

Analýza hlavních komponent

Ekonometrie

Jiří Neubauer

Katedra ekonometrie FVL UO Brno
kancelář 69a, tel. 973 442029
email: Jiri.Neubauer@unob.cz

Metoda hlavních komponent

Metoda hlavních komponent (PCA – „Principal component analysis“) je metodou, pomocí níž lze vysvětlit variančně kovarianční strukturu množiny proměnných pomocí několika lineárních kombinací těchto proměnných. K základním úkolům metody patří redukce dat (snížení dimenze) a interpretace. Máme-li p proměnných, pomocí nichž lze vysvětlit celkovou variabilitu, v řadě případů je možné většinu této variability popsat pomocí mnohem menšího počtu k hlavních komponent. Těchto k hlavních komponent potom může nahradit původních p proměnných, tedy původní datový soubor obsahující n měření p proměnných je redukován na n měření k hlavních komponent.

Metoda hlavních komponent

- Algebraicky jsou hlavní komponenty jisté lineární kombinace původních p náhodných veličin X_1, X_2, \dots, X_p .
- Z geometrického hlediska reprezentují tyto lineární kombinace výběr nového systému souřadnic získaného rotací původního souřadnicového systému s osami X_1, X_2, \dots, X_p . Nové osy reprezentují směry s největší variabilitou a nabízí jednodušší a více parsimonní popis kovarianční struktury.

Příklad

The weekly rates of return for five stocks (JPMorgan, Citibank, WellsFargo, RoyalDutchShell, ExxonMobil) listed on the New York Stock Exchange were determined for the period January 2004 through December 2005. The weekly rates of return are defined

$$\frac{\text{current week closing price} - \text{previous week closing price}}{\text{previous week closing price}},$$

adjusted for stock splits and dividends. The observations in 103 successive weeks appear to be independently distributed, but the rates of return across stock are correlated, because as one expects, stocks tend to move together in response to general economic conditions.

Johnson, R. and Wichern, D. *Applied Multivariate Statistical Analysis*. 6th ed. Pearson 2014.

Metoda hlavních komponent

Hlavní komponenty závisí výhradně na varianční matici Σ (nebo na korelační matici ρ) proměnných X_1, X_2, \dots, X_p . Nechť náhodný vektor $\mathbf{X} = (X_1, X_2, \dots, X_p)'$ má varianční matici Σ s vlastními čísly $\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_p \geq 0$. Uvažujme lineární kombinace

$$Y_1 = \mathbf{a}'_1 \mathbf{X} = a_{11}X_1 + a_{12}X_2 + \dots + a_{1p}X_p$$

$$Y_2 = \mathbf{a}'_2 \mathbf{X} = a_{21}X_1 + a_{22}X_2 + \dots + a_{2p}X_p$$

$$\vdots$$

$$Y_p = \mathbf{a}'_p \mathbf{X} = a_{p1}X_1 + a_{p2}X_2 + \dots + a_{pp}X_p.$$

Pro proměnné Y_1, Y_2, \dots, Y_p dostáváme

$$D(Y_i) = \mathbf{a}'_i \Sigma \mathbf{a}_i, \quad i = 1, 2, \dots, p \quad (1)$$

$$C(Y_i, Y_k) = \mathbf{a}'_i \Sigma \mathbf{a}_k, \quad i, k = 1, 2, \dots, p \quad (2)$$

Hlavní komponenty jsou takové nekorelované kombinace Y_1, Y_2, \dots, Y_p , jejichž rozptyly (4) jsou co možná největší.

Metoda hlavních komponent

První hlavní komponenta je lineární kombinace s největším rozptylem, tedy maximalizující výraz $D(Y_1) = \mathbf{a}'_1 \Sigma \mathbf{a}_1$. Je zřejmé, že vynásobením \mathbf{a}_1 nějakou konstantou lze hodnotu rozptylu měnit. Je třeba přidat podmínku jednotkové délky vektoru \mathbf{a}_1 , tedy $\mathbf{a}'_1 \mathbf{a}_1 = 1$.

Druhá hlavní komponenta je lineární kombinace $\mathbf{a}'_2 \mathbf{X}$ jež maximalizuje výraz $\mathbf{a}'_2 \Sigma \mathbf{a}_2$ při platnosti $\mathbf{a}'_1 \mathbf{a}_1 = 1$ a $C(\mathbf{a}'_1 \mathbf{X}, \mathbf{a}'_2 \mathbf{X}) = 0$. Podobně se postupuje dále, což znamená, že i -tá hlavní komponenta je lineární kombinace $\mathbf{a}'_i \mathbf{X}$ jež maximalizuje výraz $\mathbf{a}'_i \Sigma \mathbf{a}_i$ při platnosti $\mathbf{a}'_i \mathbf{a}_i = 1$ a $C(\mathbf{a}'_i \mathbf{X}, \mathbf{a}'_k \mathbf{X}) = 0$ pro $k < i$.

Určení hlavních komponent souvisí s vlastními čísly a vektory matice Σ .

Pozn. Nechť $\mathbf{A} = (a_{ij})$ je matice řádu n . Číslo λ se nazývá vlastní nebo charakteristické číslo matice \mathbf{A} , jestliže existuje nenulový vektor \mathbf{u} tak, že $\mathbf{A}\mathbf{u} = \lambda\mathbf{u}$. Vektor \mathbf{u} se nazývá vlastní nebo charakteristický vektor příslušný k λ .

Rovnici $\mathbf{A}\mathbf{u} = \lambda\mathbf{u}$ lze přepsat do tvaru $(\mathbf{A} - \lambda\mathbf{I})\mathbf{u} = \mathbf{0}$. Tato soustava má netriviální řešení, právě když $|\mathbf{A} - \lambda\mathbf{I}| = 0$.

Příklad: Určete vlastní čísla a vlastní vektory matice

$$\mathbf{A} = \begin{pmatrix} 4 & -2 \\ 1 & 1 \end{pmatrix}$$

Metoda hlavních komponent

Mějme náhodný vektor $\mathbf{X} = (X_1, X_2, \dots, X_p)'$ s varianční maticí Σ a dvojicemi vlastních čísel a vlastních vektorů $(\lambda_1, \mathbf{e}_1), (\lambda_2, \mathbf{e}_2), \dots, (\lambda_p, \mathbf{e}_p)$ splňující $\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_p \geq 0$. Potom i -tá hlavní komponenta je dána vztahem

$$Y_i = \mathbf{e}_i \mathbf{X} = e_{i1}X_1 + e_{i2}X_2 + \dots + e_{ip}X_p, \quad i = 1, 2, \dots, p, \quad (3)$$

přičemž

$$D(Y_i) = \mathbf{e}_i' \Sigma \mathbf{e}_i = \lambda_i, \quad i = 1, 2, \dots, p, \quad (4)$$

$$C(Y_i, Y_k) = \mathbf{e}_i' \Sigma \mathbf{e}_k = 0, \quad i \neq k. \quad (5)$$

Označíme-li prvky matice Σ jako σ_{ij} , $i, j = 1, 2, \dots, p$, potom platí

$$\sigma_{11} + \sigma_{22} + \dots + \sigma_{pp} = \sum_{i=1}^p D(X_i) = \lambda_1 + \lambda_2 + \dots + \lambda_p = \sum_{i=1}^p D(Y_i).$$

Metoda hlavních komponent

Podíl k -té hlavní komponenty na celkové variabilitě je dán podílem

$$\frac{\lambda_k}{\sum_{i=1}^p \lambda_i}.$$

Pokud lze většinu celkové variability (např. 80–90, %) pro velké p popsat několika hlavními komponentami (např. jednou, dvěma nebo třemi), potom lze původní proměnné „nahradit“ těmito komponentami bez velké ztráty informace.

Kromě vlastních čísel matice Σ si pozornost zaslouží i koeficienty vlastních vektorů $\mathbf{e} = (e_{i1}, \dots, e_{ik}, \dots, e_{ip})'$. Velikost hodnoty e_{ik} určuje význam k -té proměnné v i -té hlavní komponentě. Hodnota e_{ik} je úměrná korelačnímu koeficientu mezi Y_i a X_k , platí

$$\rho_{Y_i, X_k} = \frac{e_{ik} \sqrt{\lambda_i}}{\sqrt{\sigma_{ii}}} \quad i, k = 1, 2, \dots, p.$$

Metoda hlavních komponent pro standardizované proměnné

Standardizované proměnné získáme následovně

$$Z_1 = \frac{X_1 - \mu_1}{\sqrt{\sigma_{11}}}$$

$$Z_2 = \frac{X_2 - \mu_2}{\sqrt{\sigma_{22}}}$$

⋮

$$Z_p = \frac{X_p - \mu_p}{\sqrt{\sigma_{pp}}}$$

v maticovém vyjádření

$$\mathbf{Z} = (\mathbf{V}^{1/2})^{-1}(\mathbf{X} - \boldsymbol{\mu}), \quad (6)$$

kde $\mathbf{V}^{1/2}$ je diagonální matice směrodatných odchylek. Je zřejmé, že $E(\mathbf{Z}) = \mathbf{0}$ a

$$\text{var}(\mathbf{Z}) = (\mathbf{V}^{1/2})^{-1}\boldsymbol{\Sigma}(\mathbf{V}^{1/2})^{-1} = \boldsymbol{\rho}.$$

Metoda hlavních komponent pro standardizované proměnné

Hlavní komponenty lze získat pomocí vlastních vektorů korelační matice ρ . Uvažujme dvojice vlastních čísel a vektorů $(\lambda_1, \mathbf{e}_1), (\lambda_2, \mathbf{e}_2), \dots, (\lambda_p, \mathbf{e}_p)$ korelační matice ρ , kde $\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_p \geq 0$. Pro proměnné $\mathbf{Z} = (Z_1, Z_2, \dots, Z_p)'$ s varianční maticí $\text{var}(\mathbf{Z}) = \rho$, jsou hlavní komponenty dána vztahem

$$Y_i = \mathbf{e}_i' \mathbf{Z} = \mathbf{e}_i' (\mathbf{V}^{1/2})^{-1} (\mathbf{X} - \boldsymbol{\mu}), \quad i = 1, 2, \dots, p$$

Dále platí

$$\sum_{i=1}^p D(Y_i) = \sum_{i=1}^p D(Z_i) = p$$

a

$$\rho_{Y_i, Z_k} = \mathbf{e}_{ik} \sqrt{\lambda_i}, \quad i, k = 1, 2, \dots, p.$$

Podíl k -té hlavní komponenty na celkové variabilitě je dán vztahem

$$\frac{\lambda_k}{p}, \quad k = 1, 2, \dots, p.$$

Příklad

Mějme kovarianční matici

$$\Sigma = \begin{pmatrix} 1 & 1,5 \\ 1,5 & 4 \end{pmatrix}.$$

Odpovídající korelační matice je

$$\rho = \begin{pmatrix} 1 & 0,75 \\ 0,75 & 1 \end{pmatrix}.$$

Vlastní čísla a vlastní vektory matice Σ jsou

$$\lambda_1 = 4,621, \quad \mathbf{e}_1 = (0,383; 0,924)'$$

$$\lambda_2 = 0,379, \quad \mathbf{e}_2 = (-0,924; 0,383)'$$

Podobně vlastní čísla a vlastní vektory matice ρ jsou

$$\lambda_1^* = 1 + \rho = 1,75, \quad \mathbf{e}_1^* = (0,707; 0,707)'$$

$$\lambda_2^* = 1 - \rho = 0,25, \quad \mathbf{e}_2^* = (-0,707; 0,707)'$$

Příklad

Hlavní komponenty odpovídající matici Σ potom jsou

$$Y_1 = 0,383X_1 + 0,924X_2$$

$$Y_2 = -0,924X_1 + 0,383X_2$$

a pro matici ρ

$$Y_1^* = 0,707Z_1 + 0,707Z_2 = 0,707 \left(\frac{X_1 - \mu_1}{1} \right) + 0,707 \left(\frac{X_2 - \mu_1}{2} \right)$$

$$Y_2^* = -0,707Z_1 + 0,707Z_2 = -0,707 \left(\frac{X_1 - \mu_1}{1} \right) + 0,707 \left(\frac{X_2 - \mu_1}{2} \right)$$

Rozptyl proměnné X_2 je větší než rozptyl X_1 , k první hlavní komponentě přispívá více. Navíc tato první hlavní komponenta vysvětluje podíl

$$\frac{\lambda_1}{\lambda_1 + \lambda_2} = \frac{4,621}{4,621 + 0,379} = 0,924$$

na celkovém rozptylu.

Příklad

Podíváme-li se na standardizované proměnné, pak příspěvky obou proměnných Z_1 a Z_2 k hlavním komponentám jsou stejné. Dále platí

$$\rho_{Y_1, Z_1} = e_{11}^* \sqrt{\lambda_1^*} = 0,707 \sqrt{1,75} = 0,935$$

$$\rho_{Y_1, Z_2} = e_{21}^* \sqrt{\lambda_1^*} = -0,707 \sqrt{1,75} = -0,935$$

První hlavní komponenta určená ze standardizovaných veličin vysvětluje podíl

$$\frac{\lambda_1}{p} = \frac{1,75}{2} = 0,875$$

na celkovém rozptylu. Z uvedeného příkladu je vidět, že standardizace ovlivňuje výsledky metody hlavních komponent, hlavní komponenty získané z varianční matice Σ se liší od těch získaných z korelační matice ρ .

Výběrové hlavní komponenty

Předpokládejme, že máme náhodný výběr x_1, x_2, \dots, x_n ze základního souboru se střední hodnotou μ a varianční maticí Σ . Určíme výběrový průměr \bar{x} , výběrovou varianční matici S a výběrovou korelační matici R .

Cílem je najít takové nekorelované lineární kombinace, které budou schopny popsat většinu variability ve výběru.

- První výběrová hlavní komponenta je taková lineární kombinace $a'_1 x_j$, která maximalizuje výběrový rozptyl $a'_1 x_j$ za podmínky $a'_1 a_1 = 1$.
- Druhá výběrová hlavní komponenta je lineární kombinace $a'_2 x_j$, která maximalizuje výběrový rozptyl $a'_2 x_j$ za podmínek $a'_2 a_2 = 1$ a nulové kovariance dvojice $(a'_1 x_j, a'_2 x_j)$.
- Podobným způsobem se postupuje dále, tudíž i -tá výběrová hlavní komponenta je lineární kombinace $a'_i x_j$ maximalizující výběrový rozptyl $a'_i x_j$ při splnění podmínek $a'_i a_i = 1$ a nulové kovariance párů $(a'_1 x_j, a'_k x_j)$, $k < i$.

Výpočet výběrových hlavních komponent je úzce spojen s vlastními čísly a vektory výběrové varianční matice S .

Výběrové hlavní komponenty

Je-li $\mathbf{S} = \{s_{ik}\}$ výběrový varianční matice typu $p \times p$ s vlastními čísly a vektory tvořící dvojice $(\hat{\lambda}_1, \hat{\mathbf{e}}_1), (\hat{\lambda}_2, \hat{\mathbf{e}}_2), \dots, (\hat{\lambda}_p, \hat{\mathbf{e}}_p)$, $\hat{\lambda}_1 \geq \hat{\lambda}_2 \geq \dots \geq \hat{\lambda}_p \geq 0$, pak i -tá výběrová hlavní komponenta je dána vztahem

$$\hat{y}_i = \hat{\mathbf{e}}_i' \mathbf{x} = \hat{e}_{i1}x_1 + \hat{e}_{i2}x_2 + \dots + \hat{e}_{ip}x_p, \quad i = 1, 2, \dots, p,$$

kde \mathbf{x} je libovolné pozorování proměnných X_1, X_2, \dots, X_p . Dále platí, že výběrový rozptyl \hat{y}_k je roven $\hat{\lambda}_k$, $k = 1, 2, \dots, p$, výběrová kovariance dvojic (\hat{y}_i, \hat{y}_k) je pro $i \neq k$ rovna nule. Celkový výběrový rozptyl je potom dán vztahem

$$\sum_{i=1}^p s_{ii} = \hat{\lambda}_1 + \hat{\lambda}_2 + \dots + \hat{\lambda}_p$$

a korelační koeficienty výrazem

$$r_{\hat{y}_i, \hat{y}_k} = \frac{\hat{e}_{ik} \sqrt{\hat{\lambda}_i}}{\sqrt{s_{kk}}}, \quad i, k = 1, 2, \dots, p.$$

Při analýze dat je často používáno „centrování“ odečtením výběrového průměru $\bar{\mathbf{x}}$. Toto centrování nemá vliv na výběrovou varianční matici. Pro libovolný vektor pozorování \mathbf{x} je i -tá hlavní komponenta dána

$$\hat{y}_i = \hat{\mathbf{e}}_i' (\mathbf{x} - \bar{\mathbf{x}}), \quad i = 1, 2, \dots, p.$$

Grafické zobrazení hlavních komponent

S využitím jednoduchých grafických nástrojů lze nalézt „podezřelá“ pozorování, lze také provést ověření předpokladu normality. Obvykle se provádí ověření normality pro několik prvních hlavních komponent. To lze udělat například pomocí QQ-plotu.

Na druhé straně poslední hlavní komponenty mohou pomoci při odhalení neobvyklých pozorování. Každé pozorování lze vyjádřit jako lineární kombinaci

$$\mathbf{x}_j = (\mathbf{x}'_j \hat{\mathbf{e}}_1) \hat{\mathbf{e}}_1 + (\mathbf{x}'_j \hat{\mathbf{e}}_2) \hat{\mathbf{e}}_2 + \cdots + (\mathbf{x}'_j \hat{\mathbf{e}}_p) \hat{\mathbf{e}}_p = \hat{y}_{j1} \hat{\mathbf{e}}_1 + \hat{y}_{j2} \hat{\mathbf{e}}_2 + \cdots + \hat{y}_{jp} \hat{\mathbf{e}}_p$$

všech vlastních vektorů $\hat{\mathbf{e}}_1, \hat{\mathbf{e}}_2, \dots, \hat{\mathbf{e}}_p$ matice \mathbf{S} . Velikost posledních hlavních komponent poukazuje na to, jak dobře prvních několik hlavních komponent propisují data. Takže

$\hat{y}_{j1} \hat{\mathbf{e}}_1 + \hat{y}_{j2} \hat{\mathbf{e}}_2 + \cdots + \hat{y}_{j,q-1} \hat{\mathbf{e}}_{q-1}$ se liší od \mathbf{x}_j o $\hat{y}_{jq} \hat{\mathbf{e}}_q + \cdots + \hat{y}_{jp} \hat{\mathbf{e}}_p$, tedy o délku

$\sqrt{\hat{y}_{jq}^2 + \cdots + \hat{y}_{jp}^2}$. Neobvyklé pozorování obvykle takové, jehož souřadnice $\hat{y}_{jq}, \dots, \hat{y}_{jp}$ jsou velké.

Grafické zobrazení hlavních komponent

Dalším typem grafu, který se obvykle v analýze hlavních komponent používá, je bodový graf výběrových hlavních komponent (např. na vodorovnou osu vyneseme souřadnice první výběrové hlavní komponenty určené z analyzovaných dat, na osu svislou potom druhou výběrovou hlavní komponentu).

Předpokládejme, že analyzovaná data jsou výběrem z vícerozměrného náhodného rozdělení se střední hodnotou μ a varianční maticí Σ . Uvažujme první dvě hlavní komponenty. Výběrový rozptyl první hlavní komponenty \hat{y}_1 je $\hat{\lambda}_1$, výběrový rozptyl druhé hlavní komponenty \hat{y}_2 je $\hat{\lambda}_2$. Vykreslíme-li bodový graf hodnot $(\hat{y}_{j1}), \hat{y}_{j2}$ pro $j = 1, 2, \dots, n$, měly tyto body ležet uvnitř elipsy o rovnici

$$\frac{\hat{y}_1^2}{\hat{\lambda}_1} + \frac{\hat{y}_2^2}{\hat{\lambda}_2} \leq \chi_{1-\alpha}^2(2).$$

Vynesené body by měly mít přibližně eliptický tvar.