

METODE CORELATIONALE

Regresia si Corelatia

Stud. Master - AMP

ISAIC- MANIU ALEXANDRU
web www.amaniu.ase.ro
e-mail AL.ISAIC-MANIU@CSIE.ASE.RO

9.XII.2013

1

Context

- Intre metodele cantitative de cercatare utile sunt si cele de studiere a **dependențelor** dintre două sau mai multe variabile, incluse în "analiza de regresie și corelație".
- În cadrul acesteia se studiază dependența dintre o variabilă (caracteristică) rezultativă (**y**) și una sau mai multe variabile (caracteristici) independente (**x**).
- Caracteristica rezultativă se mai numește caracteristica *dependentă*, *endogenă* sau *efect*, iar caracteristica independentă se mai numește caracteristica *factorială*, *exogenă* sau *cauză*.
- **Regresia** ne arată **cum** (ca formă analitică) o variabilă este dependentă de altă variabilă (sau de alte variabile), iar **corelația** ne arată **gradul** (cat) în care o variabilă este dependentă de o altă variabilă (sau alte variabile).

9.XII.2013

2

I. REGRESIA

Clasificări

a) după numărul caracteristicilor independente:

- legături simple
- legături multiple

b) după direcția legăturilor, acestea pot fi:

- legături directe
- legături inverse

c) după expresia analitică a legăturilor:

- legături liniare
- legături neliniare

d) după tipul legăturii

- legături funcționale
- legături stohastice

e) după modul de manifestare în timp:

- legături concomitente
- legături cu decalaj

9.XII.2013

3

METODE ELEMENTARE DE CARACTERIZARE A LEGĂTURILOR DINTRE VARIABILE

- 1. metoda seriilor paralele independente
- 2. metoda grupărilor
- 3. metoda tabelului de corelație
- 4. metoda grafică

9.XII.2013

4

1. Metoda seriilor paralele interdependente

Se ordonează observațiile în funcție de caracteristica independentă **X** (crescător sau descrescător) și se urmează modul în care se aranjează valorile lui **Y**.

Concluzii:

- caracteristica y se ordonează aproximativ crescător - rezultă că putem aprecia că între cele două variabile există o legătură directă;
- caracteristica y se ordonează aproximativ descrescător - rezultă că putem aprecia că între cele două variabile există o legătură inversă;
- caracteristica y nu înregistrează o tendință de ordonare (crescător sau descrescător) - rezultă că putem aprecia că între cele două variabile nu există legătură.

9.XII.2013

5

2. Metoda grupărilor

- Se repartizează unitățile în grupe omogene în funcție de o caracteristică independentă.
- Pentru fiecare grupă astfel constituită se centralizează datele numerice referitoare la caracteristica rezultativă și se calculează medii pe fiecare grupă și mărimi relative.
- Prin comparația variației caracteristicii independente cu indicatorii calculați pentru caracteristica rezultativă se poate aprecia existența și forma legăturilor dintre cele două variabile.

9.XII.2013

6

3. Metoda tabelului de contingență

- Tabelul de contingență este un tabel cu dublă intrare și prezintă o grupare a unităților unei colectivității în funcție de două caracteristici: una dependentă și alta independentă.
- Se folosește în special în cadrul unui număr mare de observații.
- Dacă considerăm două variabile, de exemplu "naționalitate" și "religie", atunci tabelul poate fi de forma:

9.XII.2013

7

Exemplu - Tabel de contingență

Religia/ Naționalitatea	Română	Maghiară	.	.	General	.	.	Slovaci	Altele	Total
Ortodoxă	n_{11}	n_{12}	.	.	n_{1j}	.	.	n_{1p-1}	n_{1p}	$n_{1.}$
Romano-catolică	n_{21}	n_{22}	.	.	n_{2j}	.	.	n_{2p-1}	n_{2p}	$n_{2.}$
.
.
General	$n_{.1}$	$n_{.2}$.	.	$n_{.j}$.	.	$n_{.p-1}$	$n_{.p}$	$n_{.}$
.
.
Musulmană	n_{r-11}	n_{r-12}	.	.	n_{r-1j}	.	.	n_{r-1p-1}	n_{r-1p}	$n_{r-1.}$
Altele	$n_{r.1}$	$n_{r.2}$.	.	$n_{r.j}$.	.	$n_{r.p-1}$	$n_{r.p}$	$n_{r.}$

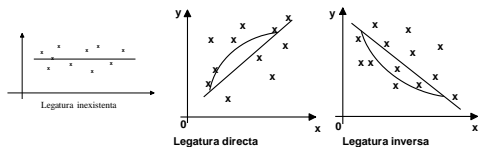
9.XII.2013

8

4. Metoda grafică

Graficul se construiește pornind de la perechile de valori observate (x , y) care se reprezintă în sistemul de axe rectangulare – *corelogramă* (nor de puncte).

Pe axa OX se reprezintă variabila independentă x , iar pe axa OY variabila dependentă y .



9.XII.2013

9

5. MODELUL UNIFACTORIAL LINIAR

Modelul probabilist la nivelul colectivității generale:

$$y_i = \alpha + \beta \cdot x_i + \varepsilon_i$$

- (x_i, y_i) reprezintă valorile numerice ale variabilelor cauză și efect înregistrate la nivelul unității statistice „i”;
- α, β = parametri constanți
- α = punctul de intersecție al dreptei de regresie cu axa Oy;
- β = panta dreptei, se mai numește și „coeficient de regresie” și arată cu câte unități de măsură se modifică Y dacă X se modifică cu o unitate
- ε_i = componenta reziduală (eroare aleatoare) pentru unitatea statistică „i”.

9.XII.2013

10

5.1. MODELUL UNIFACTORIAL LINIAR

Valoarea reală y_i a caracteristicii Y din modelul probabilistic cuprinde:

- *componenta teoretică, deterministă* (\hat{y}_i), indică partea din valoarea reală y_i care se poate determina pe baza modelului pentru o anumită valoare x_i :

$$\hat{y}_i = \alpha + \beta \cdot x_i$$

- *componenta aleatoare (reziduală)*, numită și *eroarea aleatoare* (ε_i) reprezentând acea parte din valoarea reală a lui Y care nu se poate cuantifica:

$$y_i = \hat{y}_i + \varepsilon_i$$

9.XII.2013

11

5.2. MODELUL UNIFACTORIAL LINIAR

- Dacă datele disponibile provin dintr-un eșantion, avem n perechi de observații reale:

$(x_1, y_1), (x_2, y_2), \dots, (x_n, y_n)$, pe baza cărora se vor estima parametrii ecuației, α și β .

- Modelul de regresie în eșantion va fi $y_i = a + b \cdot x_i + e_i$

cu $\hat{y}_i = a + b \cdot x_i$

unde:

- a = estimatorul parametrului α din colectiv. generală;
- b = estimatorul parametrului β din colectiv. generală;
- e_i = valoarea reziduală pt. unitatea “i” în eșantion.

9.XII.2013

12

5.3. Estimarea parametrilor modelului unifactorial liniar

Metoda celor mai mici pătrate presupune maximizarea similitudinii, a gradului de asemănare a valorilor teoretice cu valorile reale, deci minimizarea erorilor.

Cum erorile se pot produce într-un sens sau în altul față de valorile reale, ea presupune *minimizarea sumei pătratelor reziduurilor*:

9.XII.2013

13

5.3.-b Estimarea parametrilor modelului unifactorial liniar

$$S = \sum e_i^2 = \sum (y_i - \hat{y}_i)^2 = \sum (y_i - a - b \cdot x_i)^2 \rightarrow \text{m\u00e2n\u0152m}$$

Condi\u021bia de ordin 1 de minimizare a func\u021biei sunt:

$$\begin{cases} \frac{dS}{da} = 0 \\ \frac{dS}{db} = 0 \end{cases} \Rightarrow \begin{cases} \sum 2(y_i - a - b \cdot x_i)(-1) = 0 \\ \sum 2(y_i - a - b \cdot x_i)(-x_i) = 0 \end{cases} \Rightarrow \begin{cases} \sum y_i - na - b \sum x_i = 0 \\ \sum x_i y_i - a \sum x_i - b \sum x_i^2 = 0 \end{cases}$$

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

9.XII.2013

14

5.3.-c. Estimarea parametrilor modelului unifactorial liniar

Aplic\u00e2nd metoda determinan\u021bilor, se ob\u021bine:

$$\Delta a = \begin{vmatrix} \sum y_i & \sum x_i \\ \sum x_i y_i & \sum x_i^2 \end{vmatrix} \quad \Delta b = \begin{vmatrix} n & \sum y_i \\ \sum x_i & \sum x_i y_i \end{vmatrix}$$

$$\Delta = \begin{vmatrix} n & \sum x_i \\ \sum x_i & \sum x_i^2 \end{vmatrix}$$

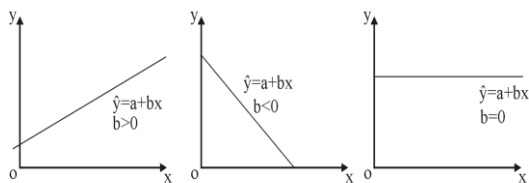
$$a = \frac{\Delta a}{\Delta} = \frac{\sum y_i \cdot \sum x_i^2 - \sum x_i \cdot \sum x_i y_i}{n \sum x_i^2 - (\sum x_i)^2}$$

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

9.XII.2013

15

5.3.-f. Estimarea parametrilor modelului unifactorial linear



Linii de regresie cu a) pantă pozitivă b) pantă negativă c) pantă egală cu zero

9.XII.2013

16

6. Modele neliniare de regresie

Modelul exponențial transformat al ecuației exponențiale are la bază ecuația:

$$y = a b^x \quad \text{care se estimează folosind modelul:}$$

$$Y = a b^x + \varepsilon$$

Prin logaritmare, modelul se poate transforma într-un model linear de forma:

$$\lg Y = \lg a + x \lg b$$

Făcând următoarele înlocuiri:

$$Y' = \lg Y; \quad a' = \lg a; \quad b' = \lg b, \text{ rezultă ecuația unei drepte,}$$

respectiv:

$$y' = a' + b' x$$

9.XII.2013

17

Interpretări

- **Coefficientul „a”**, care poate lua atât valori pozitive cât și negative, reprezintă ordonata la origine, respectiv este valoarea lui „y” când „x” este egal cu zero.
- **Coefficientul „b”** - denumit coeficient de regresie - arată măsura în care variază caracteristica dependentă în cazul în care caracteristica independentă se modifică cu o unitate.
- În funcție de semnul coeficientului de regresie, putem aprecia tipul de legătură:
 - o în cazul corelației directe, coeficientul are o valoare pozitivă;
 - o în cazul corelației inverse, valoarea lui este negativă;
 - o în cazul în care $b = 0$, se apreciază că variabilele (y și x) sunt independente.
- În graficul de corelație coeficientul „b” indică panta liniei drepte.

9.XII.2013

18

II . Metoda corelației

- A. Corelația parametrică (variabile măsurate pe scala de raport)
- B. Corelația neparametrică (variabile măsurate pe scala nominală, ordinală sau de interval)

9.XII.2013

19

A. Corelația parametrică

Metoda corelației prezintă avantajul că oferă o măsură sintetică a legăturilor dintre variabilele statistice.

Indicatorii care măsoară intensitatea legăturii sunt: covarianța, coeficientul de corelație și raportul de corelație.

1. COVARIANȚA

Covarianța se calculează sub forma mediei aritmetice simple a produselor abaterilor celor două variabile corelate, x și y , de la mediile lor aritmetice \bar{x} și \bar{y} , conform relației:

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

9.XII.2013

20

1. Covarianța

Covariația **nulă** - variabilele sunt **independente** (lipsa legăturii de corelație)

- Valoarea sa absolută cov (x,y) **nu are limită superioară**. Pe măsură ce intensitatea corelației crește și covariația crește.
- Indicatorul reprezintă **avantajul** că se calculează destul de ușor. În același timp, prezintă și **dezavantajul** că depinde de unitățile în care se măsoară variabilele aleatoare. Deci nu este comparabil de la o variabilă la alta.
- Indicatorul ia **valori pozitive** dacă legătura dintre variabile este directă și **valori negative** în caz contrar.
- Valori apropiate de zero semnifică lipsa oricărei legături între x și y; valori ridicate ale indicatorului arată o legătură puternică.

9.XII.2013

21

2. COEFICIENTUL DE CORELAȚIE LINIARĂ SIMPLĂ

Măsoară numai intensitatea legăturii de **tip liniar** dintre două variabile x și y . Se calculează ca o medie aritmetică a produsului abaterilor normale ale celor două variabile.

Notând **abaterea normală** ale variabilelor x și y :

$$z_x = \frac{x - \bar{x}}{s_x}; \quad z_y = \frac{y - \bar{y}}{s_y}$$

rezultă următoarea relație de calcul (în care „ n ” este numărul observațiilor-perechi)

$$r_{xy} = \frac{\sum (x_i - \bar{x})(y_i - \bar{y})}{n \cdot s_x \cdot s_y} = \frac{\sum (x_i - \bar{x})(y_i - \bar{y})}{\sqrt{\sum (x_i - \bar{x})^2} \sqrt{\sum (y_i - \bar{y})^2}}$$

Față de covarianță rezultă că relația:

$$r_{xy} = \frac{\text{cov}(x, y)}{s_x \cdot s_y} = \frac{\sum (x_i - \bar{x})(y_i - \bar{y})}{n \cdot s_x \cdot s_y}$$

se transformă în coeficientul de corelație liniară simplă.

9.XII.2013

22

2. COEFICIENTUL DE CORELAȚIE LINIARĂ SIMPLĂ

În practică se utilizează relația:

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

Coeficientul de corelație simplă se mai poate calcula și cu relația:

$$r_{xy} = b \frac{s_x}{s_y},$$

în care:

- b - este coeficientul de regresie simplă;
- s_x - abaterea medie pătratică a caracteristicilor factoriale;
- s_y - abaterea medie pătratică a caracteristicilor rezultative.

9.XII.2013

23

2.a. COEFICIENTUL DE CORELAȚIE LINIARĂ SIMPLĂ

- Coeficientul de corelație poate lua **valori** cuprinse între -1 și +1, adică satisface inegalitățile: $-1 \leq r_{yx} \leq 1$, iar semnul său, ca și cel al coeficientului de regresie, semnifică tipul de legătură: semnul minus indică legătura inversă, semnul plus indică legătura directă.
- Cu cât coeficientul de corelație are valori mai apropiate de 1 sau -1, cu atât corelația liniară dintre variabilele x și y este mai puternică.
- Pe măsură ce coeficientul de corelație se apropie de 0, scade și **intensitatea** legăturii dintre cele două variabile.
- În cazul în care $r_{yx}=0$, variabilele sunt independente ori necorelate liniar, iar pentru $r_{yx}=1$ rezultă dependența funcțională între cele două variabile.

9.XII.2013

24

3. - Raportul de corelație

Denumit și **coeficientul de corelație** Pearson, acest indicator măsoară atât intensitatea legăturilor liniare, cât și neliniare. Se definește cu relația:

$$R_{y,x} = \sqrt{1 - \frac{\sum (y_i - \hat{y}_i)^2}{\sum (y_i - \bar{y})^2}}$$

Coeficientul de determinare

$$R^2$$

9.XII.2013

25

3.1. - Raportul de corelație

Dispersiile au următoarele semnificații:

- măsoară acțiunea tuturor factorilor care au influențat asupra variabilei rezultative;
- măsoară variația valorilor y sub influența tuturor celorlalți factori necuprinși în model, a căror acțiune e considerată constantă; este denumită și dispersia reziduală;
- măsoară numai influența variabilei independente sau factoriale x asupra variabilei y . Cu cât ponderea acestei dispersii în cadrul dispersiei generale va fi mai mare, cu atât legătura dintre cele două variabile va fi mai puternică.

9.XII.2013

26

3.2. - Raportul de corelație

Interpretare

Raportul de corelație poate lua valori între 0 și 1. Cu cât valoarea raportului este mai apropiată de 1 cu atât legătura de corelație este mai puternică și invers.

- Dacă $R \rightarrow 1$ legătura dintre X și Y este puternică.
- Dacă $R \rightarrow 0$ legătura dintre X și Y este slabă.

În cazul corelației liniare, raportul de corelație este egal cu coeficientul de corelație luat în valoare absolută și această relație poate fi considerată ca un test de verificare a liniarității legăturii.

- În cazul legăturilor liniare:

$$R = |r_{xy}|$$

9.XII.2013

27

B. Corelația neparametrică

Dacă metodele de analiză a corelației nu mai folosesc parametrii distribuțiilor, se poate vorbi despre o corelație neparametrică (sau liberă de distribuție), iar acest lucru se întâmplă în următoarele **situații**:

- Variabilele sunt de natură **calitativă** (măsurate pe scală nominală sau ordinală);
- Datele **nu** provin dintr-o populație cu *distribuție normală* de probabilitate sau aproximativ normală;
- Sunt *informații insuficiente* pentru a putea presupune normalitatea distribuției (de exp, provin din eșantioane de volum redus)

9.XII.2013

28

1. Coeficientul de asociere Q

Distribuția persoanelor în funcție de naționalitate și religie

Religia/ Naționalitatea	Român	Non-român	Total
Ortodox	n_{11}	n_{12}	$n_{1.}$
Non-ortodox	n_{21}	n_{22}	$n_{2.}$
Total	$n_{.1}$	$n_{.2}$	n

Coeficientul Q

A fost introdus de statisticianul englez **Yule**. Formula sa de calcul este:

$$Q = \frac{n_{11} \cdot n_{22} - n_{12} \cdot n_{21}}{n_{11} \cdot n_{22} + n_{12} \cdot n_{21}}$$

Acești ai valori tot între -1 și 1 .

În cazul unor variabile independente valoarea coeficientului Q este nulă.

9.XII.2013

29

2. Coeficienții de corelație ai rangurilor

- a. Coeficientul de corelație a rangurilor introdus de **Spearman** se bazează pe analiza concordanței rangurilor acordate pentru fiecare din cele n unități statistice, după variabila X și după variabila Y .

$$C_s = 1 - \frac{6 \sum d_i^2}{n(n^2 - 1)}$$

- b. Coeficientul rangurilor **Kendall** necesită ordonarea unităților după rangurile acordate variabilei X și înscrierea în paralel a rangurilor acordate variabilei Y .

$$C_K = \frac{2 \sum S_i}{n(n-1)}$$

9.XII.2013

30

Exemplu

Pe baza datelor din anuarul statistic, s-au înregistrat datele următoare pentru 10 județe.

Nr.	Jud.	Supraf. (km ²)	Nr. Comunelor
1	AB	6242	66
2	AG	6826	93
3	AR	7754	67
4	BC	6621	79
5	BH	7544	86
6	BN	5355	53
7	BR	4766	39
8	BT	4986	68
9	BV	5363	43
10	BZ	6103	81

Să se stabilească dacă există o legătură între suprafața totală și numărul comunelor, utilizând coeficienții de corelație a rangurilor a lui Spearman și Kendall.
