

因子分析的精确模型及其解

林海明

(广东商学院 经济贸易与统计学院 广州 510320)

摘要 :本文从找因子分析精确解的角度 ,以主成份分析理论为基础 ,应用矩阵运算方法 ,建立了新的因子分析模型 ,消除了理论假设的误差 ,给出了因子分析模型的精确解 ,找到了因子分析与主成份分析的关系式。

关键词 :因子分析模型 L ;精确解 ;主成分分析 ;关系式

中图分类号 :O151 文献标识码 :A 文章编号 :1002-6487(2006)07-0004-02

1 问题的提出

多元统计方法的因子分析自 1904 年 Charles Spearman 提出以来 ,它的应用越来越广泛 ,涉及到经济、管理、金融、保险、房地产、生物、医疗、环保、体育等众多领域 ,“但是因子分析的模型和理论还是很不完美的 ,从数学上看 ,还存在许多问题。”^[1]能否进一步完善因子分析的理论模型呢 ?本文从找因子分析精确解的角度 ,以主成份分析理论为基础 ,并应用矩阵运算方法 ,对此进行了研究。

2 因子分析原模型和理论的一些缺陷

设 $(X_1, \dots, X_p)'$ 为标准化随机向量 ($p \geq 2$) R 为相关系数矩阵。

因子分析原模型为^[2] :

$$X = B_m Z_m + \varepsilon$$

$\text{Var} Z_m = I_m$ (单位矩阵 ,以下同) , $\text{Var} \varepsilon = \psi = \text{diag}(\psi_1, \dots, \psi_p)$, $\text{cov}(Z_m, \varepsilon) = 0$ 。这里 $B_m = (b_{ij})_{p \times m}$ 为因子载荷阵 $Z_m = (Z_1, \dots, Z_m)'$ 为公共因子向量 ($m < p$) $\varepsilon = (\varepsilon_1, \dots, \varepsilon_p)'$ 为特殊因子向量 $\psi_i (> 0)$ 为个性方差 [$\text{diag}(c_1, \dots, c_k)$ 是以 c_1, \dots, c_k 为对角元素的对角矩阵 (以下同)]。

因子分析原模型的缺陷有 :因子载荷阵 B_m 是估计的 ,因子 Z_m 是估计的 , $\text{Var} \varepsilon$ 为对角阵总是不能实现^[1] $R - \psi = B_m B_m'$ 有时出现负特征值^[2] !Thompson、Bartlett 因子得分哪一个更好尚无定论^[3]等。这些缺陷导致对因子分析模型和理论的认识有时模糊不清 ,如相当数量的论文和专著将因子分析与主成份分析混为一谈。

3 因子分析的精确模型及其解

因子分析模型 L :设秩 $(R) = r (\leq p)$,求 $B = (b_{ij})_{p \times r}$ 因子向量 $Z = (Z_1, \dots, Z_r)'$,使 :

$$X = BZ \quad \text{a.e.} \quad (r=p \text{ 时,无 a.e.}) \quad (1)$$

$$\text{Var} Z = I_r \quad (2)$$

Z 的前 $m (\leq r)$ 个因子对的方差贡献和达到最大 ,因子载荷阵

B 的前 m 列达到方差最大化^[2] (m 通常以前 m 个因子 Z_m 对 X 的方差贡献率 $\geq 85\%$ 、变量 X 不出现丢失、 $|b_{ij}|$ 差异大^[6]确定)。

为了与原模型比较 ,记 $B = (B_m, B_e)$, $Z = (Z_m', \varepsilon_1')'$, $B_m = (b_{ij})_{p \times m}$, $B_e = (b_{ij})_{p \times (r-m)}$, $Z_m = (Z_1, \dots, Z_m)'$, $\varepsilon_1 = (\varepsilon_{m+1}, \dots, \varepsilon_r)'$,于是因子分析模型 L 可等价写为 :

$$X = B_m Z_m + B_e \varepsilon_1 \quad \text{a.e.} \quad (r=p \text{ 时,无 a.e.}) \quad (1)'$$

$$\text{Var} Z_m = I_m, \quad \text{Var} \varepsilon_1 = I_{r-m}, \quad \text{cov}(Z_m, \varepsilon_1) = 0 \quad (2)'$$

这个模型有如下特点 :

① 待定因子 $Z = (Z_m', \varepsilon_1')'$ 只有 r 个。

② 前 m 个因子个数和相应因子 Z_m 能根据近似精度的需求确定 (数的近似精度也不过如此) ,达到了降维的目的 ,命名清晰性达到最大化。

③ 后 $r-m$ 个因子 ε_1 视为特殊因子 ,其在每个变量上是否有载荷由问题本身确定 (不存在与 x_i 专有的特殊因子 ε_i)。

④ 误差项 $B_e \varepsilon_1$ 的 $\text{Var} B_e \varepsilon_1 = B_e B_e' \neq$ 对角矩阵。

⑤ 该模型不象原模型会出现缺陷。同时 ,因子分析模型 L 包含了 Thompson 因子得分 (见推论 2) ,所以常用的就是这个因子分析模型。

本文从找因子分析精确解的角度 ,以主成份分析理论为基础 ,应用矩阵运算方法 ,给出相应结论。

由实对称矩阵中的定理有 : R 的特征值为 $\lambda_1, \lambda_2, \dots, \lambda_r, 0, \lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_r > 0$,存在实矩阵 $A = (a_{ij})_{p \times p} = (\alpha_1, \dots, \alpha_p)$,使

$$R = A \text{diag}(\lambda_1, \dots, \lambda_r, 0, \dots, 0) A' = A_r \text{diag}(\lambda_1, \dots, \lambda_r) A_r' \quad (3)$$

这里 $R \alpha_i = \lambda_i \alpha_i$, $R \alpha_k = 0$, $k = r+1, \dots, p$, $A' A = A A' = I_p$;

$$A_r = (a_{ij})_{p \times r} = (\alpha_1, \alpha_2, \dots, \alpha_r)。$$

设主成分 $F = (F_1, \dots, F_r)'$,则主成分分析 (1933 年 Hotelling 提出) 的解^[1] : $F = A' X$

$$\text{Var} F = \text{diag}(\lambda_1, \dots, \lambda_r, 0, \dots, 0)。$$

设 $A_m = (a_{ij})_{m \times m}$, $m \leq r$, $C = (c_{ij})_{m \times m}$ 为初始因子载荷矩阵 B_0 的方差最大化正交旋转矩阵^[2]。这里初始因子载荷矩阵

$$B_0 = A_m \text{diag}(\sqrt{\lambda_1}, \sqrt{\lambda_2}, \dots, \sqrt{\lambda_m}) = (\alpha_1 \sqrt{\lambda_1}, \alpha_2 \sqrt{\lambda_2}, \dots, \alpha_m \sqrt{\lambda_m})$$

记 $F_r = (F_1, \dots, F_r)'$, $F_{p-r} = (F_{r+1}, \dots, F_p)'$, $A = (A_r, A_{p-r})$, $A_{p-r} = (\alpha_{r+1}, \dots, \alpha_p)$ 。

定理 :因子分析模型 L 的解为

$$B = A_r \text{diag}(\sqrt{\lambda_1}, \sqrt{\lambda_2}, \dots, \sqrt{\lambda_m}) \text{diag}(C, I_{r-m}) \quad (6)$$

$$Z=(Z_1, \dots, Z_p)'=(Z_m' F_{m+1}/\sqrt{\lambda_{m+1}}, \dots, F_p/\sqrt{\lambda_p})', m \leq r \leq p \quad (7)$$

$$\text{这里 } Z_m=C'(F_1/\sqrt{\lambda_1}, \dots, F_m/\sqrt{\lambda_m})'=C' \begin{pmatrix} \lambda_1^{-1} & & \\ & O & \\ & & \lambda_m^{-1} \end{pmatrix} B_0' X \quad (7)'$$

在 B 确定的前提下 Z 唯一。

意义: ①式(7)'给出了因子与主成份的关系式, 其一般性(没有 R^{-1} 存在的要求)取代了现行文献的因子得分函数公式 $Z_m=B_m'R^{-1}X$ (R^{-1} 不存在时 Z_m 无意义)。②由式(7)'可计算样本因子得分值矩阵 $(Z_{ij})_{n \times m}$, 由此可对**样本进行排名与综合评价**, 也可与样本主成份值矩阵比较得出更好的结果, 这里

$$(Z_{ij})_{n \times m}=(F_{ij})_{n \times m} \text{diag}(\sqrt{\lambda_1^{-1}}, \sqrt{\lambda_2^{-1}}, \dots, \sqrt{\lambda_m^{-1}})C$$

其中 $(F_{ij})_{n \times m}$ 为样本主成份值矩阵。

证明: 因为 $A'A=I_p$, 所以由式(4)有 $X=AF$, 即

$$X=A_p F_p + A_{p-r} F_{p-r} \quad (4)'$$

取 $B=A_p \text{diag}(\sqrt{\lambda_1}, \sqrt{\lambda_2}, \dots, \sqrt{\lambda_r}) \text{diag}(C, I_{r-m})$ 得式(6);

$$Z=(Z_m' F_{m+1}/\sqrt{\lambda_{m+1}}, \dots, F_p/\sqrt{\lambda_p})' Z_m=C'(F_1/\sqrt{\lambda_1}, \dots, F_m/\sqrt{\lambda_m}) \text{得式(7)} \quad \varepsilon_0=A_{p-r} F_{p-r}$$

可直接验证 $X=BZ+\varepsilon_0$, $\text{Var} Z=I_r$ 得式(2)。

由式(5)有 $\text{Var} F_{p-r}=0$, 所以 $\text{Var} \varepsilon_0=A_{p-r} \text{Var} F_{p-r} A_{p-r}'=0$, 因为 $X=(X_1, \dots, X_p)'$ 为标准化随机向量, 所以 $E \varepsilon_0=0$, 由 $\text{Var} \varepsilon_0=0$, 得 $\varepsilon_0=0$ a.e., 得式(1) ($\varepsilon_0=0$ a.e. 结论的证明由张涤新教授给出)。

现证 Z 的前 $m(\leq r)$ 个因子 Z_m 对 X 的方差贡献和达到最大。

因子 Z_i 对变量 X 的方差贡献 $v_i=B$ 的第 i 列元素的平方和 $=\sum_{k=1}^p b_{ki}^2$, 所以前 $m(\leq P)$ 个因子 $Z_m=(Z_1, \dots, Z_m)'$ 对 X 的方差贡献和

$$\sum_{k=1}^m v_k=\sum_{i=1}^m \sum_{k=1}^p b_{ki}^2=\text{tr} B_m' B_m (\text{tr 为取矩阵对角线元素的和}),$$

这里 $B_m=(b_{ij})_{p \times m}=A_m \text{diag}(\sqrt{\lambda_1}, \dots, \sqrt{\lambda_m})C$ (8)

所以 $\text{tr}(B_m' B_m)=\text{tr}[C' \text{diag}(\lambda_1, \dots, \lambda_m)C]=\text{tr}[\text{diag}(\lambda_1, \dots, \lambda_m)]=\sum_{k=1}^m \lambda_k$ 。

因为 λ_i 是相关系数矩阵 R 的特征值, $\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_p > 0$, 故这是最大化的降序排列, 所以前 $m(\leq P)$ 个因子 $Z_m=(Z_1, \dots,$

$Z_m)'$ 对 X 的方差贡献和 $\sum_{k=1}^m v_k=\sum_{k=1}^m \lambda_k$ 达到最大。

因为 $B_m=(b_{ij})_{p \times m}=A_m \text{diag}(\sqrt{\lambda_1}, \dots, \sqrt{\lambda_m})C$ 是 B 的前 m 列组成的矩阵, 而 C 为矩阵 $A_m \text{diag}(\sqrt{\lambda_1}, \dots, \sqrt{\lambda_m})$ 的方差最大化正交旋转矩阵, 故 B 的前 m 列达到方差最大化。得证。

推论 1 方差最大化正交旋转矩阵 C 的统计意义为: 因子 Z_j 与主成分 F_i 的相关系数

$$r_{F_i Z_j}=c_{ij}, \quad 1 \leq i, j \leq m$$

$$r_{F_i Z_j}=\delta_{ik}, \quad i=1, \dots, r, k=m+1, \dots, p, \delta_{ik}=\begin{cases} 0 & i \neq k \\ 1 & i=k \end{cases} \quad (\text{证明略})$$

推论 2 $|R| \neq 0$ 时, Thompson 因子得分函数 $Z_m^*=Z_m$ 。

证明: $|R| \neq 0$ 时, R^{-1} 存在 $Z_m^*=B_m'R^{-1}X$ [2]

由式(8)、式(3)、式(7)'得:

$$Z_m^*=C' \text{diag}(\sqrt{\lambda_1}, \dots, \sqrt{\lambda_m}) A_m' A_m \text{diag}(\lambda_1^{-1}, \dots, \lambda_p^{-1}) A' X$$

由 $AA'=I_p$ 和式(4)得 $Z_m^*=Z_m$ 。

推论 2 表明因子分析中, 通过主成分法提取初始因子进

行方差最大化正交旋转, 求 Thompson 因子得分函数得到因子估计 Z_m^* 是因子分析模型 L 的前 m 个因子解 Z_m , 所以因子分析模型 L 是常用的因子分析模型, 这也是此类因子分析总能通过的原因。

对于因子分析模型 L, 式(7)的因子解 Z 是精确解, 故式(7)的因子解 Z 是无偏的, 平均预报误差为零。即因子解 Z 是精度最好的结果, 也是 X 的线性函数。

由式(1)'知 B_m, Z_m 可经各种正交、等价变换(约束)得到各种不同需求的结果。

$$\text{记 } B=(B_m, B_r), B_m=A_m \text{diag}(\sqrt{\lambda_1}, \dots, \sqrt{\lambda_m})C$$

$$B_r=A_{r-m} \text{diag}(\sqrt{\lambda_{m+1}}, \dots, \sqrt{\lambda_r})$$

有 $R=B_m B_m' \cong \Psi_1=B_r B_r' \neq$ 对角矩阵[但 $B_r' B_r=\text{diag}(\lambda_{m+1}, \dots, \lambda_r)$], $R-\Psi_1=B_m B_m'$ 特征值为 $\lambda_1, \lambda_2, \dots, \lambda_m, 0$, 且非负。

4 因子分析与主成分分析的异同

由式(1)式(4)式(7): $X=BZ=AF$, Z、F 可相互表示, 故主成分分析与因子分析本质上都是对数据向量 X 精确描述的线性等价降维技术。

从过程上看, 由式(7), 因子分析模型 L 的因子是对主成分施行了两个过程的结果: 过程 1 对主成分伸缩了 $\sqrt{\lambda_i}^{-1}$ 个单位后得到初始因子。过程 2 再对初始因子进行方差最大化正交旋转使得因子命名清晰性达到最大化方向。

从计量上看, 只要 $\lambda_i \neq 1$, 即 $\text{Var} F_i \neq \text{Var} Z_i$, 则 F_i 与 Z_i 的取值范围不同, $F_i \neq Z_i$ 永远成立, $i=1, \dots, m$ 。

从应用侧重上看, 因子分析应用侧重命名清晰性达到最大化方向, 由式(7)因子 Z 是经 X 非标准正交线性变换(有伸缩)的结果; 主成分分析应用侧重信息贡献依次达到最大化方向, 由式(4)F 是经 X 标准正交线性变换的结果。

综合以上讨论可得以下结论:

①主成分分析与因子分析本质上都是对数据向量精确描述的线性等价降维技术, 但过程、计量值、应用侧重上不同, 两者绝对不可相互替代、混淆。

②主成分分析有时命名清晰, 既达到信息贡献影响力综合评价效果, 又达到命名清晰性的综合评价效果, 此时主成分分析的结果多数优于因子分析的结果。

至此, 本文从找因子分析精确解的角度, 以主成份分析理论为基础, 应用矩阵运算方法, 进一步完善了因子分析模型理论和方法。

参考文献:

- [1] 张尧庭, 方开泰著. 多元统计分析引论[M]. 北京: 科学出版社, 1997.
- [2] Richard A. Johnson 等著. 实用多元统计分析(第四版)[M]. 陆璇译. 北京: 清华大学出版社, 2001.
- [3] Mardia, K.V., Kent, J.T. and Bibby, J.M. Multivariate Analysis[M]. Academic Press, New York, 1979.
- [4] 林海明, 张文霖. 主成分分析与因子分析详细的异同和 SPSS 软件[J]. 统计研究, 2005, (3).
- [5] 林海明, 林敏子, 丁洁花. 主成分分析法与因子分析法应用辨析[J]. 数量经济技术经济研究, 2004, (9): 155-161.
- [6] 林海明. 运用因子分析法综合评价广东烟草工业经济效益[J]. 数学的实践与认识, 2005, (6): 61-65.

(责任编辑/李友平)