

## Proprietatile descriptorilor statistici pentru serii univariate

Prof.dr. Vergil VOINEAGU, conf.dr. Tudorel ANDREI  
 Catedra de Statistica si Previziune Economica, A.S.E. Bucuresti

Studierea unui fenomen se realizeaza prin considerarea unui set de variabile masurate pentru unitatile unei populatii  $G$  constituite din  $n$  unitati elementare. Vom considera ca variabilele sunt în numar de  $p$ . Caracteristicile cuprinse în programul observarii statistice sunt definite prin urmatoarele elemente: 1. Spatiul observatiilor, notat prin  $W = W_1 \times W_2 \times \dots \times W_p$ . 2. Structura algebrica introdusa în spatiul de observatii, notata prin  $S = S_1 \times S_2 \times \dots \times S_p$ . 3. Aplicatia  $v: G \rightarrow W$ , care atribuie fiecarei unitati din cadrul populatiei valori pentru caracteristicile din planul observarii, sub forma vectorului  $v(i) = (v_1(i), v_2(i), \dots, v_p(i))$ . S-a considerat  $i \in G$ .

**Cuvinte cheie:** unitate, scala, variabila, observatie, descriptor statistic, serie univariata.

În functie de caracteristicile multimii  $\Omega$  si de structura  $S$ , se disting urmatoarele tipuri de variabile:

1. *Calitativ nominale.* În acest caz vom avea o multime formata dintr-un numar finit de coduri numerice. Aceste variabile sunt masurate pe o *scala nominala*. Pe aceasta scala nu se introduce structura de ordine din corpul numerelor reale si nici operatiile obisnuite de adunare, scadere etc. Scala este folosita numai pentru efectuarea operatiilor de clasificare si grupare a unitatilor populatiei.

2. *Calitativ ordinale.* În cadrul acestei variabile, masurata pe o scala ordinala, este permisa introducerea operatiei de ordine de pe axa numerelor reale. Nici în acest caz nu are sens efectuarea de calcule directe asupra valorilor variabilei, întrucât nu este permisa calcularea distantei dintre valori.

3. *Cantitative masurate pe o scala de interval.* În acest caz se introduce distanta între valori, dar nu are sens raportul dintre acestea, întrucât valoarea nula este arbitrara.

4. *Cantitative masurate pe o scala de raport.* În acest caz  $\Omega = \mathbb{R}$  este o multime înzestrata cu o structura de corp ordonat. Structura asociata multimii  $\Omega$  este întotdeauna indusa de semantica subiacenta parametrului analizat, motiv pentru care o numim structura paradigmatica. De exemplu, parametrului "vârsta"  $i$  se asociaza: variabila calitativa ordinala (clase de vârsta

exprimate prin cuvinte); o variabila cantitativa masurata pe o scala de raport (vârsta exacta a unei persoane masurata în ani întregi).

Pot fi utilizate doua tehnici pentru transformarea variabilelor: prin schimbarea structurii si prin diverse operatii de codificare. Fie variabila definita prin aplicatia  $v: \Gamma \rightarrow \Omega$ , care este înzestrata cu o structura  $S$ . În acest caz vom spune ca variabila este supusa unei transformari prin schimbarea structurii, daca aceasta se înlocuieste prin variabila  $v': \Gamma \rightarrow \Omega$  care are structura  $S'$ , astfel încât  $v'(i) = v(i)$  pentru orice  $i \in \Gamma$ .

Pentru a defini schimbarea de variabila prin codificare definim un alt spatiu al observatiilor  $\Omega'$ , înzestrat cu structura  $S'$ . În aceste conditii, definim aplicatia  $c: \Omega \rightarrow \Omega'$ . Astfel, noua variabila este în fapt obtinuta prin compunerea functiilor  $v$  si  $c$ .

Valorile celor  $p$  variabile observate sunt ordonate pentru cele  $n$  unitati ale populatiei într-o matrice de forma:

$$X = \begin{pmatrix} x_1^1, \dots, x_1^j, \dots, x_1^p \\ \dots \\ x_i^1, \dots, x_i^j, \dots, x_i^p \\ \dots \\ x_n^1, \dots, x_n^j, \dots, x_n^p \end{pmatrix}$$

unde:

- $\mathbf{x}^j$  este un vector coloana de dimensiune  $n$ , continând valorile unei caracteristici din planul de observare pentru unitatile populatiei;
- $\mathbf{x}_i$  este vector linie de dimensiune  $p$ , cu valori ale caracteristicilor din planul observarii pentru o unitate statistica din cadrul populatiei  $\Gamma$ .

Daca  $p_i$  sunt ponderi (frecvente relative), cu  $\sum_{i=1}^n p_i = 1$ , atunci vom defini matricea ponderilor astfel:

$$\mathbf{D} = \text{diag}(\mathbf{p}_i) = \begin{pmatrix} p_1 & 0 \dots 0 \dots 0 \\ 0 & p_2 \dots 0 \dots 0 \\ \dots & \dots & \dots & \dots \\ 0 & 0 \dots p_i \dots 0 \\ \dots & \dots & \dots & \dots \\ 0 & 0 \dots 0 \dots p_n \end{pmatrix}$$

Se observa ca în cazul în care unitatile populatiei au aceasi pondere  $p_i = \frac{1}{n}$ , atunci

$$\mathbf{D} = \frac{1}{n} \mathbf{I}, \text{ unde } \mathbf{I} \text{ este matricea unitate.}$$

Prin metode adecvate, vom prelucra seriile de date pentru fiecare caracteristica în parte, cât si pentru fiecare unitate statistica. Determinam în egala masura indicatori pentru masurarea dependentelor dintre caracteristici, a similaritatilor unitatilor statistice etc. Pentru prelucrarea seriilor de date vom introduce:

- descriptori si metrici în spatiul unitatilor;
- descriptori si metrici în spatiul variabilelor.

#### Proprietati ale descriptorului $m(\mathbf{x})$

Vom prezenta în continuare câteva proprietati care pot fi luate în considerare pentru alegerea unei masuri pentru tendinta unei serii univariate.

Definim masura statistica a tendintei centrale prin aplicatia:

$$m: \Omega_j^n \rightarrow \mathbb{R} \quad j = \overline{1, p}$$

unde  $m(\mathbf{x}) \leftarrow (x_1, x_2, \dots, x_n)$ .

Demonstram pentru fiecare masura statistica, ca tendinta centrala satisface una sau mai multe din aceste proprietati.

“Norul de puncte” asociat multimii unitatilor din  $\Gamma$  pentru o caracteristica din planul de observare se noteaza cu:

$$N = \left\{ (x_i, p_i) \mid i = \overline{1, n} \right\}$$

Pentru ansamblul caracteristicilor din planul observarii definim “norul de puncte” prin:

$$\mathbf{N} = \{ (\mathbf{x}_i, \mathbf{p}_i) \mid i = \overline{1, n} \}.$$

În cazul în care doi vectori  $(\mathbf{x}, \mathbf{y})$  sunt comparati sau sunt folositi în cadrul unei expresii algebrice, vom considera ca acest lucru este posibil atât din punct de vedere algebric, cât si economic. Pentru a asigura comparabilitatea economica pentru cei doi vectori vom folosi fie doua serii de date pentru aceeasi caracteristica, fie vom recurge la diverse transformari. Prin transformari de origine si intensitate, vom obtine un vector de valori normalizate din intervalul  $[0, 1]$ .

În cele ce urmeaza vom prezenta proprietatile mai importante pe care o masura statistica a tendintei centrale le îndeplinesc. Precizam ca aceste proprietati sunt puse în discutie pentru cazul seriilor de distributie unimodale.

Proprietatea 1 [P<sub>1</sub>]. Masura statistica  $m(\mathbf{x})$  aplicata vectorilor ansamblului  $\Omega^n$  îndeplinesc *proprietatea de intermediaritate* daca si numai daca  $\forall \mathbf{x} \in \Omega^n$  avem ca aceasta se încadreaza între cele doua valori extreme:

$$\min(x) \leq m(\mathbf{x}) \leq \max(x).$$

Toate masurile statistice ale tendintei centrale folosite în statistica descriptiva satisfac aceasta proprietate. Daca  $x \in \Omega^n$ , atunci nu este implicit si faptul ca  $m(\mathbf{x}) \in \Omega$ . Aceasta este o restrictie robusta pentru o masura a tendintei centrale a unei serii de date.

Proprietatea 2 [P<sub>2</sub>]. O masura statistica satisface *proprietatea de monotonicitate* daca, doua serii de date ale aceleasi caracteris-

tici,  $\forall \mathbf{x}, \mathbf{y} \in \Omega^n$  cu  $x_i \geq y_i$ , atunci  $m(\mathbf{x}) \geq m(\mathbf{y})$ .

Daca prima proprietate este satisfacuta de toti indicatorii tendintei centrale folositi în statistica descriptiva, nu acelasi lucru se întâmpla si în cazul acestei proprietati. Este cazul, de exemplu, al valorii modale, care nu satisface proprietatea pentru orice serie de valori.

Daca pentru doi vectori  $\mathbf{x}, \mathbf{y} \in \Omega^n$ ,  $x_i \geq y_i$   $\forall i = \overline{1, n}$  si în plus exista valori în  $\mathbf{x}$  si  $\mathbf{y}$  astfel încât  $x_i \neq y_i$ , atunci  $\mathbf{x} > \mathbf{y}$ .

Proprietatea 3 [P<sub>3</sub>]. O masura statistica satisface *proprietatea de regularitate* daca si numai daca  $\forall \mathbf{x}, \mathbf{y} \in \Omega^n$ , cu  $\mathbf{x} > \mathbf{y}$ , se obtine  $m(\mathbf{x}) > m(\mathbf{y})$ .

Aceasta proprietate este mult mai robusta decât P<sub>2</sub>. Astfel, majoritatea indicatorilor calculati pe baza structurilor de ordine satisfac P<sub>2</sub>, dar nu si proprietatea P<sub>3</sub>. Este cazul, spre exemplu, cel al medianei care verifica P<sub>2</sub> dar nu si proprietatea de regularitate.

Proprietatea 4 [P<sub>4</sub>]. Masura statistica a tendintei centrale îndeplineste *proprietatea de omogenitate* daca si numai daca  $\forall \mathbf{x} \in \Omega^n$  si  $\forall \lambda \in \mathbb{R}^*$  este satisfacuta inegalitatea:

$$m(\lambda \mathbf{x}) = \lambda m(\mathbf{x})$$

Se demonstreaza fara dificultate ca aceasta proprietate este îndeplinita de media aritmetica si de indicatorii medii de pozitie.

De regula pentru a asigura  $\lambda \mathbf{x} \hat{\in} \Omega^n$  vom considera cazul particular  $\lambda \hat{\in} \mathbb{R}_+^*$ .

Proprietatea 5 [P<sub>5</sub>]. Pentru un ansamblu  $\Omega^n$  vom spune ca masura statistica este *simitrica* daca  $\forall \mathbf{x} \in \Omega^n$ , atunci:

$$m(-\mathbf{x}) = -m(\mathbf{x}).$$

În fapt, aceasta proprietate este un caz particular al proprietatii P<sub>4</sub>, situatie în care vom lua  $\lambda = -1$ . Proprietatea este deosebit de restrictiva daca, pe lângă egalitatea de mai sus, se adauga si conditia de apartenenta  $-\mathbf{x} \in \Omega^n$ .

Proprietatea 6 [P<sub>6</sub>]. Masura  $m(\mathbf{x})$  satisface *proprietatea de aditivitate* daca  $\forall \mathbf{x}, \mathbf{y} \in \Omega^n$

doi vectori pentru care are sens sa definim operatia de adunare, atunci:

$$m(\mathbf{x} + \mathbf{y}) = m(\mathbf{x}) + m(\mathbf{y}).$$

Vom arata ca nu toate masurile care vor fi utilizate pentru analiza tendintei centrale a seriilor de date satisfac proprietatea de aditivitate. De aceea, vom introduce conceptele de masura a tendintei centrale subaditiva si supraaditiva.

O masura a tendintei centrale este subaditiva daca  $\forall \mathbf{x}, \mathbf{y} \in \Omega^n$ , atunci  $m(\mathbf{x} + \mathbf{y}) \leq m(\mathbf{x}) + m(\mathbf{y})$ . În mod asemanator definim ca o masura este supraaditiva daca  $\forall \mathbf{x}, \mathbf{y} \in \Omega^n$ , avem ca  $m(\mathbf{x} + \mathbf{y}) \geq m(\mathbf{x}) + m(\mathbf{y})$ .

Pentru a prezenta proprietatea urmatoare se noteaza  $\mathbf{u} = (1, 1, \dots, 1)$ .

Proprietatea 7 [P<sub>7</sub>]. O masura statistica accepta o *schimbare de origine*, de amplitudine  $h \in \mathbb{R}$ , daca  $\forall \mathbf{x} \in \Omega^n$  si  $\forall h \in \mathbb{R}$  avem ca  $m(\mathbf{x} + h\mathbf{u}) = h + m(\mathbf{x})$ .

Daca pentru o masura a tendintei centrale sunt verificate simultan P<sub>4</sub> si P<sub>7</sub>, atunci  $\forall \mathbf{x} \in \Omega^n$  si  $\forall h, \lambda \in \mathbb{R}$ , avem  $m(\lambda \mathbf{x} + h\mathbf{u}) = h + \lambda m(\mathbf{x})$ . De cele mai multe ori indicatorii tendintei centrale folositi în statistica descriptiva au aceasta proprietate. În cazul în care  $\lambda = 0$  vom obtine  $m(h\mathbf{u}) = h$ , deci centrul de greutate al unui sir de valori constante este însasi valoarea care constituie seria de date.

### Descriptori ai tendintei centrale

Vom prezenta în cele ce urmeaza câteva masuri ale tendintei centrale posibil de utilizat pentru o serie univariata de valori  $(x_i)_{i=1, n}$ , utilizând ponderile  $(p_i)_{i=1, n}$ .

Sistemul de ponderare, care poate fi constituit pe baza frecventelor relative sau probabilitatilor, este definit prin vectorul  $\mathbf{p} = (p_1, p_2, \dots, p_n)$  cu  $p_i \geq 0$  si  $\sum_i p_i = 1$ .

Pentru a defini un descriptor statistic peste ansamblu de valori din  $V = \Omega^n$  vom considera, în cazul în care  $\Omega \subseteq \mathbb{R}$ , functia  $d(\mathbf{x}, \mathbf{y})$  de doua variabile vectoriale  $\mathbf{x}$  si  $\mathbf{y}$  care satisface proprietatile:

i)  $d(\mathbf{x}, \mathbf{y}) > 0, \mathbf{x} \neq \mathbf{y}$ ;

ii)  $d(\mathbf{x}, \mathbf{y})=0, \mathbf{x} = \mathbf{y}$ ;

iii)  $d(\mathbf{x}, \mathbf{y})=d(\mathbf{y}, \mathbf{x})$ ;

iv)  $d(\mathbf{x}, \mathbf{z}) \leq d(\mathbf{x}, \mathbf{y}) + d(\mathbf{y}, \mathbf{z}), \mathbf{x}, \mathbf{y}, \mathbf{z} \in \Omega^n$ .

Functia cu proprietatile de mai sus se numeste metrica sau distanta euclidiană. Considerând aceasta functie vom cauta sa determinam vectorul  $\mathbf{b}=(b, b, \dots, b) \in \mathbb{R}^n$  astfel încât  $d(\mathbf{x}, \mathbf{b})$  sa fie minima. Vom arata ca  $\min_b d(x, b) = d(x, m)$ , unde prin  $\mathbf{m}=(m, m, \dots, m)$  s-a definit un vector  $n$  dimensional, cu  $m$  o valoare obtinuta prin aplicarea asupra seriei a unei masuri a tendintei centrale.

Pentru  $p \geq 1$  consideram ca spatiul vectorial  $\mathbb{R}^n$  este înzestrat cu norma:

$$L_p(x) = \|x\|_p = \left( \sum_{i=1}^n p_i |x_i|^p \right)^{\frac{1}{p}}.$$

Vom estima tendinta centrala prin interme-

$$L_p(x - tu) = \|x - tu\|_p = \left( \sum_{i=1}^n p_i |x_i - t|^p \right)^{\frac{1}{p}}$$

diul descriptorului  $m_p(x) = tu$ . Vom determina  $t \in \mathbb{R}$  din conditia de minim a functiei.

Daca  $p=1$ , atunci minimul functiei  $L_p(x-tu)$  se situeaza în intervalul median; pentru  $p>1$  valoarea minima a functiei este unica. Pentru valoarea optima a lui  $t$  definim atunci distanta dintre vectorul observatiilor si vectorul  $m_p(x)$ , prin:

$$d_p(x, m_p(x)) = \|x - u m_p(x)\|$$

În practica sunt utilizate urmatoarele cazuri particulare pentru parametrul  $p \in \mathbb{R}$ :

i)  $p=1$ , atunci valoarea minima a expresiei este mediana seriei de valori, iar indicatorul pentru masurarea dispersiei seriei de valori este:

$$e_1(x) = \sum_i p_i |x_i - x_m|.$$

ii) Pentru  $p=2$  vom obtine metrica euclidiană. În acest caz din conditia de minim a functiei  $L_2(x-tu)$  vom determina:

$$m_2 = \sum_i p_i x_i$$

iar indicatorul pentru caracterizarea gradului de dispersare este abaterea standard.

iii) Daca  $p \rightarrow \infty$  atunci, dispunem de metrica convergentei uniforme, iar:

$$m_\infty = \frac{1}{2} (\min(x) + \max(x))$$

este centrul de greutate al intervalului în care sunt plasate observatiile caracteristicii. Pentru caracterizarea dispersiei se va utiliza masura:

$$e_\infty(x) = \frac{1}{2} (\max(x) - \min(x)).$$

Pentru vectorul  $\mathbf{x}$  format din numere reale definim urmatoorii descriptori statistici:

i) media de ordinul  $p$ :

$$m_p(x) = L_p(x) = \left( \sum p_i x_i^p \right)^{\frac{1}{p}}, p \in \mathbb{R}^*;$$

Pentru  $p=0$  vom obtine formula mediei geometrice:

$$m_0 = \lim_{p \rightarrow 0} m_p(x) = \prod_i x_i^p.$$

În tabelul 1 vor fi trecute diverse particularizari ale mediei de ordinul  $p$ .

**Tabelul 1.**

$\alpha$	Media obtinuta	Formula de calcul
$-\infty$	Minimul seriei	$m_{-\infty}(x) = \min(x)$
-1	Media armonica	$m_{-1}(x) = \frac{1}{\sum_i \frac{p_i}{x_i}}$
1	Media aritmetica	$m_1(x) = \sum p_i x_i$
2	Media patratice	$m_2(x) = \left( \sum p_i x_i^2 \right)^{1/2}$
$\infty$	Maximul seriei	$m_\infty(x) = \max(x)$

Se demonstreaza ca media de ordinul  $p$ , în raport cu ordinul, este o functie monoton crescatoare. Pentru cazurile particulare din tabelul 1, se verifica relatia de ordine:

$$m_{-\infty}(x) = \min(x) \leq m_{-1}(x) \leq m_0(x) \leq m_1(x) \leq m_2(x) \leq m_\infty(x) = \max(x)$$

ii) media logaritmică de ordinul  $f$ , este definită pentru orice vector  $\mathbf{x} \in \mathbb{R}$  care are ponderile  $\mathbf{p} = (p_1, p_2, \dots, p_n)$ , prin relația:

$$l_f(x) = \frac{1}{f} \log \left( \sum_{i=1}^n p_i e^{f \cdot x_i} \right), \quad f \neq 0.$$

Dacă  $f=0$ , atunci  $l_0(x) = \lim_{f \rightarrow 0} l_f(x) = m_1(x)$

În fapt, între ultimele două medii prezentate se verifică relația:

$$l_f(x) = \log(m_f(e^x)),$$

unde  $e^e = (e^{x_1}, e^{x_2}, \dots, e^{x_n})$ .

### Descriptori generali

Pentru norul de puncte  $\mathbf{N}$  definim centrul de greutate, care coincide cu vectorul mediilor aritmetice calculate pentru cele  $p$  variabile:

$$\mathbf{G} = \sum_{i=1}^n \mathbf{p}_i \mathbf{x}_i = \sum_{i=1}^n \mathbf{p}_i \begin{bmatrix} x_i^1 \\ \vdots \\ x_i^p \end{bmatrix} = \begin{bmatrix} \bar{x}_1 \\ \vdots \\ \bar{x}_p \end{bmatrix} = \mathbf{x}^c \mathbf{D}_1.$$

Pentru a defini distanța dintre două unități înzestram spațiul cu o metrică specificată de matricea  $\mathbf{M}$ , simetrică și pozitiv definită. Produsul scalar al doi vectori se exprimă prin  $\langle \mathbf{x}_i, \mathbf{x}_j \rangle = \mathbf{x}_i' \mathbf{M} \mathbf{x}_j$ , iar distanța dintre două unități este dată de :

$$d^2(\mathbf{x}_i, \mathbf{x}_j) = \|\mathbf{x}_i - \mathbf{x}_j\|_{\mathbf{M}}^2 = (\mathbf{x}_i - \mathbf{x}_j)' \mathbf{M} (\mathbf{x}_i - \mathbf{x}_j).$$

În practică, cele mai utilizate metrici sunt:

- $\mathbf{M} = \mathbf{I}$ , situație în care se utilizează produsul scalar obișnuit;
- $\mathbf{M} = \text{diag}(1/s_j^2)$ ,  $j = \overline{1, p}$ , ceea ce revine la a diviza valorile observate ale variabilelor prin abaterile lor standard. Avantajul acestei metrici este că distanța dintre două

- unități nu depinde de unitatea de măsură, variabilele fiind echivalente ca importanță, independent de amplitudinea sau gradul de dispersare a valorilor.

Se definește astfel inerția totală a norului de puncte prin media ponderată a patratelor distanțelor punctelor față de centrul lor de greutate:

$$\mathbf{I}_g = \sum_1^n \mathbf{p}_i (\mathbf{x}_i - \mathbf{g})' \mathbf{M} (\mathbf{x}_i - \mathbf{g}).$$

În general în raport cu un punct oarecare  $\mathbf{h}$  inerția  $\mathbf{I}_g$  se definește prin:

$$\mathbf{I}_h = \sum_1^n \mathbf{p}_i (\mathbf{x}_i - \mathbf{h})' \mathbf{M} (\mathbf{x}_i - \mathbf{h}).$$

Între  $\mathbf{I}_g$  și  $\mathbf{I}_h$  există relația:

$$\mathbf{I}_h = \mathbf{I}_g + (\mathbf{g} - \mathbf{h})' \mathbf{M} (\mathbf{g} - \mathbf{h}).$$

Din ultima relație obținem că inerția în raport cu centrul de greutate este minimă.

În cazul în care  $\mathbf{g} = \mathbf{0}$ , atunci

$$\mathbf{I}_g = \sum_1^n \mathbf{p}_i (\mathbf{x}_i)' \mathbf{M} (\mathbf{x}_i).$$

### Bibliografie

- Benzecri, JP., *Histoire et prehistoire de l'analyse des données*, Dunod, Paris, 1983.
- Isaic-Maniu, Al.; Mitrut, C.; Voineagu, V., *Statistica pentru managementul afacerilor*, Ediția a II-a, Editura Economica, București, 1999.
- Saporta, G., *Probabilités, analyse des données et statistique*, Technip, Paris, 1990.
- Spircu, L, Calciu, M., Spircu, T., *Analiza datelor de marketing*, ALL, București, 1994.