

【统计理论与方法】

# 因子分析精确模型的基本思想与方法

林海明

(广东商学院 经济贸易与统计学院, 广东 广州 510320)

摘要: 文章从统计思想、等价性的方法入手, 给出了初始因子分析精确模型及解、因子分析精确模型及解、主成分分析与因子分析的关系式等结论。从基本思想、方法上完善了因子分析精确模型和理论。

关键词: 因子分析; 精确模型; 基本思想; 方法

中图分类号: O212.4 文献标识码: A 文章编号: 1007-3116(2006)05-0023-03

参考文献[1]中笔者给出了因子分析精确模型及解, 但没有给出方法的基本思想, 理解上有待深入。本文进一步从统计思想、等价性的简捷方法入手, 给出了初始因子分析精确模型及解、因子分析精确模型及解、主成分分析与因子分析的关系式等结论。从基本思想、方法上完善了因子分析精确模型和理论。

## 一、基本思想与方法

由于因子分析中常用主成分法提取初始因子载荷阵, 而主成分分析法是信息贡献最大化模型的精确解, 故本文解决问题的基本思想是从主成分分析模型和理论入手, 用等价性的方法建立因子分析精确模型和理论。具体为: 1. 从主成分分析模型精确解——变量表示主成分的表示式入手, 依据因子分析的约束要求, 建立初始因子分析精确模型和理论。2. 依据因子命名清晰性要求, 等价变换初始因子分析精确模型, 建立因子分析精确模型和理论。3. 求出因子分析精确模型的解——变量表示因子的表示式。

## 二、因子分析精确模型和解

从主成分分析模型精确解——变量表示主成分的表示式入手, 依据因子分析的约束要求, 建立初始因子分析的精确模型和理论。

设  $\mathbf{X} = (x_1, \dots, x_p)^T$  为正向化、标准化随机向量 ( $p \geq 2$ ),  $\mathbf{R}$  为相关系数矩阵, 秩  $(\mathbf{R}) = r$  ( $\mathbf{R}$  的非零特征根个数)。

主成分分析: 设  $\mathbf{R}$  的特征值为  $\lambda_1, \lambda_2, \dots, \lambda_r, 0, \lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_r > 0$ ,  $\mathbf{A} = (\alpha_1, \dots, \alpha_p)$ ,

$$\mathbf{A}^T \mathbf{A} = \mathbf{A} \mathbf{A}^T = \mathbf{I}_p, \mathbf{A}_r = (\alpha_1, \dots, \alpha_r)$$

有:

$$\mathbf{R} = \mathbf{A} \text{diag}(\lambda_1, \dots, \lambda_r, 0, \dots, 0) \mathbf{A}^T \quad (1)$$

这里  $\mathbf{R} \alpha_i = \lambda_i \alpha_i, i = 1, \dots, r, \mathbf{R} \alpha_k = 0, k = r+1, \dots, p$ 。

设主成分  $\mathbf{F} = (f_1, \dots, f_p)^T$ , 记  $\mathbf{F}_r = (f_1, \dots, f_r)^T, \mathbf{F}_m = (f_1, \dots, f_m)^T, \mathbf{F}_\epsilon = (f_{m+1}, \dots, f_r)^T$ ,

则主成分分析(1933 年 Hotelling 给出)的解:

$$\mathbf{F} = \mathbf{A}^T \mathbf{X} \quad (2)$$

$$\text{Var} \mathbf{F} = \text{diag}(\lambda_1, \dots, \lambda_r, 0, \dots, 0) \quad (3)$$

设  $\mathbf{A}_r = (\mathbf{A}_m, \mathbf{A}_\epsilon), m \leq r, \mathbf{A}_m = (\alpha_1, \dots, \alpha_m), \mathbf{A}_\epsilon = (\alpha_{m+1}, \dots, \alpha_r)$ 。

记因子载荷阵  $\mathbf{B}_m = (b_{ij})_{p \times m}$ , 特殊因子载荷阵  $\mathbf{B}_\epsilon = (b_{im+j})_{p \times (r-m)}$ , 因子向量  $\mathbf{Z}_m = (z_1, \dots, z_m)^T$ , 特殊因子向量  $\epsilon_1 = (z_{m+1}, \dots, z_r)^T$

### (一) 初始因子分析精确模型

设秩  $(\mathbf{R}) = r (\leq p)$ , 求  $\mathbf{B}_m, \mathbf{Z}_m, \mathbf{B}_\epsilon, \epsilon_1$ , 使:

$$\mathbf{X} = \mathbf{B}_m \mathbf{Z}_m + \epsilon \quad \text{a.e. } (n \geq p \text{ 时, 无 a.e.})$$

$$\epsilon \preceq \mathbf{B}_\epsilon \epsilon_1 \quad (4)$$

收稿日期: 2006-04-28

作者简介: 林海明(1959-), 男, 湖南省宁乡县人, 副教授, 研究方向: 多元统计学模型与应用。

$$\begin{aligned} \text{Var} \mathbf{Z}_m &= \mathbf{I}_m, \text{Var} \varepsilon_1 = \mathbf{I}_{r-m} \\ \text{cov}(\mathbf{Z}_m, \varepsilon_1) &= \mathbf{0} \end{aligned} \quad (5)$$

其中  $\varepsilon_1$  称为特殊因子,  $\varepsilon$  称为误差项 ( $\varepsilon$  随  $m$  取值的不同产生变异), 因子  $\mathbf{Z}_i$  对  $\mathbf{X}$  的方差贡献  $v_i (i = 1, \dots, r)$  按降序排列依次达到最大化。

$m (\leq r)$  通常以因子  $z_1, \dots, z_m$  所含变量  $\mathbf{X}$  不出现丢失确定。

现找出主成分表示变量的确定性关系式。

引理 主成分表示变量的确定性关系式为:

$$\mathbf{X} = \mathbf{A}_r \mathbf{F}_r \text{ a.e. } r = p \text{ 时, 无 a.e.} \quad (6)$$

证明 因为  $\mathbf{A} \mathbf{A}^T = \mathbf{I}_p$ , 所以将式(2)左乘  $\mathbf{A}$  有:

$$\mathbf{X} = \mathbf{A} \mathbf{F} = \mathbf{A}_r \mathbf{F}_r + \varepsilon_0$$

这里  $\varepsilon_0 = \mathbf{A}_0 \mathbf{F}_0, \mathbf{F}_0 = (f_{r+1}, \dots, f_p)^T, \mathbf{A}_0 = (\mathbf{a}_{r+1}, \dots, \mathbf{a}_p)$ 。

由式(3)有  $\text{Var} \mathbf{F}_0 = 0$ , 所以:

$$\text{Var} \varepsilon_0 = \mathbf{A}_0 \mathbf{Var} \mathbf{F}_0 \mathbf{A}_0^T = 0$$

因为  $\mathbf{X}$  为标准化向量, 所以  $E \varepsilon_0 = 0$ , 由  $\text{Var} \varepsilon_0 = 0$ , 得  $\varepsilon_0 = 0 \text{ a.e.}$  所以:  $\mathbf{X} = \mathbf{A}_r \mathbf{F}_r \text{ a.e.}$

现对主成分标准化并进行等价运算, 给出初始因子表示变量的确定性关系、变量表示初始因子的确定性关系。

定理 1 初始因子分析精确模型的解:

$$\mathbf{B}_0 = (b_{ij})_{p \times m} = \mathbf{A} \text{diag}(\sqrt{\lambda_1}, \sqrt{\lambda_2}, \dots, \sqrt{\lambda_m})$$

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

$$\mathbf{Z}_m^0 = (z_1^0, \dots, z_m^0)^T$$

$$= \text{diag}(\sqrt{\lambda_1^{-1}}, \sqrt{\lambda_2^{-1}}, \dots, \sqrt{\lambda_m^{-1}}) \mathbf{F}_m$$

$$= \text{diag}(\lambda_1^{-1}, \lambda_2^{-1}, \dots, \lambda_m^{-1}) \mathbf{B}_0^T \mathbf{X}$$

$$\mathbf{B}_\varepsilon = \mathbf{A}_\varepsilon \text{diag}(\sqrt{\lambda_{m+1}}, \dots, \sqrt{\lambda_r})$$

$$= (\alpha_{m+1} \sqrt{\lambda_{m+1}}, \dots, \alpha_r \sqrt{\lambda_r})$$

$$\varepsilon_1 = (\mathbf{Z}_{m+1}, \dots, \mathbf{Z}_r)^T$$

$$= \text{diag}(\sqrt{\lambda_{m+1}^{-1}}, \dots, \sqrt{\lambda_r^{-1}}) \mathbf{F}_\varepsilon$$

$$= \text{diag}(\lambda_{m+1}^{-1}, \dots, \lambda_r^{-1}) \mathbf{B}_\varepsilon^T \mathbf{X}$$

证明 为了得到式(4)、式(5), 将主成分  $\mathbf{F}_r$  标准化, 记其为:  $[(\mathbf{Z}_m^0)^T, \varepsilon_1^T]^T$ , 由式(2)有  $E \mathbf{F}_r = 0$ , 由式(3)有标准化主成分:

$$[(\mathbf{Z}_m^0)^T, \varepsilon_1^T]^T = \text{diag}(\sqrt{\lambda_1^{-1}}, \sqrt{\lambda_2^{-1}}, \dots,$$

$\sqrt{\lambda_r^{-1}}) \mathbf{F}_r$ , 同时  $\mathbf{Z}_m^0, \varepsilon_1$  是式(5)的解, 所以:

$$\mathbf{F}_r = \begin{bmatrix} \text{diag}(\sqrt{\lambda_1}, \sqrt{\lambda_2}, \dots, \sqrt{\lambda_m}) \mathbf{Z}_m^0 \\ \text{diag}(\sqrt{\lambda_{m+1}}, \dots, \sqrt{\lambda_r}) \varepsilon_1 \end{bmatrix}$$

将  $\mathbf{F}_r$  代入式(6)中, 由  $\mathbf{A}_r = (\mathbf{A}_m, \mathbf{A}_\varepsilon)$  得:

$$\begin{aligned} \mathbf{X} &= [\mathbf{A}_m \text{diag}(\sqrt{\lambda_1}, \sqrt{\lambda_2}, \dots, \sqrt{\lambda_m})] \mathbf{Z}_m^0 + \\ &[\mathbf{A}_\varepsilon \text{diag}(\sqrt{\lambda_{m+1}}, \dots, \sqrt{\lambda_r})] \varepsilon_1 \end{aligned}$$

即:

$$\mathbf{X} = \mathbf{B}_0 \mathbf{Z}_m^0 + \mathbf{B}_\varepsilon \varepsilon_1 \text{ a.e.} \quad (7)$$

因子  $\mathbf{Z}_i$  对变量  $\mathbf{X}$  的方差贡献  $v_i = \mathbf{B}_m^0$  的第  $i$  列元素的平方和  $= \sum_{k=1}^p (b_{ki}^0)^2 = \lambda_i \sum_{k=1}^p a_{ki}^2$ , 因为  $\mathbf{A}^T \mathbf{A} = \mathbf{I}_p$ , 有  $\sum_{k=1}^p a_{ki}^2 = 1$ , 所以因子  $\mathbf{Z}_i$  对变量  $\mathbf{X}$  的方差贡献  $v_i = \lambda_i, i = 1, \dots, r$ 。证毕。

式(7)是初始因子描述变量的确定性关系等式。

$\mathbf{B}_0$  的列有时达不到命名清晰(方差最大化), 故称  $\mathbf{B}_0$  为初始因子载荷矩阵,  $\mathbf{Z}_m^0$  为初始因子向量。

现依据因子命名清晰性要求, 等价变换初始因子分析模型, 建立因子分析的模型和理论。

(二) 因子分析精确模型

设秩  $(\mathbf{R}) = r (\leq p)$ , 求  $\mathbf{B}_m, \mathbf{Z}_m, \mathbf{B}_\varepsilon, \varepsilon_1$ , 使:

$$\mathbf{X} = \mathbf{B}_m \mathbf{Z}_m + \varepsilon$$

$$\text{a.e. } (n \geq p \text{ 时, 无 a.e.}), \varepsilon \cong \mathbf{B}_\varepsilon \varepsilon_1 \quad (4)$$

$$\text{Var} \mathbf{Z}_m = \mathbf{I}_m \quad \text{Var} \varepsilon_1 = \mathbf{I}_{r-m}$$

$$\text{cov}(\mathbf{Z}_m, \varepsilon_1) = \mathbf{0} \quad (5)$$

因子  $z_1, \dots, z_m$  对  $\mathbf{X}$  的累积方差贡献达到最大, 因子载荷阵  $\mathbf{B}_m$  的各列达到方差最大化。

$m (\leq r)$  通常以因子  $z_1, \dots, z_m$  所含变量  $\mathbf{X}$  不出现丢失,  $\mathbf{B}_m = (b_{ij})_{p \times m}$  中的  $|b_{ij}|$  差异大确定。

定理 2 设  $\mathbf{C} = (c_{ij})_{m \times m}$  为初始因子载荷矩阵  $\mathbf{B}_0$  的方差最大化正交旋转矩阵<sup>[1]</sup>, 则因子分析模型  $L$  中  $\mathbf{B}_m, \mathbf{Z}_m, \mathbf{B}_\varepsilon, \varepsilon_1$  的解为:

$$\mathbf{B}_m = \mathbf{B}_0 \mathbf{C} = \mathbf{A}_m [\text{diag}(\sqrt{\lambda_1}, \sqrt{\lambda_2}, \dots, \sqrt{\lambda_m})] \mathbf{C}$$

$$\begin{aligned} \mathbf{Z}_m &= \mathbf{C}^T \mathbf{Z}_m^0 = \mathbf{C}^T \text{diag}(\lambda_1^{-1}, \lambda_2^{-1}, \dots, \lambda_m^{-1}) \mathbf{B}_0^T \mathbf{X} \\ &= \mathbf{C}^T (f_1 / \sqrt{\lambda_1}, \dots, f_m / \sqrt{\lambda_m})^T \end{aligned} \quad (8)$$

$$\mathbf{B}_\varepsilon = \mathbf{A}_\varepsilon \text{diag}(\sqrt{\lambda_{m+1}}, \dots, \sqrt{\lambda_r})$$

$$= (\alpha_{m+1} \sqrt{\lambda_{m+1}}, \dots, \alpha_r \sqrt{\lambda_r})$$

$$\varepsilon_1 = (\mathbf{Z}_{m+1}, \dots, \mathbf{Z}_r)^T$$

$$= \text{diag}(\sqrt{\lambda_{m+1}^{-1}}, \dots, \sqrt{\lambda_r^{-1}}) \mathbf{F}_\varepsilon$$

$$= \text{diag}(\lambda_{m+1}^{-1}, \dots, \lambda_r^{-1}) \mathbf{B}_\varepsilon^T \mathbf{X}$$

在  $\mathbf{B}_m, \mathbf{B}_\varepsilon$  确定的前提下,  $\mathbf{Z}_m, \varepsilon_1$  唯一。

证明 对初始因子载荷矩阵  $\mathbf{B}_0$  的列用  $\mathbf{C}$  进行方差最大化正交旋转<sup>[2]</sup>, 所得矩阵为  $\mathbf{B}_m$  (命名清晰), 于是  $\mathbf{B}_m = \mathbf{B}_0 \mathbf{C}$  因为  $\mathbf{C}$  为正交矩阵, 有  $\mathbf{C} \mathbf{C}^T = \mathbf{I}_m$ , 所以  $\mathbf{B}_0 = \mathbf{B}_m \mathbf{C}^T$ , 将  $\mathbf{B}_0$  带入标准主成分表示变量的表达式(7)中, 有:

$$\mathbf{X} = \mathbf{B}_m [\mathbf{C}^T \mathbf{Z}_m^0] + \mathbf{B}_\varepsilon \varepsilon_1 = \mathbf{B}_m \mathbf{Z}_m + \mathbf{B}_\varepsilon \varepsilon_1 \text{ a.e.} \text{ 得}$$

式(4)], 即:

$$\mathbf{Z}_m = \mathbf{C}^T \mathbf{Z}_m^0 = \mathbf{C}^T \text{diag}(\lambda_1^{-1}, \lambda_2^{-1}, \dots, \lambda_m^{-1}) \mathbf{B}_0^T \mathbf{X},$$

因为  $\mathbf{C}$  为正交矩阵, 有  $\mathbf{C}\mathbf{C}^T = \mathbf{I}_m$ , 所以  $\text{Var}\mathbf{Z}_m = \mathbf{I}_m$  成立。

$\mathbf{B}_\varepsilon, \varepsilon_1$  的结论见定理 1 证明。

因子  $z_i$  对变量  $x$  的方差贡献  $v_i = \mathbf{B}_m$  的第  $i$  列元素的平方和  $= \sum_{k=1}^p b_{ki}^2$ , 所以  $z_1, \dots, z_m$  对  $\mathbf{X}$  的方差贡献和  $\sum_{i=1}^m v_i = \sum_{i=1}^m \sum_{k=1}^p b_{ki}^2 = \text{tr}\mathbf{B}_m^T \mathbf{B}_m$  ( $\text{tr}$  为取矩阵对角线元素的和), 由  $\mathbf{B}_m$  的表达式,  $\mathbf{A}\mathbf{A}^T = \mathbf{I}_p, \mathbf{C}^T \mathbf{C} = \mathbf{I}_m$  有  $\text{tr}\mathbf{B}_m^T \mathbf{B}_m = \text{tr}[\mathbf{C}^T \text{diag}(\lambda_1, \dots, \lambda_m) \mathbf{C}] = \text{tr}[\text{diag}(\lambda_1, \dots, \lambda_m)] = \sum_{i=1}^m \lambda_i$ , 因为  $\lambda_i$  是相关系数矩阵  $\mathbf{R}$  的特征值,  $\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_r > 0$ , 故这是最大化的降序排列, 所以因子  $z_1, \dots, z_m$  对  $\mathbf{X}$  的方差贡献和  $\sum_{i=1}^m v_i = \sum_{i=1}^m \lambda_i$  达到最大。证毕。

推论 1 设因子分析原模型 Thompson(1939 年) 因子得分函数为  $\mathbf{Z}_m^*$ , 则  $\mathbf{Z}_m^*$  与因子分析精确模型因子解  $\mathbf{Z}_m$  相等, 即:

$$\mathbf{Z}_m^* = \mathbf{Z}_m$$

证明  $|\mathbf{R}| \neq 0$  时,  $\mathbf{R}^{-1}$  存在,  $r = p, \mathbf{Z}_m^* = \mathbf{B}_m^T \mathbf{R}^{-1} \mathbf{X}^T$ , 由定理 2 的  $\mathbf{B}_m$  式(1)得:

$$\mathbf{Z}_m^* = \mathbf{C}^T \text{diag}(\sqrt{\lambda_1}, \sqrt{\lambda_2}, \dots, \sqrt{\lambda_m}) \mathbf{A}_m^T \mathbf{A} \text{diag}(\lambda_1^{-1}, \dots, \lambda_p^{-1}) \mathbf{A}^T \mathbf{X}, \text{ 由 } \mathbf{A}^T \mathbf{A} = \mathbf{I}_p \text{ 式(2) 得:}$$

$$\mathbf{Z}_m^* = \mathbf{Z}_m$$

$|\mathbf{R}| = 0$  时, 由  $\mathbf{R}$  的广义逆矩阵  $\mathbf{R}^+$  同理得:

参考文献:

- [1] 林海明. 因子分析精确模型及解[J]. 统计与决策, 2006(7).
- [2] RICHARD A JOHNSON, REAN W WICHERN. 实用多元统计分析[M]. 陆璇. 译. 北京: 清华大学出版社, 2001.
- [3] 林海明. 主成分分析与初始因子分析的异同[J]. 统计与决策, 2006(4).
- [4] 林海明, 张文霖. 主成分分析与因子分析的异同和 SPSS 软件[J]. 统计研究, 2005(3).

(责任编辑: 马 慧)

### The Precise Model of Factor Analysis and Its Thought

LIN Hai-ming

(School of Economics, Trade and Statistics, Guangdong University of Business Studies, Guangzhou 510320, China)

**Abstract:** This paper makes further improvement on the factor analysis from the models and thought, and then finds the correct models and their solutions, and what's more, gives equation concerning the relationship between factor analysis and principal component analysis. The main methods employed in this paper are the constructive proof and matrix operation.

**Key words:** factor analysis; precise model; principal component analysis; solutions

$$\mathbf{Z}_m^* = \mathbf{Z}_m$$

即因子分析精确模型的因子解是因子分析原模型中 Thompson 因子得分函数。

但因子分析原模型误差项就没有相应结论了。

### 三、主成分分析与因子分析的关系

由定理 1 中初始因子  $\mathbf{Z}_m^0$  的表达式得出:

推论 2 初始因子向量  $\mathbf{Z}_m^0$  是主成份向量  $\mathbf{F}_m$  的标准化向量, 初始因子与主成份的异同是方向一致、方差不等的等价关系, 即:

$$\begin{aligned} \mathbf{Z}_m^0 &= (\mathbf{z}_1^0, \dots, \mathbf{z}_m^0)^T \\ &= \text{diag}(\sqrt{\lambda_1^{-1}}, \sqrt{\lambda_2^{-1}}, \dots, \sqrt{\lambda_m^{-1}}) \mathbf{F}_m \\ &= (f_1/\sqrt{\lambda_1}, \dots, f_m/\sqrt{\lambda_m})^T \end{aligned}$$

初始因子与主成分的具体异同见参考文献[3]。

由定理 2 式(8)得出:

推论 3 因子与初始因子的异同是方差相等、方向有旋转的等价关系, 即:

$$\mathbf{Z}_m = \mathbf{C}^T \mathbf{Z}_m^0$$

当  $\mathbf{C} = \mathbf{I}_m$ , 因子分析精确模型变为初始因子分析精确模型。

定理 3: 因子与主成分的异同是因子是主成分进行标准化后施行方差最大化正交旋转得到的结果, 它们是方差不相等、方向有旋转的等价关系, 关系为式(8), 即:

$$\mathbf{Z}_m = (z_1, \dots, z_m)^T = \mathbf{C}^T (f_1/\sqrt{\lambda_1}, \dots, f_m/\sqrt{\lambda_m})^T$$

因子与主成分的具体异同见参考文献[4]。