

4

STUDIUL STATISTIC AL LEGĂTURILOR DINTRE FENOMENELE ȘI PROCESELE ECONOMICO-SOCIALE

Procesele și fenomenele economice (și nu numai) apar și se dezvoltă ca urmare a unor cauze variate, care pot acționa în același sens sau în sensuri opuse și cu diferite grade de intensitate. Drept urmare acestea sunt legate între ele prin conexiuni, uneori foarte complexe, care nu sunt cunoscute sau observate de la bun început, ci, de regulă, sunt descoperite pe măsura studierii lor. Manifestarea unui sau altuia dintre procese sau fenomene generează efecte care pot provoca apariția, modificarea sau încetarea altora determinând astfel relații de interdependență sau cauzalitate. Complexitatea interacțiunii dintre fenomene este cu atât mai mare cu cât acestea aparțin unor colectivități mai numeroase. De aici deducem că fenomenele și procesele economico-sociale nu sunt univoc determinate fiind rezultatul conjugării influenței mai multor fenomene-cauză, iar în sistemul acesta de conexiuni nu toate raporturile de dependență prezintă aceeași importanță întrucât există factori ce se compensează reciproc.

Studiul statistic al raporturilor de dependență dintre procese și fenomene se concentrează pe identificarea relației care există între două sau mai multe caracteristici. Importanța cunoașterii legăturii dintre un fenomen sau proces și a cauzelor care-l generează și determină este deosebită întrucât numai în acest fel se creează posibilitatea reală de control și influențare a acestuia.

În context, devine necesară utilizarea unor metode, tehnici și instrumente care să poată:

- indica existența sau absența legăturii;
- măsura intensitatea acesteia;
- preciza sensul în care acționează;
- descrie, eventual, forma legăturii.

Statistica pune la dispoziție astfel de metode, tehnici și instrumente, unele simple, altele extrem de laborioase și complexe.

4.1. Tipuri de legături dintre fenomenele și procesele economice

Așa cum am amintit, formele de manifestare a relațiilor de interdependență dintre procese și fenomene sunt extrem de variate și cel mai adesea dificil de sesizat. Problema esențială care trebuie rezolvată în analiza legăturii dintre o *variabilă dependentă* (rezultativă, efect, explicată) notată de regulă cu y și una sau mai multe *variabile independente* (factoriale, cauzale, explicative) notate de regulă cu x_i se referă la răspunsul

întrebării: există o legătură între variabile sau modificarea variabilei rezultative-efect y este influențată de modificarea variabilei (variabilelor) cauză (x_i).

În practică se întâlnesc următoarele situații:

- variabila independentă x determină modificarea variabilei dependente y , caz în care între cele două există o *legătură univocă*;
- cele două variabile se influențează reciproc (*legătură reciprocă*);
- variabilele evoluează similar independent, influențate însă de o altă variabilă simultan;
- variabilele au o evoluție similară fără să existe vreo legătură între ele.

Din această cauză, pentru studiul sistematic al relațiilor dintre cele două tipuri de variabile este necesară clasificarea lor în funcție de anumite criterii:

a) după natura relației de interdependență (de cauzalitate), distingem:

- *legături funcționale (determinate);*
- *legături statistice sau stochastice;*

Legăturile funcționale sunt univoce și se realizează direct între un fenomen-cauză și un fenomen-efect. În această situație fenomenul-efect depinde de o singură cauză, ce poate fi identificată ori de câte ori se produce. De aici concluzia că dacă se mențin aceleași condiții, atunci unei valori a caracteristicii fenomenului-cauză îi corespunde o singură valoare a caracteristicii fenomenului-efect (caracteristicii rezultative). Aceste legături se mai numesc și legături de tip determinist și relația matematică prin care putem descrie un astfel de tip de legătură este:

$$y = f(x),$$

unde: y – fenomenul-efect;
 x – fenomenul-cauză;

Legăturile funcționale de tipul $y = f(x)$ se întâlnesc rar în activitatea economică deoarece, de cele mai multe ori, modificarea variabilei efect y este rezultatul influenței simultane a mai multor variabile-cauză (x_i).

Legăturile statistice sau stochastice, cel mai frecvent întâlnite în realitate, se caracterizează prin faptul că variabila rezultativă y este influențată de una sau mai multe variabile cauză x_i , considerate ca esențiale, dar pe lângă acestea există și acționează și alte variabile neînregistrate sau nespecificate.

Influența variabilelor nespecificate este luată în calcul sub forma variației reziduale (e) numită și eroare aleatoare astfel că relația matematică ce descrie o astfel de legătură este:

$$y = f(x) + e \quad \text{– cazul unei singure variabile cauză,}$$

$$y = f(x_i) + e \quad \text{– cazul mai multor variabile cauză.}$$

b) după numărul variabilelor factoriale, distingem:

- *legături simple* – o variabilă efect y și o singură variabilă cauză x . Exemplu: profitul (y), cifra de afaceri (x);
- *legături multiple* – o variabilă efect y și două sau mai multe variabile cauză x_i . Exemplu: salariul (y), numărul de ore lucrate (x_1), vechimea în muncă (x_2) și nivelul calificării (x_3);

c) după natura caracteristicilor, distingem:

- *legături de asociere* – se referă la raporturile de interdependență dintre caracteristicile calitative sau dintre o caracteristică numerică și una calitativă. Exemplu: ramura de activitate – salariul mediu; calificare – productivitate; domeniul de activitate și dimensiunea întreprinderii. Studiul statistic al legăturilor de asociere este posibil numai în situația în care variantele pot

- exprimate numeric. De exemplu, clasele de calitate ale produselor: 0 – produse inferioare, 1 – produse medii, 2 – produse superioare;
- *legături de corelație* – corelația statistică intervine numai în cazul legăturilor de tip cauză-efect dintre două sau mai multe *variabile cantitative*;
- d) *după direcția legăturii*, distingem:
- *legături directe* – modificarea într-un sens a variabilei (variabilelor) cauză este însoțită de modificarea în același sens a variabilei efect. Exemplu: salariu-productivitatea muncii, ofertă-preț;
 - *legături inverse* – modificarea într-un sens a variabilei (variabilelor) cauză este însoțită de modificarea în sens opus a variabilei efect. Exemplu: cerere-preț;
- e) *după forma funcției sau expresia analitică prin care se descrie legătura*, distingem:
- *legături liniare* – când legătura este pusă în evidență printr-o funcție liniară;
 - *legături neliniare* – când legătura este pusă în evidență printr-o funcție neliniară (parabolă, hiperbolă etc.);
- f) *după timpul realizării legăturii*, distingem:
- *legături sincrone* – când modificarea variabilei efect se produce aproape în același timp cu cea a variabilei (variabilelor) cauză. Exemplu: modificarea prețurilor și a cererii;
 - *legături asincrone* – când modificarea variabilei efect se produce la un anumit timp (defazat) de la modificarea variabilei (variabilelor) cauză. Exemplu: modificarea investițiilor în economie și modificarea produsului intern brut (PIB).

4.2. Metode statistice utilizate în studiul legăturii dintre fenomenele și procesele economice

Metodele statistice utilizate pentru studiul legăturii dintre două sau mai multe fenomene (puse în evidență de anumite caracteristici) pot fi grupate în:

- metode elementare;
- metode analitice.

Metodele elementare sunt cele prin care se poate determina existența legăturii dintre fenomene, a tăriei, a sensului și a formei acesteia dar nu cu o precizie foarte mare, ele fiind de obicei folosite pentru orientarea către metode de altă natură, mai rafinate, pentru determinarea elementelor de mai sus foarte precis.

Metodele analitice sunt cele prin care se pot determina aceleași elemente ca și prin metodele elementare, dar cu o precizie mult mai mare, ele permițând, de asemenea, și studiul legăturii dintre un fenomen efect și mai multe fenomene cauză simultan.

4.2.1. Metode elementare utilizate în studiul legăturii dintre fenomenele și procesele economice

În categoria metodelor elementare se includ:

- metoda seriilor paralele interdependente;
- metoda grupărilor;
- metoda tabelului de corelație;

- metoda grafică.

- **Metoda seriilor paralele interdependente**

Este o metodă relativ simplă ce se recomandă a fi aplicată în cazul existenței unui număr redus de valori pentru variabilele y și x .

Aplicarea acestei metode presupune parcurgerea următoarelor etape:

- ordonarea crescătoare a datelor ce caracterizează variabila independentă (cauză) x ;
- atașarea corespunzătoare a valorilor variabilei dependente (efect) y ;
- desprinderea concluziilor referitoare la forma și direcția legăturii în raport de reacția lui y la modificările x astfel:
 - datele se modifică în același sens → corelație directă;
 - datele se modifică în sensuri diferite → corelație inversă.

Mărima modificării lui y funcție de modificările lui x permite o apreciere empirică a intensității legăturii.

- **Metoda grupărilor**

Este în fapt o variantă a metodei precedente. Potrivit acestei metode se grupează în prealabil unitățile colectivității după caracteristica factorială x . Pentru fiecare grupă se calculează media caracteristicii dependente y . În coloane paralele, se înscriu grupările ordonate ale caracteristicii x și mediile corespunzătoare ale lui y . Prin compararea variației celor două caracteristici x și y obținem informații ce permit formularea de concluzii privind existența, sensul și intensitatea legăturii.

În cazul în care se analizează o singură variabilă rezultativă în raport de mai multe variabile factoriale (corelație multiplă) se înregistrează pe grupe valorile caracteristicilor factoriale înscriindu-se valorile respective în coloane distincte, în ordinea importanței lor pentru caracteristica rezultativă. Pentru caracteristica rezultativă se calculează valorile medii condiționate pe grupe.

- **Metoda tabelului de corelație**

Permite evidențierea tuturor elementelor necesare pentru confirmarea existenței unei legături dintre două fenomene, pe baza observației modului de manifestare. Pentru aceasta se utilizează măsurătorile unor variabile care caracterizează fenomenele supuse studiului.

Pentru utilizarea acestei metode este necesară distribuția bidimensională obținută prin prelucrarea perechilor de valori determinate prin măsurarea celor două variabile care caracterizează fenomenul cauză, respectiv fenomenul efect.

Modul în care se distribuie frecvențele în interiorul acestei distribuții (tabelul 4.1.) oferă toate elementele pentru evidențierea unei eventuale legături între cele două fenomene.

Tabelul 4.1.

$X \backslash Y$	y_1	y_2	y_3	$\dots y_j \dots$	y_n	F_X
x_1	f_{11}	f_{12}	f_{13}	$\dots f_{1j} \dots$	f_{1n}	F_{X_1}
x_2	f_{21}	f_{22}	f_{23}	$\dots f_{2j} \dots$	f_{2n}	F_{X_2}
x_3	f_{31}	f_{32}	f_{33}	$\dots f_{3j} \dots$	f_{3n}	F_{X_3}
\vdots	\vdots	\vdots	\vdots	\dots	\vdots	\vdots
x_i	f_{i1}	f_{i2}	f_{i3}	f_{ij}	f_{in}	F_{X_i}
\vdots	\vdots	\vdots	\vdots	\dots	\vdots	\vdots
x_n	f_{n1}	f_{n2}	f_{n3}	$\dots f_{nj} \dots$	f_{nn}	F_{X_n}
F_Y	F_{Y_1}	F_{Y_2}	F_{Y_3}	$\dots F_{Y_j} \dots$	F_{Y_n}	F

unde: X - variabila cauză;
 Y - variabila efect;
 $x_1 \dots x_n$ - valorile variabilei cauză;
 $y_1 \dots y_n$ - valorile variabilei efect;
 f_{ij} - frecvența de apariție a perechii de valori (x_i, y_j) ;
 F_{X_i} - frecvența de apariție a valorii x_i ;
 F_{Y_i} - frecvența de apariție a valorii y_i ;
 F - numărul total de perechi de valori (x_i, y_j) .

Atunci când este posibil se recomandă folosirea intervalelor egale de grupare, un număr suficient de grupe și același număr de grupe pentru ambele caracteristici.

Elementele care pot fi evidențiate cu ajutorul acestei metode:

1. *Existența legăturii dintre variabila X factorială și Y rezultativă:*

Dacă frecvențele f_{ij} se distribuie într-o bandă grupată de-a lungul unei diagonale a tabelului (figurile 4.1., 4.2.);

2. *Sensul legăturii:*

Dacă banda în care sunt grupate frecvențele f_{ij} se află pe diagonala tabelului care corespunde aceleiași sens de variație a valorilor corespunzătoare celor două variabile X, Y înseamnă că între cele două variabile există o legătură directă (figura 4.1.). Dacă se află pe cealaltă diagonală care corespunde sensului diferit de variație a celor două variabile X, Y atunci legătura dintre cele două variabile este inversă (figura 4.2.).

3. *Intensitatea legăturii:*

Este dată de lățimea benzii în care sunt grupate frecvențele f_{ij} . Cu cât banda este mai îngustă cu atât intensitatea legăturii crește (figurile 4.3., 4.4.).

4. *Forma legăturii:*

Este dată de forma benzii, putând fi liniară dacă forma benzii este liniară (figura 4.5.) sau neliniară dacă banda are altă formă decât cea liniară (figura 4.6.).

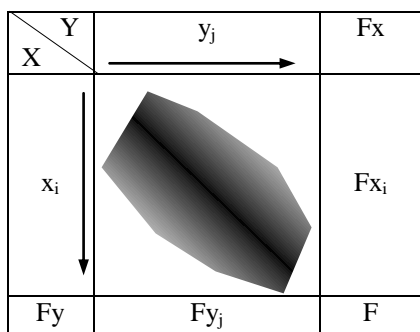


Figura 4.1. Corelație directă

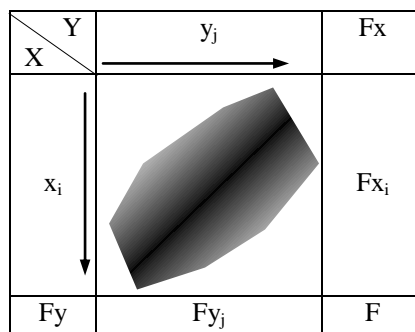


Figura 4.2. Corelație inversă

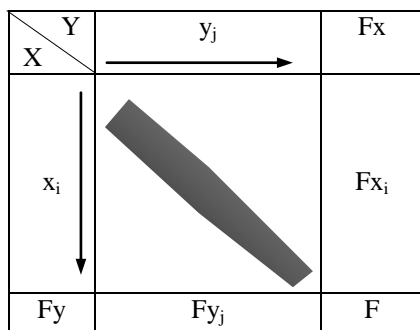


Figura 4.3. Corelație puternică

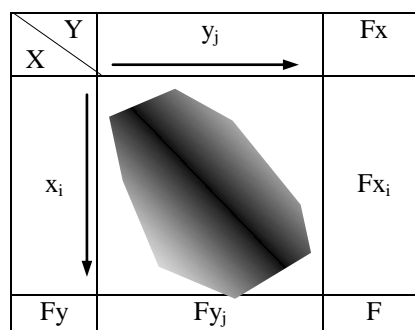


Figura 4.4. Corelație slabă

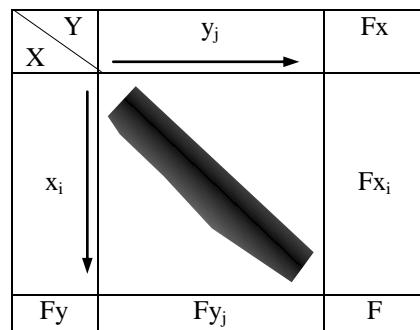


Figura 4.5. Corelație liniară

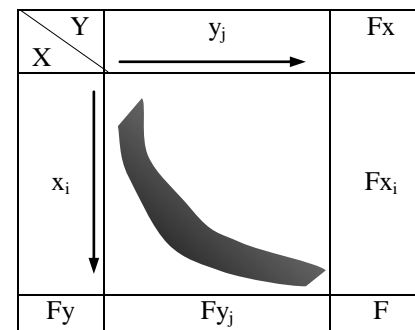


Figura 4.6. Corelație neliniară

Exemplul 4.1.

Pentru muncitorii unei firme se cunosc următoarele date (tabelul 4.2.):

Tabelul 4.2.

vechimea (ani) \ vârsta (ani)	1-5	5-10	10-15	15-20	20-25	25-30	Total
18-25	1	-	-	-	-	-	1
25-32	-	2	-	-	-	-	2
32-39	-	-	3	2	-	-	5
39-46	-	-	-	3	1	-	4
46-53	-	-	-	-	2	4	6
53-60	-	-	-	-	2	2	4
Total	1	2	3	6	5	5	22

În tabelul anterior se poate observa existența unei corelații directe, între cele două variabile, de intensitate mare și formă aproape liniară.

• **Metoda grafică**

Ca și metoda precedentă permite evidențierea prin apreciere vizuală a elementelor ce caracterizează legătura dintre două variabile.

În acest caz este necesară construirea corelogramei. Pe abscisă se trec valorile scării de reprezentare corespunzătoare variabilei cauză X, iar pe ordonată, valorile scării de reprezentare corespunzătoare variabilei Y. Prin unirea cu segmente de dreaptă a punctelor obținute reprezentând grafic perechile de valori (x_i, y_j) se obține corelograma (figura 4.7.).

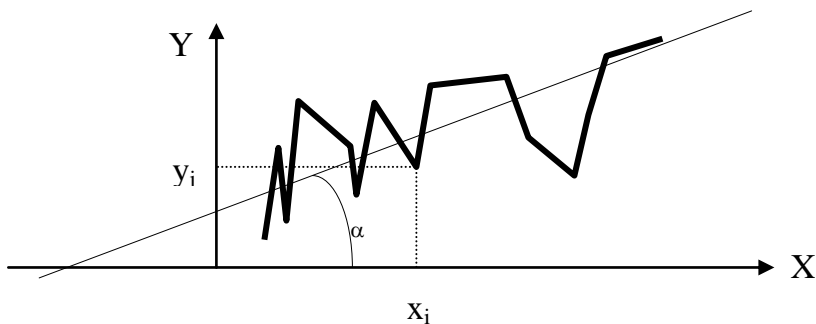


Figura 4.7. Corelograma.

Cu ajutorul acestei metode se pot evidenția:

1. *Existența legăturii:*

Se determină prin existența unghiului α (diferit de 0) realizat de linia de tendință cu orizontala .

2. *Sensul legăturii:*

- legătură directă - atunci când linia de tendință este ascendentă (figura 4.8.);
- legătură inversă - atunci când linia de tendință este descendentă (figura 4.9.);

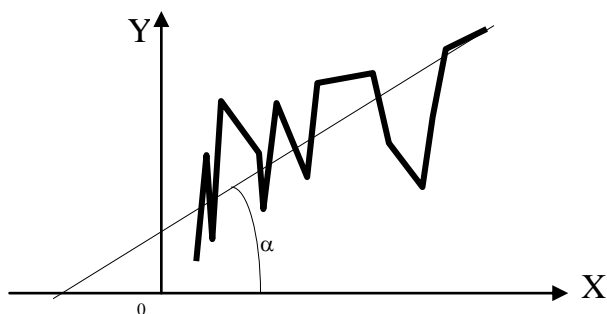


Figura 4.8. Corelație directă.

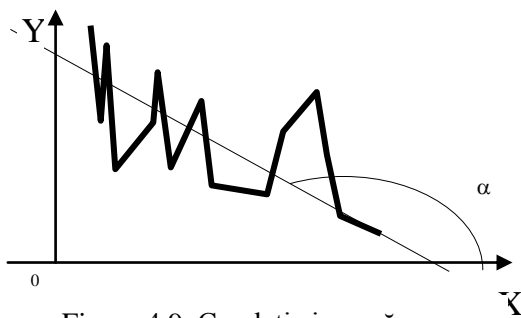


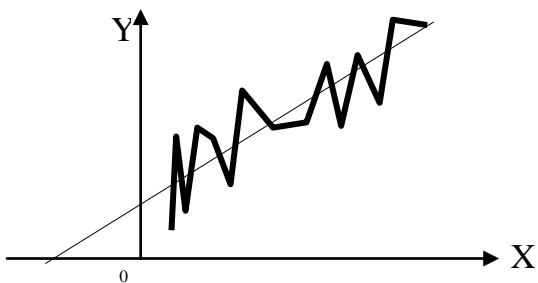
Figura 4.9. Corelație inversă.

3. Intensitatea legăturii:

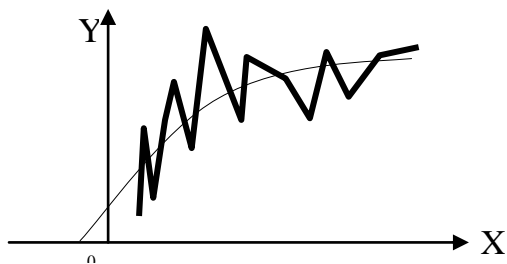
Este dată de mărimea unghiului α . Cu cât α este mai mare cu atât legătura este mai intensă (și invers).

4. Forma legăturii:

Este dată de forma corelogramei (figura 4.10.).



Corelație liniară



Corelație neliniară

Figura 4.10.

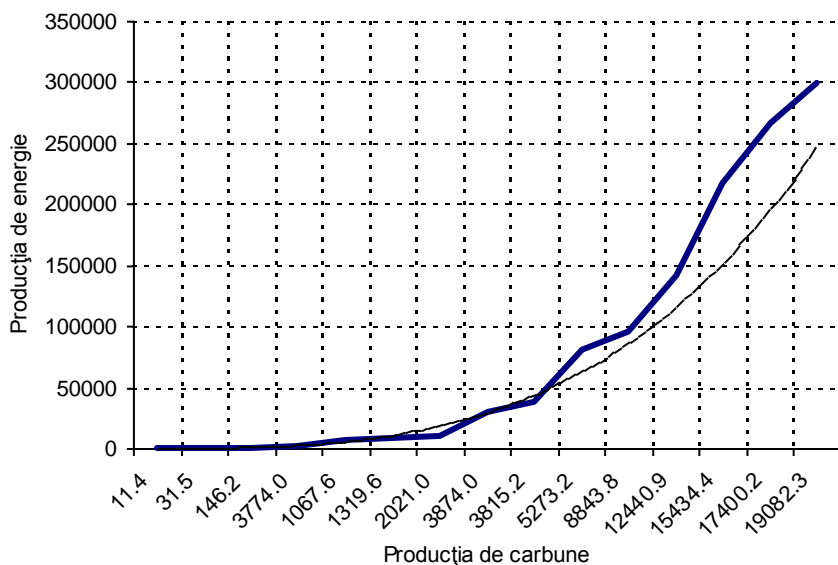
Exemplul 4.2.

În urma unui studiu efectuat pentru dependența dintre producția de energie electrică și termică și cea de cărbune, în perioada 1990-2004, s-au determinat datele (tabelul 4.3.):

Tabelul 4.3.
miliarde lei prețuri curente

Anul	Producția de cărbune	Producția de energie electrică și termică, gaze și apă
1990	11,4	61,2
1991	31,5	210,4
1992	146,2	735,3
1993	377,4	1553,2
1994	1067,6	5877,0
1995	1319,6	7581,3
1996	2021,0	10506,5
1997	3874,0	28763,9
1998	3815,2	37689,4
1999	5273,2	79889,1
2000	8843,8	94826,7
2001	12440,9	141103,0
2002	15434,4	217040,6
2003	17400,2	266203,2
2004	19082,3	299615,9

Sursa datelor: Anuarul statistic al României 2005



Din corelogramă rezultă că între cele două variabile (producția de cărbune și producția de energie) există o legătură directă, neliniară (exponențială) cu o intensitate destul de ridicată.

4.2.2. Metode analitice (parametrice) de măsurare a legăturilor dintre fenomenele și procesele economice

Metodele parametrice sunt cele care permit determinarea precisă atât a legăturii dintre două sau mai multe variabile cât și a intensității acesteia. Metodele parametrice sunt:

- metoda regresiei;
- metoda coeficientului (raportului) de corelație.

4.2.2.1. Metoda regresiei

Se bazează pe utilizarea funcțiilor matematice pentru descrierea formei legăturii dintre variabile. Funcția de regresie are forma generală:

$$y = f(x_1, x_2, \dots, x_n) + \varepsilon,$$

unde: y – variabila dependentă (efect);

x_1, x_2, \dots, x_n – variabilele independente (factorii de influență);

n – numărul factorilor de influență (variabilelor independente);

ε – variabila aleatoare (perturbatoare) sau eroarea ce sintetizează influența factorilor nespecificați (de regulă greu de cuantificat sau nesemnificativi).

În raport de numărul factorilor de influență înregistrați avem:

- regresie simplă (unifactorială);
- regresie multiplă (multifactorială).

Regresia simplă

Se bazează pe funcția:

$$y = f(x) + \varepsilon$$

și studiază variația unei caracteristici rezultative (dependente) y în raport cu un singur factor de influență x ceilalți factori fiind considerați neglijabili și cu acțiune constantă.

Alegerea funcției se face cu ajutorul graficului de corelație. Cele mai frecvent utilizate funcții de corelație simplă sunt:

- funcția liniară: $y = a + bx$;

- funcția parabolică: $y = a + bx + cx^2$;

- funcția exponențială: $y = ab^x$;

- funcția hiperbolică: $y = \frac{I}{a + bx}$ sau $y = a + b \frac{I}{x}$;

- funcția logaritmică: $y = a \cdot \lg x$ sau $y = a + b \cdot \lg x$,

unde: y – variabila dependentă sau rezultativă;

x – variabila independentă sau factorială;

a, b, c – parametrii ce urmează a fi determinați.

Pentru determinarea concretă a valorilor numerice a parametrilor se utilizează, de obicei, metoda celor mai mici pătrate, conform căreia pentru ca funcția de regresie aleasă să fie cu adevărat semnificativă trebuie să avem:

$$S = \sum_{i=1}^n (y_i - y_{x_i})^2 \rightarrow \min, \quad (*)$$

unde: $i = 1, 2, \dots, n$ – numărul unităților statistice observate;

y_i - valorile empirice (observate) ale variabilei dependente;

y_{x_i} – valorile teoretice exprimate prin ecuația de regresie.

Din condiția de mai sus rezultă că suma pătratelor abaterilor valorilor reale observate (y_i) de la valorile exprimate prin ecuația de regresie y_{x_i} trebuie să fie minimă.

Pentru calculul concret al parametrilor funcțiilor se anulează derivatele parțiale în raport cu fiecare parametru (a, b, c etc.) al expresiei (*), obținându-se un sistem cu un număr de ecuații egal cu numărul parametrilor funcției. Prin rezolvarea acestui sistem se obțin valorile concrete ale parametrilor.

În continuare vom ilustra metoda pentru funcția liniară

$$y = a + bx.$$

Vom avea condiția:

$$S = \sum [y_i - a - b \cdot x_i]^2 \rightarrow \min$$

$$\frac{\partial S}{\partial a} = -2 \sum (y_i - a - b \cdot x_i) \cdot 1 = 0$$

$$\frac{\partial S}{\partial b} = -2 \sum (y_i - a - b \cdot x_i) \cdot x_i = 0$$

După calcule simple și ținând cont de faptul că a și b sunt constante, rezultă sistemul:

$$\begin{cases} na + b \sum x_i = \sum y_i \\ a \sum x_i + b \sum x_i^2 = \sum x_i y_i \end{cases}$$

care va avea soluția:

$$b = \frac{n \sum x_i y_i - \sum x_i \cdot \sum y_i}{n \sum x_i^2 - (\sum x_i)^2}$$

$$a = \frac{\sum y_i \cdot \sum x_i^2 - \sum x_i \cdot \sum x_i y_i}{n \sum x_i^2 - (\sum x_i)^2} \text{ sau}$$

$$a = \bar{y} - b\bar{x}$$

Sistemele de ecuații normale specifice celor mai uzuale funcții de regresie sunt prezentate în tabelul 4.4.

Tabelul 4.4.

Sistemele de ecuații normale ale principalelor funcții de extrapolare

Tipul funcției	Funcția	Sistemul de ecuații normale corespunzător
Liniară	$y = a \pm bx$	$\begin{cases} na \pm b \sum x = \sum y \\ a \sum x \pm b \sum x^2 = \sum xy \end{cases}$
Parabolică	$y = a + bx + cx^2$	$\begin{cases} na + b \sum x + c \sum x^2 = \sum y \\ a \sum x + b \sum x^2 + c \sum x^3 = \sum xy \\ a \sum x^2 + b \sum x^3 + c \sum x^4 = \sum x^2 y \end{cases}$
Exponențială	$y = ab^x (*)$	$\begin{cases} n \lg a + \lg b \sum x = \sum \lg y \\ \lg a \sum x + \lg b \sum x^2 = \sum x \lg y \end{cases}$

Tipul funcției	Funcția	Sistemul de ecuații normale corespunzător
Putere	$y = ax^b$ (**)	$\begin{cases} n \lg a + b \sum \lg x = \sum \lg y \\ \lg a \sum \lg x + b \sum (\lg x)^2 = \sum \lg x \cdot \lg y \end{cases}$
Hiperbolică	$y = \frac{1}{a+bx}$ (***)	$\begin{cases} na + b \sum x = \sum \frac{1}{y} \\ a \sum x + b \sum x^2 = \sum \frac{1}{y} \cdot x \end{cases}$
Logistică clasică	$y = \frac{k}{1+be^{-cx}}$ (****)	$\begin{cases} (n-1)A + B \sum y = \sum \frac{\Delta y}{y} \\ A \sum y + B \sum y^2 = \sum \Delta y \end{cases}$
Törnquist	$y = \frac{kx}{x+a}$ (*****)	$\begin{cases} n \cdot \frac{1}{k} + \frac{a}{k} \sum \frac{1}{x} = \sum \frac{1}{y} \\ \frac{1}{k} \sum \frac{1}{x} + \frac{a}{k} \sum \frac{1}{x^2} = \sum \frac{1}{xy} \end{cases}$

(*) se liniarizează mai întâi, prin logaritmare, și $\Rightarrow \lg y = \lg a + x \cdot \lg b$.

(**) se liniarizează mai întâi, prin logaritmare, și $\Rightarrow \lg y = \lg a + b \lg x$.

(***) se pornește de la inversa sa: $\frac{1}{y} = a + bx$.

(****) Se scrie, în primul rând, forma transformată a acesteia $\frac{\Delta y}{y} = c - \frac{c}{k} \cdot y$ și se parcurg două etape de lucru. În prima etapă se calculează parametrii c și k , unde k reprezintă nivelul de saturație. În acest scop se notează $c=A$ și $-c/k = B$ și rezultă $\frac{\Delta y}{y} = A + By$. De

aici obținem sistemul de ecuații normale inserat în tabel. Observăm că parametrul A se înmulțește cu numărul termenilor seriei diminuat cu 1, deoarece Δy_i reprezintă diferența dintre y_i și y_{i-1} , și numărul termenilor Δy_i este mai mic cu 1 decât numărul termenilor y_i . În etapa a doua se calculează parametrul b , pornind de la relația

$$\frac{k}{y} - 1 = be^{-cx}.$$

Prin logaritmare (în acest caz folosim logaritmi naturali, pentru că modelul matematic conține numărul e , adică baza acestor logaritmi) se obține:

$$\ln\left(\frac{k}{y} - 1\right) = \ln b - cx \Leftrightarrow \ln b = \ln\left(\frac{k}{y} - 1\right) + cx.$$

Fiind vorba de o serie statistică cu n variabile x_i și y_i relația devine:

$$\ln b = \frac{1}{n} \left[\sum \ln\left(\frac{k}{y} - 1\right) + c \sum x \right].$$

Când variabila independentă este timpul, b se calculează după relația:

$$\ln b = \frac{1}{n} \left[\sum \ln \left(\frac{k}{y} - 1 \right) + c \cdot \frac{n(n+1)}{2} \right],$$

unde n reprezintă numărul anilor din perioada de analiză retrospectivă.

(*****) Ca și la funcția hiperbolică, se pornește de la inversa sa:

$$\frac{1}{y} = \frac{1+a}{kx} = \frac{1}{k} \cdot \frac{x+a}{x} = \frac{1}{k} \left(1 + \frac{a}{x} \right) = \frac{1}{k} + \frac{a}{k} \cdot \frac{1}{x}.$$

De regulă aceste sisteme se rezolvă cu ajutorul calculatorului utilizând programe adecvate.

Dacă totuși se apelează la rezolvarea manuală se construiește un tabel pentru determinarea constantelor ecuațiilor respective. De exemplu, pentru funcția parabolică de ordinul 2, tabelul va avea următorul conținut (tabelul 4.5.):

Tabelul 4.5.

Anii	y	x	x^2	x^3	x^4	xy	x^2y
t_1	y_1	x_1	x_1^2	x_1^3	x_1^4	$x_1 y_1$	$x_1^2 y_1$
t_2	y_2	x_2	x_2^2	x_2^3	x_2^4	$x_2 y_2$	$x_2^2 y_2$
:	:	:	:	:	:	:	:
t_i	y_i	x_i	x_i^2	x_i^3	x_i^4	$x_i y_i$	$x_i^2 y_i$
:	:	:	:	:	:	:	:
t_n	y_n	x_n	x_n^2	x_n^3	x_n^4	$x_n y_n$	$x_n^2 y_n$
Σ	Σy	Σx	Σx^2	Σx^3	Σx^4	Σxy	Σx^2y

În cazul rezolvării manuale, pentru ușurarea calculelor, se poate proceda astfel: în locul valorilor x_i se consideră $x_i = x_i - \bar{x}$ unde:

$$\bar{x} = \frac{\sum x_i}{n} \text{ - media aritmetică simplă a valorilor } x_i \text{ observate}$$

În acest fel vom avea:

$$x_1' = x_1 - \frac{x_1 + x_2 + \dots + x_n}{n}$$

...

$$x_i' = x_i - \frac{x_1 + x_2 + \dots + x_n}{n}$$

...

$$x_n' = x_n - \frac{x_1 + x_2 + \dots + x_n}{n}$$

Sumând aceste relații vom obține:

$$\sum x_i' = \sum x_i - n \cdot \frac{\sum x_i}{n} = 0$$

În acest fel, în locul variabilelor x_i se lucrează cu variabilele x_i' a căror sumă este 0.

De aici rezultă că și sumele de puteri impare ale lui x_i' devin tot 0.

Prin acest procedeu se simplifică mult calculele astfel:

a) în cazul funcției de regresie liniară $y = a + bx$, sistemul de ecuații normale devine:

$$\begin{cases} na = \sum y \\ b \sum x^2 = \sum x' y \end{cases} \Rightarrow \begin{cases} a = \frac{\sum y}{n} = \bar{y} \\ b = \frac{\sum x' y}{\sum x^2} = \frac{\sum (x_i - \bar{x})(y_i - \bar{y})}{n \sigma_x^2} = \frac{\sum (x_i - \bar{x})(y_i - \bar{y})}{n} \cdot \frac{1}{\sigma_x^2} = \text{cov}(x, y) \cdot \frac{1}{\sigma_x^2} \end{cases}$$

b) în cazul funcției parabolice de ordinul II, $y = a + bx + cx^2$, vom avea:

$$\begin{cases} na + c \sum (x')^2 = \sum y \\ b \sum (x')^2 = \sum x' y \\ a \sum (x')^2 + c \sum (x')^4 = \sum (x')^2 y \end{cases}$$

După aflarea parametrilor funcției de regresie se calculează valorile teoretice (ajustate) ale variabilei y pe baza ecuației explicitate (ecuația în care a și b au valorile numerice rezultate din calcule).

Pentru verificarea calculului parametrilor funcției de regresie se utilizează relația $\sum_{i=1}^n y_i = \sum_{i=1}^n y_{x_i}$ ceea ce arată că prin ajustare nu se face decât o redistribuire a influenței factorilor.

Exemplul 4.3.

În ultimii ani o firmă a obținut următoarele rezultate economice (tabelul 4.6.):

Tabelul 4.6.

Anii	t_1	t_2	t_3	t_4
Cifra de afaceri (mii \$)	2200	3400	3800	4700
Bugetul de publicitate (mii \$)	45	54	76	81

Conducerea întreprinderii dorește să știe care va fi cifra de afaceri dacă bugetul de publicitate va fi majorat la 150 mii \$.

Pentru rezolvarea acestei probleme ne propunem utilizarea unei funcții de corelație simplă de forma:

$y = a + bx$, unde

y = cifra de afaceri (în mii \$)

x = bugetul de publicitate (în mii \$)

Aplicând metoda celor mai mici pătrate simplificată vom scrie

$\min S = \sum_{i=1}^n (y_i - a - bx')^2$ de unde anulând derivatele parțiale în raport de a și b

obținem:

$$\frac{\partial S}{\partial a} = -2 \sum_{i=1}^n (y_i - a - bx') \cdot 1 = 0$$

$$\frac{\partial S}{\partial b} = -2 \sum_{i=1}^n (y_i - a - bx') \cdot x' = 0$$

Dezvoltând va rezulta sistemul

$$\begin{cases} na + b\sum x' = \sum y \\ a\sum x' + b\sum (x')^2 = \sum x'y \end{cases} \text{ unde } x' = x_i - \bar{x}$$

Datele necesare rezolvării sistemului se vor obține din tabelul 4.7.

Tabelul 4.7.

Anii	y	x	x' = x _i - \bar{x}	(x') ²	x'y	y' _i	y _i - y'	$\frac{y_i - y'}{y_i} \cdot 100$	$\left(\frac{y_i - y'}{y_i} \cdot 100\right)^2$	y ²
t ₁	2200	45	-19	361	-41800	2468,79	-268,79	-12,22	149,33	4840000
t ₂	3400	54	-10	100	-34000	2969,10	430,90	12,67	160,53	11560000
t ₃	3800	76	12	144	45600	4192,08	-392,08	10,32	106,50	14440000
t ₄	4700	81	17	289	79900	4470,03	229,97	4,89	23,91	22090000
Σ	14100	256	0	894	49700	*	*	*	440,27	52930000

$$\bar{x} = \frac{256}{4} = 64$$

Rezultă sistemul:

$$\begin{cases} 4a = 14110 \\ 894b = 49700 \end{cases}$$

De unde a=3525 și b=55,59.

Rezultă ecuația dreptei de regresie:

$$Y = 3525 + 55,59x' \quad \text{și}$$

$$y'_1 = 2468,79$$

$$y'_2 = 2969,1$$

$$y'_3 = 4192,08$$

$$y'_4 = 4470,03$$

Calculăm abaterea medie pătratică procentuală cu relația:

$$\sigma_{\%} = \sqrt{\frac{\sum \left(\frac{y_i - y'}{y_i} \cdot 100\right)^2}{n}} = \sqrt{\frac{440,27}{4}} = 10,50\%$$

Coeficientul de corelație care măsoară legătura dintre y și x se poate calcula cu relația:

$$r = \frac{n\sum x'y - \sum x' \cdot \sum y}{\sqrt{[n\sum y^2 - (\sum y)^2][n\sum (x')^2 - (\sum x')^2]}}$$

$$r = \frac{4 \cdot 49700}{\sqrt{[4 \cdot 52930000 - 198810000][4 \cdot 894 - 0^2]}}$$

$$r = \frac{198800}{\sqrt{4616616000}} = 0,9252$$

Cum coeficientul de corelație este apropiat de 1 putem utiliza ecuația

$$y = 3552 + 55,59x'$$

pentru estimarea cifrei de afaceri.

Cifra de afaceri y pentru un buget de publicitate x de 150 mii \$ va fi
 $y = 3552 + 55,59(150 - 64) = 8305,74$ mii \$.

Exemplul 4.4.

Pentru ultimele cinci luni se cunosc datele statistice privind prețurile unitare practicate (x) și cantitatea corespunzătoare de produse vândute (y), conform tabelului 4.8. (coloanele 1 și 2). Să se studieze relația dintre cele două variabile și să se estimeze vânzările pentru un preț $x = 9,5$ cu o probabilitate de 95 %.

Tabelul 4.8.

Luna	x	y	$x - \bar{x}$	$(x - \bar{x})^2$	$y - \bar{y}$	$(x - \bar{x})(y - \bar{y})$	x^2	xy	y^2	$(y - \bar{y})^2$
1	13	25	1,5	2,25	-14	-21	169	325	625	196
2	12	30	0,5	0,25	-9	-4,5	144	360	900	81
3	11,5	45	0	0	6	0	132,25	517,5	2025	36
4	11	45	-0,5	0,25	6	-3	121	495	2025	36
5	10	50	-1,5	2,25	11	-16,5	100	500	2500	121
Total	57,5	195	0	5	0	-45	666,25	2197,5	8075	470
Media	11,5	39	-	-	-	-	-	-	-	-

În acest scop vom parcurge următoarele etape:

- 1) Se reprezintă grafic datele statistice x și y (vezi figura 4.11., linie discontinuă)

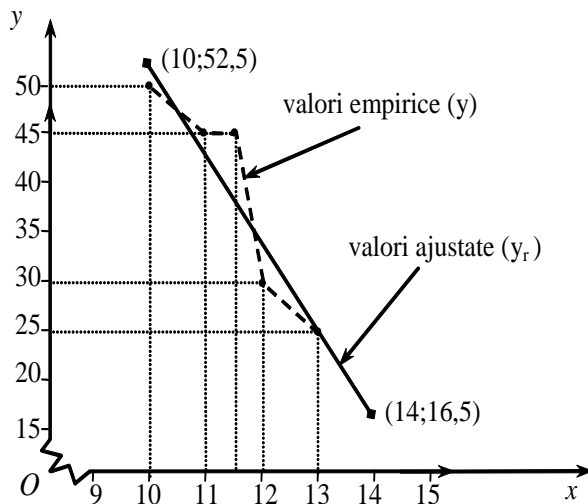


Figura 4.11. Evoluția vânzărilor funcție de preț.

Dacă punctele marcate sunt relativ aliniate înseamnă că între x și y există o legătură liniară. Dacă punctele rezultate au un mare grad de împrăștiere sau au o alură curbilinie înseamnă că între cele două variabile există alt tip de dependență și, în consecință, nu se mai parcurg etapele cerute de regresia simplă.

- 2) Se calculează coeficientul de corelație (r), astfel:

$$r = \frac{\sum(x - \bar{x})(y - \bar{y})}{\sqrt{\sum(x - \bar{x})^2} \cdot \sqrt{\sum(y - \bar{y})^2}}$$

Cu datele din tabelul 4.8. rezultă $r = \frac{-45}{\sqrt{5} \cdot \sqrt{470}} = -0,9282791$.

Interpretare:

- dacă r este pozitiv, între cele două variabile există o legătură directă, iar dacă r este negativ (cazul de mai sus), variabilele se află într-o legătură inversă (când x crește, y scade și invers);
- valoarea absolută a lui r constituie un indiciu al corelației (legăturii) dintre x și y : legătura este foarte strânsă când r tinde spre 1 și foarte redusă când r tinde spre zero.

3) Se calculează coeficientul de determinare (R^2), astfel:

$$R^2 = r^2$$

în cazul dat, $R^2 = (-0,9282791)^2 = 0,8617021$

Interpretare:

- dacă R^2 ia valoarea minimă $R^2 = 0$ înseamnă că între cele două variabile x și y nu există nici o legătură liniară, iar dacă $R^2 = 1$ – maxim, înseamnă că variabilele sunt perfect legate una de alta.

Complementul lui R^2 , respectiv $(1-R^2)$, se numește coeficient de nedeterminare și arată proporția în care y nu este explicat de x ci de alți factori neluați în considerare.

În cazul de mai sus se poate afirma că factorul preț explică volumul vânzărilor în proporție de $0,8617021 = 86,17\%$, în timp ce restul $(1-R^2) = 0,1382979 = 13,83\%$ se datorează influenței altor factori.

4) Ajustarea datelor empirice

Convinși fiind de mărimea lui R^2 că între cele două variabile există o indiscutabilă legătură liniară, se caută acea dreaptă numită „dreaptă de regresie” de forma $y_r = a + bx$ care să reprezinte cât mai fidel datele reale care au o evoluție greu de exprimat matematic.

În acest scop, pe baza metodei celor mai mici pătrate care minimizează suma pătratelor diferențelor dintre y și y_r , coeficienții a și b rezultă din sistemul:

$$\begin{cases} na + b\Sigma x = \Sigma y \\ a\Sigma x + b\Sigma x^2 = \Sigma xy \end{cases}$$

unde: n = numărul datelor empirice

Cu datele problemei (vezi și tabelul 4.8.), avem:

$$\begin{cases} 5a + 57,5b = 195 & \Rightarrow & a = 142,5 \\ 57,5a + 666,25b = 2197,5 & & b = -9 \end{cases}$$

Coeficienții a și b mai pot fi calculați și direct, astfel:

$$b = \frac{\sum(x-\bar{x})(y-\bar{y})}{\sum(x-\bar{x})^2} = \frac{-45}{5} = -9$$

$$a = \bar{y} - b\bar{x} = 39 - (-9 \cdot 11,5) = 39 + 103,5 = 142,5$$

Ca urmare, pentru cazul dat,

$$y_r = 142,5 - 9x \quad (\text{vezi și reprezentarea grafică în figura 4.11.)}$$

Observații:

- coeficientul b (panta dreptei de regresie) are următoarea semnificație: o scădere cu o unitate a lui x (prețul) conduce la o creștere a lui y (vânzări) cu b unități și invers.
- dacă coeficienții a și b au fost bine calculați trebuie să existe relația:

$$\bar{y} = a + b\bar{x}$$

5) Calculul erorii dintre y și y_r

Datorită ajustării efectuate, între valorile reale ale lui y și valorile corespunzătoare de pe dreapta $y_r = 142,5 - 9x$, vor exista în mod firesc deosebiri (vezi figura 4.11. și tabelul 4.9., coloanele 2 și 3).

Tabelul 4.9.

Luna	x	y	y_r	$(y_r - \bar{y}_r)^2$	$(y - y_r)^2$
1	13	25	25,5	182,25	0,25
2	12	30	34,5	20,25	20,25
3	11,5	45	39	0	36
4	11	45	43,5	20,25	2,25
5	10	50	52,5	182,25	6,25
Total	57,5	195	195	405	65
Medie	11,5	39	39	-	-

Important este ca aceste abateri să fie cât mai mici, respectiv, diferențele (erorile) să fie cuantificate pentru a fi luate în considerare atunci când se va face previziunea. Deoarece întotdeauna $\bar{y} = \bar{y}_r$, rezultă că abaterile dintre valorile lui y și y_r nu pot fi măsurate decât de dispersie, în termenii analizei dispersionale, se pot face următoarele asocieri:

- variația totală a datelor y reale = $\sum (y - \bar{y})^2 = 470$ (vezi col. 10, tabelul 4.8.);
- variația totală a datelor y_r estimate = $\sum (y_r - \bar{y}_r)^2 = 405$ (vezi col.4, tabelul 4.9.);
- variația dintre datele reale și cele estimate = $\sum (y - y_r)^2 = 65$ (vezi col.5, tabelul 4.9.).

Cele de mai sus se interpretează astfel: din cele 470 de unități ale variației datelor reale, 450 sunt explicate de funcția y_r , iar restul de 65 rămân neexplicate.

În aceeași optică, R^2 și $(1 - R^2)$ se pot calcula astfel:

$$R^2 = \frac{\sum (y_r - \bar{y}_r)^2}{\sum (y - \bar{y})^2} = \frac{\sum (y_r - \bar{y})^2}{\sum (y - \bar{y})^2}$$

$$(1 - R^2) = \frac{\sum (y - y_r)^2}{\sum (y - \bar{y})^2}$$

Cu datele problemei, $R^2 = \frac{405}{470} = 0,8617021$ și $(1 - R^2) = \frac{65}{470} = 0,1382979$, adică

exact rezultatele obținute anterior.

Variația dintre datele reale și cele estimate $\sum (y - y_r)^2$ se mai numește și variație reziduală sau variație neexplicată sau eroare de estimăție. Se exprimă sub forma abaterii medii pătratice (e) astfel:

$$e = \sqrt{\frac{\sum (y - y_r)^2}{\nu}}, \quad \text{unde } \nu = \text{numărul gradelor de libertate}$$

Observații:

- deoarece se cunosc doi parametri (coeficienții a și b), în cazul regresiei simple $\nu = n - 2$;
- pentru un anumit nivel de semnificație (α) și un număr dat de grade de libertate (ν), R^2 trebuie să aibă o valoare minimă teoretică. Pentru ca dreapta de regresie să poată fi

considerată ca reprezentativă (în cazul de față, unde $\nu = (n - 2) = (5 - 2) = 3$, pentru $\alpha = 0,05$, $R^2_{calculat} = 0,8617021 > R^2_{teoretic} = 0,7714$).

În situația concretă presupusă,

$$e = \sqrt{\frac{65}{5-2}} = 4,65$$

6) Se estimează cererea (volumul vânzărilor)

Prin similitudine cu eșantionajul, unde $m = \bar{x} \pm t \cdot S_{\bar{x}}$

$$P = (a+bx) \pm t \cdot e$$

Eșantionul fiind mic, pentru $\alpha = 0,05$ și $\nu = n-1 = 5-1 = 4$, $t = 2,778$ și, ca urmare:

$$P = (142,5-9-9,5) \pm 12,9, \text{ respectiv } 44,1 < y < 69,9$$

Concluzie. Cu o probabilitate de 95%, în cazul practicării unui preț de 9,5 unități bănești / produs, în luna următoare cererea (vânzările) se va situa între 44 și 70 bucăți.

Regresia multifactorială

În viața economică reală se înregistrează influențe multiple chiar și asupra celor mai simple fenomene și procese. De aceea, modelul regresiei unifactoriale este un model simplificat. Mult mai apropiat de realitate este modelul regresiei multiple (multifactoriale) bazat pe ecuația de tipul:

$$y_{x_i} = f(x_1, x_2, \dots, x_n) + \varepsilon.$$

În acest model variația fenomenului efect (y_{x_i}) se datorează influenței unei multitudini de factori cauză (independenți) explicați prin intermediul variabilelor x_1, x_2, \dots, x_n și aleatori explicați prin valoarea reziduală ε .

Având în vedere complexitatea abordării multifactoriale cel mai accesibil și utilizat model este cel liniar de forma:

$$y_{x_i} = a_0 + a_1x_1 + a_2x_2 + \dots + a_nx_n,$$

unde: y_{x_i} – valorile ajustate (teoretice) ale variabilei efect (dependente);

x_1, x_2, \dots, x_n – factorii de influență înregistrați;

a_0 – parametru ce exprimă influența factorilor neînregistrați;

a_1, a_2, \dots, a_n – parametri, coeficienți parțiali de regresie care arată cu cât se modifică y_{x_i} , atunci factorii de influență înregistrați x_1, x_2, \dots, x_n se modifică cu o

unitate, iar toate celelalte variabile rămân constante (a_1 arată cu cât se modifică y_{x_i} dacă x_1 se modifică cu o unitate ceilalți factori rămânând neschimbați

ș.a.m.d.).

Pentru determinarea parametrilor funcției se utilizează metoda celor mai mici pătrate pornind de la relația:

$$S = \sum_{i=1}^n (y_i - a_0 - a_1x_1 - a_2x_2 - \dots - a_nx_n)^2 \rightarrow \min$$

Anulând derivatele parțiale ale expresiei de mai sus în raport cu parametrii $a_0, a_1, a_2, \dots, a_n$ se obține un sistem de ecuații normale care rezolvat conduce la valorile numerice ale parametrilor $a_0, a_1, a_2, \dots, a_n$. Acest sistem este:

$$\begin{cases} na_0 + a_1 \sum x_1 + a_2 \sum x_2 + \dots + a_n \sum x_n = \sum y \\ a_0 \sum x_1 + a_1 \sum x_1^2 + a_2 \sum x_1 x_2 + \dots + a_n \sum x_1 x_n = \sum x_1 y \\ a_0 \sum x_2 + a_1 \sum x_1 x_2 + a_2 \sum x_2^2 + \dots + a_n \sum x_2 x_n = \sum x_2 y \\ \vdots \\ a_0 \sum x_i + a_1 \sum x_1 x_i + a_2 \sum x_2 x_i + \dots + a_n \sum x_i x_n = \sum x_i y \\ \vdots \\ a_0 \sum x_n + a_1 \sum x_1 x_n + a_2 \sum x_2 x_n + \dots + a_n \sum x_n^2 = \sum x_n y \end{cases}$$

Dacă avem în vedere modelul liniar cu numai doi factori de influență $y_{x_i} = a_0 + a_1 x_1 + a_2 x_2$ va rezulta sistemul

$$\begin{cases} na_0 + a_1 \sum x_1 + a_2 \sum x_2 = \sum y \\ a_0 \sum x_1 + a_1 \sum x_1^2 + a_2 \sum x_1 x_2 = \sum x_1 y \\ a_0 \sum x_2 + a_1 \sum x_1 x_2 + a_2 \sum x_2^2 = \sum x_2 y \end{cases}$$

În cazul regresiei multifactoriale este posibilă existența unor interdependențe între factorii de influență (fenomenul de multicolaritate) caz în care efectele acestuia afectează și influențează concluziile analizei.

Calitatea ajustării prin intermediul funcției de regresie se apreciază cu ajutorul indicatorilor:

- a) eroarea standard calculată prin intermediul relației:

$$S_{y_i / y_{x_i}} = \sqrt{\frac{\sum (y_i - y_{x_i})^2}{n}}$$

Relația de mai sus este în fapt abaterea medie pătratică a valorilor reale y_i față de cele teoretice y_{x_i} .

- b) coeficientul de eroare e (sau abaterea medie pătratică procentuală $\sigma_{\%}$) calculat pe baza relației:

$$e = \frac{S_{y_i / y_{x_i}}}{\bar{y}} \cdot 100$$

$$\sigma_{\%} = \sqrt{\frac{\sum \left(\frac{y_i - y_{x_i}}{y_i} \cdot 100 \right)^2}{n}}$$

Funcția aleasă este cu atât mai reprezentativă cu cât valorile celor doi indicatori sunt mai apropiate de 0.

- c) coeficientul de determinație calculat pe baza relației:

$$R = \left(1 - \frac{\sum (y_i - y_{x_i})^2}{\sum (y_i - \bar{y})^2} \right) \cdot 100$$

care arată proporția în care variabila (variabilele) independentă x (x_i) explică variația caracteristicii variabilei dependente y_{x_i} .

Variația totală a lui y față de media sa \bar{y} se poate scrie:

$$y_i - \bar{y} = (y_{x_i} - \bar{y}) + (y_i - y_{x_i}),$$

unde: $y_i - \bar{y}$ - variația totală;

$(y_{x_i} - \bar{y})$ - variația explicată de regresie;

$(y_i - y_{x_i})$ - variația neexplicată de regresie.

Pornind de la relația de mai sus și calculând:

- dispersia totală a lui y : $\sigma_y^2 = \frac{\sum (y_i - \bar{y})^2}{n}$ care exprimă influența tuturor factorilor

asupra variabilei efect y putem descompune această dispersie (influență) în două:

- dispersie explicată de factorii cuprinși în ecuația de regresie

$$\sigma_{y/x_i}^2 = \frac{\sum (y_{x_i} - \bar{y})^2}{n}$$
 și care exprimă influența tuturor factorilor explicați

asupra variației lui y ;

- dispersie neexplicată prin modelul ecuației de regresie (reziduală)

$$\sigma_{y/r}^2 = \frac{\sum (y_i - y_{x_i})^2}{n}$$
 care exprimă influența factorilor reziduali

(neexplicați în model).

Pentru validarea modelului de regresie se utilizează testul Fisher-Snedecor (testul F) conform căruia

$$F_{calc} = \frac{\frac{\sum (y_{x_i} - \bar{y})^2}{k-1}}{\frac{\sum (y_i - y_{x_i})^2}{n-k}}$$

unde: k – numărul parametrilor funcției de regresie (modelului);

n – numărul perechilor de valori $x_i \cdot y_i$.

Valoarea calculată F_{calc} se compară cu valoarea teoretică a lui F obținută din tabel

$F_{\alpha, k-1, n-k}$ pentru un prag de semnificație (probabilitate) α și $k-1$, $n-k$ grade de libertate.

Modelul de regresie se validează dacă:

$$F_{calc} > F_{\alpha, k-1, n-k}$$

Modelul bazat pe regresie constituie numai o ipoteză statistică prin care se exprimă tendința medie a legăturii dintre variabila dependentă y și variabila (variabilele) independentă x (x_i) și reprezintă primul pas pentru măsurarea intensității legăturii, lucru ce se realizează prin metoda corelației.

4.2.2.2. Metoda corelației

Pentru măsurarea intensității legăturii dintre variabila dependentă y și variabila (variabilele) independentă x (x_i) se utilizează metoda corelației. În funcție de natura legăturii dintre variabila dependentă y și variabila (variabilele) independentă x (x_i) – legătura directă sau inversă – corelația poate să fie pozitivă (în cazul legăturii directe) sau negativă (în cazul legăturii inverse). În cadrul acestei metode se utilizează indicatorii: covarianța, coeficientul de corelație și raportul de corelație.

Covarianța surprinde existența și direcția legăturii dintre variabila dependentă y și o variabilă independentă (x). Se calculează sub forma mediei aritmetice simple a produselor

abaterilor celor două variabile corelate y și x de la mediile lor aritmetice \bar{y} și \bar{x} cu ajutorul relației:

$$\text{cov}(x, y) = \frac{1}{n} \cdot \sum (x_i - \bar{x})(y_i - \bar{y}).$$

Valorile pozitive ale acestui indicator reflectă o legătură directă, iar cele negative o legătură inversă.

Valorile mari ale indicatorului arată o legătură puternică, în timp ce valorile apropiate de zero semnifică lipsa de legătură între variabilele y și x .

Coeficientul de corelație simplă măsoară intensitatea *legăturii liniare* dintre două variabile y – rezultativă (endogenă) și x – factorială.

Deoarece fenomenele aflate în relații de interdependență prezintă, în mod normal, o poziție similară a valorilor individuale (y_i și x_i) față de media corespunzătoare (\bar{y} și \bar{x}) va rezulta că și abaterile normalizate $\frac{(y_i - \bar{y})}{\sigma_y}$, respectiv $\frac{(x_i - \bar{x})}{\sigma_x}$ au mărimi apropiate pentru valorile perechi (y_i, x_i).

Pentru a obține mărimea sintetică a abaterilor normalizate la nivelul întregii colectivități se calculează *coeficientul de corelație* $r_{y/x}$ cu ajutorul relației:

$$r_{y/x} = \frac{\sum \left(\frac{x_i - \bar{x}}{\sigma_x} \right) \left(\frac{y_i - \bar{y}}{\sigma_y} \right)}{n}.$$

numit și *coeficientul de corelație liniară al lui Pearson*.

$$\text{Se observă că } r_{y/x} = \frac{\sum (x_i - \bar{x})(y_i - \bar{y})}{n \sigma_x \sigma_y}.$$

Ținând cont de relațiile de calcul pentru σ_x și σ_y se obține relația de calcul simplificat:

$$r = \frac{n \sum x_i y_i - (\sum x_i)(\sum y_i)}{\sqrt{[n \sum x_i^2 - (\sum x_i)^2][n \sum y_i^2 - (\sum y_i)^2]}}.$$

Examinând elementele din formula de calcul simplificat se observă că exceptând $\sum y_i^2$ toate celelalte se găsesc în tabelul de calcul al parametrilor funcției de regresie liniară. De aceea, în tabelele de lucru pentru calculul parametrilor funcției de regresie se recomandă includerea coloanei y_i^2 .

În cazul cunoașterii sumei abaterilor de la medie pentru cele două variabile se recomandă utilizarea relației echivalente:

$$r_{y/x} = \frac{\sum (x_i - \bar{x})(y_i - \bar{y})}{\sqrt{\sum (x_i - \bar{x})^2 \cdot \sum (y_i - \bar{y})^2}}.$$

Între coeficientul de corelație liniară simplă $r_{y/x}$, coeficientul de regresie (b) al funcției liniare $y = a + bx$ și abaterile medii pătratice ale lui x și y există relația:

$$b = r_{y/x} \cdot \frac{\sigma_y}{\sigma_x}$$

$$b = \frac{\sum (x_i - \bar{x})(y_i - \bar{y})}{\sum (x_i - \bar{x})^2}$$

Deoarece $\sum y_i = \sum y_{x_i}$ se poate înlocui \bar{y} cu \bar{y}_{x_i} și cum cele două medii sunt cunoscute (\bar{y} și \bar{x}) rezultă:

$$a = \bar{y} - b\bar{x}.$$

Cum legătura dintre cele două variabile este liniară rezultă:

$$r_{y/x} = \sqrt{b_{y/x} \cdot b_{x/y}},$$

unde: $b_{y/x}$ – coeficientul de regresie corespunzător dependenței lui y față de x (coeficientul b din ecuația $y=a+bx$);

$b_{x/y}$ – coeficientul de regresie corespunzător dependenței lui x față de y (coeficientul b din ecuația $x=a+by$).

Coeficientul de corelație poate lua valori între -1 și 1 astfel că:

$$r_{y/x} \leq |1|.$$

Dacă: $r_{y/x} \in [-1, 0)$ – legătură inversă;

$r_{y/x} \in (0, 1]$ – legătură directă;

$r_{y/x} = 0$ – cele două variabile nu se corelează liniar.

În practică, funcție de valorile lui $r_{y/x}$ avem:

- $|r_{y/x}| \in (0, 0,2)$ – nu există legătură;
- $|r_{y/x}| \in [0,2; 0,5)$ – există o legătură slabă;
- $|r_{y/x}| \in [0,5; 0,75)$ – legătură de intensitate medie;
- $|r_{y/x}| \in [0,75; 0,95)$ – legătură puternică;
- $|r_{y/x}| \in [0,95; 1]$ – relație (legătură) deterministă.

Formula anterioară este utilizată în cazul unui număr relativ redus de valori individuale pentru variabilele y_i și x_i .

În situația unor observații mai ample datele statistice pot fi sistematizate prin grupări simple sau combinate.

Dacă s-a utilizat gruparea simplă și variabilele y_i și x_i au frecvențe comune f_i formula de calcul a coeficientului de corelație simplă devine:

$$r_{y/x} = \frac{\sum f_i \cdot \sum x_i y_i f_i - \sum x_i f_i \cdot \sum y_i f_i}{\sqrt{\left[\sum f_i \cdot \sum x_i^2 f_i - \left(\sum x_i f_i \right)^2 \right] \left[\sum f_i \cdot \sum y_i^2 f_i - \left(\sum y_i f_i \right)^2 \right]}}$$

În situația unei distribuții bidimensionale variabilele y_i și x_i au frecvențe distincte f_j cât și frecvențe comune f_{ij} . Din acest considerent formula de calcul a coeficientului de corelație simplă devine:

$$r_{y/x} = \frac{\left(\sum_i \sum_j f_{ij} \cdot \sum_i \sum_j x_i y_j f_{ij} \right) - \sum_i x_i f_i \cdot \sum_j y_j f_j}{\sqrt{\left[\sum_i f_i \cdot \sum_i x_i^2 f_i - \left(\sum_i x_i f_i \right)^2 \right] \left[\sum_j f_j \cdot \sum_j y_j^2 f_j - \left(\sum_j y_j f_j \right)^2 \right]}}.$$

Verificarea semnificației coeficientului de corelație liniară simplă se face cu ajutorul testului „t” (STUDENT) parcurgând următorii pași:

- 1) Se determină $t_{calculat}$ cu relația:

$$t_{calc} = \frac{|r_{y/x}|}{\sqrt{1 - r_{y/x}^2}} \cdot \sqrt{n-2},$$

unde: $r_{y/x}$ – coeficientul de corelație liniară simplă;
 n – numărul observațiilor;
 $n-2$ – numărul gradelor de libertate.

- 2) Se compară valoarea rezultată din calcul (t_{calc}) cu valoarea teoretică din tabelul repartiției Student (t_{tab}) – $t(\alpha, f)$ în raport cu α (probabilitatea cu care se dorește garantarea rezultatului) și $f=n-2$ (numărul gradelor de libertate).
- 3) Dacă $t_{calc} \geq t_{tab}$ – coeficientul de corelație liniară simplă este semnificativ;
 Dacă $t_{calc} < t_{tab}$ – influența caracteristicii factoriale x asupra lui y nu este reală sau nu este garantată cu probabilitatea dorită.

În cazul corelației multiple de tip liniar se recomandă calculul coeficienților de corelație liniară simplă luând pe rând câte un factor x_i pentru a măsura intensitatea legăturii sale cu variabila dependentă y . Se va obține astfel câte un coeficient de corelație liniară simplă pentru fiecare caracteristică factorială înregistrată: $r_{y/x_1}, r_{y/x_2}, \dots, r_{y/x_n}$. Aceștia vor putea fi utilizați pentru calculul coeficientului de corelație multiplă $R_{y/x_1, x_2, \dots, x_n}$.

Raportul de corelație

Pentru măsurarea intensității legăturii dintre variabila dependentă y și variabila independentă x în cazul funcțiilor de regresie neliniare se utilizează raportul de corelație simplă. În scopul determinării modelului de calcul al raportului de corelație se pornește de la ideea că variația totală a caracteristicii rezultative y are două componente:

- 1) o componentă esențială (determinantă) explicată prin influența caracteristicii factoriale x (variabila cauză esențială);
- 2) o componentă neesențială (reziduală) explicată prin influența factorilor aleatori (neînregistrați).

Adâncind analiza, putem pune în evidență trei feluri de abatere:

- abaterea valorilor empirice ale lui y_i de la medie: $(y_i - \bar{y})$ sintetizată la nivelul seriei în *dispersia totală* $\sigma_y^2 = \frac{\sum (y_i - \bar{y})^2}{n}$, care reflectă influența tuturor factorilor esențiali și neesențiali (întâmplători);
- abaterea valorilor calculate pe baza funcției de regresie (valori teoretice) $(y_i - y_{x_i})$ exprimată pe total prin *dispersia reziduală* $\sigma_{y/r}^2 = \frac{\sum (y_i - y_{x_i})^2}{n}$ și care măsoară influența factorilor aleatori;
- abaterea valorilor teoretice de la medie $(y_{x_i} - \bar{y})$ sintetizată la nivelul seriei de *dispersia sistematică* $\sigma_{y/x}^2 = \frac{\sum (y_{x_i} - \bar{y})^2}{n}$ și care arată influența variabilei independente x , considerată ca factor determinant al variației y .

În acest fel putem scrie:

$y_i - \bar{y} = (y_i - y_{x_i}) + (y_{x_i} - \bar{y})$ la nivelul fiecărui y_i și ținând cont de câteva transformări elementare și de relațiile de calcul ale dispersiilor obținem:

$$\sigma_y^2 = \sigma_{y/r}^2 + \sigma_{y/x}^2.$$

Împărțim prin σ_y^2 și găsim:

$$I = \frac{\sigma_{y/r}^2}{\sigma_y^2} + \frac{\sigma_{y/x}^2}{\sigma_y^2} \Rightarrow \text{notând } N_{y/x}^2 = \frac{\sigma_{y/r}^2}{\sigma_y^2} \text{ și } R_{y/x}^2 = \frac{\sigma_{y/x}^2}{\sigma_y^2}$$

$$I = N_{y/x}^2 + R_{y/x}^2 \quad (**)$$

$$N_{y/x}^2 = \frac{\sigma_{y/r}^2}{\sigma_y^2} - \text{coeficient de nedeterminație};$$

$$R_{y/x}^2 = \frac{\sigma_{y/x}^2}{\sigma_y^2} - \text{coeficient de determinație}.$$

Coefficientul de determinație arată ponderea (proporția) din variația totală a fenomenului reprezentat de y explicată prin variația factorului înregistrat x .

Coefficientul de nedeterminație arată ponderea (proporția) factorilor aleatori în variația totală a fenomenului reprezentat de y .

Din relația (**) deducem:

$$R_{y/x}^2 = I - N_{y/x}^2 = I - \frac{\sigma_{y/r}^2}{\sigma_y^2}.$$

De aici deducem relația de calcul a *raportului de corelație* extrăgând rădăcina pătrată din coeficient de determinație:

$$R_{y/x} = \sqrt{I - \frac{\sigma_{y/r}^2}{\sigma_y^2}} = \sqrt{I - \frac{\frac{\sum (y_i - y_{x_i})^2}{n}}{\frac{\sum (y_i - \bar{y})^2}{n}}} = \sqrt{I - \frac{\sum (y_i - y_{x_i})^2}{\sum (y_i - \bar{y})^2}}$$

Din relația de calcul observăm că valorile raportului de corelație sunt totdeauna pozitive și cuprinse între 0 și 1.

În cazul legăturilor de tip liniar raportul de corelație trebuie să fie egal cu coeficientul de corelație.

Relația de mai sus se utilizează pentru un volum mic de date negrupate. În situația unei grupări simple în care x și respectiv y au frecvențe egale vom avea:

$$R_{y/x_i} = \sqrt{I - \frac{\frac{\sum (y_i - y_{x_i})^2 \cdot f_i}{n}}{\frac{\sum (y_i - \bar{y})^2 \cdot f_i}{n}}},$$

iar în situația unei distribuții bidimensionale

$$R_{y/x} = \sqrt{I - \frac{\sum_i \sum_j (y_j - y_{x_i})^2 \cdot f_{ij}}{\sum_j (y_j - \bar{y})^2 \cdot f_j}}.$$

Atunci când este analizată modificarea variabilei dependente y în raport de variația mai multor factori de influență intensitatea legăturii se măsoară cu ajutorul *raportului de corelație multiplă* determinat cu relația:

$$R_{y/x_1, x_2, \dots, x_n} = \sqrt{1 - \frac{\sum (y_i - y_{x_1, x_2, \dots, x_n})^2}{\sum (y_i - \bar{y})^2}}$$

Coeficientul de corelație multiplă

Coeficientul de corelație multiplă, simbolizat cu R_{y, x_i} sau simplu R , măsoară intensitatea legăturii dintre variabila dependentă y și două sau mai multe variabile independente x_i . În acest caz, trebuie să se țină seama de faptul că fiecărei variabile independente îi revine numai o fracțiune din influența totală asupra variabilei dependente. Relația de calcul este următoarea:

$$R_{y, x_i} = R = \sqrt{1 - \frac{\sum (y_i - y_{x_i})^2}{\sum (y_i - \bar{y})^2}}, \text{ unde } \bar{y} \text{ reprezintă media aritmetică simplă a valorilor}$$

empirice ale variabilei dependente pe perioada de analiză statistică, sau, dacă se ține seama de funcția de corelație multiplă, relația devine:

$$R = \sqrt{\frac{a_0 \sum y + a_1 \sum y_{x_1} + a_2 \sum y_{x_2} + \dots + a_n \sum y_{x_n} - \frac{(\sum y)^2}{n}}{\sum y^2 - \frac{(\sum y)^2}{n}}}$$

Acest indicator are întotdeauna valoare pozitivă și este mai mare decât oricare coeficient de corelație simplă dintre variabila dependentă și cele independente, luat în valoare absolută. Pătratul coeficientului de corelație multiplă este cunoscut în literatura de specialitate sub denumirea de coeficient de determinație multiplă (R^2). Acesta exprimă ponderea cu care variabilele independente influențează concomitent asupra variabilei dependente. Ponderea influenței celorlalți factori, neincluși în model, se calculează ca diferență între unitate și R^2 , adică $1 - R^2$.

Coeficientul de corelație multiplă în cazul a doi factori de influență în ipoteza unei legături liniare de forma:

$$y = a + b_1 x_1 + b_2 x_2$$

se poate determina cu relația:

$$R_{y/x_1, x_2} = \sqrt{\frac{r_{yx_1}^2 + r_{yx_2}^2 - 2r_{yx_1} r_{yx_2} r_{x_1 x_2}}{1 - r_{x_1 x_2}^2}}$$

unde: r_{yx_1} - coeficientul de corelație simplă dintre y și x_1 ;

r_{yx_2} - coeficientul de corelație simplă dintre y și x_2 ;

$r_{x_1 x_2}$ - coeficientul de corelație simplă dintre x_1 și x_2 .

Pornind de la coeficienții de corelație simplă se pot calcula coeficienții de corelație parțială dintre variabila dependentă y și un factor de influență cu excluderea influenței celuilalt factor pe baza relațiilor:

$$R_{yx_1/x_2} = \frac{r_{yx_1} - r_{yx_2} r_{x_1 x_2}}{\sqrt{(1 - r_{yx_2}^2)(1 - r_{x_1 x_2}^2)}}$$

$$R_{yx_2/x_1} = \frac{r_{yx_2} - r_{yx_1} r_{x_1 x_2}}{\sqrt{(1 - r_{yx_1}^2)(1 - r_{x_1 x_2}^2)}}$$

unde: R_{yx_1/x_2} – coeficientul de corelație parțială dintre variabila dependentă y și variabila independentă x_1 cu excluderea influenței lui x_2 ;

R_{yx_2/x_1} – coeficientul de corelație parțială dintre variabila dependentă y și variabila independentă x_2 cu excluderea influenței lui x_1 .

Valorile coeficienților de corelație parțială sunt mai mici decât valoarea coeficientului de corelație multiplă dar mai mari decât valoarea coeficientului de corelație simplă.

Testarea semnificației raportului de corelație calculat se poate face cu ajutorul testului „F” de analiză dispersională.

Mărimea lui F_{calc} se determină cu ajutorul relației:

$$F_{calc} = \frac{\frac{\sum (y_x - y)^2}{r-1}}{\frac{\sum (y_i - y_x)^2}{n-r}}$$

unde: n – numărul unităților statistice;

r – numărul grupelor.

Se compară valoarea lui F_{calc} cu F_{tab} identificată funcție de nivelul de semnificație ales α (probabilitatea cu care se dorește obținerea rezultatului) și de numărul gradelor de libertate $f_1=r-1$ și $f_2=n-r$.

Dacă $F_{calc} \geq F_{tab}$ – raportul de corelație este semnificativ.

Dacă $F_{calc} \leq F_{tab}$ – raportul de corelație este respins.

Exemplul 4.5.

Regresia multiplă cercetează și stabilește existența unei legături liniare între o variabilă dependentă (y) și mai multe variabile explicative (x_1, x_2, \dots, x_n), de tipul:

$$y = a + b_1x_1 + b_2x_2 + \dots + b_nx_n + e$$

unde: a = media variabilei y când $x_1=x_2= \dots =x_n=0$

b_1, b_2, \dots, b_n = variația (\pm) a lui y când $x_1, x_2 \dots x_n$ variază cu o unitate

e = eroarea care rezultă independent de $x_1, x_2 \dots x_n$ și care are o distribuție normală.

Modul de lucru se prezintă pe baza exemplului următor care este o extensie a exemplului precedent.

Pentru ultimele 5 luni se cunosc datele statistice privind prețurile unitare (x_1) bugetul de publicitate alocat (x_2) și cantitatea corespunzătoare (y) de produse vândute (vezi tabelul 4.10., col.1,2 și 3). Să se studieze relația dintre cele trei variabile și să se estimeze vânzările pentru $x_1=9,5$ și $x_2=60$ cu o probabilitate de 95 %.

Tabelul 4.10.

Luna	x_1	x_2	y	x_1^2	$x_1 \cdot x_2$	x_2^2	$x_1 \cdot y$	$x_2 \cdot y$	$(y - \bar{y})(x_1 - \bar{x}_1)$	$(y - \bar{y})(x_2 - \bar{x}_2)$
0	1	2	3	4	5	6	7	8	9	10
1	13	4	25	169	52	16	325	100	-21	154
2	12	6	30	144	72	36	360	180	-4,5	81
3	11,5	15	45	132,25	172,5	225	517,5	675	0	0
4	11	20	45	121	220	400	495	900	-3	30
5	10	30	50	100	300	900	500	1550	-16,5	165
Total	57,5	75	195	666,25	816,5	1577	2197,5	3.355	-45	430
Media	11,5	15	39	-	-	-	-	-	-	-

Pentru rezolvarea acestei aplicații vom parcurge următoarele etape:

1) Ajustarea datelor empirice printr-o dreaptă de regresie

$$y_r = a + b_1 x_1 + b_2 x_2$$

Aplicându-se metoda celor mai mici pătrate, se demonstrează că parametrii a , b_1 , b_2, \dots, b_n se determină rezolvând sistemul de ecuații normale corespunzător, conform modelului următor:

$$\begin{cases}
 a \cdot n + b_1 \sum x_1 + b_2 \sum x_2 + b_3 \sum x_3 + \dots + b_m \sum x_m = \sum y \\
 a \sum x_1 + b_1 \sum x_1 x_1 + b_2 \sum x_1 x_2 + b_3 \sum x_1 x_3 + \dots + b_m \sum x_1 x_m = \sum y x_1 \\
 a \sum x_2 + b_1 \sum x_2 x_1 + b_2 \sum x_2 x_2 + b_3 \sum x_2 x_3 + \dots + b_m \sum x_2 x_m = \sum y x_2 \\
 a \sum x_3 + b_1 \sum x_3 x_1 + b_2 \sum x_3 x_2 + b_3 \sum x_3 x_3 + \dots + b_m \sum x_3 x_m = \sum y x_3 \\
 \dots \dots \dots \\
 a \sum x_m + b_1 \sum x_m x_1 + b_2 \sum x_m x_2 + b_3 \sum x_m x_3 + \dots + b_m \sum x_m x_m = \sum y x_m
 \end{cases}$$

unde: n = numărul datelor empirice, iar m = numărul necunoscutelor ($b_1, b_2, b_3, \dots, b_m$)

Notă: x_1, x_2, \dots, x_m pot fi și $\frac{1}{x}, x^2, \log x, \dots$

În cazul considerat, sistemul cu trei necunoscute este următorul:

$$\begin{cases}
 a \cdot n + b_1 \sum x_1 + b_2 \sum x_2 = \sum y \\
 a \sum x_1 + b_1 \sum x_1 x_1 + b_2 \sum x_1 x_2 = \sum y x_1 \\
 a \sum x_2 + b_1 \sum x_2 x_1 + b_2 \sum x_2 x_2 = \sum y x_2
 \end{cases}$$

Cu datele concrete, acesta are forma (vezi și tabelul 4.10., col. 1-8)

$$\begin{cases}
 5a + 57,5b_1 + 75b_2 = 195 \\
 57,5a + 666,25b_1 + 816,5b_2 = 2197,5 \\
 75a + 816,5b_1 + 1577b_2 = 3355
 \end{cases} \Rightarrow \begin{cases}
 a = 75,388... \\
 b_1 = -3,888... \\
 b_2 = 0,555...
 \end{cases}$$

Rezolvând sistemul, ecuația de regresie căutată are forma:

$$y_r = 75,388 - 3,888x_1 + 0,555x_2$$

Verificare: $\bar{y} = a + b_1 \bar{x}_1 + b_2 \bar{x}_2$

$$39 = 75,388 - 3,888 \cdot 11,5 + 0,555 \cdot 15$$

2) Se calculează abaterile dintre y și $y_r = 75,3888 - 3,888x_1 + 0,555x_2$

Pentru ușurința calculului datele sunt sintetizate în tabelul 4.11.

Tabelul 4.11.

Luna	x_1	x_2	y	y_r	$(y - \bar{y})^2$	$(y_r - \bar{y}_r)^2$	$(y - y_r)^2$
1	13	4	25	27,055..	196	142,67	4,23
2	12	6	30	22,055..	81	48,23	4,23
3	11,5	15	45	39	36	0	36
4	11	20	45	43,722..	36	22,30	1,63
5	10	30	50	53,166..	121	200,68	10,02
Total	57,5	75	195	195	470	413,88..	56,11..
Medie	11,5	15	39	39	-	-	-

- variația totală a datelor $y = \sum (y - \bar{y})^2 = 470$ (vezi col.5, tabelul 4.11.)
- variația totală a datelor $y_r = \sum (y_r - \bar{y}_r)^2 = 413,88$ (vezi col.6, tabelul 4.11.)
- variația dintre y și $y_r = \sum (y - y_r)^2 = 56,11$ (vezi col. 7, tabelul 4.11.)

Interpretare: din cele 470 de unități ale variației totale, 413,88 sunt explicate de funcția y_r (adică de x_1 și x_2), iar restul de 56,11 rămân neexplicate.

La rândul său, variația explicată se poate descompune pe cei doi factori, astfel:

- variația explicată de $x_1 = b_1 \sum (y - \bar{y})(x_1 - \bar{x}_1)$
- variația explicată de $x_2 = b_2 \sum (y - \bar{y})(x_2 - \bar{x}_2)$

Cu datele problemei avem:

- $b_1 \sum (y - \bar{y})(x_1 - \bar{x}_1) = (-3,888)(-45) = 175$ (vezi și tabelul 4.10., col.9)
- $b_2 \sum (y - \bar{y})(x_2 - \bar{x}_2) = (0,555)(430) = \frac{238,88}{413,88}$ (vezi și tabelul 4.10., col. 10)

3) Se calculează eroarea standard (e)

În acest scop se utilizează relația (*) în care $\nu = n-3$ deoarece, în acest caz se cunosc trei parametri (a , b_1 și b_2)

$$e = \sqrt{\frac{\sum (y - y_r)^2}{n-3}} = \sqrt{\frac{56,11}{5-3}} = 5,3$$

4) Se compară R^2 calculat cu R^2 teoretic

$$R^2 = \frac{\sum (y_r - \bar{y}_r)^2}{\sum (y - \bar{y})^2} = \frac{413,88}{470} = 0,8806146$$

Pentru $\alpha = 0,05$, R^2 minim trebuie să fie egal cu 0,9025 pentru $\nu = n-3 = 5-3 = 2$ grade de libertate. Deoarece R^2 observat $< R^2$ teoretic deducem că funcția $y_r = 75,388 - 3,888x_1 + 0,555x_2$ nu poate fi socotită acceptabilă decât dacă ne mulțumim cu o probabilitate a previziunii $< 95\%$.

5) Se estimează vânzările.

Ca și în cazul regresiei simple,

$$P = y_r \pm t \cdot e = (a + b_1 x_1 + b_2 x_2) + t \cdot e$$

Pentru cazul dat, unde $x_1 = 9,5$, $x_2 = 60$ și $t = 2,78$,

$$P = (75,388 - 3,888 \cdot 9,5 + 0,555 \cdot 60) \pm 2,78 \cdot 5,3 = 71,777 \pm 14,73, \text{ respectiv } 57,65 < y_r < 86,51$$

Concluzie: în cazul practicării unui preț $x_1 = 9,5$ și a alocării unui buget publicitar $x_2 = 60$, volumul vânzărilor se va situa între 57 și 87 unități, dar cu o probabilitate $< 95\%$.

Exemplul 4.6.

O întreprindere a observat pe un eșantion de 10 magazine că cifra de afaceri variază funcție de suprafața magazinului și numărul de salariați, conform datelor din tabelul 4.12.

Tabelul 4.12.

Suprafața magazinului (m ²)	100	600	600	700	700	500	800	300	200	200
Numărul de salariați	24	28	20	23	26	21	28	28	20	25
Cifra de afaceri a magazinului (mil \$)	11	23	20	21	21	13	30	18	7	18

Întreprinderea dorește să construiască un magazin cu o suprafață de 600 m² și un număr de 24 salariați. Care va fi cifra de afaceri previzibilă a acestui magazin?

Pentru rezolvarea acestei probleme vom folosi un model liniar de regresie multiplă de forma:

$$y = a + b_1x_1 + b_2x_2, \text{ unde:}$$

y = cifra de afaceri a magazinului (în mil. \$)

x_1 = suprafața magazinului (m²)

x_2 = numărul de salariați

Aplicând metoda celor mai mici pătrate

$$\min S = \sum_{i=1}^n (y_i - a - b_1x_{i1} - b_2x_{i2})^2 \text{ și anulând derivatele parțiale în raport cu } a, b_1 \text{ și}$$

b_2 vom obține sistemul de ecuații normale, care va permite calculul parametrilor.

$$\frac{\partial S}{\partial a} = -2 \sum_{i=1}^{10} (y_i - a - b_1x_{i1} - b_2x_{i2}) \cdot 1 = 0$$

$$\frac{\partial S}{\partial b_1} = -2 \sum_{i=1}^{10} (y_i - a - b_1x_{i1} - b_2x_{i2}) \cdot x_{i1} = 0$$

$$\frac{\partial S}{\partial b_2} = -2 \sum_{i=1}^{10} (y_i - a - b_1x_{i1} - b_2x_{i2}) \cdot x_{i2} = 0$$

de unde:

$$\begin{cases} na + b_1 \sum x_1 + b_2 \sum x_2 = \sum y \\ a \sum x_1 + b_1 \sum x_1^2 + b_2 \sum x_1 x_2 = \sum x_1 y \\ a \sum x_2 + b_1 \sum x_1 x_2 + b_2 \sum x_2^2 = \sum x_2 y \end{cases}$$

Pentru rezolvarea sistemului vom organiza datele ca în tabelul 4.13.

Tabelul 4.13.

Număr magazin	y (mil. lei)	x_1	x_2	x_1^2	x_2^2	x_1x_2	x_1y	x_2y	y^2
1	11	100	24	10000	576	2400	1100	264	121
2	23	600	28	360000	784	16800	13800	644	529
3	20	600	20	360000	400	12000	12000	400	400
4	21	700	23	490000	529	16100	14700	483	441
5	21	700	26	490000	676	18200	14700	546	441
6	13	500	21	250000	441	10500	6500	273	169
7	30	800	28	640000	784	22400	24000	840	900
8	18	300	28	90000	784	8400	5400	504	324
9	7	200	20	40000	400	4000	1400	140	49
10	18	200	25	40000	625	5000	3600	450	324
Σ	182	4700	243	2770000	5999	115800	97200	4544	3698

Rezultă sistemul:

$$\begin{cases} 10a + 4700b_1 + 243b_2 = 182 \\ 4700a + 277000b_1 + 115800b_2 = 97200 \\ 243a + 115800b_1 + 5999b_2 = 4544 \end{cases}$$

Prin rezolvarea sistemului obținem:

$$a = -14,21$$

$$b_1 = 0,0179$$

$$b_2 = 0,99$$

de unde ecuația modelului liniar de corelație:

$$y = -14,21 + 0,0179x_1 + 0,99x_2$$

Fiind vorba de o corelație multiplă liniară vom calcula coeficientul de corelație cu relația:

$$R_{y/x_1, x_2} = \sqrt{\frac{r_{yx_1}^2 + r_{yx_2}^2 - 2r_{yx_1} \cdot r_{yx_2} \cdot r_{x_1x_2}}{1 - r_{x_1x_2}^2}}$$

unde: r_{yx_1} = coeficientul de corelație simplă dintre y și x_1

r_{yx_2} = coeficientul de corelație simplă dintre y și x_2

$r_{x_1x_2}$ = coeficientul de corelație simplă dintre x_1 și x_2

$$\begin{aligned} r_{yx_1} &= \frac{n \sum x_1 y - \sum x_1 \cdot \sum y}{\sqrt{[n \sum x_1^2 - (\sum x_1)^2][n \sum y^2 - (\sum y)^2]}} = \\ &= \frac{10 \cdot 97200 - 4700 \cdot 182}{\sqrt{[10 \cdot 2770000 - 22090000][10 \cdot 3698 - 33124]}} = \\ &= \frac{116600}{\sqrt{5610000 \cdot 3856}} = \frac{116600}{147079} = 0,793 \end{aligned}$$

$$r_{yx_2} = \frac{n \sum x_2 y - \sum x_2 \cdot \sum y}{\sqrt{[n \sum x_2^2 - (\sum x_2)^2][n \sum y^2 - (\sum y)^2]}}$$

$$\begin{aligned}
 &= \frac{10 \cdot 4544 - 243 \cdot 182}{\sqrt{[10 \cdot 5999 - 59049][10 \cdot 3698 - 33124]}} = \\
 &= \frac{1214}{\sqrt{941 \cdot 3856}} = \frac{1214}{1905} = 0,637 \\
 r_{x_1 x_2} &= \frac{n \sum x_1 x_2 - \sum x_1 \cdot \sum x_2}{\sqrt{[n \sum x_1^2 - (\sum x_1)^2][n \sum x_2^2 - (\sum x_2)^2]}} = \\
 &= \frac{10 \cdot 115800 - 4700 \cdot 243}{\sqrt{5610000 \cdot 941}} = \frac{1158000 - 1142100}{72657} = 0,219
 \end{aligned}$$

Rezultă

$$\begin{aligned}
 R_{y/x_1, x_2} &= \sqrt{\frac{0,628849 + 0,405769 - 2 \cdot 0,793 \cdot 0,637 - 0,219}{1 - 0,047961}} = \\
 &= \sqrt{\frac{1,034618 - 0,221252}{0,952039}} = 0,9243
 \end{aligned}$$

Rezultă că variația volumului cifrei de afaceri se datorează în proporție de circa 92,43% suprafeței magazinului și numărului de salariați.

Calculând coeficienții de corelație parțială cu excluderea influenței celuilalt factor rezultă:

$$\begin{aligned}
 R_{y x_1 / x_2} &= \frac{r_{y x_1} - r_{y x_2} \cdot r_{x_1 x_2}}{\sqrt{(1 - r_{y x_2}^2)(1 - r_{x_1 x_2}^2)}} = \frac{0,793 - 0,637 \cdot 0,219}{\sqrt{(1 - 0,637^2)(1 - 0,219^2)}} = 0,869 \\
 R_{y x_2 / x_1} &= \frac{r_{y x_2} - r_{y x_1} \cdot r_{x_1 x_2}}{\sqrt{(1 - r_{y x_1}^2)(1 - r_{x_1 x_2}^2)}} = \frac{0,637 - 0,793 \cdot 0,219}{\sqrt{(1 - 0,793^2)(1 - 0,219^2)}} = 0,779
 \end{aligned}$$

În consecință construirea unui magazin cu o suprafață de 600 m² și 24 salariați va conduce la obținerea de către acesta a unei cifre de afaceri

$$Y = -14,21 + 0,0179 \cdot 600 + 0,99 \cdot 24 = 20,29 \text{ mil. \$}$$

4.2.3. Metode neparametrice de măsurare a legăturilor

Sunt metode care se utilizează atunci când:

- nu se cunoaște forma legăturii;
- caracteristicile sunt calitative și nu se pot exprima numeric dar este posibilă ierarhizarea lor;
- în cazul distribuțiilor asimetrice;
- când dispunem de un număr mic de observații.

Cele mai utilizate metode neparametrice sunt:

- metoda tabelului de asociere și a coeficientului de asociere;
- metoda corelației rangurilor.

A. Metoda tabelului de asociere se utilizează în cazul caracteristicilor alternative care admit numai două forme de manifestare sau valori. Tabelul de asociere este de fapt un caz particular al tabelului cu dublă intrare. Pe linii se înscrie variația caracteristicii factoriale x (variabila independentă), iar pe coloane variația caracteristicii rezultative y

(variabila efect). La intersecția liniilor cu coloanele se trec frecvențele cu care unitățile colectivității se înscriu în cele patru grupe formate prin intersecția variației caracteristicilor x și y (tabelul 4.14.).

Tabelul 4.14.

	<i>Y</i>		
<i>X</i> \	y_1	y_2	<i>Total</i>
x_1	a	b	$a+b$
x_2	c	d	$c+d$
<i>Total</i>	$a+c$	$b+d$	$a+b+c+d$

Pentru măsurarea intensității legăturii dintre x și y se utilizează *coeficientul de asociere* propus de Yulle:

$$Q = \frac{ad - bc}{ad + bc}.$$

Valorile acestui indicator sunt cuprinse între -1 și $+1$.

Dacă $Q > 0$ – asociere directă.

Dacă $Q < 0$ – asociere inversă.

Dacă $Q = 0$ – nu există asociere.

Dacă Q tinde la ± 1 avem asociere foarte puternică.

Dacă o frecvență din cele patru (a, b, c, d) este nulă, atunci asocierea este completă și $Q = \pm 1$.

B. Metoda corelației rangurilor

După cum s-a văzut până aici, în calculul coeficientului de corelație liniară sau a raportului de corelație, se folosesc valorile variabilelor care alcătuiesc cuplul corelativ fapt care evidențiază oarecum că nivelul acestor indicatori depinde de nivelul variabilelor pentru care se realizează studiul.

O posibilă soluție pentru eliminarea acestei dependențe este utilizarea unor metode de studiu a corelației care nu folosesc direct în calculul respectivilor coeficienți valorile variabilelor din cuplul corelativ.

Corelația rangurilor este o metodă de studiu a corelației care presupune folosirea în calculul indicatorilor pentru măsurarea intensității corelației, a rangurilor corespunzătoare variabilelor din cuplul corelativ.

Rang – poziția pe care o ocupă valorile variabilelor X și Y din cuplul corelativ în cadrul șirului din care fac parte, ordonat crescător sau descrescător.

Pentru ranguri se folosesc următoarele notații:

- u_i - rangurile valorilor x_i , din cadrul șirului ordonat x_1, x_2, \dots, x_n ;
- w_i - rangurile valorilor y_i , din cadrul șirului ordonat y_1, y_2, \dots, y_n .

Coeficientului lui Spearman

Folosind în relația de calcul a coeficientului de corelație liniară (r) rangurile corespunzătoare, în locul valorilor variabilelor X respectiv Y , utilizând o serie de proprietăți a seriilor de ranguri și după câteva calcule elementare se ajunge la:

$$\theta = 1 - \frac{6 \sum d^2}{n(n^2 - 1)}$$

unde: d - diferența dintre rangurile u_i și w_i ;
 n - numărul termenilor seriei.

Relația dă rezultate corecte atâta timp cât sunt îndeplinite premisele folosite la obținerea ei, și anume:

1. $\sum u_i = \sum w_i$;
2. $\bar{u}_i = \frac{\sum u_i}{n} = \frac{\sum w_i}{n} = \frac{n+1}{2} = \bar{w}_i$;

3. rangurile u_i și w_i , sunt unice, nu se repetă în cadrul șirului din care fac parte. Dacă această ultimă condiție nu este îndeplinită atunci se poate proceda astfel: valoarea lui x_i sau y_i care se repetă se va trece o singură dată în șirul din care face parte, iar ca valoare corespondentă se va trece media valorilor celeilalte variabile, corespunzătoare valorii care se repeta.

Coeficientul lui Spearman ia valori în intervalul $[-1;1]$.

Semnificația acestui coeficient este similară cu cea a coeficientului de corelație liniară (r).

Coeficientul lui Kendall

Presupune îndeplinirea aceluiași condiții ca și coeficientul lui Spearman și se poate calcula cu ajutorul relației:

$$\tau = \frac{2S}{n(n-1)} = \frac{2(P-Q)}{n(n-1)},$$

unde: P - suma rangurilor w_i mai mari decât rangul curent.

Q - suma rangurilor w_i mai mici decât rangul curent.

Pentru determinarea lui P și lui Q se procedează astfel:

- se ordonează crescător perechile de valori (x_i, y_i) ;
- se elimină eventualele repetiții de valori pentru variabila X și variabila Y și dacă este cazul, se reordonează crescător perechile de valori (x_i, y_i) ;
- se determină rangurile u_i atașate valorilor variabilei X ;
- se determină rangurile w_i atașate valorilor variabilei Y ;
- pornind de la prima valoare spre sfârșitul șirului rangurilor w_i , se determină succesiv pentru fiecare rang:
 - câte ranguri w_i sunt mai mari de cât rangul curent;
 - câte ranguri w_i sunt mai mici de cât rangul curent;

Notă: Numărul rangurilor mai mari sau mai mici decât rangul curent se face pornind numărătoarea de la următoarea poziție față de rangul curent, către sfârșitul șirului.

- se determină P și Q .

Și coeficientul lui Kendall $\in [-1;1]$ și are o interpretare similară cu cea de la coeficientul lui Spearman.

Dacă pentru același set de date se calculează ambii coeficienți se va observa că:

$$\theta > \tau.$$

Exemplul 4.3.

Presupunând că avem un set de perechi de valori (x_i, y_i) corespunzătoare a două variabile X și Y , pentru care se dorește evidențierea existenței unei legături cauzale, procedăm astfel (tabelul 4.15.):

Tabelul 4.15.

Inițial		După eliminare repetiții		După ordonare	
x_i	y_i	x_i	y_i	x_i	y_i
100	28	100	27	51	9
98	24	98	24	54	8
51	9	51	9	69,5	12
100	26	86	14	86	14
86	14	69,5	12	98	24
70	12	54	8	100	27
69	12	108	31	102	35
54	8	102	35	108	31
108	31				
102	35				
n = 10		n = 8		n = 8	

Pentru calculul celor doi coeficienți determinăm mai întâi rangurile u_i și w_i (tabelul 4.16.).

Tabelul 4.16.

Valori ordonate		Ranguri		Indicatori			
x_i	y_i	u_i	w_i	d	d^2	P	Q
51	9	1	2	-1	1	6	1
54	8	2	1	1	1	6	0
69,5	12	3	3	0	0	5	0
86	14	4	4	0	0	4	0
98	24	5	5	0	0	3	0
100	27	6	6	0	0	2	0
102	35	7	8	-1	1	0	1
108	31	8	7	1	1	0	0
n = 10		36	36	0	4	26	2

$$\theta = 1 - \frac{6 \cdot 4}{8(8^2 - 1)} = 0,994$$

$$\tau = \frac{2 \cdot (26 - 2)}{8(8 - 1)} = 0,857$$

Între cele două variabile există o corelație directă foarte intensă.

Tot în studiul corelației rangurilor poate fi utilizată și corelograma rangurilor.

Modul de construcție al acesteia este similar cu cel al corelogramei prezentate la metodele elementare de studiu a corelației cu deosebirea că la corelograma rangurilor se reprezintă perechile de ranguri (u_i, w_i).

Diagonalele rețelei reprezintă corelație maximă directă respectiv inversă. Interpretarea corelogramei rangurilor este similară cu cea a corelogramei prezentate la metode elementare în studiul corelației.

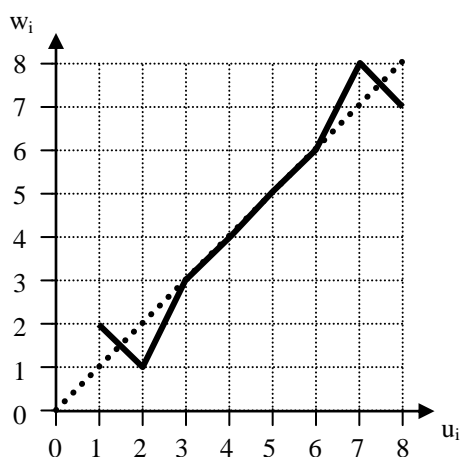


Figura 4.13. Corelograma rangurilor

Exemplul 4.4.

Exporturile (FOB) și importurile (CIF) de mărfuri ale României în / din țările Uniunii Europene în anul 2006 au fost (tabelul 4.17.).

Tabelul 4.17.

Nr. crt.	Țara	Export (x)	Import (y)	Rang după		Diferența de rang
				Export (x)	Import (y)	
1.	Austria	685,2	1535,3	20	22	-2
2.	Belgia	426,3	626,5	15	15	0
3.	Cehia	294,5	1007,4	14	19	-5
4.	Cipru	50,9	15,2	8	4	4
5.	Danemarca	63,6	164,9	9	8	1
6.	Estonia	6,3	6,7	3	3	0
7.	Franța	1938,3	2664,7	23	23	0
8.	Finlanda	31,2	185,1	5	9	-4
9.	Germania	4060,2	6176,8	24	25	-1
10.	Grecia	507,5	481,4	17	14	3
11.	Irlanda	42,1	240,4	7	11	-4
12.	Italia	4637,2	5954,9	25	24	1
13.	Letonia	5,8	3,8	2	2	0
14.	Lituania	15,5	23,8	4	5	-1
15.	Luxemburg	5,4	29,9	1	6	-5
16.	Malta	76,6	3,7	10	1	9
17.	Olanda	637,4	739,6	19	16	3
18.	Polonia	468,3	1128,4	16	18	-2
19.	Portugalia	33,5	128,7	6	7	-1

Nr. crt.	Țara	Export (x)	Import (y)	Rang după		Diferența de rang
				Export (x)	Import (y)	
20.	Spania	605,8	953,1	18	17	-1
21.	Suedia	142,3	419,4	12	12	0
22.	Regat Unit M. Britanii	1216,8	1011,7	21	20	1
23.	Slovacia	180,4	461,7	13	13	0
24.	Slovenia	92,6	192,0	11	10	1
25.	Ungaria	1276,0	1331,6	22	21	1
	Total	17499,7	25486,7	-	-	-

Sursa: *Buletin statistic lunar*, nr.12, 2006, Institut Național de Statistică.

Să se analizeze legătura dintre cele două variabile.

Rezolvare

$$\theta = 1 - \frac{6 \cdot 223}{25(25^2 - 1)} = 0,914$$

Pentru calculul coeficientului lui Kendall (tabelul 4.18.) se ordonează crescător țările UE după variabila x (export) înscriind în coloana alăturată rangurile corespunzătoare după y (import). Se determină apoi pentru fiecare țară, pe baza rangurilor y (import):

- P_i – numărul de țări (de la rândul i până la sfârșitul seriei) având la importuri ranguri superioare rangului țării i ;
- Q_i – numărul de țări (de la rândul i până la sfârșitul seriei) având la importuri ranguri inferioare rangului țării i .

Tabelul 4.18.

Nr. crt.	Țara	Export (x)	Import (y)	Rang după		P_i	Q_i	$P_i - Q_i$
				x	y			
1.	Austria	5,4	29,9	1	6	19	5	14
2.	Belgia	5,8	3,8	2	2	22	1	21
3.	Cehia	6,3	6,7	3	3	21	1	20
4.	Cipru	15,5	23,8	4	5	19	2	17
5.	Danemarca	31,2	185,1	5	9	16	4	12
6.	Estonia	33,5	128,7	6	7	17	2	15
7.	Franța	42,1	240,4	7	11	14	4	10
8.	Finlanda	50,9	15,2	8	4	16	1	15
9.	Germania	63,6	164,9	9	8	15	1	14
10.	Grecia	76,6	3,7	10	1	15	0	15
11.	Irlanda	92,6	192	11	10	14	0	14
12.	Italia	142,3	419,4	12	12	13	0	13
13.	Letonia	180,4	461,7	13	13	12	0	12
14.	Lituania	294,5	1007,4	14	19	6	4	2
15.	Luxemburg	426,3	626,5	15	15	9	1	8
16.	Malta	468,3	1128,4	16	18	6	3	3
17.	Olanda	507,5	481,4	17	14	8	0	8
18.	Polonia	605,8	953,1	18	17	6	1	5
19.	Portugalia	637,4	739,6	19	16	6	0	6
20.	Spania	685,2	1535,3	20	22	3	2	1

Nr. crt.	Țara	Export (x)	Import (y)	Rang după		P _i	Q _i	P _i - Q _i
				x	y			
21.	Suedia	1216,8	1011,7	21	20	4	0	4
22.	Regat Unit M. Britanii	1276	1331,6	22	21	3	0	3
23.	Slovacia	1938,3	2664,7	23	23	2	0	2
24.	Slovenia	4060,2	6176,8	24	25	0	1	-1
25.	Ungaria	4637,2	5954,9	25	24	0	0	0
	Total	17499,7	25486,7	-	-	266	33	233

$$\tau = \frac{2 \cdot (266 - 33)}{25(25 - 1)} = 0,777$$

Valorile celor doi coeficienți indică o legătură directă puternică între exportul și importul României în / din țările Uniunii Europene.

STUDIUL STATISTIC AL LEGĂTURILOR DINTRE FENOMENELE ȘI PROCESELE ECONOMICO-SOCIALE	110
4.1. Tipuri de legături dintre fenomenele și procesele economice	110
4.2. Metode statistice utilizate în studiul legăturii dintre fenomenele și procesele economice ...	112
4.2.1. Metode elementare utilizate în studiul legăturii dintre fenomenele și procesele economice	112
4.2.2. Metode analitice (parametrice) de măsurare a legăturilor dintre fenomenele și procesele economice	119
4.2.2.1. Metoda regresiei.....	119
4.2.2.2. Metoda corelației	130
4.2.3. Metode neparametrice de măsurare a legăturilor	141