

## CURSUL VI

### ANALIZA LEGĂTURILOR DINTRE PROCESELE ȘI FENOMENELE ECONOMICE

6.1 Tipologia legaturilor statistice

6.2. Metoda regresiei

6.3 Coeficientul de determinație

6.4 Coeficientul de corelație

#### 6.1 Tipologia legaturilor statistice

Este cunoscut faptul că variabilitatea fenomenelor social-economice este, în majoritatea cazurilor, determinată de acțiunea simultană a mai multor factori; o parte dintre acești factori favorizează evoluția unui fenomen, alții o frânează sau acționează chiar în sens invers.

Sensurile și intensitățile influențelor diferiților factori se schimbă în condiții timp și spațiu, astfel că evoluția fenomenelor dependente înregistrează și ea tendințe diferite față de cele anterioare.

Raporturile de cauzalitate dintre fenomenele social-economice pot fi cuantificate și analizate cu ajutorul corelației. Informațiile obținute sunt deosebit de utile, mai ales pentru faptul că, metodele specifice pe care statistica le pune la dispoziția cercetătorului, oferă posibilitatea cunoașterii, în principal, a următoarelor aspecte:

- existența raporturilor de cauzalitate dintre fenomene;
- contribuția fiecărui factor la variabilitatea globală a fenomenelor efect;
- intensitatea legăturilor cauzale dintre fenomenele și procesele social-economice;
- tendințele evolutive ale corelației dintre fenomene.

Se remarcă faptul că analiza corelației oferă o paletă mult mai largă de informații, fiind preferată deși determinarea indicatorilor specifici corelației este mult mai dificilă.

*Corelația este expresia sintetică a intensității legăturilor cauzale dintre fenomene.*

Cuplul corelativ poate cuprinde două sau mai multe variabile, din care una este variabila efect, cunoscută sub numele de variabilă **rezultativă**, simbolizată în general cu  $Y$ , iar celelalte sunt variabile cauză, cunoscute sub numele de variabile **factoriale** - fiind simbolizate cu  $X_1, X_2, \dots, X_n$ .

Se deosebesc mai multe tipuri de corelație.

a) După numărul de variabile din cuplul corelativ, se disting:

- **corelația simplă**, în care cuplul corelativ cuprinde două variabile, din care una este variabila rezultativă  $Y$  și cealaltă - variabila factorială  $X$ .

- **corelație multiplă**, în care cuplul corelativ cuprinde trei sau mai multe variabile, din care una este variabila rezultativă  $Y$  și celelalte sunt variabilele factoriale  $X_1, X_2, \dots$ , etc.

b) După sensul legăturilor factoriale, se deosebesc:

- **corelația directă**, în care variabilitatea rezultativei  $Y$  se produce în același sens cu variabilitatea factorialei  $X$  (sau factorialelor).

- **corelația inversă**, în care variabilitatea rezultativei  $Y$  se produce în sens invers, comparativ cu variabilitatea factorilor determinanți.

c) După forma legăturilor cauzale, se disting :

- **Corelație lineară**. Constă în aceea că variabila rezultativă  $Y$  înregistrează o tendință lineară, ca urmare a influenței factorilor determinanți.

- **Corelație nelineară**. În acest caz, variabilitatea rezultativei  $Y$  se integrează într-o tendință de tip nelinear (parabolic, exponențial, etc).

Cunoașterea formei corelațiilor prezintă un interes deosebit în estimarea tendințelor evolutive ale fenomenelor efect în strânsă legătură cu variabilitatea factorilor determinanți. Evident, o tendință lineară a rezultativei Y, sugerează faptul că nivelurile acestei variabile cresc sau descresc aproximativ în progresie aritmetică; tendința nelineară, evidențiază modificarea acestor niveluri în progresie geometrică, exponențială, etc. Aceste informații sunt deosebit de utile în calculele de previziune precum și în modelarea tendințelor evolutive.

Dintre **metodele elementare** de analiză a corelației, prezentate în literatura de specialitate, se folosesc mai frecvent: tabelul de corelație și metoda grafică.

• **Tabelul de corelație.** Este de fapt tabelul cu dublă intrare folosit în prezentarea seriilor de distribuție bidimensionale; rolul variabilei factoriale este preluat de caracteristica principală X și cel al variabilei rezultative este preluat de caracteristica secundară Y.

Cu ajutorul tabelului de corelație pot fi evidențiate următoarele aspecte: existența corelației, sensul corelației, forma corelației și intensitatea corelației.

Tabelul de corelație servește, în majoritatea cazurilor în analizele profesionale, ca bază preliminară a folosirii metodelor analitice în studiul corelației.

• **Metoda grafică.** Presupune construirea și folosirea **corelogramei** în analiza corelației. Pentru a construi corelograma, pe abscisă se trec valorile scării de reprezentare a variabilei factoriale X iar pe ordonată valorile variabilei rezultative Y. Prin unirea punctelor corespunzătoare coordonatelor XY se obține corelograma.

Cu ajutorul acestei metode se pot evidenția următoarele aspecte: existența corelației, sensul corelației, forma corelației și intensitatea corelației.

Și metoda grafică este frecvent folosită în analiza corelației, fiind utilizată în foarte multe cazuri ca element esențial în alegerea funcțiilor statistico-matematice pentru analiza regresiei și intensității corelației.

În continuare va fi abordată metodologia aplicării funcțiilor statistico-matematice în studiul regresiei simple, precum și metodologia cuantificării intensității corelației simple.

## 6.2 Metoda regresiei

În cadrul analizei corelației se au în vedere, în principal, două aspecte esențiale :

- **regresia** - cu ajutorul căreia, prin folosirea și interpretarea coeficienților de regresie ai diferitelor funcții statistico-matematice, se determină contribuțiile factorilor determinanți la variabilitatea fenomenelor efect;

- **intensitatea corelației** - sintetizată cu ajutorul coeficienților de corelație.

Pentru corelația simplă, primul aspect poate fi evidențiat cu ajutorul funcțiilor : lineară - pentru legăturile cauzale de tip linear; parabolică de ordin superior, hiperbolică, exponențială, logaritmică, semilogaritmică, logistică, etc - pentru legături cauzale de tip nelinear.

Corelația simplă, după forma pe care o au tendințele legăturilor cauzale, poate fi **lineară și nelineară.**

Modelul **linear** de ordinul 1 al regresiei simple este dat de expresia:

$$y = \beta_0 + \beta_1 x + \varepsilon$$

în care

$y$  = variabila dependentă (rezultativă);

$x$  = variabila independentă (factorială);

$\beta_0$  = valoarea inițială a lui  $y$ ;

$\beta_1$  = modificarea lui  $y$  ca urmare a modificării lui  $x$  egală cu unitatea;

$\varepsilon$  = eroarea variabilei.

Pentru determinarea parametrilor  $\beta_0$  și  $\beta_1$  se folosesc ca estimatori parametrii funcției lineare:

$$\hat{y} = b_0 + b_1 x$$

in care

$b_0$ -valorile variabilei rezultative Y determinate în afara influenței variabilei factoriale X;

$b_1$ - coeficientul de regresie; el sintetizează creșterea (când are semnul +) sau descreșterea (când are semnul -) înregistrată de variabila rezultativă Y corespunzătoare unei creșteri sau descreșteri a variabilei X egală cu unitatea;

$x$ - valorile  $x_1, x_2, x_3, \dots, x_n$  ale variabilei factoriale X.

Parametrul  $b_1$  are un rol esențial în cadrul analizei regresiei, deoarece cu ajutorul său se determină contribuția factorialei X la variabilitatea rezultativei Y.

Aplicând metoda celor mai mici pătrate,  $\sum (y - \hat{y})^2 \Rightarrow \min$ , cei doi parametri  $b_0$  și  $b_1$  se pot determina astfel:

$$b_1 = \frac{\text{cov}(X, Y)}{s_x^2} ; \quad b_0 = \bar{y} - b_1 \bar{x}$$

Elementele folosite ca bază de calcul pentru cei doi parametri se determină astfel:

$$\bar{x} = \frac{\sum_{i=1}^n x_i}{n}$$

$$\bar{y} = \frac{\sum_{i=1}^n y_i}{n}$$

$$s_x^2 = \frac{\sum_{i=1}^n (x_i - \bar{x})^2}{n - 1}$$

$$\text{cov}(X, Y) = \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{n - 1}$$

În analiza fenomenelor social-economice sunt foarte frecvente cazurile în care valorile empirice  $y_i$  și valorile teoretice  $\hat{y}_i$  există diferențe; acestea sunt determinate de acțiunea factorilor întâmplători și constituie valoarea reziduală a estimării, cunoscută și sub numele de eroare ( $y_i - \hat{y}_i$ ).

Eroarea medie (sau eroarea standard) a estimației, se calculează cu formula:

$$S_e = \sqrt{\frac{SSE}{n - 2}}$$

unde  $SEE$  reprezintă suma pătratelor erorilor de estimație, și are următoarea relație de calcul:

$$SSE = \sum_{i=1}^n (y_i - \hat{y}_i)^2$$

Pentru a stabili că  $b_1$  estimează corect pe  $\beta_1$ , trebuie verificată ipoteza nulă:  $H_0: \beta_1 = 0$ . Alternativa ei este:  $H_0 : \beta_1 \neq 0$ . În acest scop se poate folosi testul t, determinat cu ajutorul relației:

$$t = \frac{b_1 - \beta_1}{S_{b_1}}$$

unde

$S_{b_1}$  = eroarea standard a lui  $b_1$ .

$$S_{b_1} = \frac{S_e}{\sqrt{(n-1)s_x^2}}$$

Dacă erorile variabilei sunt normal distribuite, testul statistic t (Student), considerat ca nivel referențial, este cel căruia îi corespund  $n-2$  grade de libertate.

Regiunea de respingere a ipotezei nule este:

$$t > t_{\alpha/2, n-2} \text{ sau } t < -t_{\alpha/2, n-2}$$

Spre exemplu, dacă  $n = 100$  și  $\alpha = 0,05$  atunci zona de respingere este:

$$t > t_{0,025,98} \cong 1.984 \text{ sau } t < -t_{0,025,98} \cong -1.984$$

### 6.3 Coeficientul de determinație

În practică, pentru determinarea contribuției unuia dintre factorii de influență la variabilitatea fenomenului dependent, se folosește coeficientul de determinație, a cărui relație de calcul este:

$$R^2 = \frac{[\text{cov}(X, Y)]^2}{s_x^2 s_y^2}$$

sau

$$R^2 = 1 - \frac{\sum_{i=1}^n (y_i - \hat{y}_i)^2}{\sum_{i=1}^n (y_i - \bar{y})^2} = 1 - \frac{SSE}{\sum_{i=1}^n (y_i - \bar{y})^2}$$

Având în vedere faptul că abaterile  $(y_i - \bar{y})$  evidențiază variația totală a variabilei Y, determinată de toți factorii care au influențat-o și știind că o parte din această variație este preluată de funcția de regresie și evidențiată prin abaterile  $(\hat{y}_i - \bar{y})$ , iar cealaltă parte - care constituie variația reziduală sau eroarea estimației - este evidențiată prin abaterile  $(y_i - \hat{y}_i)$ , înseamnă că:

$$(y_i - \bar{y}) = (y_i - \hat{y}_i) + (\hat{y}_i - \bar{y})$$

Dar și suma pătratelor acestor abateri verifică această relație. Deci:

$$\sum_{i=1}^n (y_i - \bar{y})^2 = \sum_{i=1}^n (y_i - \hat{y}_i)^2 + \sum_{i=1}^n (\hat{y}_i - \bar{y})^2$$

Considerând:

$$SS_y = \sum_{i=1}^n (y_i - \bar{y})^2$$

$$SSE = \sum_{i=1}^n (y_i - \hat{y}_i)^2$$

$$SSR = \sum_{i=1}^n (\hat{y}_i - \bar{y})^2$$

rezultă că pentru calculul coeficientului de determinație se pot folosi mai multe variante de calcul:

$$R^2 = 1 - \frac{SSE}{\sum_{i=1}^n (y_i - \bar{y})^2}$$

$$R^2 = \frac{SS_y - SSE}{\sum_{i=1}^n (y_i - \bar{y})^2}$$

$$R^2 = \frac{SSR}{\sum_{i=1}^n (y_i - \bar{y})^2}$$

În practică  $R^2$  se exprimă și în procente.

#### 6.4 Coeficientul de corelație al lui Pearson.

Coeficientul de determinație  $R^2$  evidențiază proporția variației variabilei dependente  $Y$  determinată de influența variației variabilei independente  $X$ .

Coeficientul de corelație  $r$  evidențiază intensitatea legăturii cauzale dintre cele două variabile. Ambii coeficienți au la bază relația lui Pearson; coeficientul de determinație este pătratul coeficientului de corelație. Deci:

$$r = \sqrt{R^2} = \frac{\text{cov}(X, Y)}{s_x s_y}$$

Coeficientul de corelație  $r$  ia valori cuprinse între -1 și 1 cu următoarea semnificație:

- valorile cuprinse între 0 și 1 evidențiază o corelație directă, din ce în ce mai intensă pe măsură ce acestea se apropie de 1;
- valorile cuprinse între 0 și -1 evidențiază o corelație inversă, din ce în ce mai intensă pe măsură ce se apropie de -1;
- valoarea zero evidențiază faptul că între cele două variabile nu există nici o legătură.

De regulă, în practică se rafinează intervalul dintre -1 și 1 astfel:

- dacă  $0 \leq r < 0,2$  nu avem o legătură semnificativă între variabile;
- dacă  $0,20 \leq r < 0,50$  legătura dintre variabile este slabă;
- dacă  $0,50 \leq r < 0,75$  legătura dintre variabile este medie;
- dacă  $0,75 \leq r < 0,95$  legătura dintre variabile este puternică;
- dacă  $0,95 \leq r \leq 1$  între cele două variabile există o legătură de tip funcțional.

Pentru testarea semnificației coeficientului de corelație se are în vedere verificarea ipotezei nule ( $H_0$ ) - ce se face cu ajutorul *testului t*- potrivit căreia între cele două variabile nu este nici o legătură lineară.

*Testul t (Student)* folosit în acest caz are următoarea relație de calcul:

$$t = r \sqrt{\frac{n-2}{1-r^2}}$$

unde

$r$  - coeficientul de corelație;

$n$  - numărul perechilor de valori  $x$  și  $y$ .

Valoarea calculată a lui  $t$  se compară cu valoarea teoretică obținută din tabelul  $t$  (*Student*), pentru  $n-2$  grade de libertate și nivelul de semnificație stabilit.

Regiunea de respingere a ipotezei  $H_0$ , (pentru  $n-2$  grade de libertate) se stabilește în același mod ca cel de la parametrii funcției de regresie, adică:

$$t > t_{\alpha/2, n-2} \text{ sau } t < -t_{\alpha/2, n-2}$$

Dacă  $H_0$  se respinge, se trage concluzia, cu un risc asumat (de regulă 5%), că valoarea coeficientului de corelație nu este egală cu 0, adică între cele două variabile există o legătură semnificativă sau altfel spus, coeficientul de corelație este semnificativ statistic.

În cazul că  $H_0$  se acceptă, coeficientul de corelație nu este semnificativ, și deci între variabile există o legătură întâmplătoare.

În practică este foarte des întâlnit și coeficientul de corelație a rangurilor al lui Spearman, care are la bază relația coeficientului lui Pearson. În locul valorilor  $x, y$  Spearman a folosit rangurile corespunzătoare acestor valori, ajungând la următoarea relație:

$$r_s = \frac{\text{cov}(a, b)}{s_a s_b}$$

unde

$a$  = rangurile valorilor  $x$

$b$  = rangurile valorilor  $y$

Coeficientul lui Spearman ia valori cuprinse între -1 și 1, având aceeași semnificație ca cel a lui Pearson.

Pentru verificarea ipotezei nule  $H_0 : \rho = 0$  se folosește testul statistic  $z$ , care are următoarea relație:

$$z = \frac{r_s - 0}{1/\sqrt{n-1}} = r_s \sqrt{n-1}$$

Se compară  $z$  calculat cu valoarea tabelară corespunzătoare unui anumit nivel de semnificație și dacă valoarea calculată este mai mare decât  $z_{tab}$  sau mai mică decât  $-z_{tab}$ , atunci ipoteza  $H_0$  se respinge.