

Cursul 10

Covarianța și coeficientul de corelație

Știind să identificăm variabile aleatoare independente, ne întrebăm în mod natural cum caracterizăm intensitatea legăturii dintre două variabile ce nu sunt independente. Această intensitate este măsurată de covarianța, respectiv coeficientul lor de corelație.

10.1 Covarianța a două variabile aleatoare

Definiția 10.1.1 *Covarianța* variabilelor aleatoare X și Y , ce au mediile $m_X = M(X)$, $m_Y = M(Y)$ finite, este definită prin

$$\text{cov}(X, Y) = M((X - m_X)(Y - m_Y)). \quad (10.1)$$

În cele ce urmează vom presupune că toți indicatorii (media, dispersia, covarianța) ce apar în formule există și sunt numere reale (finite).

Observația 10.1.1 Covarianța unei variabile cu ea însăși coincide cu dispersia sa, căci

$$\text{cov}(X, X) = M((X - m_X)(X - m_X)) = M((X - m_X)^2) = \sigma^2(X).$$

Propoziția 10.1.1 *Covarianța variabilelor aleatoare X, Y se poate calcula astfel:*

$$\text{cov}(X, Y) = M(XY) - M(X) M(Y). \quad (10.2)$$

Demonstrație: Efectuând produsul

$$(X - m_X)(Y - m_Y) = XY - m_X Y - m_Y X + m_X m_Y$$

și aplicând proprietățile mediei, obținem

$$\begin{aligned} \text{cov}(X, Y) &= M(XY) - m_X M(Y) - m_Y M(X) + m_X m_Y \\ &= M(XY) - M(X) M(Y). \end{aligned}$$

□

Definiția 10.1.2 Două variabile aleatoare X, Y ce au covarianța zero se numesc *variabile aleatoare necorelate*.

Deoarece am introdus covarianța ca o măsură a intensității dependenței dintre două variabile aleatoare, este natural să ne așteptăm ca două variabile aleatoare independente să aibă covarianța zero, adică să fie necorelate.

Propoziția 10.1.2 Dacă X, Y sunt independente, atunci X și Y sunt necorelate, adică dacă X, Y sunt independente, atunci $cov(X, Y) = 0$.

Demonstrație: Rezultă din faptul că dacă X și Y sunt variabile aleatoare independente, atunci $M(XY) = M(X)M(Y)$. \square

Observația 10.1.2 Reciproca propoziției de mai sus nu este adevărată. Două variabile aleatoare necorelate nu sunt în mod necesar independente, adică dacă $cov(X, Y) = 0$, **nu rezultă** în general că X și Y sunt independente.

Folosind definiția covarianței sau relația (10.2) și proprietățile operatorului mediei, au loc următoarele reguli de calcul:

- 1) $cov(X, X) = \sigma^2(X)$;
- 2) $cov(X, Y) = cov(Y, X)$;
- 3) $cov(aX, Y) = a cov(X, Y)$, $a \in \mathbb{R}$;
- 4) $cov(X + Y, Z) = cov(X, Z) + cov(Y, Z)$;
- 5) $\sigma^2(X + Y) = \sigma^2(X) + \sigma^2(Y) + 2cov(X, Y)$;

Dacă în plus X și Y sunt variabile aleatoare independente, atunci dispersia sumei lor este egală cu suma dispersiilor, adică

$$\sigma^2(X + Y) = \sigma^2(X) + \sigma^2(Y).$$

Demonstrație: O să arătăm doar relația 5), celelalte fiind evidente. Fie m_X, m_Y mediile celor două variabile aleatoare. Conform definiției dispersiei, avem:

$$\begin{aligned} \sigma^2(X + Y) &= M([(X + Y) - M(X + Y)]^2) \\ &= M([(X - m_X) + (Y - m_Y)]^2) \\ &= M((X - m_X)^2 + (Y - m_Y)^2 + 2(X - m_X)(Y - m_Y)) \\ &= M((X - m_X)^2) + M((Y - m_Y)^2) + 2M((X - m_X)(Y - m_Y)) \\ &= \sigma^2(X) + \sigma^2(Y) + 2cov(X, Y). \end{aligned}$$

\square

Ținând seama că $\sigma^2(aX) = a^2\sigma^2(X)$, $\forall a \in \mathbb{R}$, al doilea rezultat din relația 5) de mai sus se poate generaliza astfel:

Propoziția 10.1.3 *Dacă X_1, X_2, \dots, X_n sunt variabile aleatoare independente, atunci are loc:*

$$\sigma^2(a_1X_1 + a_2X_2 + \dots + a_nX_n) = a_1^2\sigma^2(X_1) + a_2^2\sigma^2(X_2) + \dots + a_n^2\sigma^2(X_n), \quad (10.3)$$

pentru orice $a_1, a_2, \dots, a_n \in \mathbb{R}$.

10.2 Coeficientul de corelație a două variabile aleatoare

Deoarece covarianța a două variabile aleatoare este un număr real oarecare, deci dificil de interpretat datorită nemărginirii mulțimii valorilor posibile, definim o altă măsură a dependenței lor, care ia valori într-un interval mărginit.

Definiția 10.2.1 *Coeficientul de corelație a două variabile aleatoare X și Y , de abateri standard nenule, este un număr real, notat cu $\rho(X, Y)$, definit prin*

$$\rho(X, Y) = \frac{\text{cov}(X, Y)}{\sigma_X \sigma_Y}, \quad (10.4)$$

unde σ_X, σ_Y sunt abaterile standard ale variabilelor aleatoare X , respectiv Y .

Observația 10.2.1 Coeficientul de corelație a două variabile aleatoare X și Y este, de fapt, covarianța standardizată a variabilelor X și Y , adică a variabilelor aleatoare

$$Z_1 = \frac{X - m_X}{\sigma_X}, \quad Z_2 = \frac{Y - m_Y}{\sigma_Y}.$$

Într-adevăr, ținând seama că $M(Z_1) = M(Z_2) = 0$, avem

$$\begin{aligned} \text{cov}(Z_1, Z_2) &= M(Z_1 Z_2) = M\left(\frac{X - m_X}{\sigma_X} \cdot \frac{Y - m_Y}{\sigma_Y}\right) = \frac{1}{\sigma_X \sigma_Y} M((X - m_X)(Y - m_Y)) \\ &= \frac{\text{cov}(X, Y)}{\sigma_X \sigma_Y} = \rho(X, Y). \end{aligned}$$

Propoziția 10.2.1 *Coeficientul de corelație a două variabile aleatoare X, Y are valoarea absolută subunitară, adică*

$$\rho(X, Y) \in [-1, 1]. \quad (10.5)$$

Demonstrație: Are loc

$$\sigma^2(tX + Y) \geq 0, \forall t \in \mathbb{R}.$$

Dar

$$\begin{aligned} \sigma^2(tX + Y) &= \sigma^2(tX) + \sigma^2(Y) + 2\text{cov}(tX, Y) \\ &= t^2\sigma^2(X) + 2t\text{cov}(X, Y) + \sigma^2(Y), \end{aligned}$$

deci $\sigma^2(tX + Y) \geq 0$ pentru orice $t \in \mathbb{R}$ dacă și numai dacă discriminantul

$$\Delta = 4[\text{cov}(X, Y)]^2 - 4\sigma^2(X)\sigma^2(Y) \leq 0,$$

ceea ce este echivalent cu

$$\left| \frac{\text{cov}(X, Y)}{\sigma(X)\sigma(Y)} \right| \leq 1,$$

adică ceea ce trebuia demonstrat. □

Observăm că **pentru două variabile aleatoare independente, coeficientul de corelație este 0**. Este natural să ne întrebăm în ce caz coeficientul de corelație a două variabile aleatoare ia valorile extreme ± 1 . Răspunsul este dat de următorul rezultat:

Propoziția 10.2.2 *Dacă între variabilele aleatoare X și Y există o relație liniară de forma*

$$Y = aX + b, \quad a, b \in \mathbb{R}, \quad a \neq 0,$$

atunci coeficientul de corelație al variabilelor aleatoare X și Y este ± 1 și anume:

$$\rho(X, Y) = \begin{cases} -1, & \text{dacă } a < 0, \\ 1, & \text{dacă } a > 0. \end{cases}$$

Reciproc, dacă modulul coeficientului de corelație a două variabile aleatoare X, Y este 1, atunci între ele există o relație liniară, $Y = aX + b$, $a \neq 0$.

Demonstrație: Dacă $Y = aX + b$, atunci

$$\rho(X, Y) = \frac{\text{cov}(X, aX + b)}{\sqrt{\sigma^2(X)\sigma^2(aX + b)}}. \quad (10.6)$$

Dar $\text{cov}(X, aX + b) = a\text{cov}(X, X) + \text{cov}(X, b)$, iar

$$\text{cov}(X, b) = M(bX) - M(b)M(X) = bM(X) - bM(X) = 0,$$

deci $\text{cov}(X, aX + b) = a\sigma^2(X)$. Pe de altă parte, $\sigma^2(aX + b) = a^2\sigma^2(X)$. Astfel, rezultă

$$\rho(X, Y) = \frac{a\sigma^2(X)}{\sqrt{\sigma^2(X)a^2\sigma^2(X)}} = \frac{a}{|a|}. \quad (10.7)$$

Prin urmare, pentru $a > 0$, $\rho(X, Y) = 1$, iar pentru $a < 0$, $\rho(X, Y) = -1$.

Reciproc, considerăm funcția $g(a, b) = M((Y - aX - b)^2)$ (variabilele aleatoare X și Y sunt fixate). Determinăm $a, b \in \mathbb{R}$ astfel încât g să fie minimă, adică determinăm parametrii a, b astfel încât media abaterii la pătrat a lui Y față de o funcție de gradul întâi după X să fie minimă.

Pentru a arăta că funcția g are un punct de minim local o descompunem astfel:

$$\begin{aligned} g(a, b) &= M((Y - aX - b)^2) \\ &= M(Y^2 + a^2X^2 + b^2 - 2aXY - 2bY + 2abX) \\ &= M(Y^2) + a^2M(X^2) + b^2 - 2aM(XY) - 2bM(Y) + 2abM(X). \end{aligned}$$

Evident că $M(Y^2), M(X^2), M(XY), M(Y), M(X)$ sunt constante. Rezolvând sistemul

$$\begin{cases} \frac{\partial g}{\partial a} = 0, \\ \frac{\partial g}{\partial b} = 0, \end{cases}$$

obținem punctul critic (a_0, b_0) , unde

$$a_0 = \frac{\text{cov}(X, Y)}{\sigma^2(X)}, \quad b_0 = M(Y) - a_0M(X).$$

Se arată că (a_0, b_0) este punct de minim local al lui g , adică matricea

$$\begin{bmatrix} \frac{\partial^2 g}{\partial a^2}(a_0, b_0) & \frac{\partial^2 g}{\partial a \partial b}(a_0, b_0) \\ \frac{\partial^2 g}{\partial a \partial b}(a_0, b_0) & \frac{\partial^2 g}{\partial b^2}(a_0, b_0) \end{bmatrix} = \begin{bmatrix} 2M(X^2) & 2M(X) \\ 2M(X) & 2 \end{bmatrix}$$

este pozitiv definită. Într-adevăr, avem $\Delta_1 = 2M(X^2) > 0$ (dacă $M(X^2) = 0$, atunci, cum $M(X^2) \geq M(X)^2$, rezultă că $M(X) = 0$, deci $\sigma^2(X) = 0$, ceea ce este fals), respectiv

$$\Delta_2 = \begin{vmatrix} 2M(X^2) & 2M(X) \\ 2M(X) & 2 \end{vmatrix} = 4[M(X^2) - M(X)^2] = 4\sigma^2(X) > 0.$$

Punctul (a_0, b_0) este chiar punct de minim global al funcției g . Calculând $g(a_0, b_0)$, obținem valoarea minimă a funcției g . Astfel,

$$g(a_0, b_0) = \min_{a, b} g(a, b) = \min_{a, b} M((Y - aX - b)^2) = \sigma^2(Y)(1 - \rho^2(X, Y)).$$

Dar cum $\rho(X, Y) = \pm 1$, obținem

$$g(a_0, b_0) = M((Y - a_0X - b_0)^2) = 0.$$

Deoarece $(Y - a_0X - b_0)^2 \geq 0$, media sa este zero dacă și numai dacă $Y - a_0X - b_0 = 0$ sau, echivalent, $Y = a_0X + b_0$. \square

În concluzie:

- când coeficientul de corelație a două variabile aleatoare este apropiat de zero, variabilele sunt slab corelate (intensitatea legăturii dintre ele este redusă);
- când valoarea absolută a coeficientului de corelație este apropiată de 1, relația dintre variabilele aleatoare este "aproape liniară", adică valorile (x, y) ale vectorului aleator (X, Y) sunt ușor dispersate în jurul unei drepte de ecuație $y = ax + b$.

Exemplul 1. Fie X o variabilă aleatoare ce are media $M(X) = 3$ și dispersia $\sigma^2(X) = 1$, iar $Y = -2X + 5$. Să se calculeze covarianța și coeficientul de corelație pentru variabilele X, Y .

Rezolvare: Deoarece între X și Y există o relație liniară de forma $Y = aX + b$ cu $a < 0$, coeficientul de corelație este

$$\rho(X, Y) = -1.$$

Dar cum

$$\rho(X, Y) = \frac{\text{cov}(X, Y)}{\sigma(X) \sigma(Y)},$$

calculând $\sigma^2(Y) = \sigma^2(-2X + 5) = 4\sigma^2(X) = 4$, rezultă că $-1 = \frac{\text{cov}(X, Y)}{2}$, deci $\text{cov}(X, Y) = -2$.

□

10.3 Matricea de covarianță a unui vector aleator

Fie $\mathbf{X} = (X_1, X_2, \dots, X_n)^T$ un vector aleator ale cărui componente nu sunt neapărat variabile aleatoare independente. Pentru a cuantifica intensitatea legăturii dintre două câte două componente, se asociază vectorului \mathbf{X} *matricea de covarianță*. Înainte de a o defini, notăm prin

$$M(\mathbf{X}) = (M(X_1), M(X_2), \dots, M(X_n))^T$$

vectorul ce are drept coordonate mediile corespunzătoare variabilelor X_i , $i = \overline{1, n}$, sau $\mathbf{m} = (m_1, m_2, \dots, m_n)^T$, $m_i = M(X_i)$.

Definiția 10.3.1 *Matricea de covarianță* a vectorului aleator \mathbf{X} este matricea notată cu Σ , ale cărei elemente sunt $\sigma_{ij} = \text{cov}(X_i, X_j)$, $i, j = \overline{1, n}$.

Remarcăm că $\sigma_{ii} = \text{cov}(X_i, X_i) = \sigma^2(X_i)$. Cu alte cuvinte, elementele de pe diagonala principală ale matricei de covarianță a unui vector aleator sunt dispersiile componentelor vectorului.

Notând cu $\mathbf{Y} = \mathbf{X} - \mathbf{m} = (X_1 - m_1, X_2 - m_2, \dots, X_n - m_n)^T$, matricea de covarianță se poate exprima astfel:

$$\Sigma = M(\mathbf{Y}\mathbf{Y}^T),$$

unde

$$\mathbf{Y}\mathbf{Y}^T = \begin{bmatrix} Y_1 \\ Y_2 \\ \vdots \\ Y_n \end{bmatrix} \begin{bmatrix} Y_1 & Y_2 & \dots & Y_n \end{bmatrix} = \begin{bmatrix} Y_1Y_1 & Y_1Y_2 & \dots & Y_1Y_n \\ Y_2Y_1 & Y_2Y_2 & \dots & Y_2Y_n \\ \vdots & \vdots & \dots & \vdots \\ Y_nY_1 & Y_nY_2 & \dots & Y_nY_n \end{bmatrix},$$

iar media matricei $\mathbf{Y}\mathbf{Y}^T$ este matricea mediilor elementelor sale, $M(Y_iY_j) = \text{cov}(X_i, X_j)$, $i, j = \overline{1, n}$.

Propoziția 10.3.1 *Matricea de covarianță a unui vector aleator $\mathbf{X} = (X_i)$, $i = \overline{1, n}$, este simetrică și semipozitiv definită.*

În recunoasterea formelor se studiază intensitatea legăturii dintre un număr imens de variabile X_1, X_2, \dots, X_n . Matricea de covarianță Σ este supusă analizei PCA (Principal Component Analysis) din care se extrage informație valoroasă despre corelațiile dintre variabile. Informația se extrage din descompunerea $\Sigma = QDQ^T$ a matricei simetrice Σ , unde D este matricea diagonală a valorilor proprii $\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_n \geq 0$ ale lui Σ , iar Q este matricea ortogonală (i.e. $Q^TQ = I_n$) ce are pe coloane coordonatele vectorilor proprii ortonormați u_1, u_2, \dots, u_n , corespunzători valorilor proprii, $\Sigma u_i = \lambda_i u_i$, $i = \overline{1, n}$.