

Faktorová analýza

(přednáška v rámci předmětu Statistika a spolehlivost v lékařství)

Mirko Navara

Centrum strojového vnímání

katedra kybernetiky FEL ČVUT

Karlovo náměstí, budova G, místnost 104a

<https://cw.felk.cvut.cz/doku.php/courses/a6m33ssl>

25. prosince 2011

Cíl: Popsat podstatné závislosti mnoha náhodných veličin pomocí menšího počtu společných faktorů.

(Co nejjednodušší popis s malou ztrátou informace.)

$\mathbf{X} = (X_1, \dots, X_p)^\top$ = náhodný vektor tvořený p normovanými náhodnými veličinami, $E X_i = 0$, $D X_i = 1$,

$\Sigma_{\mathbf{X}} = E(\mathbf{X} \mathbf{X}^\top)$ = kovarianční matice.

Lineární závislosti (=korelace) složek náhodného vektoru \mathbf{X} chceme vysvětlit jako důsledek následujících náhodných veličin:

F_1, \dots, F_m , $m < p$, ... **společné faktory**, které vysvětlují korelace

R_1, \dots, R_p ... **specifické faktory** (=chyběvě faktory=reziduální složky), které vysvětlují zbytkový rozptyl
Model:

$$\begin{aligned} X_1 &= w_{11} F_1 + w_{12} F_2 + \dots + w_{1m} F_m + R_1, \\ X_2 &= w_{21} F_1 + w_{22} F_2 + \dots + w_{2m} F_m + R_2, \\ &\dots \\ X_p &= w_{p1} F_1 + w_{p2} F_2 + \dots + w_{pm} F_m + R_p, \end{aligned}$$

kde F_j jsou normované, $E F_j = 0$, $D F_j = 1$, a nekorelované, $E(F_j F_k) = 0$,

R_i jsou centrované, $E R_i = 0$, a nezávislé navzájem i na všech F_j , $E(E_j E_k) = E(E_j F_k) = 0$.

Maticový zápis pomocí $\mathbf{F} = (F_1, \dots, F_m)^\top$, $\mathbf{R} = (R_1, \dots, R_p)^\top$ a matice $\mathbf{W} \in \mathbb{R}^{p \times m}$:

$$\mathbf{X} = \mathbf{W} \mathbf{F} + \mathbf{R},$$

kde $E \mathbf{R} = \mathbf{O}_p = (0, \dots, 0) \in \mathbb{R}^p$,

$\Sigma_{\mathbf{R}} = E(\mathbf{R} \mathbf{R}^\top)$ je diagonální kovarianční matice,

$E \mathbf{F} = \mathbf{O}_m = (0, \dots, 0) \in \mathbb{R}^m$,

$\Sigma_{\mathbf{F}} = E(\mathbf{F} \mathbf{F}^\top) = \mathbf{I}_m \in \mathbb{R}^{m \times m}$ je jednotková matice,

$E(\mathbf{F} \mathbf{R}^\top) = \mathbf{O}_{m,p} \in \mathbb{R}^{m \times p}$ je nulová matice.

Vztah mezi kovariančními maticemi:

$$\begin{aligned} \Sigma_{\mathbf{X}} &= E(\mathbf{X} \mathbf{X}^\top) = E((\mathbf{W} \mathbf{F} + \mathbf{R})(\mathbf{W} \mathbf{F} + \mathbf{R})^\top) = E((\mathbf{W} \mathbf{F} \mathbf{F}^\top \mathbf{W}^\top + \mathbf{W} \mathbf{F} \mathbf{R}^\top + \mathbf{R} \mathbf{F}^\top \mathbf{W}^\top + \mathbf{R} \mathbf{R}^\top) = \\ &= E(\mathbf{W} \mathbf{F} \mathbf{F}^\top \mathbf{W}^\top) + E(\mathbf{W} \mathbf{F} \mathbf{R}^\top) + E(\mathbf{R} \mathbf{F}^\top \mathbf{W}^\top) + E(\mathbf{R} \mathbf{R}^\top) = \\ &= \underbrace{\mathbf{W} \mathbf{F} \mathbf{F}^\top \mathbf{W}^\top}_{\mathbf{I}_m} + \underbrace{\mathbf{W} \mathbf{F} \mathbf{R}^\top}_{\mathbf{O}_{m,p}} + \underbrace{\mathbf{R} \mathbf{F}^\top \mathbf{W}^\top}_{\mathbf{O}_{p,m}} + \underbrace{\mathbf{R} \mathbf{R}^\top}_{\Sigma_{\mathbf{R}}} = \\ &= \mathbf{W} \mathbf{W}^\top + \Sigma_{\mathbf{R}}. \end{aligned}$$

První sčítanec vysvětlují společné faktory; má být co největší.

Druhý sčítanec vysvětlují specifické faktory; má být co nejmenší.

Kovarianční matici $\Sigma_{\mathbf{X}}$ tvoří prvky na diagonále,

$$(\Sigma_{\mathbf{X}})_{ii} = D X_i = \sum_{j=1}^m w_{ij}^2 + D R_i,$$

a prvky mimo diagonálu ($i \neq k$),

$$(\Sigma_{\mathbf{X}})_{ik} = \text{Cov}(X_i, X_k) = \sum_{j=1}^m w_{ij} w_{kj}.$$

$$\mathbf{W} \mathbf{W}^\top = \Sigma_{\mathbf{X}} - \Sigma_{\mathbf{R}}$$

je **redukovaná kovarianční matici**. Její i -tý prvek na diagonále,

$$(\mathbf{W} \mathbf{W}^\top)_{ii} = \sum_{j=1}^m w_{ij}^2 = DX_i - DR_i,$$

se nazývá **komunalita** proměnné X_i ; představuje část rozptylu DX_i , vysvětlenou *společnými* faktory; zbytek, DR_i , je **specifický rozptyl**.

Význam koeficientů w_{ij} :

$$w_{ij} = \text{Cov}(X_i, F_j)$$

Důkaz:

A. Po složkách:

$$\begin{aligned} \text{Cov}(X_i, F_j) &= E(X_i F_j) = E((w_{i1} F_1 + w_{i2} F_2 + \dots + w_{im} F_m + R_i) F_j) = \\ &= E(w_{i1} F_1 F_j) + E(w_{i2} F_2 F_j) + \dots + E(w_{im} F_m F_j) + E(R_i F_j) = \\ &= w_{i1} E(F_1 F_j) + w_{i2} E(F_2 F_j) + \dots + w_{im} E(F_m F_j) + E(R_i F_j) = \\ &= w_{ij}. \end{aligned}$$

B. Maticový:

$$\begin{aligned} E(\mathbf{X} \mathbf{F}^\top) &= E((\mathbf{W} \mathbf{F} + \mathbf{R}) \mathbf{F}^\top) = E(\mathbf{W} \mathbf{F} \mathbf{F}^\top + \mathbf{R} \mathbf{F}^\top) = \\ &= \mathbf{W} \underbrace{E(\mathbf{F} \mathbf{F}^\top)}_{\mathbf{I}_m} + \underbrace{E(\mathbf{R} \mathbf{F}^\top)}_{\mathbf{O}_{p,m}} = \mathbf{W} \end{aligned}$$

Úloha: K dané (symetrické pozitivně definitní) kovarianční matici $\Sigma_{\mathbf{X}}$ hledáme diagonální matici $\Sigma_{\mathbf{R}}$ (s nezápornými prvky DR_i na diagonále) takovou, aby rozdíl $\Sigma_{\mathbf{X}} - \Sigma_{\mathbf{R}}$ byl tvaru $\mathbf{W} \mathbf{W}^\top$, kde matice \mathbf{W} je typu $p \times m$ (takže $\mathbf{W} \mathbf{W}^\top$ má hodnost nejvýše m).

Řešitelnost: Pomocí

$p m$ koeficientů matice \mathbf{W} a

p koeficientů diagonální matice $\Sigma_{\mathbf{R}}$

chceme vysvětlit

$\frac{p(p-1)}{2}$ korelací (resp. koeficientů symetrické matice $\Sigma_{\mathbf{X}}$).

Pro $m = p$ řešení existuje a lze volit $\Sigma_{\mathbf{R}} = \mathbf{0}_{p,p}$.

Nás ale zajímá případ $m \ll p$; pro řešitelnost známe podmínky

- nutné
- nebo postačující,
- nikoli však nutné a postačující (otevřený problém).

Jednoznačnost: Pro $m > 1$ není řešení jednoznačné.

Je určen jen lineární podprostor (prostoru \mathcal{N} všech náhodných veličin, které mají nulovou střední hodnotu a konečný rozptyl), který je lineárním obalem vektorů F_1, \dots, F_m .

Za společné faktory můžeme volit i jakoukoli jinou jeho bázi (kterou dostaneme **rotací faktorů**).

Proto interpretace faktorů je problematická.

Jednoznačnost může obvykle zajistit **kanonický tvar faktorového modelu**: Faktory vyčerpávají postupně maximum celkového rozptylu (první faktor nejvíce atd.).

Problémy:

- Numerické řešení (např. metoda maximální věrohodnosti za předpokladu normálního rozdělení náhodných veličin).
- Počet faktorů (existují testy posuzující vhodnou složitost modelu podle jeho přesnosti).
- Interpretace faktorů (případně jejich rotace na lépe interpretovatelné veličiny).

Literatura

[Hebák] Hebák, P., Hustopecký, J.: *Vícerozměrné statistické metody s aplikacemi*. SNTL, Praha, 1987.

[Navara: PMS] Navara, M.: *Pravděpodobnost a matematická statistika*. Skriptum ČVUT, Praha, 2007.