

Legăturile dintre fenomenele colective

Statistica studiază *legăturile cauzale* dintre fenomenele colective. Legile naturale îmbracă întotdeauna forma enunțării unei legături fie între manifestările a două sau mai multor fenomene, fie între prezența și absența anumitor caractere. Spunem, de exemplu, că electricitățile de sens contrar se atrag, iar cele de același sens se resping, că temperatura descrește odată cu creșterea altitudinii, că infracționalitatea crește dacă scade nivelul de trai. Acestea sunt legi cunoscute. Prima este o lege rigidă, se verifică, prin urmare, aplicată oricărui caz. Celelalte sunt legi *stochastice*..

Forma sub care am enunțat legea stochastică de mai sus constă numai în existența legăturii dintre fenomene și prin aceasta suntem privați de cunoașterea unui element destul de interesant în această direcție: *gradul de asociație între cele două fenomene sau caractere*.

Prezentul capitol este consacrat mijloacelor de care dispune statistica matematică pentru a măsura acestui nou element.

Este clar că o definiție perfect circumscrisă a gradului de legătură între două variabile nu poate fi atinsă. În procesul de cercetare a legăturilor cauzale existente între fenomenele naturale avem de-a face cu contopirea acțiunilor unei multitudini de factori (cauze), dintre care unii esențiali, alții neesențiali, unii pot fi determinați, alții nu. Tocmai de aceea, în cercetarea legăturii reciproce dintre două fenomene apar dificultăți mari deoarece pot să existe cauze necunoscute. În astfel de situații este util să determinăm gradul de corelare și apoi să analizăm separat unele din aceste cauze. Cu alte cuvinte, trebuie să definim relații posibile între diferiți factori, evenimente, atribute sau caracteristici care ar putea avea o influență, cel puțin parțială asupra datelor experimentale. În acest mod este posibilă stabilirea unui tablou al condițiilor în care se desfășoară un anumit fenomen natural, fapt care duce la modelarea matematică a sa. Din această cauză, diferite puncte de vedere au putut sta alături pe această chestiune, fără ca să se decidă care este cel mai adecvat. Și natural, fiecare din ele a condus la un mijloc diferit de măsură a acestei legături.

În multe procese naturale, pe lângă complexa întrepătrundere cu alte fenomene (proces), acestea mai sunt supuse unor evoluții care la prima vedere pot fi considerate probabiliste (aleatoare). De aceea, pentru cunoașterea modului de evoluție probabilă în viitor - prognoza unui anumit fenomen - trebuie să ne bazăm pe *cunoașterea evoluției trecute*, precum și pe situația prezentă. Experiența unui mare număr de observații (probe) au dus la concluzia că între diferitele mărimi variabile pot exista următoarele tipuri de relații:

a) ***Relația de dependență***, Y depinde de X sau invers. O modificare a unei variabile duce la o modificare a celei de-a doua. În cazul unei astfel de relații s-ar putea aminti existența relației cauzale în care o variabilă este cauza, iar cealaltă este efectul, cauza fiind un fenomen sau un complex de fenomene care provoacă, generează sau determină un alt fenomen - efectul. Operația logică prin care efectul este dedus din cauză se numește *inferență cauzală*.

b) ***Relația de interdependență***, Y depinde de X și X depinde de Y. În acest caz modificarea unei variabile provoacă modificarea celei de-a doua variabile, iar modificarea acesteia din urmă are influență asupra primei variabile.

c) ***Relația de tranziție***, X se transformă parțial sau total în Y și invers.

d) **Corelația statistică** sau **covarianța**. Pentru X există întotdeauna Y și invers. Este o relație reciprocă dintre două variabile, dintre care una în mod logic apelează la alta și pe baza analizei datelor experimentale se poate pune în evidență o asociere între ele.

e) **Relația stochastică**; dacă se realizează X atunci cu o anumită probabilitate se realizează și Y, sau invers.

Datorită caracterului complex al fenomenelor și datorită multitudinii de factori esențiali și întâmplători care intervin, aceste legături se manifestă *sub formă de tendință*; ele pot fi identificate în condițiile acțiunii legii numerelor mari, în colectivitățile de volum ridicat.

O problemă importantă în cadrul analizei seriilor interdependente o reprezintă **identificarea legăturilor cu adevărat semnificative și stabile**.

În cazul în care s-a stabilit că între două variabile există un raport de dependență, se ridică o altă problemă, mult mai dificilă: găsirea unei **măsuri a acestei legături**, prin intermediul unui indicator sintetic de corelație.

Există însă și situații înșelătoare, în care variațiile (creșterea sau scăderea) celor două fenomene aparent interdependente sunt similare, **dar nu există nici o legătură logică între ele**. De exemplu, în ultimii ani în România cresc simultan atât rata sărăciei, cât și înzestrarea populației cu telefoane mobile. Alături explicația *tendinței de asociere* a celor două variabile este dată de **existența unei cauze comune**. De exemplu, creșterea simultană și semnificativă a vânzărilor la pulovere de lână și la medicamente antigripale are aceeași cauză: venirea iernii. Analiza calitativă a datelor statistice rezolvă probleme de acest tip.

SERIILE INTERDEPENDENTE

TIPURI DE LEGĂTURI

Analiza seriilor interdependente urmărește

1. **verificarea existenței** și
2. **măsurarea intensității legăturilor dintre fenomene**.

Cu cât fenomenele sunt mai complexe, cu atât numărul factorilor care le influențează este mai mare, ceea ce face legăturile cauzale dificil de evidențiat. Analiza dependențelor se complică atunci când factorii de influență intercondiționează determinând apariția de cauzalități în lanț.

Pentru a putea determina intensitatea relațiilor de corelație, este necesar să se analizeze, în primul rând, conținutul și forma acestor relații.

Vom nota cu:

- y fenomenul a cărui variație este influențată (*variabilă dependentă sau rezultativă*) și cu
- x factorul de influență (*variabilă independentă sau factorială*).

După **natura legăturilor de interdependență** deosebim două categorii de relații:

a) **funcționale (determinate)** și

b) **stochastice (statistice)**.

a) **Funcționale**: $y = f(x)$.

Fenomenul cauză x determină în mod univoc fenomenul efect y, astfel încât *unei valori a variabilei x îi corespunde o valoare unică a variabilei y*.

b) **Stochastice**: $y = f(x_1, x_2, \dots, x_n)$.

Fenomenul efect este influențat de o multitudine de factori, esențiali și întâmplători. Unei valori a fenomenului cauză x_i îi pot corespunde mai multe valori diferite ale fenomenului efect y , în funcție de acțiunea combinată a celorlalți factori de influență. De exemplu, infracționalitatea unei persoane depinde de: profilul psihologic, sex, gradul de educație, gradul de sărăcie, situația familială, starea de sănătate etc.

Fiecare caracteristică x_i determină numai o parte a variației fenomenului y , restul variației fiind explicat prin alte caracteristici, care din punct de vedere al legăturii x_i - y sunt întâmplătoare.

Variația fenomenului y poate fi analizată în funcție de unul sau mai mulți factori de influență (x_1, x_2, \dots, x_n), dar întotdeauna va rămâne o variație neexplicată, determinată de factorii neînregistrați.

În cazul în care se identifică și se analizează factorii de influență esențiali, componenta aleatoare, care sintetizează acțiunea factorilor întâmplători, va avea o pondere redusă și nu va influența semnificativ veridicitatea rezultatelor.

Legăturile stochastice sunt specifice fenomenelor din societate și economie.

Extrema diversitate a legăturilor stochastice impune sistematizarea lor după mai multe criterii.

După **numărul caracteristicilor** analizate, legăturile stochastice pot fi simple sau multiple.

a) Legături simple - se alege o singură caracteristică determinantă pentru variația fenomenului y , toate celelalte caracteristici care îl influențează, fie că sunt esențiale sau întâmplătoare, sunt considerate cu acțiune constantă.

De exemplu, analiza legăturii dintre recolta totală și suprafața cultivată.

b) Legături multiple - se analizează variația fenomenului y în funcție de mai multe caracteristici esențiale x_1, x_2, \dots, x_n . Rămâne și în acest caz o componentă aleatoare, care sintetizează acțiunea, presupusă constantă, a celorlalți factori de influență.

De exemplu, analiza variației salariului într-o colectivitate în funcție de productivitate, vechime și calificare.

După **natura caracteristicii** pot exista legături stochastice de asociere sau de corelație.

a) Asocierea statistică se referă la raporturile de interdependență dintre caracteristicile calitative, sau dintre o caracteristică numerică și una calitativă. De exemplu, legătura dintre locul de muncă și studii, calificare și productivitate, între zona geografică și clasa de fertilitate a terenurilor agricole etc.

Analiza statistică a raporturilor de asociere este posibilă doar dacă se găsește o modalitate de exprimare numerică a variantelor. De exemplu, clasele de calitate ale produselor pot fi codificate și ierarhizate: 0 – produs inferior, 1 - produs mediu, 2 - produs superior.

b) Corelația statistică exprimă raporturile de cauzalitate dintre două sau mai multe caracteristici exprimate cantitativ. De exemplu, analiza cifrei de afaceri în funcție de numărul salariaților și valoarea capitalului fix, analiza gradului de poluare în funcție de producția de substanțe chimice.

După **sensul relației de cauzalitate**, legăturile stochastice pot fi directe sau inverse.

a) Legături directe există atunci când modificarea într-un anumit sens (creștere sau scădere) a fenomenului cauză x determină modificarea în același sens a fenomenului efect y . De exemplu, legătura dintre numărul salariaților și volumul producției, dintre mărimea creditului și masa dobânzii, dintre costul unitar și costul total etc.

b) Legăturile inverse există atunci când modificarea într-un anumit sens a lui x determină o modificare în sens contrar a lui y . De exemplu, legătura dintre profitul unitar și costul unitar de producție, dintre impozitul pe profit și profitul net, dintre mărimea dividendelor și profitul reinvestit etc.

După **forma matematică** a legăturilor, acestea pot fi liniare sau neliniare.

a) Liniare: legătura se realizează după ecuația dreptei.

b) Neliniare: exponențiale, hiperbolice, parabolice, logaritmice.

Forma legăturii este, de regulă, vizibilă pe grafic. Atunci când legătura grafică nu este clară, se poate continua analiza pe variantele sugerate de grafic folosind metode analitice și folosind anumite criterii pentru a alege varianta cea mai bună.

După **momentul** producerii lor deosebim legături sincrone și asincrone.

- a) **Sincrone** - modificarea lui x determină modificarea imediată a lui y. De exemplu, creșterea veniturilor populației determină mărirea imediată a cererii de consum, creșterea producției se obține concomitent cu creșterea cheltuielilor etc.
- b) **Asincrone** - fenomenul x determină variația fenomenului efect y după o perioadă de timp. De exemplu, legătura dintre investiții și creșterea producției sau legătura dintre rata dobânzii și volumul masei monetare.

METODE ELEMENTARE DE VERIFICARE A EXISTENȚEI LEGĂTURILOR

Studierea legăturilor dintre fenomenele economice presupune parcurgerea mai multor etape:

1. **depistarea factorilor** care influențează variația fenomenului analizat și **ierarhizarea** acestora;
2. **alegerea factorului sau factorilor esențiali** a căror influență asupra fenomenului dependent urmează să fie analizată;
3. culegerea și sistematizarea datelor referitoare la variabilele studiate; verificarea gradului de cuprindere a datelor înregistrate (dacă datele provin dintr-un sondaj interpretarea rezultatelor se va face în sens probabilistic);
4. **verificarea existenței și formei legăturii** dintre caracteristicile corelate în vederea alegerii corecte a procedurilor statistico -matematice de măsurare a dependenței statistice;
Verificarea existenței legăturilor se poate face cu ajutorul unor metode simple:
 - metoda seriilor paralele interdependente;
 - metoda grupărilor;
 - metoda tabelului de corelație;
 - metoda grafică;
 - analiza dispersională.Metodele elementare evidențiază **direcția legăturii**, iar unele dintre ele pot indica și **forma** acesteia. Aplicarea acestor metode trebuie completată cu o analiză calitativă a fenomenelor, bazată pe conținutul lor, pe legătura logică dintre ele
5. **măsurarea intensității legăturii** cu ajutorul indicatorilor de corelație selectați în funcție de forma de legătură și de natura informațiilor de care dispunem;
6. **testarea semnificației indicatorilor de corelație** calculați când datele au provenit dintr-un sondaj.

Corelația statistică

Eficiența aplicării metodei corelației depinde de punerea (enunțarea) corectă a problemei în studiu precum și de aplicarea corectă a statisticii matematice.

Caracterul complex al dependenței statistice pune pe primul plan problema identificării existenței legăturilor. Calculul indicatorilor de corelație este admis cu condiția stabilirii anticipate a unei legături cauzale reale între fenomenele cercetate. Statistica nu poate să rezolve o astfel de problemă fără ajutorul științei din domeniul căreia face parte fenomenul studiat. Cu alte cuvinte, specialistul din domeniul respectiv trebuie să cunoască temeinic noțiunile analizei statistice implicate pentru a da o interpretare corectă a rezultatelor. Pentru a asigura deducții suficiente de

întemeiate, este necesar includerea în cercetare, dacă este posibil, a tuturor factorilor cu acțiune esențială.

La fenomenele simple, unde cauzele acționează separat, relația dintre fenomenul-efect și fenomenul-cauză se reprezintă sub forma:

$$y=f(x)$$

unde x reprezintă cauza, iar y efectul.

La fenomenele complexe, dependența se exprimă sub forma generală:

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

Fenomenul y este generat de acțiunea comună a factorilor x_1, x_2, \dots, x_n (cauze), din care luăm însă în calcul numai o parte.

Să admitem că am luat în calcul factorul x_1 . Întrebarea care se pune este următoarea: în ce condiții indicatorii corelației obținuți exprimă măsura reală a influenței variabilei x_1 asupra variabilei y ? Numai cu condiția ca factorul x_1 să fie hotărâtor în determinarea lui y , ceilalți fiind neesențiali. În cazul în care fenomenul este sub acțiunea unui complex de factori esențiali și aceasta este situația obișnuită, pentru a exprima influența și gradul de intensitate a legăturilor în raport cu un singur factor trebuie să eliminăm influența celorlalți.

Să considerăm o colectivitate statistică caracterizată prin mărimile X și Y . Efectuând o serie de determinări experimentale (sau observații) asupra acestei colectivități, putem întocmi tabela datelor respective:

$$X | x_1, x_2, \dots, x_n$$

$$Y | y_1, y_2, \dots, y_n$$

Repartiția empirică a celor două variabile se poate afișa grafic, într-un sistem de axe XOY , unde vom reprezenta punctele de coordonate x_i și y_i . Un ansamblu de astfel de puncte se numește **câmp de corelație, tabel de corelație** sau **nor statistic**.

În Excel, acest lucru se poate face utilizând **diagrama XY prin puncte (XY-scatter)**

Analiza vizuală a organizării și formei norului de puncte obținut poate oferi indicii importante asupra relației dintre variabile.

Datele de sondaj vor susține ipoteza asocierii între variabile dacă forma norului de puncte se apropie de o curbă funcțională.

Dacă punctele $M_i(x_i, y_i)$ sunt distribuite de-a lungul unei fâșii, care în general, urmează o curbă determinată, spunem că între mărimile respective există o dependență funcțională. Dacă punctele $M_i(x_i, y_i)$ nu arată o dependență funcțională strictă, dar există o tendință ca valorile lui Y să depindă de cele ale lui X , deși nu în mod riguros, între mărimile X și Y există o corelație. Aceasta

poate să fie liniară (fig. 1) sau neliniară (fig. 2). În cazul când între X și Y nu există nici un fel de dependență, câmpul de distribuție se va prezenta asemănător cu acela arătat în fig. 3, 4. Cele două caracteristici sunt independente. Într-un caz particular, dependența corelațională se poate transforma într-o dependență funcțională, dar cu un anumit grad de certitudine. Apare problema de stabili cantitativ (numeric) în ce măsură dependența corelațională se apropie sau se depărtează de dependența funcțională.

În foarte multe cazuri, din observarea fenomenelor naturale sau a proceselor sociale, fără a cunoaște natura exactă a acestora și nici cauzele prin care este pusă în evidență o anumită caracteristică, se pot trage concluzii foarte importante prin examinarea corelației dintre aceste trăsături și alte evenimente. În acest mod se poate aprecia existența unei relații statistice între două sau mai multe variabile, adică, în astfel de cazuri se vorbește despre corelații dintre mărimile care indică o dependență reciprocă.

Corelația este o metodă statistică de determinare a relației dintre două variabile existente.

Regresia este o metodă statistică utilizată pentru a descrie natura relației dintre variabile dacă este pozitivă, liniară sau neliniară

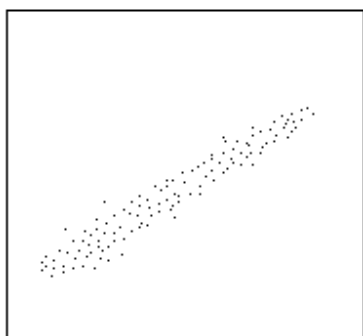


Figura 1 Distribuție liniară

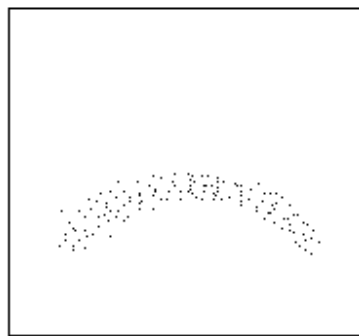


Figura 2 Distribuție neliniară

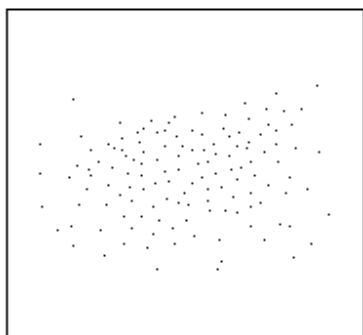


Figura 3 Distribuție aleatoare

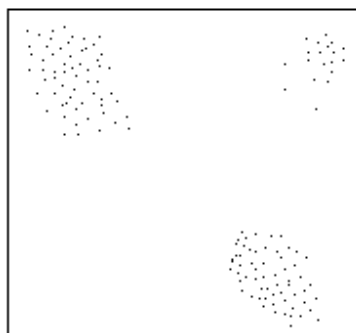


Figura 4 Distribuție grupată

Fiind date două variabile X și Y , se pune problema dacă între cele două variabile, respectiv între fenomenele descrise de acestea, există o anumită dependență numită și **corelație**. O primă concluzie se poate obține reprezentând grafic într-un sistem de coordonate XY , cele două șiruri de date observate pentru cele două variabile. Dacă punctele graficului se împrăștie pe toată suprafața fără a urma o anumită regulă, atunci vom spune că cele două variabile nu sunt corelate. Dacă în schimb punctele descriu o anumită curbă, numită și **curbă de regresie**, atunci vom spune că există corelație și ea este cu atât mai intensă, cu cât domeniul pe care se întind punctele este mai îngust. Mai mult, dacă punctele se așează pe o curbă care poate fi aproximată de o curbă clasică (dreaptă, parabolă, exponențială, etc.) atunci vom spune că legătura dintre cele două variabile este una liniară sau parabolică sau exponențială, etc. și vom folosi ecuația acelei curbe clasice pentru prognoză.

Exemplu:

Un exemplu de studiu a legăturii specifice dintre :

X = producția totală de substanțe chimice din Europa înregistrată între anii 1996 până în 2007

$X=(x_i, i=1..12)$ și

Y =gradul de poluare a aerului prin măsurarea cantității de gaze emise $Y=(y_i, i=1..12)$

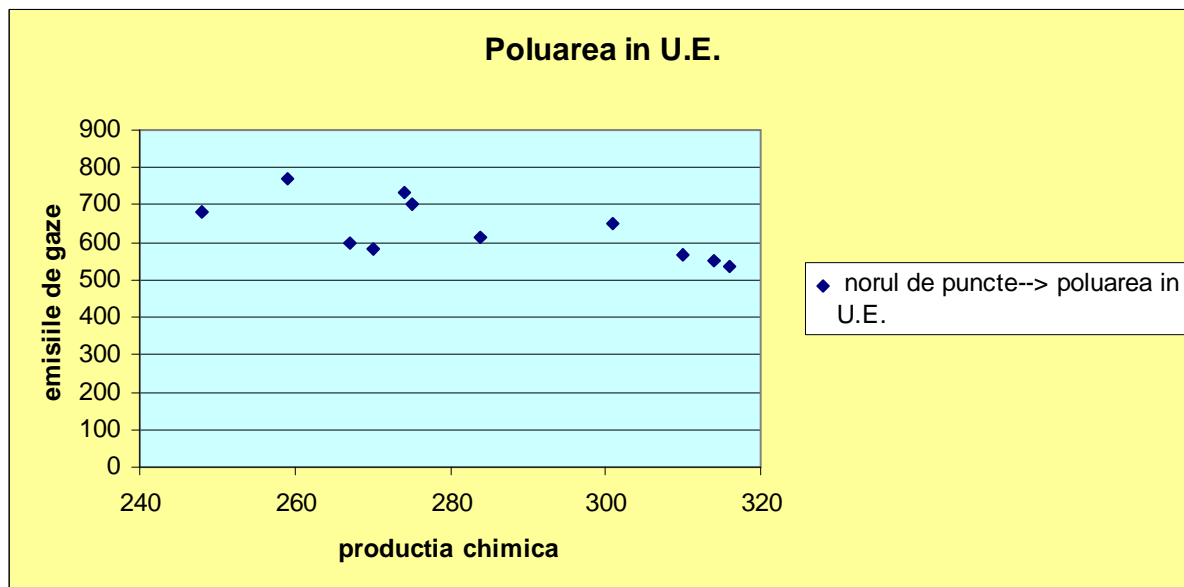
În Excel :

$x = (259, 274, 275, 248, 301, 284, 267, 270, 310, 314, 316, 317)$ și

$y = (769.75, 732.58, 702.34, 682.04, 648.38, 614.57, 600.43, 583.63, 565.42, 552.28, 533.93, 521.32)$

Datele au fost obținute de pe site-ul www.eurostat.eu

anii	1996	1997	1998	1999	2000	2001	2002	2003	2004	2005	2006	2007
producția chimică	259	274	275	248	301	284	267	270	310	314	316	317
volumul emisiilor de gaze	769.75	732.58	702.34	682.04	648.38	614.57	600.43	583.63	565.42	552.28	533.93	521.32



Metoda celor mai mici pătrate.

Dependența funcțională a unei variabile y față de altă variabilă x poate fi studiată empiric, pe cale experimentală, efectuându-se o serie de măsurători asupra variabilei y pentru diferite valori ale lui x . Rezultatele se pot prezenta sub formă de tabel sau grafic.

Problema care se pune în acest caz este de a găsi reprezentarea analitică a dependenței funcționale căutate, adică de a alege o formulă care să descrie rezultatele experimentului.

Formula se alege dintr-o mulțime de formule de tip determinat, de exemplu.

$$y = ax + b, \quad y = ax^2 + bx + c, \quad y = ae^{bx} + c, \quad y = a + h \sin(\omega t + \varphi)$$

Cu alte cuvinte, problema constă în a determina parametrii a, b, c , ai formulei, în timp ce tipul formulei este cunoscut dinainte ca urmare a unor considerente teoretice sau după forma prezentării grafice a materialului empiric.

Să notăm, la modul general când avem n parametrii, dependența funcțională prin:

$$y = f(x; a_0, a_1, \dots, a_n)$$

Parametrii a_0, a_1, \dots, a_n nu se pot determina exact pe baza valorilor empirice y_1, y_2, \dots, y_n ale funcției, deoarece acestea din urmă conțin erori aleatoare. Este vorba de obținerea unei estimății "suficient de bune".

Formularea problemei

Dacă toate măsurătorile valorilor funcției sunt y_1, y_2, \dots, y_n atunci estimațiile parametrilor a_0, a_1, \dots, a_n se determină din condiția ca suma pătratelor abaterilor valorilor măsurate y_k de la cele calculate $f(x_k; a_0, a_1, \dots, a_n)$, adică expresia

$$S = \sum_{k=1}^n [y_k - f(x_k; a_0, a_1, \dots, a_n)]^2$$

să ia valoarea minimă.

Aflarea valorilor parametrilor a_0, a_1, \dots, a_n , care conduc la cea mai mică valoare a funcției

$$s = s(a_0, a_1, \dots, a_n)$$

revine la rezolvarea sistemului de ecuații

$$\frac{\partial S}{\partial a_0} = 0, \frac{\partial S}{\partial a_1} = 0, \dots, \frac{\partial S}{\partial a_n} = 0$$

Dacă formula empirică depinde liniar de parametri necunoscuți atunci sistemul de mai sus va fi de asemenea liniar.

Dreapta de regresie.

În cazul cel mai simplu se studiază numai două variabile X, Y și se dorește găsirea dependenței:

$$Y = aX + b$$

în ipoteza că X este cauza și Y este efectul.

În urma celor n probe se cunosc datele (x_i, y_i) , $i=1, \dots, n$ și trebuie să determinăm coeficienții a și b astfel încât suma

$$S(a, b) = \sum_{i=1}^n (ax_i + b - y_i)^2$$

să fie minimă. Se obține

$$a = \frac{\bar{c}_{xy}}{\sigma_x^2} = \frac{\sigma_y}{\sigma_x} \cdot \frac{\bar{c}_{xy}}{\sigma_x \sigma_y} = \frac{\sigma_y}{\sigma_x} \cdot r_{xy} \quad b = \bar{y} - a \cdot \bar{x}$$

unde σ_x^2 este dispersia variabilei x , iar σ_y^2 este dispersia variabilei y .

Mărimea

$$\bar{c}_{xy} = \overline{xy} - \bar{x} \cdot \bar{y}$$

se numește **covarianța variabilelor X și Y**. sau notația echivalentă :

$$\text{cov}(x, y) = M(x \cdot y) - M(x) \cdot M(y)$$

$$M(x \cdot y) = \frac{\sum_{i=1}^N x_i \cdot y_i}{N}$$

Pentru a analiza dacă între variabilele X și Y există o legătură liniară, se calculează **coeficientul de corelație liniară**, dat de formulele echivalente:

$$r_{xy} = \bar{c}_{xy} / \sigma_x \sigma_y$$

$$\text{Correl}(X, Y) = \frac{\sum (x - \bar{x})(y - \bar{y})}{\sqrt{\sum (x - \bar{x})^2 \sum (y - \bar{y})^2}}$$

Funcția Excel corespunzătoare este CORREL.

Sintaxa funcției :

CORREL (X,Y)= coeficientul de corelație

-parametrii de intrare reprezintă doi vectori de aceeași dimensiune care conțin valorile celor două variabile pentru care dorim să calculăm coeficientul de corelație.

Interpretare:

Valorile coeficientului de corelație sunt în intervalul $[-1, 1]$.

Dacă **r = 0** între cele două variabile **nu există corelație**.

Dacă **r = 1**, corelația între cele două variabile este **maximă și directă**.

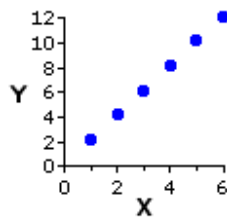
Dacă **r = -1**, corelația între cele două variabile este **maximă și inversă**.

Cu cât avem o valoare mai apropiată de 1 sau -1 cu atât corelația e mai puternică (directă pentru valori pozitive și inversă pentru valori negative), cu cât avem o valoare apropiată de 0 corelația este mai slabă.

În ambele cazuri (r=-1, sau 1) , sintagma "*tendința de a fi asociat*" este un alt mod de a spune că variabilitatea lui X tinde să fie asociată cu variabilitate în Y, și vice-versa, sau, pe scurt, că *X și Y au tendința de a varia împreună*.

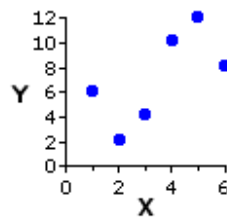
Example I.

$r = +1.0, r^2 = 1.0$



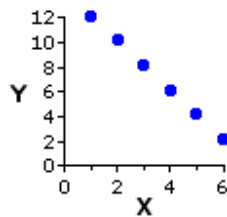
Example II.

$r = +0.66, r^2 = 0.44$



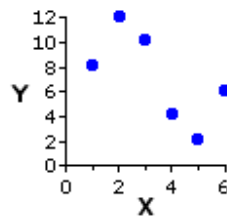
Example III.

$r = -1.0, r^2 = 1.0$



Example IV.

$r = -0.66, r^2 = 0.44$

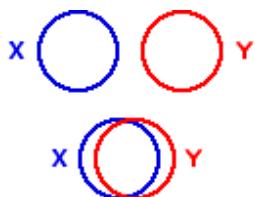


Coeficient de determinare

r^2 = *coeficient de determinare*

Dacă reprezentăm variația lui X ,respectiv Y prin suprafața a două cercuri (suprafața cercului reprezintă 100% din variația fie a lui X sau Y). În caz de corelație zero nu există nici o tendință pentru X și Y pentru a co-varia; și, prin urmare, după cum este ilustrat ,avem două cercuri separate în partea de sus;

corelația nu este zero,atunci cele două cercuri se suprapun(o parte din variabilitatea lui X este explicată de variabilitatea lui Y). și anume r^2 (care reprezintă procentul din variația lui y determinată de variația lui x)



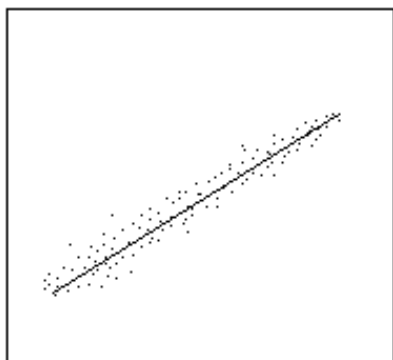
În acest caz cele două cercuri se suprapun ($r = -0.86, r^2 = 0.74$) Zona de suprapunere reprezintă 74% din variația lui Y (Yeste cuplat în variabilitate cuX, și invers), precum și faptul că 26% din variația lui Y este fără nici o legătură cu variabilitatea lui X, la fel ca din variația de X este fără nici o legătură cu variabilitate în Y.(zonele de ne-suprapunere reprezintă 26%) Aceasta zonă de ne-suprapunere sau porțiune de dezacord, între X șiY se numește **variație reziduală** = $1 - r^2$.

Dar faptul că relațiile de cauzalitate între variabile pot produce corelații nu implică faptul că o relație de cauzalitate se află în spatele fiecărui exemplu de corespondență.

În final se va obține ecuația de regresie:

$$Y - \bar{y} = \frac{\sigma_y}{\sigma_x} r_{xy} (X - \bar{x})$$

Această dependență reprezintă o dreaptă numită dreaptă de regresie a variabilei Y în raport cu variabila X.



Dreapta de regresie

Observație. Se poate vorbi și de dependența variabilei X în funcție de Y. Urmând un calcul asemănător se ajunge la dreapta de regresie a variabilei X în raport cu Y:

$$X - \bar{x} = \frac{\sigma_x}{\sigma_y} r_{xy} (Y - \bar{y})$$

Se observă că cele două drepte de regresie coincid dacă și numai dacă $r_{xy}^2 = 1$.

Observații.

1. Trebuie să facem observația că, indiferent de gradul de împrăștiere al punctelor, întotdeauna se poate găsi o dreaptă de regresie, dar în cazul unei dispersii mari aceasta devine inutilă. De aceea un studiu preliminar a distribuției punctelor în plan sau spațiu se impune cu necesitate.

2. Coeficientul de corelație este o mărime foarte importantă în cadrul regresiei liniare. El măsoară gradul de dependență liniară între cauză și efect și are o valoare cuprinsă între -1 și 1 . Aproximarea de 1 implică o dependență liniară puternică între mărimi, iar apropierea de zero indică o lipsă a corelației. Valorile negative semnifică o corelație inversă.

EXAMPLE

