

MA012 Statistika II

11. Analýza hlavních komponent (PCA)

Ondřej Pokora (pokora@math.muni.cz)

Ústav matematiky a statistiky, Přírodovědecká fakulta, Masarykova univerzita, Brno

(podzim 2015)



Analýza hlavních komponent

Analýza hlavních komponent (Principal component analysis, PCA) je statistická metoda pro snížení dimenze (počtu) proměnných, které bývají často korelované. Připomeňme např. velmi častý problém multikolinearity ve vícenásobné lineární regresi.

PCA využívá ortogonalizace báze vektorového prostoru původních náhodných veličin. Nové báze náhodné veličiny, nazývané hlavní komponenty, jsou nekorelované. Současně je volba nové báze optimální v tom, že pro vysvětlení zadaného podílu celkové variability dat stačí uvažovat nižší počet hlavních komponent, než byl počet (dimenze) původních náhodných veličin. Jednotlivé hlavní komponenty jsou postupně hledány tak, aby vysvětlily co největší část celkové variability dat pomocí co nejžšího počtu nových náhodných veličin.

Nevýhrou PCA může být nemožnost interpretace hlavních komponent, např. z důvodu nekompatibility použitých jednotek. V praxi toto často řešíme uvažováním bezrozměrných veličin.

Pozorování ve standardní bázi

Data n pozorování k veličin máme v matici \mathbf{X} rozměru $n \times k$, $n > k$. Můžeme to chápat jako n pozorování v k -rozměrném prostoru. Uvažujeme-li standardní bázi jednotkových vektorů v \mathbb{R}^k ,

$$e_1, \dots, e_k,$$

jednotlivé řádky matice \mathbf{X} jsou souřadnice jednotlivých pozorování v této standardní bázi, která je ortonormální.

Předpokládejme nyní, že máme obecnou ortonormální bázi

$$u_1, \dots, u_k \in \mathbb{R}^k, \quad u_i' u_j = \begin{cases} 1, & i = j \\ 0, & i \neq j \end{cases}.$$

Vytvořme matici \mathbf{U} tak, že tyto bázevé vektory umístíme do sloupců. Taková matice \mathbf{U} je potom ortogonální, tzn. platí

$$\mathbf{U}'\mathbf{U} = \mathbf{I}_k, \quad \mathbf{U}^{-1} = \mathbf{U}'.$$

Pro každé $x \in \mathbb{R}^k$ pak platí

$$x = a_1 u_1 + \dots + a_k u_k = \mathbf{U}a, \quad a = \mathbf{U}'x,$$

kde $a = (a_1, \dots, a_k)'$ jsou souřadnice x v bázi sloupců matice \mathbf{U} .

Volba optimální báze

Bázi $\mathbf{u}_1, \dots, \mathbf{u}_k$ chceme zvolit tak, aby pro $\mathbf{x} \in \mathbb{R}^k$ platilo

$$\mathbf{x} \approx \mathbf{x}^*,$$

kde \mathbf{x}^* je lineární kombinací pouze prvních $r < k$ báзовých vektorů,

$$\mathbf{x}^* = a_1 \mathbf{u}_1 + \dots + a_r \mathbf{u}_r = \mathbf{U}^* \mathbf{a}^*,$$

s redukovanou maticí \mathbf{U}^* tvořenou prvními r sloupci matice \mathbf{U} . Máme tedy

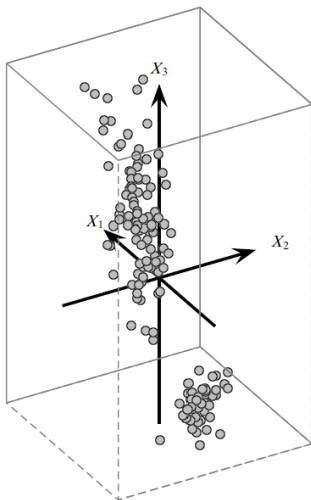
$$\mathbf{x}^* = \mathbf{U}^* \mathbf{a}^* = \underbrace{\mathbf{U}^* (\mathbf{U}^*)'}_P \mathbf{x}, \quad \text{neboť} \quad \mathbf{a}^* = (\mathbf{U}^*)' \mathbf{x}.$$

Aproximace \mathbf{x}^* je tedy projekcí \mathbf{x} do podprostoru generovaného sloupci redukované matice \mathbf{U}^* , a to pomocí projekční matice P .

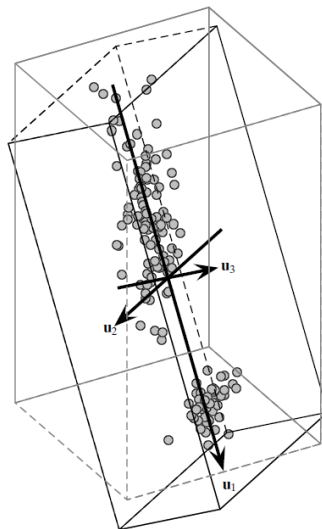
Dopustíme se přitom chyby

$$\boldsymbol{\varepsilon} = \mathbf{x} - \mathbf{x}^* = a_{r+1} \mathbf{u}_{r+1} + \dots + a_k \mathbf{u}_k.$$

Různé volby bází



(a) Original Basis



(b) Optimal Basis

Analýza hlavních komponent

Od začátku: chceme nalézt první vektor \mathbf{u} optimální ortonormální báze, pomocí něhož dokážeme \mathbf{x} aproximovat pomocí \mathbf{x}^* nejlépe z pohledu vysvětleného rozptylu.

Dále předpokládáme, že data v matici \mathbf{X} jsou centrována; v opačném případě je centrujeme, tzn. od každého sloupce odečteme jeho průměr.

Projekce i -tého pozorování \mathbf{x}_i , do směru bázevého vektoru \mathbf{u} je tvaru

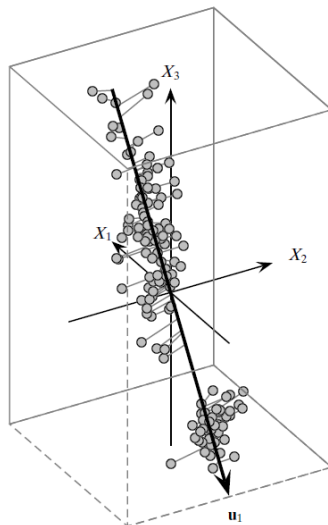
$$\mathbf{x}_i^* = a_i \mathbf{u}, \quad i = 1, \dots, n.$$

Rozptyl projektovaných pozorování $\mathbf{x}_1^*, \dots, \mathbf{x}_n^*$ je rovný

$$\begin{aligned} \sigma_{\mathbf{u}}^2 &= \frac{1}{n} \sum_{i=1}^n (a_i - \underbrace{\mu_i}_{=0})^2 = \frac{1}{n} \sum_{i=1}^n a_i^2 = \frac{1}{n} \sum_{i=1}^n (\mathbf{u}' \mathbf{x}_i)^2 = \\ &= \frac{1}{n} \sum_{i=1}^n \mathbf{u}' \mathbf{x}_i \mathbf{x}_i' \mathbf{u} = \mathbf{u}' \underbrace{\left(\frac{1}{n} \sum_{i=1}^n \mathbf{x}_i \mathbf{x}_i' \right)}_{\Sigma} \mathbf{u} = \mathbf{u}' \Sigma \mathbf{u}, \end{aligned}$$

kde Σ je kovarianční matice ($k \times k$) centrovanych veličin.

Volba optimálního směru



Obrázky: Zaki & Meira: Data mining and analysis, 2014.

Analýza hlavních komponent

Optimální vektor \mathbf{u} , $\mathbf{u}'\mathbf{u} = 1$, hledáme tak, abychom maximalizovali rozptyl $\sigma_{\mathbf{u}}^2$ pozorování projektovaných do směru \mathbf{u} . Řešíme tedy úlohu na vázaný extrém

$$\sigma_{\mathbf{u}}^2 = \mathbf{u}'\Sigma\mathbf{u} \rightarrow \max, \quad \text{za podmínky} \quad \mathbf{u}'\mathbf{u} = 1.$$

Použijeme Lagrangeovu metodu. Lagrangeova funkce L s Lagrangeovým multiplikátorem λ je

$$L(\mathbf{u}) = \mathbf{u}'\Sigma\mathbf{u} - \lambda(\mathbf{u}'\mathbf{u} - 1),$$

pro niž nalezneme pomocí stacionárního bodu $\frac{\partial L(\mathbf{u})}{\partial \mathbf{u}} = \mathbf{0}$ řešení zadané vztahem

$$\Sigma\mathbf{u} = \lambda\mathbf{u}.$$

To znamená, že λ je vlastní číslo kovarianční matice Σ a \mathbf{u} je jemu odpovídající vlastní vektor.

Přitom hledaný rozptyl je rovný $\sigma_{\mathbf{u}}^2 = \lambda$. Má-li být $\sigma_{\mathbf{u}}^2$ maximální, je třeba vzít za λ největší vlastní číslo kovariační matice Σ a za \mathbf{u} jemu odpovídající vlastní vektor, který označujeme jako první hlavní komponentu \mathbf{u}_1 .

Analýza hlavních komponent

Na uvedeném principu je založena metoda hlavních komponent. Postupně hledáme vektory ortonormální báze, které s co nejmenším dimenzí generovaného podprostoru vysvětlují co nejvíce variability dat.

Ve druhém kroku tedy hledáme druhou hlavní komponentu u_2 ,

$$u_2' u_2 = 1, \quad u_1' u_2 = 0,$$

tak, aby se maximalizovala variabilita v datech projektovaných do podprostoru generovaného vektory u_1 a u_2 .

Analýza hlavních komponent

Kovarianční matice Σ datové matice X rozměru $n \times k$ je symetrická a pozitivně semidefinitní čtvercová matice řádu k , tzn. její vlastní čísla jsou vždy nezáporná, a počet kladných vlastních čísel odpovídá hodnotě matice Σ . Uspořádejme tato vlastní čísla do nerostoucí posloupnosti

$$\lambda_1 \geq \dots \geq \lambda_k$$

a označme k nim příslušející vlastní vektory $\mathbf{u}_1, \dots, \mathbf{u}_k$.

Věta 1

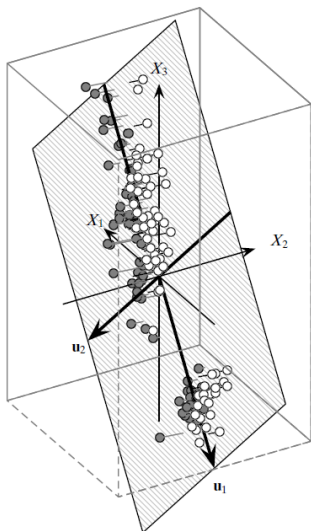
$$\text{Platí} \quad \Sigma = \mathbf{U}\Lambda\mathbf{U}',$$

kde sloupce matice \mathbf{U} jsou vlastní vektory a $\Lambda = \text{diag}(\lambda_1, \dots, \lambda_k)$.

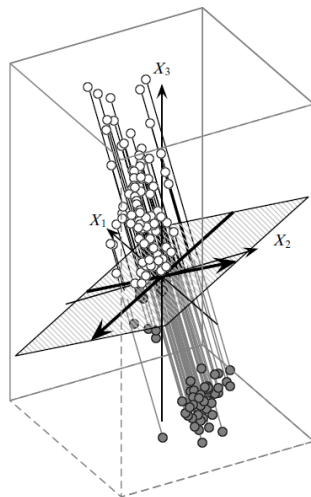
Celková variabilita dat v matici X je rovna

$$\sum_{j=1}^k \lambda_j = \text{Tr } \Sigma.$$

Variabilita v různých bázích



(a) Optimal basis



(b) Nonoptimal basis

Algoritmus PCA I

- 1 Pro datovou matici \mathbf{Y} typu $n \times k$ spočítáme výběrové průměry pro každou veličinu (tj. sloupcové průměry)

$$\boldsymbol{\mu}' = \frac{1}{n} \sum_{i=1}^n \mathbf{x}'_i = \left(\frac{1}{n} \sum_{i=1}^n x_{i1}, \dots, \frac{1}{n} \sum_{i=1}^n x_{ik} \right).$$

- 2 Vytvoříme matici cetrovaných veličin rozměru $n \times k$,

$$\mathbf{X} = (1, \dots, 1)' \boldsymbol{\mu}',$$

složenou z n řádků $\mathbf{x}'_1, \dots, \mathbf{x}'_n$.

- 3 Označíme $\boldsymbol{\Sigma}$ kovariační matici matice \mathbf{X} ,

$$\boldsymbol{\Sigma} = \text{cov} \mathbf{X} = \frac{1}{n} \sum_{i=1}^n \mathbf{x}_i \mathbf{x}'_i.$$

Algoritmus PCA II

- 4 Spočítáme vlastní čísla kovarianční matice Σ , seřadíme je do nerostoucí posloupnosti a odpovídající vlastní vektory uspořádáme (ve stejném pořadí) do sloupců matice \mathbf{U} ,

$$\lambda_1 \geq \dots \geq \lambda_k, \quad \mathbf{U} = (\mathbf{u}_1 \mid \dots \mid \mathbf{u}_k).$$

- 5 Nalezneme vhodnou dimenzi $r < k$, např. s podmínkou na minimální hodnotu podílu rozptylu

$$\frac{\sum_{j=1}^r \lambda_j}{\sum_{j=1}^k \lambda_j} \geq \alpha$$

pro stanovenou hodnotu $\alpha \in (0; 1)$, často se uvádí $\alpha = 0,8$.

- 6 Vytvoříme matici redukované báze hlavních komponent \mathbf{U}_r tak, že z matice \mathbf{U} vezmeme pouze prvních r sloupců.
- 7 Datová matice v bázi hlavních komponent \mathbf{A} je rozměru $n \times r$ a je tvořena řádky \mathbf{a}'_i , kde

$$\mathbf{a}'_i = \mathbf{U}_r \mathbf{x}'_i, \quad i = 1, \dots, n.$$

Příklad

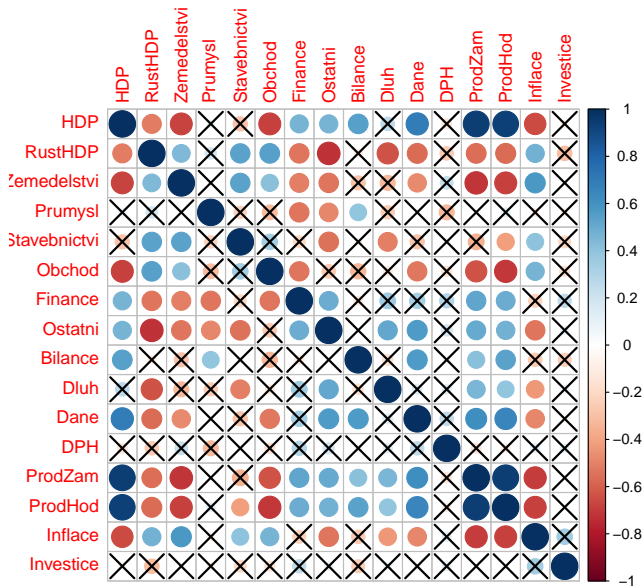
V *R* je klasická PCA naprogramována ve funkci `prcomp`, která pro hledání vlastních čísel a vektorů využívá klasický spektrální rozklad kovarianční matice $\Sigma = \mathbf{U}\Lambda\mathbf{U}'$.

V praxi se však obvykle používá tzv. *singular value decomposition (SVD)* datové matice $\mathbf{X} = \mathbf{L}\mathbf{\Delta}\mathbf{R}'$, která je výpočetně výhodnější. V *R* tento postup využívá funkce `princomp`.

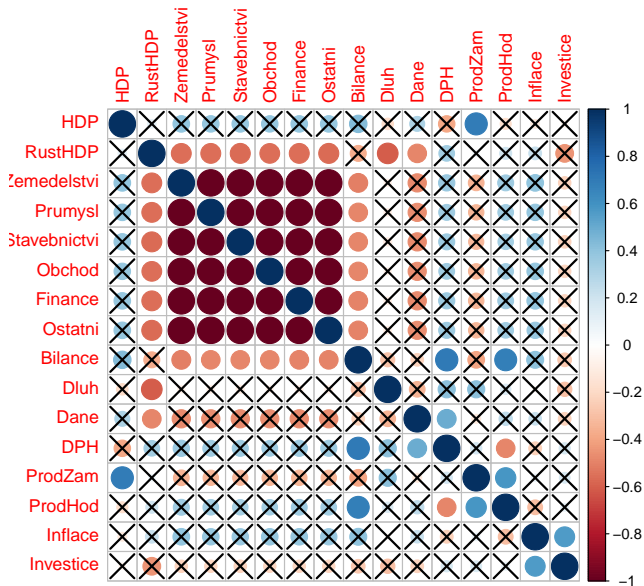
Pro grafické reprezentace výsledků PCA pak používáme mj. `screeplot` a `biplot`.

Příklad 1

V souboru `ukazatele.csv` je 16 vybraných ekonomických a finančních ukazatelů 29 evropských zemí z roku 2007, které publikoval EUROSTAT v roce 2009. Proved'te analýzu hlavních komponent a identifikujte hlavní komponenty ukazatelů.

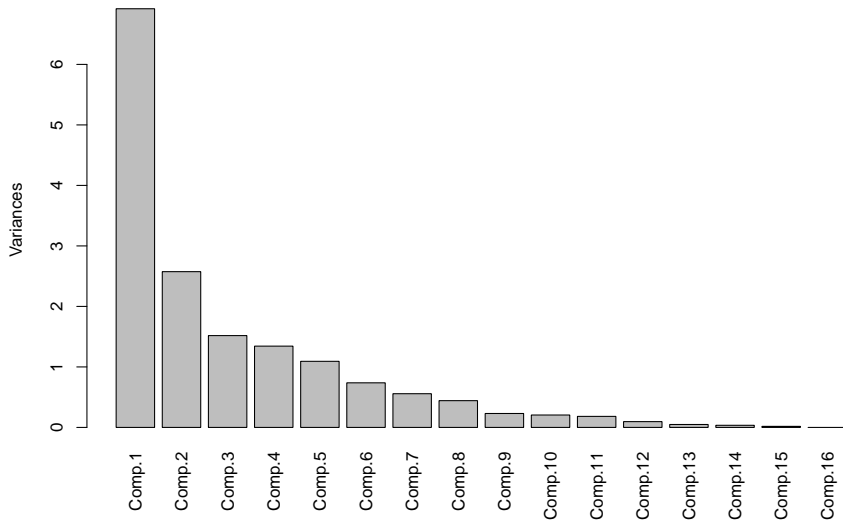


Korelogram

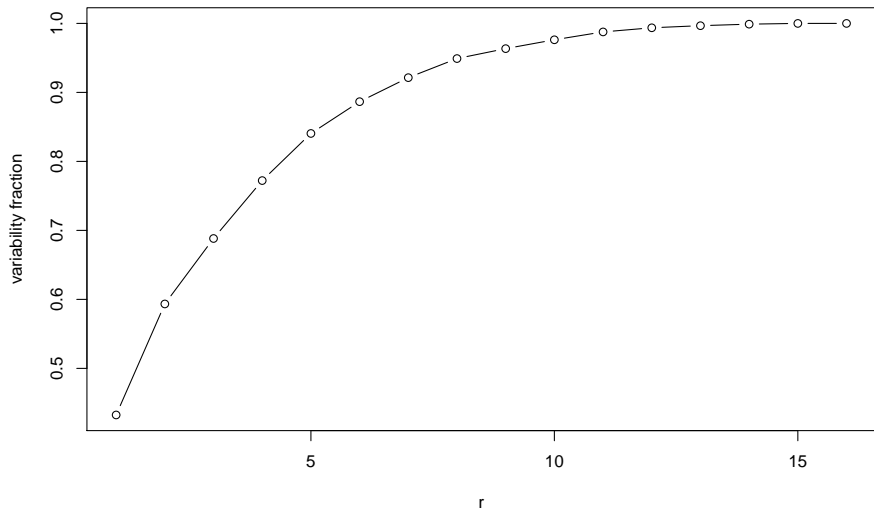


Korelogram parciálních korelací

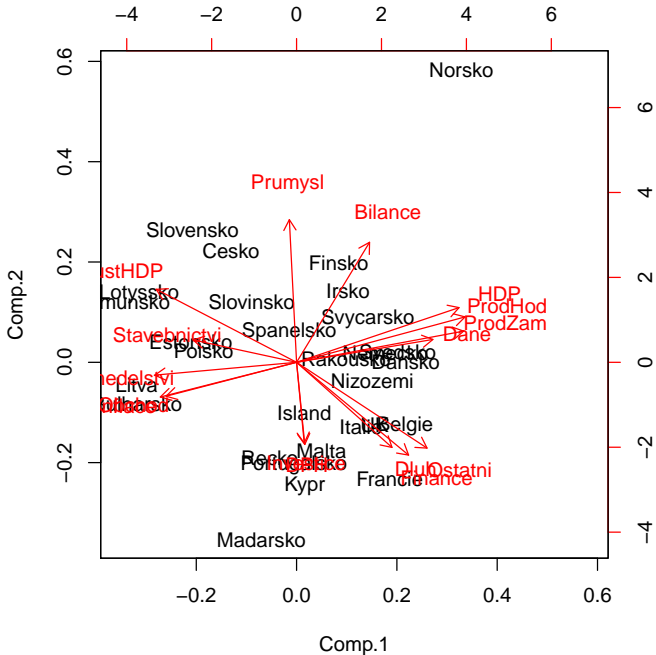
Scree plot



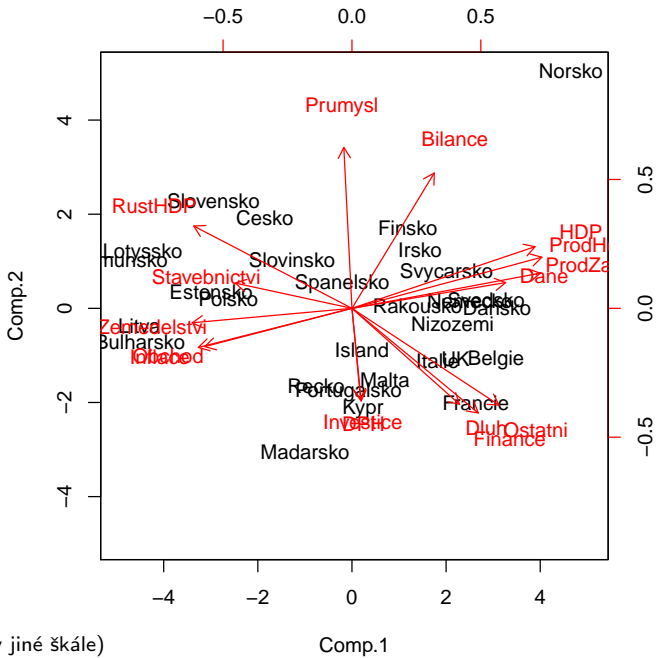
Tzv. suťový graf



Graf podílu vysvětlené variability



Tzv. biplot



Tzv. biplot (v jiné škále)