

ECONOMETRIE

CURS 2

- note de curs -

IAȘI
- 2023-

C2. REGRESIA LINIARĂ SIMPLĂ (1)

TEMATICA CURS 02

Prezentarea modelului - regresia empirică și cea teoretică

Estimarea punctuală și prin interval de încredere a parametrilor - MCMMP

Testarea parametrilor

Probleme specifice utilizând SPSS și Excel

1. Scurt istoric al regresiei

1886 - Francis Galton, "Family Likeness in Stature", *Proceedings of Royal Society*, London, vol. 40, 1886, pp. 42-72

1897 - G. U. Yule, "On the Theory of Correlation", *Journal of the Royal Statistical Society*, pp. 812-54.

1903 - Karl Pearson, G. U. Yule, Norman Blanchard, and Alice Lee, "The Law of Ancestral Heredity", *Biometrika*

1925 - R.A. Fisher, *Statistical Methods for Research Workers*

2. Noțiuni (1)

Regresia este o **legătură statistică** între două sau mai multe variabile statistice.

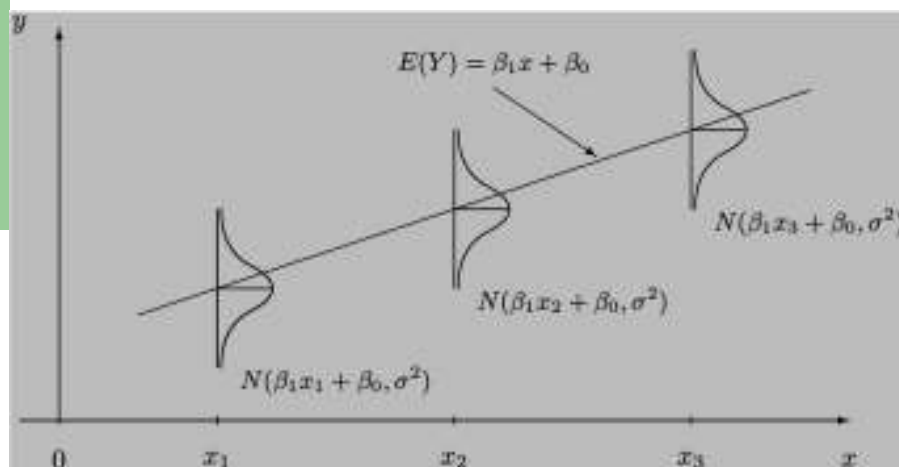
În descrierea **legăturilor statistice** pentru **variabilele dependente** utilizăm **variabile aleatoare (stohastice)**, ceea ce înseamnă că acestora le corespund **distribuții de probabilitate** (fig. 1).

$$\mathbf{x}_i \rightarrow \mathbf{Y}: \begin{pmatrix} y_1 \dots y_j \dots y_m \\ n_i \dots n_j \dots n_m \end{pmatrix} \Rightarrow \mathbf{M}(\mathbf{Y} | \mathbf{x}_i) = \mathbf{f}(\mathbf{x}_i) \Rightarrow y_i = \mathbf{f}(\mathbf{x}_i) + \epsilon_i$$

\mathbf{X} – **variabilă independentă** - variabilă nestohastică

\mathbf{Y} – **variabilă dependentă** – variabilă stohastică (prezintă distribuții pentru fiecare valoare a lui \mathbf{X})

2. Noțiuni (2) - Fig.1



În *legăturile de tip funcțional* unei valori i se asociază o altă valoare și nu o distribuție de probabilitate.

$$x_i \rightarrow y_i \Rightarrow f(x_i) = y_i$$

Analiza de regresie studiază *forma legăturii* dintre una sau mai multe variabile \Rightarrow **Model de regresie**

Analiza de corelație studiază *intensitatea legăturii* dintre una sau mai multe variabile puse în relație printr-un *model de regresie*.

Modele de regresie

Modele de regresie	Simple	Multiple
Liniare	$C = \beta_0 + \beta_1 X + \varepsilon$	$C = \beta_0 + \beta_1 X_1 + \beta_2 X_2 + \varepsilon$
Neliniare	$C = e^{\beta_0 + \beta_1 X + \varepsilon}$	$C = e^{\beta_0 + \beta_1 X_1 + \beta_2 X_2 + \varepsilon}$

3. Modelul de regresie liniară simplă

$$\varepsilon_i = y_i - M(Y | X=x_i) \Leftrightarrow y_i = M(Y | X=x_i) + \varepsilon_i = \beta_0 + \beta_1 x_i + \varepsilon_i$$

> $M(Y | X=x_i) = \beta_0 + \beta_1 x_i$ – media condiționată de $X=x_i$, a variabilei stohastice Y ,

> $\beta_0 = f(0)$ – parametrul “intersecția dreptei de regresie liniară cu axa OY ” (engl. intercept)

> β_1 – parametrul “panta a dreptei” care reprezintă variația absolută **in medie** a variabilei Y atunci când variabila X crește cu o unitate ($\beta_1 = \Delta Y / \Delta X$):

- $\beta_1 > 0$: legătură **directă** între variabile, Y variază în același sens cu X

- $\beta_1 < 0$: legătură **inversă** între variabile, Y nu variază în același sens cu X

3. Modalități de scriere ale modelului liniar simplu

Dacă este scris pentru *valorile variabilelor*

$$y_i = \beta_0 + \beta_1 x_i + \varepsilon_i$$

Dacă este scris pentru *variabile în general*

$$Y = \beta_0 + \beta_1 X + \varepsilon$$

4. Componentele modelului de regresie

A. Componenta deterministă ($\beta_0 + \beta_1 x_i$)

B. Componenta aleatoare (ε_i).

Factorii care influențează componenta aleatoare:

- *natura fenomenului studiat;*
- *specificarea modelului ;*
- *erorile de măsurare s.a.*

5. Ipoteze clasice ale modelului de regresie

Ipotezele modelului de regresie vizează **variabila reziduală** și **variabila independentă**.

Cele mai importante **ipoteze cu privire la variabila reziduală** sunt:

- **normalitatea erorilor**: $\varepsilon_i \sim N(\mu_i, \sigma_i^2)$, adică variabila reziduală urmează o lege de repartiție normală de medie μ_i și varianță σ_i^2 ;
- **media erorilor de modelare este nulă**: $M(\varepsilon_i) = 0 \Rightarrow \varepsilon_i \sim N(0, \sigma_i^2)$
- **homoscedasticitate**: $V(\varepsilon_i) = M(\varepsilon_i^2) = \sigma^2$, adică varianța erorii este constantă la nivelul distribuțiilor condiționate de tipul $Y_i | X = x_i$
 $\rightarrow \varepsilon_i \sim N(0, \sigma^2)$
- **necorelarea erorilor**: $\text{cov}(\varepsilon_i, \varepsilon_j) = 0, i \neq j; i, j = \overline{1, n}$, adică erorile nu se influențează reciproc;
- **lipsa corelației dintre variabila independentă și variabila eroare**:
 $\text{cov}(\varepsilon_i, x_i) = 0$

6. Exemple de modele liniare simple

Funcția de consum

- **cererea sau consumul populației** în funcție de **venit**:

$$C_i = \beta_0 + \beta_1 V_i + \varepsilon_i$$

unde parametrul β_1 arată cu cât crește **in medie consumul** unui anumit produs (C_i) la o creștere cu o unitate a venitului și este de regulă pozitiv.

Legea cererii

- **cererea populației** pentru o gamă de produse în funcție de **prețul** acestora:

$C_i = \beta_0 + \beta_1 P_i + \varepsilon_i$, unde parametrul β_1 este de regulă negativ și arată cu cât scade **in medie cererea** la o creștere a prețului cu o unitate.

STOP

7. Estimarea punctuală parametrilor modelului de regresie prin MCMMP

MCMMP (engl. Method of Ordinary Least Squares - OLS)

Fie: $\hat{y}_i = \hat{\beta}_0 + \hat{\beta}_1 x_i$ - **valorile estimate** (teoretice ale lui Y)

și $y_i = \beta_0 + \beta_1 x_i + \varepsilon_i$ - **valorile reale, înregistrate**

Erorile estimate pot fi obținute ca diferența între **valorile reale** și **valorile teoretice**:

$$\hat{\varepsilon}_i = y_i - \hat{y}_i = y_i - (\hat{\beta}_0 + \hat{\beta}_1 x_i) = y_i - \hat{\beta}_0 - \hat{\beta}_1 x_i$$

Pentru a ușura notațiile în procesul de estimare vom utiliza direct **notațiile pentru estimății**. Astfel relația de mai sus se va scrie:

$$\sum_{i=1}^n e_i^2 = \sum_{i=1}^n (y_i - \hat{y}_i)^2 = \sum_{i=1}^n (y_i - (b_0 + b_1 x_i))^2 = \sum_{i=1}^n (y_i - b_0 - b_1 x_i)^2 - \min.$$

Notăm
$$S = \sum_{i=1}^n e_i^2 = \sum_{i=1}^n (y_i - b_0 - b_1 x_i)^2$$

Rezolvarea acestei probleme de minim presupune îndeplinirea a **două condiții**:

1. **Anularea derivatelor parțiale de ordinul I ale lui S în raport b_0 și b_1** :

$$\begin{cases} \frac{\partial S}{\partial b_0} = 2 \sum_{i=1}^n (y_i - b_0 - b_1 x_i)(-1) = 0 \\ \frac{\partial S}{\partial b_1} = 2 \sum_{i=1}^n (y_i - b_0 - b_1 x_i)(-x_i) = 0 \end{cases} \Rightarrow \begin{cases} nb_0 + b_1 \sum_{i=1}^n x_i = \sum_{i=1}^n y_i \\ b_0 \sum_{i=1}^n x_i + b_1 \sum_{i=1}^n x_i^2 = \sum_{i=1}^n y_i x_i \end{cases}$$

$$b_0 = \frac{\Delta b_0}{\Delta} = \frac{\sum y_i \sum x_i^2 - \sum x_i \sum x_i y_i}{n \sum x_i^2 - (\sum x_i)^2}$$

$$\text{sau } b_0 = \bar{y} - b_1 \bar{x}$$

$$b_1 = \frac{\Delta b_1}{\Delta} = \frac{n \sum x_i y_i - \sum x_i \sum y_i}{n \sum x_i^2 - (\sum x_i)^2}$$

2. Matricea derivatelor parțiale de ordinul doi să fie pozitiv definită:

$$\det \begin{pmatrix} \frac{\partial^2 S}{\partial^2 b_0} = 2n & \frac{\partial^2 S}{\partial b_0 \partial b_1} = 2 \sum x_i \\ \frac{\partial^2 S}{\partial b_0 \partial b_1} = 2 \sum x_i & \frac{\partial^2 S}{\partial^2 b_1} = 2 \sum x_i^2 \end{pmatrix} > 0 \Leftrightarrow \det \begin{pmatrix} 2n & 2 \sum x_i \\ 2 \sum x_i & 2 \sum x_i^2 \end{pmatrix} > 0$$

$$\Leftrightarrow \det \begin{pmatrix} n & \sum x_i \\ \sum x_i & \sum x_i^2 \end{pmatrix} > 0$$

Matricea derivatelor parțiale de ordin doi este pozitiv definită deoarece $n \sum x_i^2 - (\sum x_i)^2 = n^2 \sigma^2 > 0$

7. Proprietățile estimatorilor parametrilor modelului de regresie

Estimatorii parametrilor modelului de regresie sunt variabile de selecție care:

- urmează o **distribuție normală**: $\hat{\beta}_0 \sim N(\beta_0, \sigma_{\hat{\beta}_0}^2)$, $\hat{\beta}_1 \sim N(\beta_1, \sigma_{\hat{\beta}_1}^2)$
- sunt **nedeplasați**: $M(\hat{\beta}_0) = \beta_0$, $M(\hat{\beta}_1) = \beta_1$
- **convergenți** (in probabilitate): $(\hat{\beta}_0)_{n \in N} \xrightarrow{p} \beta_0$, $(\hat{\beta}_1)_{n \in N} \xrightarrow{p} \beta_1$
- **eficienți**: dintre toți estimatorii posibili pentru β_1 , $\hat{\beta}_1$ are varianța cea mai mică

8. Estimarea prin interval de încredere a parametrilor modelului de regresie liniară

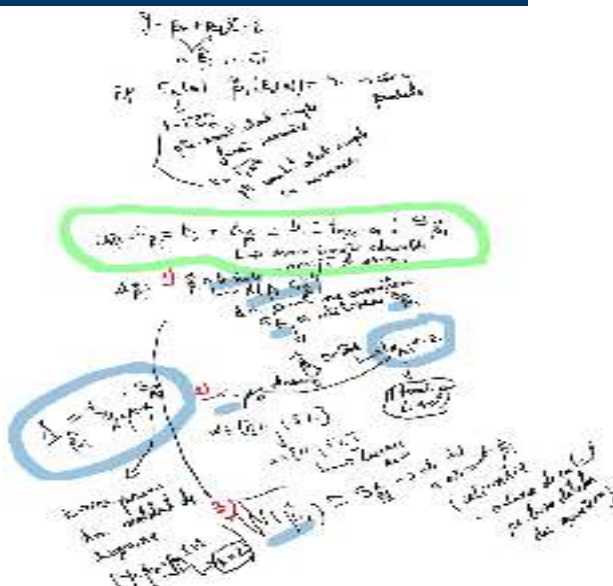
Atât pentru β_0 , cât și pentru β_1 , *intervalele de încredere pentru* se vor construi astfel:

$$\beta_0 \in [b_0 - t_{\alpha/2, n-k} s_{\hat{\beta}_0}, b_0 + t_{\alpha/2, n-k} s_{\hat{\beta}_0}]$$

$$\beta_1 \in [b_1 - t_{\alpha/2, n-k} s_{\hat{\beta}_1}, b_1 + t_{\alpha/2, n-k} s_{\hat{\beta}_1}]$$

unde k = *numărul parametrilor estimați din model* (pentru modelul liniar $k=2$)

Precizari suplimentare



Vezi exemplu excel (n=19 iar $\alpha=5\%$): Y-> greutatea (kg) si X-> inaltimea (cm) $Y=\beta_0+\beta_1X+\epsilon \rightarrow K=2$

	Coefficients (b_i)	Standard Error ($S_{\hat{\beta}_i}, i=0,1$)	t Stat (t_{calc})	P-value (sig.)	Lower 95% (L_i)	Upper 95% (U_i)
Intercept (β_0)	-151.816 (b_0)	61.808 ($S_{\hat{\beta}_0}$)	-2.456	0.025	-282.219	-21.412
Inaltimea (X) (β_1)	1.266 (b_1)	0.364 ($S_{\hat{\beta}_1}$)	3.483	0.003	0.499	2.034

$t_{\alpha/2, n-k} \rightarrow$ aceasta este o valoare teoretica si se ia din tabelul valorilor statisticii Student si nu din tabelul de mai sus!!!!

$$\beta_0 \in [b_0 - t_{\alpha/2, n-k} S_{\hat{\beta}_0}, b_0 + t_{\alpha/2, n-k} S_{\hat{\beta}_0}] = -151.816 \pm 2.11 \cdot 61.808 \rightarrow \text{vezi } L_1 \text{ si } L_s \text{ din tabelul de mai sus}$$

$$\beta_1 \in [b_1 - t_{\alpha/2, n-k} S_{\hat{\beta}_1}, b_1 + t_{\alpha/2, n-k} S_{\hat{\beta}_1}] = 1.266 \pm 2.11 \cdot 0.364 \rightarrow \text{vezi } L_1 \text{ si } L_s \text{ din tabelul de mai sus}$$

$$t_{\alpha/2, n-k} = t_{0.025; (19-2)} = t_{0.025; 17} = 2.11$$

[Student t Distribution Table \(thoughtco.com\)](http://www.thoughtco.com)

5. Testarea parametrilor modelului liniar

1. Formularea ipotezelor

pentru β_0

$$H_0: \beta_0 = 0$$

$$H_1: \beta_0 \neq 0$$

pentru β_1

$$H_0: \beta_1 = 0$$

$$H_1: \beta_1 \neq 0$$

2. Fixarea pragului de semnificație

$$\alpha = 0.05$$

3. Alegerea si calculul statisticii test

$$t_{\beta_0} = \frac{b_0}{S_{\hat{\beta}_0}} \cdot t_{teoretic} = t_{\alpha/2, n-2}$$

$$t_{\beta_1} = \frac{b_1}{S_{\hat{\beta}_1}} \cdot t_{teoretic} = t_{\alpha/2, n-2}$$

4. Criterii de decizie:

$$|t_{calc}| \leq t_{teoretic} = t_{\alpha/2, n-2} \text{ sau } sig. \geq \alpha \Rightarrow \text{A}H_0 \text{ cu o probabab. de } 1-\alpha.$$

$$|t_{calc}| > t_{teoretic} = t_{\alpha/2, n-2} \text{ sau } sig. < \alpha \Rightarrow \text{R}H_0 \text{ cu un risc asumat } \alpha.$$

	Coefficients (b_i)	Standard Error ($s_{\hat{\beta}_i}, i=0,1$)	t Stat (t_{calc})	P-value (sig.)	Lower 95% (L_i)	Upper 95% (U_i)
Intercept (β_0)	-151.816	61.808	-2.456	0.025	-282.219	-21.412
Inaltimea (X) (β_1)	1.266	0.364	3.483	0.003	0.499	2.034

1. Formularea ipotezelor

pentru β_0
 $H_0: \beta_0 = 0$
 $H_1: \beta_0 \neq 0$

pentru β_1
 $H_0: \beta_1 = 0$
 $H_1: \beta_1 \neq 0$

2. Fixarea pragului de semnificație

$\alpha = 0,05$

3. Alegerea si calcul statisticii test

$$t_{calc \beta_0} = \frac{b_0}{s_{\hat{\beta}_0}} \sim t_{\alpha/2, n-2}$$

$$t_{calc \beta_1} = \frac{b_1}{s_{\hat{\beta}_1}} \sim t_{\alpha/2, n-2}$$

4. Criterii de decizie:

$$|t_{calc \beta_0}| = |-2,456| > t_{teoretic} = t_{0,025; 17} = 2,11 \quad |t_{calc \beta_1}| = |-2,456| \leq t_{teoretic} = t_{0,025; 17} = 2,11$$

$$\text{sig.} = 0.025 < \alpha = 0,05$$

$$\text{sig.} = 0.003 < \alpha = 0,05$$

\Rightarrow RH_0 cu un risc asumat α

\Rightarrow RH_0 cu un risc asumat α

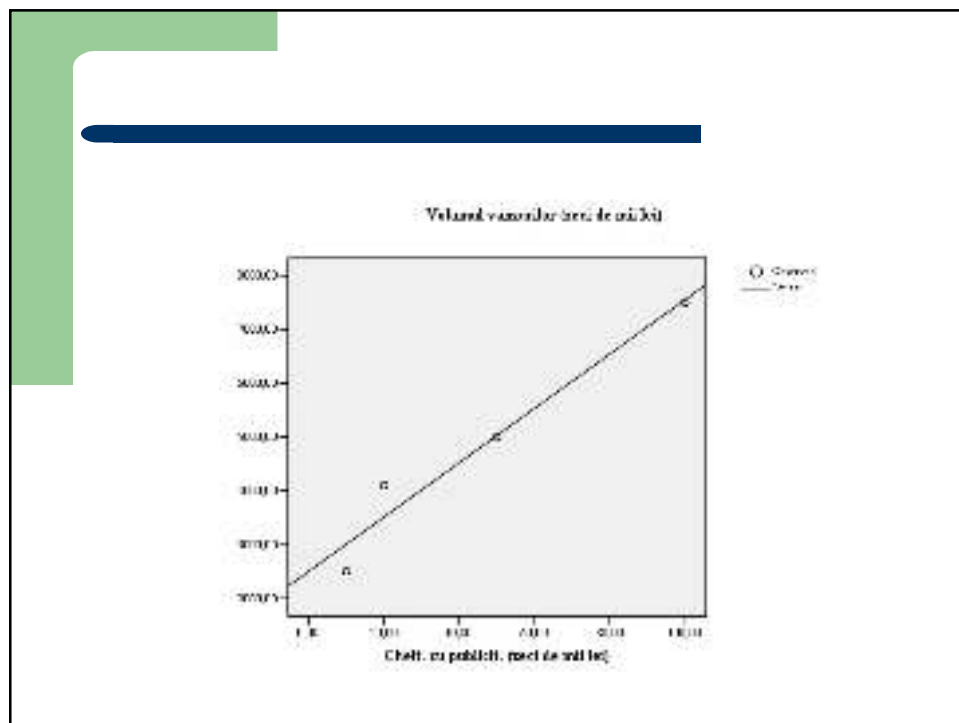
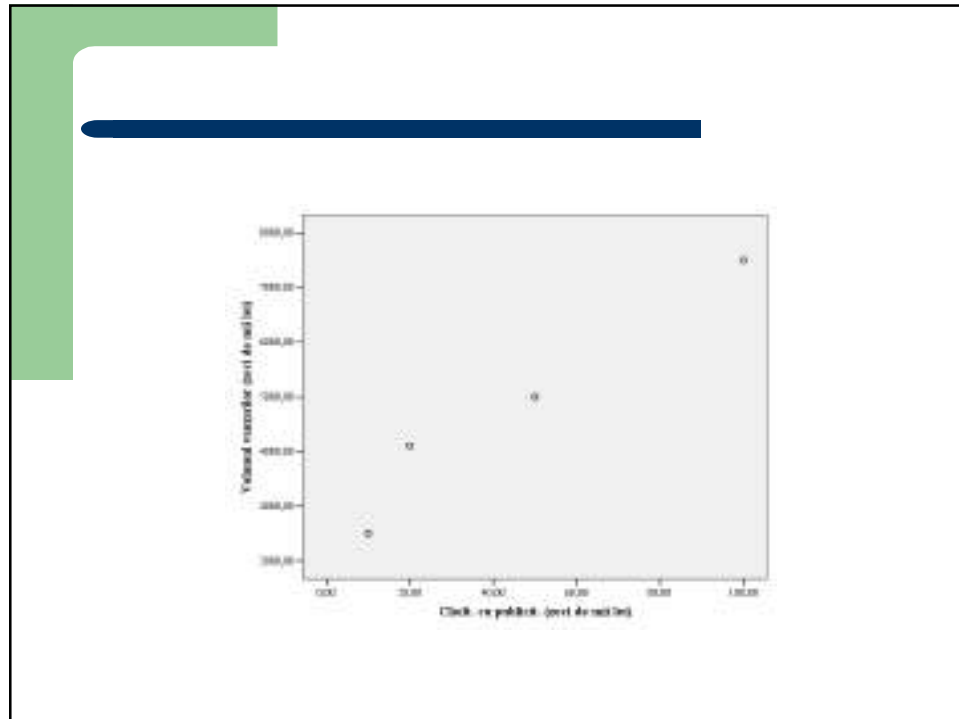
5. Probleme specifice analizei de corelație și regresie

Analiza de regresie si corelație în SPSS

Se consideră datele cu privire la *Valoarea vânzărilor* și *Cheltuielile cu publicitatea* pentru un eșantion de 4 firme.

Datele sunt prezentate în tabelul alăturat.

x_i	y_i
10	2500
20	4100
50	5000
100	7500
180	19100



Model Summary

Model	R	R Square	Adjusted R Square	Std. Error of the Estimate
1	.977 ^a	.954	.931	550.55830

a. Predictors: (Constant), X

ANOVA^a

Model		Sum of Squares	df	Mean Square	F	Sig.
1	Regression	12501275.51	1	12501275.51	41.243	.023 ^a
	Residual	606224.490	2	303112.245		
	Total	13107500.00	3			

a. Dependent Variable: Y

b. Predictors: (Constant), X

Correlations

		X	Y
X	Pearson Correlation	1	.977 [*]
	Sig. (2-tailed)		.023
	N	4	4
Y	Pearson Correlation	.977 [*]	1
	Sig. (2-tailed)	.023	
	N	4	4

*. Correlation is significant at the 0.05 level (2-tailed).

Coefficients^a

Model		Unstandardized Coefficients		Standardized Coefficients	t	Sig.	95.0% Confidence Interval for B	
		B	Std. Error	Beta			Lower Bound	Upper Bound
1	(Constant)	2502.041	449.379		5.560	.031	572.821	4431.260
	X	50.510	7.865	.977	6.422	.023	16.689	84.351

a. Dependent Variable: Y