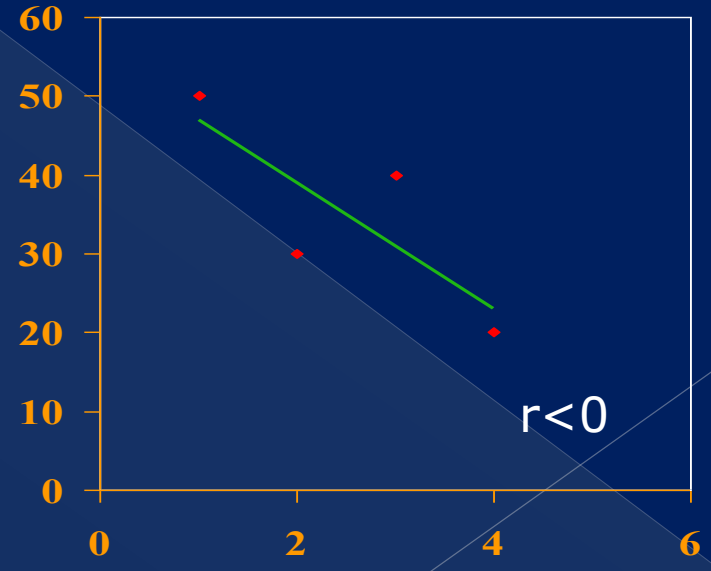
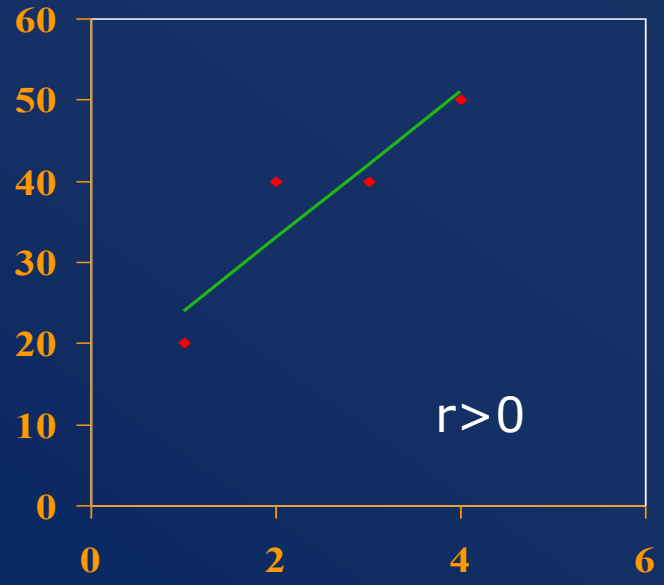
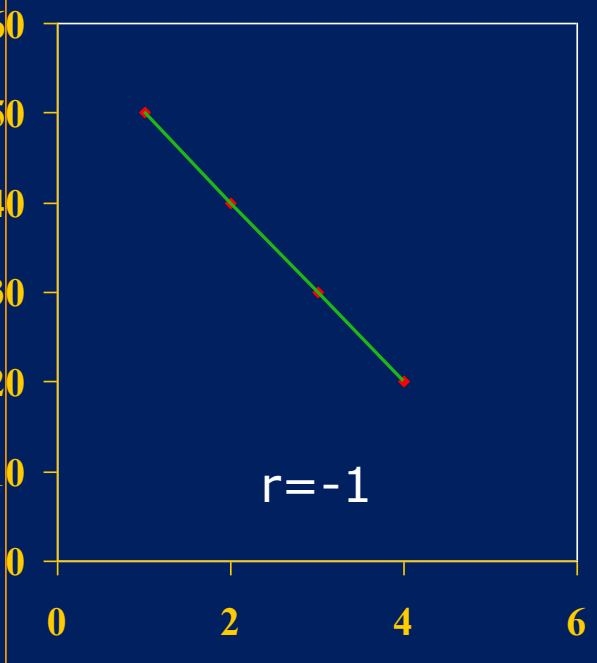
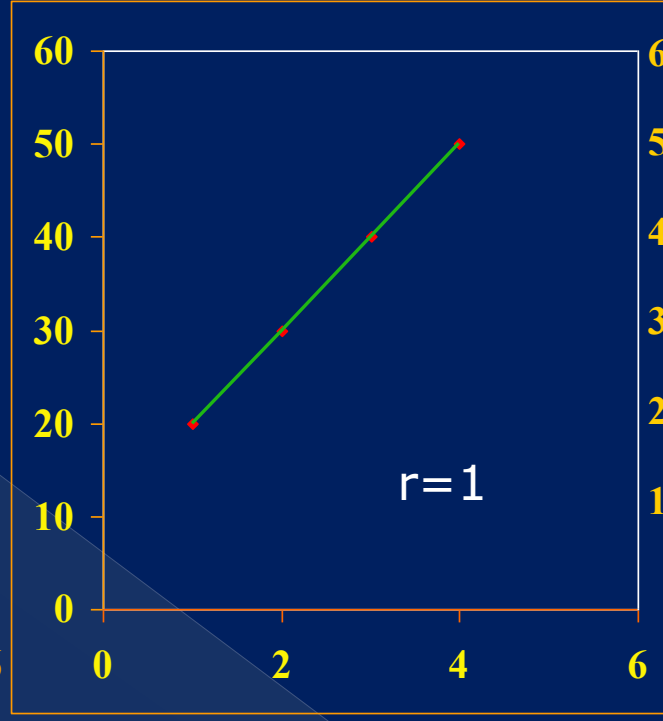
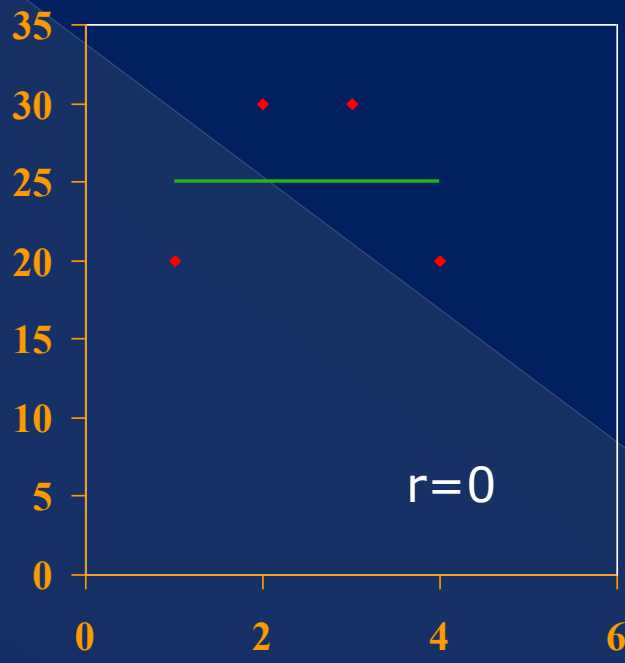
S_x și S_y reprezintă abaterile standard pentru seriile X și respectiv Y :

$$S = \sqrt{s^2}$$

$$s^2 = \frac{\sum_{i=1}^n (X_i - \bar{X})^2}{n}$$

Coeficientul de corelație -proprietăți

- $r \in [-1, 1]$
- măsoară intensitatea relației dintre variabilele X și Y
-
- dacă $0 < r < 1$, norul de puncte poate fi înlocuit (ajustat) printr-o dreaptă de pantă pozitivă .
- Dacă $-1 < r < 0$ atunci norul de puncte poate fi aproximat cu o dreaptă de pantă negativă. Dispersia punctelor față de dreaptă va fi cu atât mai mică cu cât r este mai apropiat de -1.



Testul de semnificație pentru coeficientul de corelație Pearson

- Semnificația coeficientului de corelație Pearson poate fi evaluată dacă valoarea observată a apărut datorită întâmplării (dacă este semnificativ diferită de zero).

Coeficientul de corelație -interpretare

Regulile empirice ale lui Colton (1974)

Valoarea r	$p > 0,05$	$p < 0,05$
in $(-0.25 ; 0,25)$	corelație slabă sau nulă	corelație slabă sau nulă
in $[0.25 ; 0.50)$ sau in $(-0.50 ; -0,25]$	Nu are semnificație statistica	Grad de asociere acceptabil
in $[0.50 ; 0.75)$ sau in $(-0.75 ; -0,50]$	Nu are semnificație statistica	O corelație moderată spre bună
>0.75 sau $< -0,75$	Nu are semnificație statistica	O foarte bună asociere sau corelație
>1 sau <-1	Eroare	Eroare

Interpretarea r,p

Atentie:

- > Variabilele corelate trebuie sa fie cantitative
- > Intre variabilele pentru care se calculeaza corelatia trebuie sa existe o relatie de cauzalitate

Coeficientul de corelație-interpretare

Alura norului de puncte

- $r > 0$

- > O creștere a lui X determină o creștere a lui Y (direct proporționale).

- $r < 0$

- > O creștere a lui X determină o diminuare a lui Y (invers proporționale)

Indici de corelație.

Coeficientul de determinare

- ◎ **$d = r^2$**
- ◎ Reprezintă partea din variația totală a lui Y explicată prin relația liniară existentă între X și Y.
- ◎ Cazuri particulare:
 - > $d=1$:
 - > $d=0$:

Dacă d este exprimat în procente: d reprezintă procentul în care variația lui Y este dată prin relația liniară între cele două variabile.

Două variabile ordinale (sau o variabilă ordinală și una cantitativă) Coeficientul de corelație al lui Spearman

Se procedează astfel:

- i. Se înlocuiește seria bivariată $(x_1, \dots, x_n; y_1, \dots, y_n)$ cu seria rangurilor $(R_{x_1}, \dots, R_{x_n}; R_{y_1}, \dots, R_{y_n})$, valorilor x_i și y_i după ordonarea lor în ordine crescătoare (pentru valorile egale se ia media aritmetică a rangurilor).
- ii. Pentru determinarea coeficientului r_s al lui Spearman se calculează coeficientul de corelație (Pearson) pentru seria rangurilor.

Analiza chestionarelor

Reproductibilitatea

- **Indicele Kappa** de concordanța între observatori
- Coeficientul de corelație interclase
- Coeficientul **Cronbach alfa** de consistență internă

Analiza chestionarelor

Validitatea

- Validitatea de continut (ceea ce testam (intrebarile) este in acord cu cunostintele care sunt testate)
- Validitatea de criteriu (daca e in acord cu alte masuratori cu care ar trebui sa fie in relatie)
- Validitatea de construct (masuratoarea e in acord cu caracteristica masurata si nu cu alte caracteristici)

Statistici descriptive in două dimensiuni.

Drepte de regresie

- Dreapta de regresie $Y(X)$:

$$y = a + b x$$

Dreapta de regresie $Y(X)$

$$\min_{a,b \in R} \sum_{i=1}^n (a + bX_i - Y_i)^2$$

Valorile lui a și b pentru care este atins minimul sumei sunt date prin formulele:

$$b = \frac{COV(X, Y)}{S_x}$$

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

Utilizarea funcțiilor de regresie

■ Extrapolare și interpolare

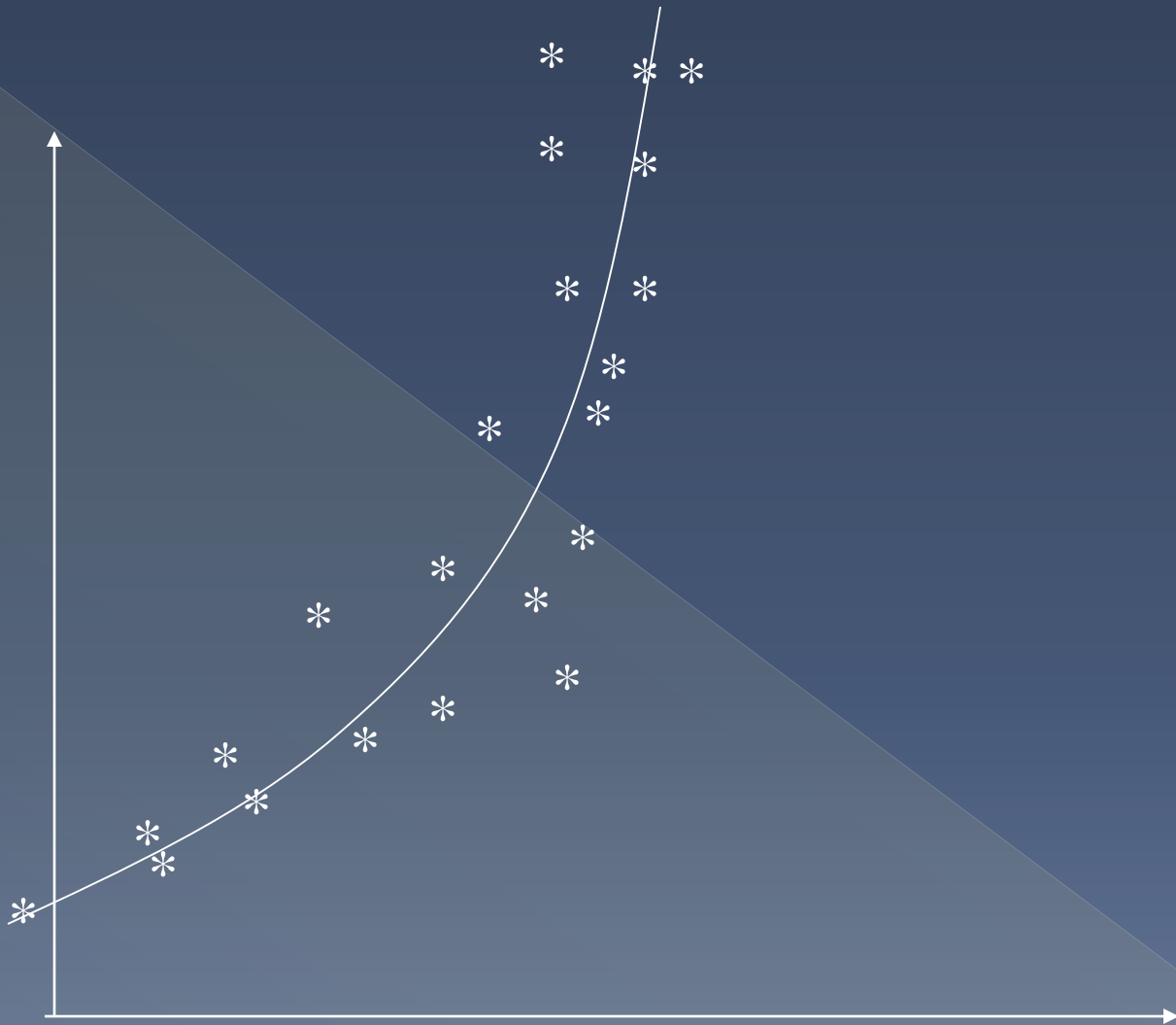
- ◆ Când se determină valoarea funcției (adică a lui Y), pentru un X cuprins intervalul $[X_{\min}, X_{\max}]$, atunci se efectuează o operație de interpolare, iar când X se află în afara intervalului se spune că este vorba de o extrapolare.

■ Prezicerea lui Y pentru un X dat

■ Simulari

Funcția de regresie. Schimbări de variabile

- In unele cazuri se constată că relația liniară pare a nu fi adecvată pentru descrierea dependenței dintre variabilele X și Y , sau că scalele utilizate nu sunt cele mai potrivite.



Regresii multidimensionale

Fiind date variabilele:

$$X_i: X_{i1}, \dots, X_{in} \quad , \quad i=1, 2, \dots, m$$

$$Y: Y_1, \dots, Y_n$$

se caută o relație de forma:

$$Y = a + b_1 X_1 + \dots + b_m X_m,$$

unde coeficienții a și b_i ($i=1, \dots, m$) se determină astfel încât să minimizeze expresia:

$$\sum_{i=1}^n (Y_i - (a + b_1 X_{1i} + \dots + b_m X_{mi}))^2 .$$

