

Faktorová analýza

Ekonometrie

Jiří Neubauer

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

Faktorová analýza – úvod

Na sledovaných objektech se často zjišťují hodnoty většího počtu statistických znaků nebo proměnných, ty jsou obvykle mezi sebou korelovány. Lze si představit, že korelace mezi jednotlivými proměnnými jsou způsobeny vlivem nějakého menšího počtu nepřímo měřitelných společných **faktorů**, které ovlivňují hodnoty sledovaných proměnných. Každý z těchto faktorů může ovlivnit hodnoty pozorování každé ze zkoumaných proměnných.

Cílem faktorové analýzy pak je tyto faktory odhadnout, odhadnout počet statisticky významných faktorů a konečně odhadnout hodnoty každého z faktorů pro každý sledovaný objekt, tedy popsat objekty pomocí nalezených faktorů. V tomto směru je faktorová analýza metodou pro snížení rozsahu dat. Navíc sledováním realizací společných faktorů u jednotlivých objektů mohou být identifikovány výjimečné objekty.

Ortogonalní faktorový model

Ve faktorovém modelu předpokládáme, že p -rozměrný sloupcový náhodný vektor sledovaných proměnných $\mathbf{X} = (X_1, \dots, X_p)'$ se střední hodnotou $\boldsymbol{\mu}$ a varianční maticí $\boldsymbol{\Sigma}$ lineárně závisí na m -rozměrném náhodném vektoru společných faktorů $\mathbf{F} = (F_1, \dots, F_m)'$ a p -rozměrném vektoru specifických faktorů $\boldsymbol{\epsilon} = (\epsilon_1, \dots, \epsilon_p)'$. Maticově lze faktorový model zapsat

$$\mathbf{X} - \boldsymbol{\mu} = \mathbf{L}\mathbf{F} + \boldsymbol{\epsilon},$$

kde \mathbf{L} je matice typu $p \times m$. Speciálně pro i -tou proměnnou X_i platí

$$X_i - \mu_i = l_{i1}F_1 + l_{i2}F_2 + \dots + l_{im}F_m + \epsilon_i.$$

V uvedeném modelu popisuje koeficient l_{ij} vliv faktoru F_j na proměnnou X_i , a proto se matice \mathbf{L} nazývá **matice faktorových zátěží**.

Ortogonalní faktorový model

Vektor společných faktorů \mathbf{F} a vektor specifických faktorů ϵ nelze přímo pozorovat. Aby bylo možné je odhadnout, jsou na ně kladeny následující předpoklady.

- 1 Střední hodnoty obou těchto vektorů jsou nulové, tj. $E(\mathbf{F}) = \mathbf{0}$, $E(\epsilon) = \mathbf{0}$.
- 2 Variační matice vektoru \mathbf{F} je jednotková, tj. $\text{var}(\mathbf{F}) = \mathbf{I}$, a varianční matice vektoru ϵ je diagonální, tj. $\text{var}(\epsilon) = \mathbf{\Psi} = \text{diag}(\psi_1, \dots, \psi_p)$.
- 3 Vektory \mathbf{F} a ϵ jsou nekorelované, tedy jejich kovarianční matice $\text{cov}(\epsilon, \mathbf{F}) = \mathbf{0}$.

Za uvedených předpokladů lze snadno odvodit, že

$$\begin{aligned} (\mathbf{X} - \boldsymbol{\mu})(\mathbf{X} - \boldsymbol{\mu})' &= (\mathbf{LF} + \epsilon)(\mathbf{LF} + \epsilon)' = (\mathbf{LF} + \epsilon)((\mathbf{LF})' + \epsilon') = \\ &= \mathbf{LF}(\mathbf{LF})' + \epsilon(\mathbf{LF})' + \mathbf{LF}\epsilon' + \epsilon\epsilon', \end{aligned}$$

odkud pro varianční matici vektoru sledovaných proměnných \mathbf{X} platí

$$\begin{aligned} \text{var}(\mathbf{X}) &= E(\mathbf{X} - \boldsymbol{\mu})(\mathbf{X} - \boldsymbol{\mu})' = \mathbf{LE}(\mathbf{FF}')\mathbf{L}' + E(\epsilon\mathbf{F}')\mathbf{L}' + \mathbf{LE}(\mathbf{F}\epsilon') + E(\epsilon\epsilon') = \\ &= \mathbf{LL}' + \mathbf{\Psi}, \end{aligned}$$

neboli

$$\begin{aligned} D(X_i) &= l_{i1}^2 + l_{i2}^2 + \dots + l_{im}^2 + \psi_i, \quad i = 1, 2, \dots, p, \\ C(X_i, X_k) &= l_{i1}l_{k1} + l_{i2}l_{k2} + \dots + l_{im}l_{km}, \quad i, k = 1, 2, \dots, p, \quad i \neq k. \end{aligned} \tag{1}$$

Ortogonalní faktorový model

Dále platí, že kovarianční matice vektorů \mathbf{X} a \mathbf{F} je rovna matici faktorových zátěží \mathbf{L} , tj.

$$\begin{aligned}\text{cov}(\mathbf{X}, \mathbf{F}) &= E[(\mathbf{X} - \boldsymbol{\mu})\mathbf{F}'] = E[(\mathbf{L}\mathbf{F} + \boldsymbol{\epsilon})\mathbf{F}'] = E[\mathbf{L}\mathbf{F}\mathbf{F}' + \boldsymbol{\epsilon}\mathbf{F}'] = \\ &= \mathbf{L}E(\mathbf{F}\mathbf{F}') + E(\boldsymbol{\epsilon}\mathbf{F}') = \mathbf{L},\end{aligned}$$

neboli

$$C(X_i, F_j) = l_{ij}, \quad i = 1, 2, \dots, p, \quad j = 1, 2, \dots, m.$$

V této souvislosti pak říkáme, že uvedená vyjádření varianční matice $\text{var}(\mathbf{X})$ a kovarianční matice $\text{cov}(\mathbf{X}, \mathbf{F})$ popisují kovarianční strukturu ortogonalního faktorového modelu.

Model $\mathbf{X} - \boldsymbol{\mu} = \mathbf{L}\mathbf{F} + \boldsymbol{\epsilon}$ je lineární vzhledem ke společným faktorům. Je-li vztah mezi \mathbf{X} a faktory jiný než lineární, popis kovarianční struktury $\mathbf{L}\mathbf{L}' + \boldsymbol{\Psi}$ uvedený v (1) nebude odpovídající.

Ortogonalní faktorový model

Ta část rozptylu proměnné X_i , která je vysvětlená pomocí m společných faktorů, se nazývá **komunalita** („communality“), příspěvek specifického faktoru k této variabilitě je označován jako **specifický rozptyl**, tedy

$$D(X_i) = \underbrace{\sigma_{ii}}_{\text{komunalita}} = \underbrace{l_{i1}^2 + l_{i2}^2 + \dots + l_{im}^2}_{\text{komunalita}} + \underbrace{\psi_i}_{\text{specifický rozptyl}} \quad (2)$$

Pro i -tou komunalitu h_i^2 lze psát

$$h_i^2 = l_{i1}^2 + l_{i2}^2 + \dots + l_{im}^2,$$

a rozptyl $D(X_i)$ vyjádřit ve tvaru

$$\sigma_{ii} = h_i^2 + \psi_i, \quad i = 1, 2, \dots, p.$$

Je-li dán ortogonalní faktorový model, je třeba pomocí opakovaných nezávislých pozorování odhadnout matici faktorových zátěží \mathbf{L} a dále najít odhady hodnot společných faktorů $F_1 \dots F_p$ pro každou statistickou jednotku.

Ortogonalní faktorový model

Matice faktorových zátěží L a vektor společných faktorů F nejsou v daném ortogonalním modelu určeny jednoznačně. Pro libovolnou ortogonální transformaci vektoru společných faktorů F získáme nový vektor společných faktorů $F^* = T'F$ (zde T je ortogonální matice typu $m \times m$, tj. $TT' = T'T = I$) a novou matici faktorových zátěží $L^* = LT$, které opět splňují předpoklady modelu, neboť

$$X - \mu = LTT'F + \epsilon = L^*F^* + \epsilon,$$

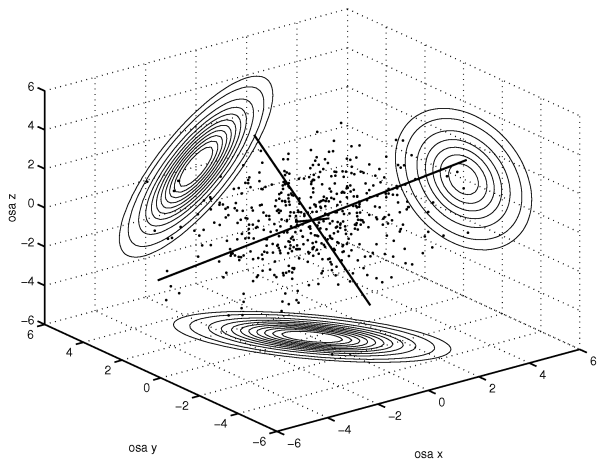
$$\text{var}(X) = LTT'L' + \Psi = L^*L^{*'} + \Psi.$$

Jak bude dále naznačeno, lze uvedenou nejednoznačnost modelu s výhodou použít při interpretaci výsledků.

Geometrická interpretace

Bylo nasimulováno 500 pozorování trojrozměrného náhodného vektoru \mathbf{X} z normálního rozdělení s nulovou střední hodnotou a nediagonální varianční maticí Σ . Na obr. 1 jsou jednotlivá pozorování znázorněna body. Z obrázku je vidět, jaká je variabilita dat. V každém směru se souřadnice jednotlivých pozorování nacházejí přibližně v rozmezí od -4 do 4 . Úsečky v 1 představují směry vlastních vektorů \mathbf{e}_1 , \mathbf{e}_2 a \mathbf{e}_3 varianční matice Σ . Jejich délky jsou úměrné vlastním číslům λ_1 , λ_2 a λ_3 varianční matice Σ , $\lambda_1 > \lambda_2 > \lambda_3$. Elipsy v rovinách os znázorňují vrstevnice hustot příslušných dvourozměrných marginálních rozdělení.

Geometrická interpretace

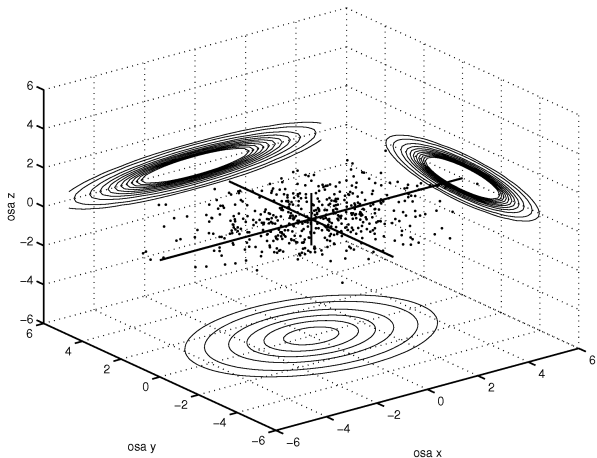


Obrázek: Simulovaná data před transformací

Geometrická interpretace

Když zvolíme novou souřadnou soustavu tak, že osa x bude ve směru vlastního vektoru e_1 , osa y ve směru vlastního vektoru e_2 a osa z ve směru vlastního vektoru e_3 a všech 500 simulovaných bodů znázorníme v této nové souřadné soustavě, dostaneme situaci znázorněnou na obr. 2. Z tohoto obrázku je dobře patrné, že největší variabilita znázorněných bodů je ve směru nové osy x . Při tom není možné najít jiný směr, v němž by byla větší variabilita dat, než ve směru osy x , tedy ve směru určeném vlastním vektorem e_1 . Vektor e_2 , který udává směr nové osy y , zároveň udává směr, kolmý na osu x , v němž je opět největší variabilita ze všech těchto kolmých směrů. Konečně se ukazuje, že ve směru osy z , tedy ve směru vlastního vektoru v_3 , je variabilita ze všech těchto tří směrů nejmenší. Variabilita znázorněných bodů ve směru vlastního vektoru v_i je úměrná vlastnímu číslu λ_i a pro celkovou variabilitu vektoru \mathbf{X} platí $\sigma_T^2(\mathbf{X}) = \sigma_1^2 + \sigma_2^2 + \sigma_3^2 = \lambda_1 + \lambda_2 + \lambda_3$, zde $\sigma_i^2 = D(X_i)$.

Geometrická interpretace



Obrázek: Simulovaná data po transformaci do báze tvořené vlastními vektory varianční matice

Odhadování parametrů modelu

Předpokládejme, že je předem znám pevný počet faktorů m . Při odhadování parametrů modelu vycházíme z náhodného výběru $\mathbf{X}_1, \dots, \mathbf{X}_n$ z rozdělení náhodného vektoru \mathbf{X} . Tento výběr je často před zpracováním standardizován, takže předpokládáme, že výsledný výběr je pak z rozdělení standardizované náhodné veličiny $\mathbf{Y} = (\text{Diag}(\boldsymbol{\Sigma}))^{-1/2}(\mathbf{X} - \boldsymbol{\mu})$, pro kterou lze faktorový model přepsat následovně

$$\mathbf{Y} = \mathbf{L}\mathbf{F} + \boldsymbol{\epsilon},$$

$$\text{var}(\mathbf{Y}) = \text{cor}(\mathbf{X}) = \mathbf{L}\mathbf{L}' + \boldsymbol{\Psi}.$$

Z tohoto náhodného výběru nejdříve stanovíme výběrový průměr $\bar{\mathbf{X}}$ jako odhad vektoru střední hodnoty $\boldsymbol{\mu}$ dále odhadneme korelační matici $\text{cor}(\mathbf{X})$ výběrovou korelační maticí \mathbf{R} . Pokud mimodiagonální prvky matice \mathbf{R} nejsou malé, tj. pokud jsou složky náhodného vektoru \mathbf{X} silně korelované, má smysl hledat společné faktory. V opačném případě by totiž hlavní roli hrály pouze specifické faktory.

Z kovarianční struktury ortogonálního faktorového modelu vyplývá, že při faktorové analýze hledáme rozklad varianční matice $\boldsymbol{\Sigma}$ náhodného vektoru \mathbf{X} na symetrickou, pozitivně definitní matici $\mathbf{L}\mathbf{L}'$ a diagonální matici $\boldsymbol{\Psi}$.

Metoda založená na hlavních komponentách

Mějme varianční matici Σ a jí odpovídající dvojice vlastních čísel a vlastních vektorů $(\lambda_i, \mathbf{e}_i)$, $\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_p \geq 0$, potom ji lze vyjádřit jako

$$\begin{aligned} \Sigma &= \lambda_1 \mathbf{e}_1 \mathbf{e}_1' + \lambda_2 \mathbf{e}_2 \mathbf{e}_2' + \dots + \lambda_p \mathbf{e}_p \mathbf{e}_p' = \\ &= \left(\sqrt{\lambda_1} \mathbf{e}_1 \mid \sqrt{\lambda_2} \mathbf{e}_2 \mid \dots \mid \sqrt{\lambda_p} \mathbf{e}_p \right) \begin{pmatrix} \sqrt{\lambda_1} \mathbf{e}_1' \\ \sqrt{\lambda_2} \mathbf{e}_2' \\ \vdots \\ \sqrt{\lambda_p} \mathbf{e}_p' \end{pmatrix} \end{aligned}$$

Uvedený rozklad se nazývá **spektrální dekompozice**.

Metoda založená na hlavních komponentách

Tímto lze popsat faktorový model mající tolik faktorů kolik proměnných ($m = p$) a specifické rozptyly $\psi_i = 0$ pro všechna i , lze tedy psát

$$\underset{(p \times p)}{\Sigma} = \underset{(p \times p)(p \times p)}{\mathbf{L} \mathbf{L}'} + \underset{(p \times p)}{\mathbf{0}} = \mathbf{L} \mathbf{L}' \quad (3)$$

Faktorové zátěže („factor loadings“) j -tého faktoru jsou až na $\sqrt{\lambda_j}$ rovny j -té hlavní komponentě.

Reprezentace matice Σ popsaná v (3) je sice přesná, nicméně obsahuje stejný počet faktorů kolik je proměnných. Cílem bude najít model, který vystihuje kovarianční strukturu pomocí několika málo faktorů. Jednou z možností je zanedbání posledních $p - m$ členů ve spektrálním rozkladu, pokud odpovídající vlastní čísla jsou malá.

$$\Sigma \doteq \left(\sqrt{\lambda_1} \mathbf{e}_1 \mid \sqrt{\lambda_2} \mathbf{e}_2 \mid \cdots \mid \sqrt{\lambda_m} \mathbf{e}_m \right) \begin{pmatrix} \sqrt{\lambda_1} \mathbf{e}'_1 \\ \sqrt{\lambda_2} \mathbf{e}'_2 \\ \vdots \\ \sqrt{\lambda_m} \mathbf{e}'_m \end{pmatrix} = \underset{(p \times m)(m \times p)}{\mathbf{L} \mathbf{L}'} \quad (4)$$

Metoda založená na hlavních komponentách

Přibližné vyjádření popsané v (4) předpokládá, že specifické faktory ϵ lze zanedbat. Pokud tyto specifické faktory do modelu zahrneme, lze psát

$$\Sigma \doteq \mathbf{L}\mathbf{L}' + \mathbf{\Psi} = \left(\sqrt{\lambda_1}\mathbf{e}_1 \mid \sqrt{\lambda_2}\mathbf{e}_2 \mid \cdots \mid \sqrt{\lambda_m}\mathbf{e}_m \right) \begin{pmatrix} \sqrt{\lambda_1}\mathbf{e}'_1 \\ \sqrt{\lambda_2}\mathbf{e}'_2 \\ \vdots \\ \sqrt{\lambda_m}\mathbf{e}'_m \end{pmatrix} + \begin{pmatrix} \psi_1 & 0 & \cdots & 0 \\ 0 & \psi_2 & \cdots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \cdots & \psi_p \end{pmatrix} \quad (5)$$

kde pro rozptyly ψ_i specifických faktorů platí $\psi_i = \sigma_{ii} - \sum_{j=1}^m l_{ij}^2$ pro $i = 1, 2, \dots, p$.

Metoda založená na hlavních komponentách

Předtím, než uvedený postup aplikujeme na data (náhodný výběr), obvykle se provádí centrování odečtením výběrového průměru, případně standardizace proměnných.

Výběrová kovarianční matice \mathbf{S} takto standardizovaných proměnných je rovna výběrové korelační matici původních dat \mathbf{R} .

Mějme výběrovou kovarianční matici \mathbf{S} , $(\hat{\lambda}_1, \hat{\mathbf{e}}_1), (\hat{\lambda}_2, \hat{\mathbf{e}}_2), \dots, (\hat{\lambda}_p, \hat{\mathbf{e}}_p)$, jsou dvojice vlastních čísel a vektorů této matice splňující $\hat{\lambda}_1 \geq \hat{\lambda}_2 \geq \dots \geq \hat{\lambda}_p$. Předpokládejme, že $m < p$ je počet společných faktorů. Odhad matice faktorových zátěží je dán výrazem

$$\hat{\mathbf{L}} = \left(\sqrt{\hat{\lambda}_1} \hat{\mathbf{e}}_1 \mid \sqrt{\hat{\lambda}_2} \hat{\mathbf{e}}_2 \mid \dots \mid \sqrt{\hat{\lambda}_m} \hat{\mathbf{e}}_m \right).$$

Odhady specifických rozptylů jsou diagonální prvky matice $\mathbf{S} - \hat{\mathbf{L}}\hat{\mathbf{L}}'$

$$\hat{\Psi} = \begin{pmatrix} \hat{\psi}_1 & 0 & \dots & 0 \\ 0 & \hat{\psi}_2 & \dots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \dots & \hat{\psi}_p \end{pmatrix}, \quad \hat{\psi}_i = s_{ii} - \sum_{j=1}^m \hat{l}_{ij}^2.$$

Komunalita lze odhadnout ze vztahu

$$\hat{h}_i^2 = \hat{l}_{i1}^2 + \hat{l}_{i2}^2 + \dots + \hat{l}_{im}^2.$$

Odhady metodou maximální věrohodnosti

Jestliže předpokládáme, že společné faktory \mathbf{F} a specifické faktory ϵ mají normální rozdělení, lze pro jejich odhady použít metodu maximální věrohodnosti. Věrohodnostní funkce má tvar

$$\begin{aligned} L(\boldsymbol{\mu}, \boldsymbol{\Sigma}) &= (2\pi)^{-\frac{np}{2}} |\boldsymbol{\Sigma}|^{-\frac{n}{2}} e^{-\frac{1}{2} \text{tr}[\boldsymbol{\Sigma}^{-1}(\sum_{j=1}^n (\mathbf{x}_j - \bar{\mathbf{x}})(\mathbf{x}_j - \bar{\mathbf{x}})' + n(\bar{\mathbf{x}} - \boldsymbol{\mu})(\bar{\mathbf{x}} - \boldsymbol{\mu})')] } = \\ &= (2\pi)^{-\frac{(n-1)p}{2}} |\boldsymbol{\Sigma}|^{-\frac{n-1}{2}} e^{-\frac{1}{2} \text{tr}[\boldsymbol{\Sigma}^{-1}(\sum_{j=1}^n (\mathbf{x}_j - \bar{\mathbf{x}})(\mathbf{x}_j - \bar{\mathbf{x}})')] } \times \\ &\quad \times (2\pi)^{-\frac{p}{2}} |\boldsymbol{\Sigma}|^{-\frac{1}{2}} e^{-\frac{n}{2}(\bar{\mathbf{x}} - \boldsymbol{\mu})' \boldsymbol{\Sigma}^{-1}(\bar{\mathbf{x}} - \boldsymbol{\mu})}. \end{aligned} \quad (6)$$

Věrohodnostní funkce (6) závisí na \mathbf{L} a $\boldsymbol{\Psi}$ prostřednictvím $\boldsymbol{\Sigma} = \mathbf{L}\mathbf{L}' + \boldsymbol{\Psi}$. Tento model není zcela jednoznačně definovaný vzhledem k možným volbám matice \mathbf{L} . K modelu se obvykle přidává podmínka

$$\mathbf{L}'\boldsymbol{\Psi}^{-1}\mathbf{L} = \boldsymbol{\Delta}$$

je diagonální matice. Maximálně věrohodné odhady $\hat{\mathbf{L}}$ a $\hat{\boldsymbol{\Psi}}$ se získají numerickou maximalizací (6). Odhady komunalit jsou tvaru

$$\hat{h}_i^2 = \hat{l}_{i1}^2 + \hat{l}_{i2}^2 + \dots + \hat{l}_{im}^2.$$

Odhady metodou maximální věrohodnosti

Pokud jsou proměnné standardizované, podobně jako tomu bylo u metody hlavních komponent, máme $\mathbf{Z} = \mathbf{V}^{-1/2}(\mathbf{X} - \boldsymbol{\mu})$, potom kovarianční matice je rovna matici korelační a platí

$$\boldsymbol{\rho} = \mathbf{V}^{-1/2} \boldsymbol{\Sigma} \mathbf{V}^{-1/2} = (\mathbf{V}^{-1/2} \mathbf{L})(\mathbf{V}^{-1/2} \mathbf{L})' + \mathbf{V}^{-1/2} \boldsymbol{\Psi} \mathbf{V}^{-1/2}.$$

Matice faktorových zátěží odpovídající $\boldsymbol{\rho}$ má tvar $\mathbf{L}_z = \mathbf{V}^{-1/2} \mathbf{L}$ a matice specifických rozptylů je rovna $\boldsymbol{\Psi}_z = \mathbf{V}^{-1/2} \boldsymbol{\Psi} \mathbf{V}^{-1/2}$. Maximálně věrohodný odhad $\boldsymbol{\rho}$ je potom dán výrazem

$$\begin{aligned} \hat{\boldsymbol{\rho}} &= (\hat{\mathbf{V}}^{-1/2} \hat{\mathbf{L}})(\hat{\mathbf{V}}^{-1/2} \hat{\mathbf{L}})' + \hat{\mathbf{V}}^{-1/2} \hat{\boldsymbol{\Psi}} \hat{\mathbf{V}}^{-1/2} = \\ &= \hat{\mathbf{L}}_z \hat{\mathbf{L}}_z' + \hat{\boldsymbol{\Psi}}_z \end{aligned}$$

Poznámka: Počítačové programy obvykle provádějí standardizaci proměnných, faktorizovaná je tedy výběrová korelační matice \mathbf{R} . Díky tomu získáme maximálně věrohodné odhady $\hat{\mathbf{L}}_z$ a $\hat{\boldsymbol{\Psi}}_z$. Maximálně věrohodné odhady matice faktorových zátěží a specifických rozptylů odpovídající matici $\frac{n-1}{n} \mathbf{S}$ jsou $\hat{\mathbf{L}} = \hat{\mathbf{V}}^{1/2} \hat{\mathbf{L}}_z$ a $\hat{\boldsymbol{\Psi}} = \hat{\mathbf{V}}^{1/2} \hat{\boldsymbol{\Psi}}_z \hat{\mathbf{V}}^{1/2}$, neboli

$$\hat{l}_{ij} = \hat{l}_{z,ij} \sqrt{\hat{\sigma}_{ii}} \quad \text{a} \quad \hat{\psi}_i = \hat{\psi}_{z,i} \hat{\sigma}_{ii},$$

kde $\hat{\sigma}_{ii}$ v tomto případě neznačí výběrový rozptyl, ale rozptyl momentový.

Stanovení počtu faktorů

Obecně počet faktorů není známý. Uvažujme **reziduální matici**

$$\mathbf{S} - (\hat{\mathbf{L}}\hat{\mathbf{L}}' + \hat{\mathbf{\Psi}})$$

získanou aproximací výběrové kovarianční matice \mathbf{S} . Diagonální prvky jsou nulové, a pokud jsou i ostatní prvky matice malé, lze usoudit, že zvolený počet faktorů m je dostatečný. Pro jeho stanovení se při metodě hlavních komponent užívá podíl variability vysvětlené pomocí zvoleného počtu faktorů a celkové variability. Podíl j -tého faktoru na celkové výběrové variabilitě je

$$\frac{\hat{\lambda}_j}{s_{11} + s_{22} + \dots + s_{pp}} = \frac{\hat{\lambda}_j}{\hat{\lambda}_1 + \hat{\lambda}_2 + \dots + \hat{\lambda}_p},$$

kde $\hat{\lambda}_1, \hat{\lambda}_1 \geq \hat{\lambda}_2 \geq \dots \geq \hat{\lambda}_p$ jsou vlastní čísla matice \mathbf{S} . Máme-li m faktorů, pak jejich relativní příspěvek k celkové variabilitě je

$$\frac{\sum_{j=1}^m \hat{\lambda}_j}{\sum_{i=1}^p s_{ii}} = \frac{\sum_{j=1}^m \hat{\lambda}_j}{\sum_{i=1}^p \hat{\lambda}_i}.$$

Snahou je nalézt takový počet faktorů m , při kterém je tento podíl dostatečně blízký 1.

Stanovení počtu faktorů

U metody maximální věrohodnosti je možné podíl j -tého faktoru na celkové variabilitě vyjádřit zlomkem

$$\frac{\hat{l}_{1j}^2 + \hat{l}_{2j}^2 + \dots + \hat{l}_{pj}^2}{s_{11} + s_{22} + \dots + s_{pp}}$$

Za předpokladu normality lze odvodit test adekvátnosti modelu s m společnými faktory. Jedná se o test věrohodnostním poměrem a je založen a testovací statistice

$$-2 \ln \left(\frac{|\hat{\Sigma}|}{|\mathbf{S}_n|} \right)^{-n/2} + n \left[\text{tr}(\hat{\Sigma}^{-1} \mathbf{S}_n) - p \right]$$

kde $\hat{\Sigma} = \hat{\mathbf{L}}\hat{\mathbf{L}}' + \hat{\Psi}$, $\hat{\mathbf{L}}$, $\hat{\Psi}$ jsou odhady matic \mathbf{L} a Ψ získané metodou maximální věrohodnosti a $\mathbf{S}_n = \frac{n-1}{n} \mathbf{S}$. Tato testovací statistika má asymptoticky rozdělení $\chi^2(\nu)$ se stupni volnosti $\nu = \frac{1}{2}[(p-m)^2 - (p+m)]$. Lze dokázat, že výraz $\text{tr}(\hat{\Sigma}^{-1} \mathbf{S}_n) - p = 0$, za předpokladu, že $\hat{\Sigma} = \hat{\mathbf{L}}\hat{\mathbf{L}}' + \hat{\Psi}$ je maximálně věrohodný odhad $\Sigma = \mathbf{L}\mathbf{L}' + \Psi$. Dostáváme tedy testovací statistiku

$$n \ln \left(\frac{|\hat{\Sigma}|}{|\mathbf{S}_n|} \right). \quad (7)$$

Stanovení počtu faktorů

Je-li realizace této statistiky větší než příslušný kvantil χ^2 rozdělení, pak zamítáme hypotézu o dostatečném počtu faktorů. Aproximaci rozdělení testovací statistiky χ^2 rozdělením lze zpřesnit nahrazením n v testovací statistice (7) hodnotou $n - 1 - (2p + 4m + 5)/6$. Vzhledem k tomu, že stupně volnosti jsou kladné, pro použití testu musí platit nerovnost

$$m < \frac{1}{2}(2p + 1 - \sqrt{8p + 1}).$$

Rotace faktorů

Lepší interpretaci faktorů je možno získat po provedení ortogonální transformace společných faktorů. Takovou ortogonální transformací (rotací) neporušíme předpoklady modelu.

Je-li $\hat{\mathbf{L}}$ odhad $p \times m$ matice faktorových zátěží, potom

$$\hat{\mathbf{L}}^* = \hat{\mathbf{L}}\mathbf{T}, \quad \text{kde } \mathbf{T}\mathbf{T}' = \mathbf{T}'\mathbf{T} = \mathbf{I},$$

je $p \times m$ matice „rotovaných“ faktorových zátěží. Navíc odhadnuté kovarianční (nebo korelační) matice zůstávají nezměněny, neboť

$$\hat{\mathbf{L}}\hat{\mathbf{L}}' + \hat{\mathbf{\Psi}} = \hat{\mathbf{L}}\mathbf{T}\mathbf{T}'\hat{\mathbf{L}}' + \hat{\mathbf{\Psi}} = \hat{\mathbf{L}}^*\hat{\mathbf{L}}^{*'} + \hat{\mathbf{\Psi}}. \quad (8)$$

Z rovnice (8) lze usoudit, že reziduální matice $\mathbf{S}_n - \hat{\mathbf{L}}\hat{\mathbf{L}}' - \hat{\mathbf{\Psi}} = \mathbf{S}_n - \hat{\mathbf{L}}^*\hat{\mathbf{L}}^{*'} - \hat{\mathbf{\Psi}}$ zůstávají nezměněny. Totéž platí pro specifické rozptyly $\hat{\psi}_i$ a komunalitu \hat{h}_i^2 .

Rotace faktorů

Původní matice faktorových zátěží nemusí být snadno interpretovatelná. V praxi je obvyklé provádět takovou rotaci, která umožní snadnější interpretaci. V ideálním případě docílit toho, aby každá proměnná byla silně zastoupena v jednom faktoru a v ostatních faktorech se již téměř nevyskytovala.

Jedním z nejčastěji používaných kritérií optimální ortogonální transformace je **varimax** kritérium. Definujme $\tilde{l}_{ij}^* = \hat{l}_{ij}^* / \hat{h}_i$ jako rotované koeficienty škálované pomocí druhých odmocnin komunalit. Cílem je najít takovou ortogonální transformaci \mathbf{T} , která maximalizuje výraz

$$V = \frac{1}{p} \sum_{j=1}^m \left[\sum_{i=1}^p \tilde{l}_{ij}^{*4} - \frac{1}{p} \left(\sum_{i=1}^p \tilde{l}_{ij}^{*2} \right)^2 \right].$$

Po nalezení matice \mathbf{T} přenásobíme získané zátěže \tilde{l}_{ij}^* konstantami \hat{h}_i a získáme potřebné zátěže \hat{l}_{ij}^* . Varimax kritérium tedy minimalizuje počet proměnných vysvětlovaných jedním faktorem.

Odhad společných faktorů

Při faktorové analýze se pozornost obvykle zaměřuje na parametry faktorového modelu, nicméně odhady hodnot společných faktorů, které se nazývají **faktorové skóry**, mohou být také užitečné. Tyto hodnoty jsou často používány pro diagnostické účely, případně jako vstupy pro následující analýzy.

Faktorové skóry nejsou odhady neznámých parametrů v obvyklém smyslu, jsou to odhady hodnot nepozorovaných náhodných faktorů F_j , $j = 1, 2, \dots, n$. Faktorové skóry \hat{f}_j tedy jsou odhady f_j získané pro F_j .

Považujeme-li nyní získané odhady matic L a Ψ za skutečné pevné hodnoty, můžeme odhadnout společné faktory buď váženou metodou nejmenších čtverců, nebo metodou regresní.

Odhad společných faktorů – vážená metoda nejmenších čtverců

Předpokládejme, že vektor středních hodnot $\boldsymbol{\mu}$, matice faktorových zátěží \mathbf{L} a matice specifických rozptylů jsou ve faktorovém modelu známe, tedy

$$\underset{(p \times 1)}{\mathbf{X}} - \underset{(p \times 1)}{\boldsymbol{\mu}} = \underset{(p \times m)}{\mathbf{L}} \underset{(m \times 1)}{\mathbf{F}} + \underset{(p \times 1)}{\boldsymbol{\epsilon}}$$

Považujme specifické faktory za chybové složky $\boldsymbol{\epsilon} = (\epsilon_1, \epsilon_2, \dots, \epsilon_p)'$. Protože rozptyly $D(\epsilon_i) = \psi_i$, $i = 1, 2, \dots, p$, nemusí být stejné, použijeme pro odhad společných faktorů váženou metodu nejmenších čtverců. Součet čtverců vážený převrácenými hodnotami rozptylů je

$$\sum_{i=1}^p \frac{\epsilon_i^2}{\psi_i} = \boldsymbol{\epsilon}' \boldsymbol{\Psi}^{-1} \boldsymbol{\epsilon} = (\mathbf{x} - \boldsymbol{\mu} - \mathbf{L}\mathbf{f})' \boldsymbol{\Psi} (\mathbf{x} - \boldsymbol{\mu} - \mathbf{L}\mathbf{f}). \quad (9)$$

Minimalizací výrazu (9) lze získat odhad

$$\hat{\mathbf{f}} = \left(\hat{\mathbf{L}}' \hat{\boldsymbol{\Psi}}^{-1} \hat{\mathbf{L}} \right)^{-1} \hat{\mathbf{L}}' \hat{\boldsymbol{\Psi}}^{-1} (\mathbf{X}_j - \boldsymbol{\mu}). \quad (10)$$

Odhad společných faktorů – vážená metoda nejmenších čtverců

Dosazením odhadů do (10) získáme faktorové skóry ve tvaru

$$\hat{f}_j = \left(\hat{\mathbf{L}}' \hat{\Psi}^{-1} \hat{\mathbf{L}} \right)^{-1} \hat{\mathbf{L}}' \hat{\Psi}^{-1} (\mathbf{X}_j - \bar{\mathbf{X}}), \quad j = 1, 2, \dots, n.$$

V případě, že byla faktorizována korelační matice, jsou faktorové skóry dány vztahem

$$\hat{f}_j = \left(\hat{\mathbf{L}}_z' \hat{\Psi}_z^{-1} \hat{\mathbf{L}}_z \right)^{-1} \hat{\mathbf{L}}_z' \hat{\Psi}_z^{-1} \mathbf{z}_j, \quad j = 1, 2, \dots, n,$$

kde \mathbf{z}_j jsou standardizované proměnné a $\hat{\rho} = \hat{\mathbf{L}}_z \hat{\mathbf{L}}_z' + \hat{\Psi}_z$. Faktorové skóry mají nulovou střední hodnotu a nulovou výběrovou kovarianci.

Odhad společných faktorů – regresní metoda

Regresní metoda vede k odhadu

$$\hat{\mathbf{f}}_j = \hat{\mathbf{L}}' \left(\hat{\mathbf{L}}\hat{\mathbf{L}}' + \hat{\mathbf{\Psi}} \right)^{-1} (\mathbf{X}_j - \bar{\mathbf{X}}), \quad j = 1, 2, \dots, n.$$

Z důvodu snížení efektu možného nesprávného určení počtu faktorů v modelu, je někdy místo matice $\hat{\mathbf{\Sigma}} = \hat{\mathbf{L}}\hat{\mathbf{L}}' + \hat{\mathbf{\Psi}}$ použita matice výběrová kovarianční matice \mathbf{S} . Odhady potom mají tvar

$$\hat{\mathbf{f}}_j = \hat{\mathbf{L}}' \mathbf{S}^{-1} (\mathbf{X}_j - \bar{\mathbf{X}}), \quad j = 1, 2, \dots, n.$$

Byla-li faktorizována korelační matice, potom

$$\hat{\mathbf{f}}_j = \hat{\mathbf{L}}' \mathbf{R}^{-1} \mathbf{z}, \quad j = 1, 2, \dots, n.$$