

2.6 Teoretické základy korelace

Nechť X a Y jsou náhodné veličiny s konečnými druhými momenty a s kladnými rozptyly. Závislost těchto veličin na sobě se často měří pomocí *korelačního koeficientu*

$$\rho = \frac{\text{cov}(X, Y)}{\sqrt{(\text{var } X)(\text{var } Y)}}.$$

Někdy místo ρ píšeme $\rho_{X, Y}$, abychom vyznačili, o které veličiny se jedná. Je zřejmé, že $\rho_{X, Y} = \rho_{Y, X}$. Nechť dále a, b, c, d jsou taková reálná čísla, že $ac \neq 0$. Snadno se dá dokázat, že platí

$$\rho_{aX+b, cY+d} = \begin{cases} \rho_{X, Y} & \text{pro } ac > 0, \\ -\rho_{X, Y} & \text{pro } ac < 0. \end{cases}$$

Při lineární transformaci se tedy korelační koeficient buď nezmění vůbec, nebo pouze změní znaménko.

Věta 2.16 *Pro korelační koeficient platí $-1 \leq \rho_{X, Y} \leq 1$. Rovnost $\rho_{X, Y} = 1$ platí právě tehdy, je-li $Y = a + bX$ s pravděpodobností 1, přičemž $b > 0$. Analogicky rovnost $\rho_{X, Y} = -1$ platí právě tehdy, je-li $Y = a + bX$ s pravděpodobností 1, přičemž $b < 0$.*

Důkaz. Schwarzova nerovnost

$$|E(X - EX)(Y - EY)| \leq \sqrt{E(X - EX)^2 E(Y - EY)^2}$$

má v našem označení podobu $|\text{cov}(X, Y)| \leq \sqrt{(\text{var } X)(\text{var } Y)}$. Z toho plyne, že $-1 \leq \rho_{X, Y} \leq 1$. Rovnosti ve Schwarzově nerovnosti je dosahováno právě tehdy, platí-li buď $X - EX = 0$ nebo $Y - EY = 0$ skoro jistě (což v našem případě nepřichází v úvahu vzhledem k předpokladu $\text{var } X > 0, \text{var } Y > 0$), nebo když platí $Y - EY = b(X - EX)$ skoro jistě pro nějaké $b \neq 0$. Výpočtem se ověří, že v tomto posledním případě je $\rho_{X, Y} = 1$ pro $b > 0$ a $\rho_{X, Y} = -1$ pro $b < 0$. \square

Při práci s korelačním koeficientem bývá užitečná následující interpretace. Uvažujme prostor náhodných veličin s konečnými druhými momenty, přičemž tyto veličiny jsou definovány na stejném pravděpodobnostním prostoru. Vytvořme třídy ekvivalence těchto veličin tak, že X a Y prohlásíme za ekvivalentní, když si jsou veličiny $X - EX$ a $Y - EY$ rovny skoro všude. Definujme na třídách ekvivalence skalární součin předpisem $(X, Y) = \text{cov}(X, Y)$ pro libovolné reprezentanty těchto tříd, jak je obvyklé ve funkcionální analýze. Prostor těchto tříd je Hilbertův prostor a označme ho \mathcal{H} . Norma veličiny X je $\|X\| = \sqrt{(X, X)}$. V tomto označení lze korelační koeficient ρ veličin $X \neq 0$ a $Y \neq 0$ vyjádřit ve tvaru

$$\rho_{X, Y} = \frac{(X, Y)}{\|X\| \cdot \|Y\|}.$$

Na základě analogie se vzorcem známým z geometrie můžeme tedy korelační koeficient interpretovat jako kosinus úhlu mezi veličinami X a Y , nebo názorněji jako kosinus úhlu mezi veličinami $X - EX$ a $Y - EY$.

Mějme náhodné vektory $\mathbf{X} = (X_1, \dots, X_n)'$ a $\mathbf{Y} = (Y_1, \dots, Y_m)'$ s konečnými druhými momenty. Nechť všechny složky těchto dvou vektorů mají kladné rozptyly. Matici typu $n \times m$, jejíž (i, j) -tý prvek je roven korelačnímu koeficientu veličin X_i a Y_j , nazveme *korelační maticí vektorů* \mathbf{X} a \mathbf{Y} a označíme ji $\text{cor}(\mathbf{X}, \mathbf{Y})$. Matice $\text{cor}(\mathbf{X}, \mathbf{X})$ se nazývá *korelační matice vektoru* \mathbf{X} a značí se $\text{cor } \mathbf{X}$.

Lemma 2.17 *Nechť $\mathbf{X} = (X_1, \dots, X_n)'$ je náhodný vektor s konečnými druhými momenty, jehož všechny složky mají kladné rozptyly. Označme $\mathbf{V} = \text{var } \mathbf{X}$, $\mathbf{P} = \text{cor } \mathbf{X}$, $\sigma_i = \sqrt{\text{var } X_i}$ pro $i = 1, \dots, n$ a*

$$\mathbf{D} = \text{Diag}\{\sigma_1, \dots, \sigma_n\} = \begin{pmatrix} \sigma_1 & 0 & \dots & 0 \\ \dots & \dots & \dots & \dots \\ 0 & 0 & \dots & \sigma_n \end{pmatrix}.$$

Pak $\mathbf{P} = \mathbf{D}^{-1} \mathbf{V} \mathbf{D}^{-1}$.

Důkaz je zřejmý. \square

Lemma 2.18 *Budiž \mathbf{V} symetrická pozitivně definitní matice typu $n \times n$. Pak pro libovolné n -rozměrné vektory \mathbf{b} a \mathbf{c} platí $(\mathbf{b}' \mathbf{V} \mathbf{c})^2 \leq (\mathbf{b}' \mathbf{V} \mathbf{b})(\mathbf{c}' \mathbf{V} \mathbf{c})$.*

Důkaz. Je známo, že za uvedených předpokladů existuje taková symetrická pozitivně semidefinitní matice $\mathbf{V}^{\frac{1}{2}}$, pro kterou platí $\mathbf{V}^{\frac{1}{2}} \mathbf{V}^{\frac{1}{2}} = \mathbf{V}$. Tato matice $\mathbf{V}^{\frac{1}{2}}$ se nazývá *odmocninová*. Ze Schwarzovy nerovnosti plyne

$$\begin{aligned} (\mathbf{b}' \mathbf{V} \mathbf{c})^2 &= (\mathbf{b}' \mathbf{V}^{\frac{1}{2}} \mathbf{V}^{\frac{1}{2}} \mathbf{c})^2 = [(\mathbf{V}^{\frac{1}{2}} \mathbf{b})' (\mathbf{V}^{\frac{1}{2}} \mathbf{c})]^2 \\ &\leq [(\mathbf{V}^{\frac{1}{2}} \mathbf{b})' (\mathbf{V}^{\frac{1}{2}} \mathbf{b})][(\mathbf{V}^{\frac{1}{2}} \mathbf{c})' (\mathbf{V}^{\frac{1}{2}} \mathbf{c})] = (\mathbf{b}' \mathbf{V} \mathbf{b})(\mathbf{c}' \mathbf{V} \mathbf{c}). \quad \square \end{aligned}$$

Mějme náhodnou veličinu Y a náhodný vektor $\mathbf{X} = (X_1, \dots, X_n)'$ s konečnými druhými momenty podobně jako v případě lineární regrese. Nechť $\text{var } Y > 0$ a nechť $\mathbf{V} = \text{var } \mathbf{X}$ je regulární matice. Závislost mezi Y a celým vektorem \mathbf{X} měříme pomocí koeficientu mnohonásobné korelace $\rho_{Y, \mathbf{X}}$, což je korelační koeficient mezi veličinou Y a její nejlepší lineární aproximací $\hat{Y} = \alpha + \beta' \mathbf{X}$, kde α a β jsou uvedeny ve větě 2.15. Je-li $\beta = \mathbf{0}$, definuje se $\rho_{Y, \mathbf{X}} = 0$. Někdy se místo $\rho_{Y, \mathbf{X}}$ píše podrobněji $\rho_{Y, X_1, \dots, X_n}$. Protože platí $\rho_{Y, \mathbf{X}} = \rho_{Y, \alpha + \beta' \mathbf{X}} = \rho_{Y, \beta' \mathbf{X}}$ a

$$\text{cov}(Y, \beta' \mathbf{X}) = \text{cov}(Y, \mathbf{X}) \beta = \text{cov}(Y, \mathbf{X}) \mathbf{V}^{-1} \text{cov}(\mathbf{X}, Y) \geq 0,$$

je koeficient mnohonásobné korelace vždy nezáporný.

Věta 2.19 *Označme $\mathbf{P} = \text{cor } \mathbf{X}$. Pak platí*

$$\rho_{Y, \mathbf{X}}^2 = \text{cor}(Y, \mathbf{X}) \mathbf{P}^{-1} \text{cor}(\mathbf{X}, Y). \quad (2.8)$$

Důkaz. Pro $\beta = 0$ je tvrzení $\rho_{Y, \alpha + \beta' \mathbf{X}} = \rho_{Y, \beta' \mathbf{X}}$. Protože

Dosadíme-li $\text{cov}(Y, \mathbf{X}) = \beta$

$$\beta = \mathbf{V}^{-1} \text{cov}(\mathbf{X}, Y) =$$

a proto

$$\begin{aligned} \beta' \mathbf{V} \beta &= \sigma_Y^2 \\ &= \sigma_Y^2 \end{aligned}$$

Dosazením do (2.9) se dost

Koeficient mnohonásobné korelace $\rho_{Y, \mathbf{X}}$ ze vzorců (2.7) a (2.10) dost

$$\sigma_{Y, \mathbf{X}}^2 =$$

Často se počítá koeficient korelace $\rho_{Y, \mathbf{X}}$ mezi $(X_1, X_2)'$, který má jen dva argumenty. Koeficient mezi X_i a X_j ($i, j = 1, 2$) v případě z věty 2.19 plyne

Korelační matice \mathbf{P} vektorů $(X_1, X_2)'$ determinant je nezáporný.

$$|\mathbf{P}| = \begin{vmatrix} 1 & \rho_{01} \\ \rho_{01} & 1 \\ \rho_{02} & \rho_{12} \end{vmatrix}$$

Předpokládejme nyní, že $\rho_{01} = 1$. Vidíme, že v případě $\rho_{01} = 1$ $|\rho_{02}|$ blízké 1.

Další zdůvodnění koeficientu korelace $\rho_{Y, \mathbf{X}}$ následující tvrzení.

Důkaz. Pro $\beta = 0$ je tvrzení zřejmé. Nechť tedy $\beta \neq 0$. Označme $V = \text{var } X$. Máme $\rho_{Y, \alpha + \beta' X} = \rho_{Y, \beta' X}$. Protože $\text{cov}(Y, \beta' X) = \text{cov}(Y, X)\beta$ a $\text{var } \beta' X = \beta' V \beta$, je

$$\rho_{Y, X}^2 = \frac{[\text{cov}(Y, X)\beta]^2}{\text{var } Y \cdot \beta' V \beta}. \quad (2.9)$$

Dosadíme-li $\text{cov}(Y, X) = \beta' V$, dostaneme odtud

$$\rho_{Y, X}^2 = \frac{\beta' V \beta}{\text{var } Y}. \quad (2.10)$$

Užitím lemmatu 2.17 dostaneme

$$\beta = V^{-1} \text{cov}(X, Y) = D^{-1} P^{-1} D^{-1} \text{cov}(X, Y) = D^{-1} P^{-1} \text{cor}(X, Y) \sigma_Y,$$

a proto

$$\begin{aligned} \beta' V \beta &= \sigma_Y^2 \text{cor}(Y, X) P^{-1} D^{-1} V D^{-1} P^{-1} \text{cor}(X, Y) \\ &= \sigma_Y^2 \text{cor}(Y, X) P^{-1} \text{cor}(X, Y). \end{aligned}$$

Dosazením do (2.9) se dostane tvrzení věty. \square

Koeficient mnohonásobné korelace těsně souvisí s reziduálním rozptylem $\sigma_{Y, X}^2$. Ze vzorců (2.7) a (2.10) dostáváme

$$\sigma_{Y, X}^2 = \sigma_Y^2 (1 - \rho_{Y, X}^2), \quad \rho_{Y, X}^2 = 1 - \frac{\sigma_{Y, X}^2}{\sigma_Y^2}.$$

Často se počítá koeficient mnohonásobné korelace mezi Y a vektorem $X = (X_1, X_2)'$, který má jen dvě složky. Položme $X_0 = Y$ a označme ρ_{ij} korelační koeficient mezi X_i a X_j ($i, j = 0, 1, 2$). Dále označme $\rho_{0.12} = \rho_{Y, X_1, X_2}$. V tomto případě z věty 2.19 plyne

$$\rho_{0.12}^2 = \frac{\rho_{01}^2 + \rho_{02}^2 - 2\rho_{01}\rho_{02}\rho_{12}}{1 - \rho_{12}^2}.$$

Korelační matice P vektoru $(X_0, X_1, X_2)'$ je pozitivně semidefinitní, a proto její determinant je nezáporný. Odtud dostáváme

$$|P| = \begin{vmatrix} 1 & \rho_{01} & \rho_{02} \\ \rho_{01} & 1 & \rho_{12} \\ \rho_{02} & \rho_{12} & 1 \end{vmatrix} = 1 + 2\rho_{01}\rho_{02}\rho_{12} - \rho_{01}^2 - \rho_{02}^2 - \rho_{12}^2 \geq 0.$$

Předpokládejme nyní, že $\rho_{01} = 0$. Pak $|P| = 1 - \rho_{02}^2 - \rho_{12}^2 \geq 0$, takže $\rho_{02}^2 / (1 - \rho_{12}^2) \leq 1$. Vidíme, že v případě $\rho_{01} = 0$ může být $\rho_{0.12}^2 = \rho_{02}^2 / (1 - \rho_{12}^2)$ blízké 1, pokud je $|\rho_{02}|$ blízké 1.

Další zdůvodnění koeficientu mnohonásobné korelace jako míry závislosti přináší následující tvrzení.

(2.8)

