

doi:10.3969/j.issn.1673-9833.2014.03.017

多属性数据聚类的一种因子分析新方法

张爱平, 陈志彬

(湖南工业大学 理学院, 湖南 株洲 412007)

摘要: 根据因子分析法的思想, 用统计学的方法, 建立多属性数据样本间的相似矩阵, 探索求因子载荷矩阵的有效方法, 实现多属性数据的样本聚类。文中的方法是因子分析法在聚类分析中的进一步推广与应用。

关键词: 多属性; 样本; 相似矩阵; 数据聚类

中图分类号: O212

文献标志码: A

文章编号: 1673-9833(2014)03-0083-05

A New Factor Analysis Method in Multiple Attribute Data Clustering

Zhang Aiping, Chen Zhibin

(School of Science, Hunan University of Technology, Zhuzhou Hunan 412007, China)

Abstract: According to the thought of factor analysis method established the similar matrix between multiple attribute data samples by statistical methods, and explored the effective solution to the factor loading matrix for the realization of sample clustering of multiple attribute data. This method is the factor analysis method further extended and applied in clustering analysis.

Keywords: multiple attribute; sample; similar matrix; data clustering

0 引言

因子分析法是通过建立数学模型, 用线性代数方法, 研究多属性变量间的内在依赖关系。用少于属性变量个数的几个抽象变量表示被观测数据的基本结构, 实现对被观测数据变量的降维, 达到简化数据结构的目的。这几个抽象的变量通常被称为因子, 每个因子综合地包含了多个属性变量的信息, 是一些异于可观测原始变量的不可观测的潜在变量。因子分析的内容较丰富, 常见的类型可以概括为两类^[1-5]: 一类是R型因子分析, 另一类是Q型因子分析, 前者是基于变量间的相关系数矩阵, 后者则基于样本间的相似矩阵, 两种矩阵均为非负定矩阵。

这两种类型选择因子分析的对象和计算的出发点不同但方法类似。在实际问题中, 由于被观测的样本数目 n 通常较大, 因此Q型因子分析中的样本相似矩阵是一个阶数较高的 n 阶方阵, 其计算量与 n^2 同阶且可能是非正定的; 而求解样本相似矩阵的特征根与特征向量的计算量与 n^3 同阶。由于计算量随阶数 n 的增大而急剧增大, 这就限制了以样本为变量的Q型因子分析法在经济学、生物学和社会学等领域中的应用。

为此, 本文根据高阶样本相似矩阵与因子载荷矩阵的关系, 通过间接地求解一个与高阶样本相似矩阵有联系的低阶矩阵的特征根与特征向量, 探讨因子载荷矩阵的计算方法。

收稿日期: 2014-03-10

基金项目: 湖南省教育科学研究基金资助项目(10C0656), 湖南省教育改革基金资助项目(288)

作者简介: 张爱平(1967-), 女, 湖南冷水江人, 湖南工业大学副教授, 主要从事应用数学方面的教学与研究,

E-mail: zaping@163.com

1 预备知识

文中所讨论的随机向量是具有 p 种属性变量的总体, 用列向量 $\mathbf{X} = (X_1, X_2, \dots, X_p)^T$ 表示。 n 个样本为源于总体 n 次观测所获得的 n 个 p 维数据组 ($p \ll n$); 第 i 个样本用列向量 $\mathbf{Y}_i = (y_{i1}, y_{i2}, \dots, y_{ip})^T$ 表示, 均值 $\bar{y}_i = \frac{1}{p} \sum_{k=1}^p y_{ik}$; n 个样本的数据矩阵表示为 $\mathbf{Y} = (y_{ij})_{n \times p}$ 或用列向量表示为 $\mathbf{Y} = (\mathbf{Y}_1, \mathbf{Y}_2, \dots, \mathbf{Y}_n)^T$ 。为了方便, 在文中引入如下记号:

样本向量的均值

$$\bar{\mathbf{Y}} = \left(\frac{1}{p} \sum_{k=1}^p y_{1k}, \frac{1}{p} \sum_{k=1}^p y_{2k}, \dots, \frac{1}{p} \sum_{k=1}^p y_{nk} \right)^T = (\bar{y}_1, \bar{y}_2, \dots, \bar{y}_n)^T。$$

样本向量的离差矩阵

$$\mathbf{S} = \sum_{i=1}^p (\mathbf{Y}_{(i)} - \bar{\mathbf{Y}})(\mathbf{Y}_{(i)} - \bar{\mathbf{Y}})^T = (s_{ij})_{n \times n},$$

式中: $\mathbf{Y}_{(i)} - \bar{\mathbf{Y}} = (y_{1i} - \bar{y}_1, y_{2i} - \bar{y}_2, \dots, y_{ni} - \bar{y}_n)^T$;

$$s_{ij} = \sum_{k=1}^p (y_{ik} - \bar{y}_i)(y_{jk} - \bar{y}_j)。$$

样本向量的协方差矩阵

$$\mathbf{V} = \frac{1}{p-1} \mathbf{S} = (v_{ij})_{n \times n},$$

式中 $v_{ij} = \frac{1}{p-1} \sum_{k=1}^p (y_{ik} - \bar{y}_i)(y_{jk} - \bar{y}_j)。$

样本向量的相似阵

$$\mathbf{R} = (r_{ij})_{n \times n},$$

式中 $r_{ij} = \frac{s_{ij}}{\sqrt{s_{ii}} \sqrt{s_{jj}}} = \frac{v_{ij}}{\sqrt{v_{ii}} \sqrt{v_{jj}}}。$

令 $z_{ij} = \frac{y_{ij} - \bar{y}_i}{\sqrt{v_{ii}}}$, 将样本数据矩阵 $\mathbf{Y} = (y_{ij})_{n \times p}$ 标准化后得到矩阵 $\mathbf{Z} = (z_{ij})_{n \times p}$, 其列向量矩阵表示为

$\mathbf{Z} = (\mathbf{Z}_1, \mathbf{Z}_2, \dots, \mathbf{Z}_n)^T$, 其中 $\mathbf{Z}_i = (z_{i1}, z_{i2}, \dots, z_{ip})^T, i=1, 2, \dots, n。$

\mathbf{Z} 的离差矩阵为 \mathbf{ZZ}^T , 协方差矩阵为 $\frac{1}{p-1} \mathbf{ZZ}^T$, 标准化后样本向量 \mathbf{Z} 的协方差矩阵等于原样本向量 \mathbf{Y} 的相似矩阵, 即

$$\frac{1}{p-1} \mathbf{ZZ}^T = \mathbf{R}。$$

引理 1^[6] 实对称矩阵的不同特征值的特征向量彼此正交。

引理 2^[6] 对于 n 阶实对称矩阵 \mathbf{B} , 必存在一个 n 阶正交矩阵 \mathbf{P} 使得 $\mathbf{P}^T \mathbf{B} \mathbf{P} = \mathbf{\Lambda}$ (其中 $\mathbf{\Lambda}$ 是以矩阵 \mathbf{B} 的

特征值为对角元素的对角矩阵), 即实对称矩阵都可以对角化。

引理 3 对于实对称矩阵 \mathbf{ZZ}^T 与 $\mathbf{Z}^T \mathbf{Z}$ 有如下结论:

1) 矩阵 \mathbf{ZZ}^T 与 $\mathbf{Z}^T \mathbf{Z}$ 矩阵的特征值为非负实数, 且 2 个矩阵具有相同的正特征值。

2) 若 ξ_i 为矩阵 \mathbf{ZZ}^T 的非零特征值 λ_i 的单位特征向量, 则 $\eta_i = \frac{1}{\sqrt{\lambda_i}} \mathbf{Z}^T \xi_i$ 为矩阵 $\mathbf{Z}^T \mathbf{Z}$ 的非零特征值 λ_i 的单位特征向量。

3) 若 η_i 为矩阵 $\mathbf{Z}^T \mathbf{Z}$ 的非零特征值 λ_i 的单位特征向量, 则 $\xi_i = \frac{1}{\sqrt{\lambda_i}} \mathbf{Z} \eta_i$ 为矩阵 \mathbf{ZZ}^T 的非零特征值 λ_i 的单位特征向量。

证明 1) 根据实对称矩阵 \mathbf{ZZ}^T 及 $\mathbf{Z}^T \mathbf{Z}$ 具有非负定性, 由引理 2 即可得它们的特征值为非负实数。不妨设 ξ_i 是矩阵 \mathbf{ZZ}^T 关于正特征值 λ_i 的单位特征向量, 即 $\mathbf{ZZ}^T \xi_i = \lambda_i \xi_i$, 则有

$$\mathbf{Z}^T \mathbf{Z} (\mathbf{Z}^T \xi_i) = \lambda_i (\mathbf{Z}^T \xi_i), \mathbf{Z}^T \xi_i \neq \mathbf{0}。 \quad (1)$$

根据式 (1), 可推得矩阵 $\mathbf{Z}^T \mathbf{Z}$ 的特征多项式为

$$|\mathbf{Z}^T \mathbf{Z} - \lambda_i \mathbf{I}| = 0, \quad (2)$$

即 λ_i 是矩阵 $\mathbf{Z}^T \mathbf{Z}$ 的特征值。

同理由 $\mathbf{Z}^T \mathbf{Z} \eta_i = \lambda_i \eta_i$, 可推得矩阵 \mathbf{ZZ}^T 的特征多项式为

$$|\mathbf{ZZ}^T - \lambda_i \mathbf{I}| = 0,$$

即 λ_i 是矩阵 \mathbf{ZZ}^T 的特征值。

2) 若 ξ_i 是矩阵 \mathbf{ZZ}^T 关于 λ_i 的单位特征向量, 则有

$$\mathbf{ZZ}^T \xi_i = \lambda_i \xi_i, \xi_i^T \xi_i = 1。$$

于是得

$$\begin{cases} \mathbf{Z}^T \mathbf{Z} \left(\frac{1}{\sqrt{\lambda_i}} \mathbf{Z}^T \xi_i \right) = \lambda_i \left(\frac{1}{\sqrt{\lambda_i}} \mathbf{Z}^T \xi_i \right), \\ \left(\frac{1}{\sqrt{\lambda_i}} \mathbf{Z}^T \xi_i \right)^T \left(\frac{1}{\sqrt{\lambda_i}} \mathbf{Z}^T \xi_i \right) = 1。 \end{cases} \quad (3)$$

取 $\eta_i = \frac{1}{\sqrt{\lambda_i}} \mathbf{Z}^T \xi_i$, 即 η_i 为矩阵 $\mathbf{Z}^T \mathbf{Z}$ 相应于特征值 λ_i 的单位特征向量。

3) 同理可证, 若 η_i 为非零特征值 λ_i 关于矩阵 $\mathbf{Z}^T \mathbf{Z}$ 的单位特征向量, 则 $\xi_i = \frac{1}{\sqrt{\lambda_i}} \mathbf{Z} \eta_i$ 为矩阵 \mathbf{ZZ}^T 相应于特征值 λ_i 的单位特征向量。

2 主要结论及证明

对于 p 种属性的 n 个样本, 首先将原始数据矩阵标准化, 得矩阵 $\mathbf{Z} = (z_{ij})_{n \times p}$, 若用列向量表示, 则记

为 $\mathbf{Z} = (\mathbf{Z}_1, \mathbf{Z}_2, \dots, \mathbf{Z}_n)^\top$, 其中 $\mathbf{Z}_i = (z_{i1}, z_{i2}, \dots, z_{ip})^\top$, $i=1, 2, \dots, n$ 。

如果观测到的 n 个样本之间具有强相似性, 则可依照样本相似性的大小将 n 个样本分组, 使得同组的样本之间相似性较高, 不同组的样本之间相似性较低, 实现样本的聚类, 并对类中样本所具有的共性进行分析和解析。

根据因子分析的思想, 引入一组抽象公共因子变量 $\mathbf{F} = (\mathbf{F}_1, \mathbf{F}_2, \dots, \mathbf{F}_m)^\top$ 和一组特殊因子变量 $\boldsymbol{\theta} = (\boldsymbol{\theta}_1, \boldsymbol{\theta}_2, \dots, \boldsymbol{\theta}_n)^\top$, 将样本分量 $\mathbf{Z}_i (i=1, 2, \dots, n)$ 近似地用因子变量 $\mathbf{F}_i (i=1, 2, \dots, m)$ 的线性关系式表示。依据样本与 \mathbf{F} 的相似程度, 近似地通过因子载荷矩阵用公共因子 $\mathbf{F}_i (i=1, 2, \dots, m)$ 来描述位于同组样本具有的基本结构。样本 $\mathbf{Z} = (\mathbf{Z}_1, \mathbf{Z}_2, \dots, \mathbf{Z}_n)^\top$ 关于公共因子 \mathbf{F} 和特殊因子变量 $\boldsymbol{\theta}$ 的数学模型为

$$\mathbf{Z}_i = a_{i1}\mathbf{F}_1 + a_{i2}\mathbf{F}_2 + \dots + a_{im}\mathbf{F}_m + \boldsymbol{\theta}_i, \quad i=1, 2, \dots, n, \quad (4)$$

且满足下列条件:

i) $m \leq n$;

ii) $\text{cov}(\mathbf{F}, \boldsymbol{\theta}) = 0$, 即协方差矩阵为零, 公共因子与特殊因子不相关;

iii) $\text{cov}(\mathbf{Z}_i, \mathbf{Z}_i) = 1, i=1, 2, \dots, n$, 即样本的方差等于 1;

iv) $\text{cov}(\mathbf{F}_i, \mathbf{F}_j) = \begin{cases} 1, & i=j, \\ 0, & i \neq j, \end{cases} i, j=1, 2, \dots, m$, 即公共因子彼此不相关, \mathbf{F} 的协方差矩阵为单位矩阵。

v) $\text{cov}(\boldsymbol{\theta}_i, \boldsymbol{\theta}_j) = \begin{cases} \sigma_i^2, & i=j, \\ 0, & i \neq j, \end{cases} i, j=1, 2, \dots, n$ 即特殊因子彼此不相关, 特殊因子 $\boldsymbol{\theta}_i$ 的方差等于 σ_i^2 。

令式 (4) 中公共因子变量的系数矩阵与特殊因子变量的列矩阵分别为

$$\mathbf{A} = (\mathbf{A}_1, \mathbf{A}_2, \dots, \mathbf{A}_m),$$

$$\boldsymbol{\theta} = (\boldsymbol{\theta}_1, \boldsymbol{\theta}_2, \dots, \boldsymbol{\theta}_n)^\top。$$

则模型 (4) 简化为

$$\mathbf{Z} = \mathbf{A}\mathbf{F} + \boldsymbol{\theta}, \quad (5)$$

其中系数矩阵 $\mathbf{A} = (a_{ij})_{n \times m}$ 称为公共因子变量 \mathbf{F} 的载荷矩阵。

定理 1 对于模型 (5), 在满足条件 i) ~ v) 时, 有如下结论:

1) $\text{cov}(\mathbf{Z}_i, \mathbf{F}_j) = a_{ij}, i=1, 2, \dots, n, j=1, 2, \dots, m$;

2) $a_{i1}^2 + a_{i2}^2 + \dots + a_{im}^2 + \sigma_i^2 = 1, i=1, 2, \dots, n$;

3) \mathbf{Z}_i 与 \mathbf{F}_j 的相关系数 $r_{ij} = a_{ij}$;

4) $\mathbf{R} = \mathbf{A}\mathbf{A}^\top + \text{diag}(\sigma_1^2, \sigma_2^2, \dots, \sigma_n^2)$, 其中

$$\mathbf{R} = \frac{1}{p-1} \mathbf{Z}\mathbf{Z}^\top, \text{cov}(\boldsymbol{\theta}, \boldsymbol{\theta}) = \text{diag}(\sigma_1^2, \sigma_2^2, \dots, \sigma_n^2)。$$

证明 1) 根据式 (4) 及条件 ii) 有

$$\begin{aligned} \text{cov}(\mathbf{Z}_i, \mathbf{F}_j) &= \text{cov}(a_{i1}\mathbf{F}_1 + a_{i2}\mathbf{F}_2 + \dots + a_{im}\mathbf{F}_m + \boldsymbol{\theta}_i, \mathbf{F}_j) = \\ &= \text{cov}[(a_{i1}, a_{i2}, \dots, a_{im})\mathbf{F}, \mathbf{F}_j] + \text{cov}(\boldsymbol{\theta}_i, \mathbf{F}_j) = \\ &= (a_{i1}, a_{i2}, \dots, a_{im})\text{cov}(\mathbf{F}, \mathbf{F})(0, \dots, 1, 0, \dots, 0)^\top = \\ &= a_{ij}。 \end{aligned} \quad (6)$$

2) 根据式 (4) 及条件 ii) ~ iv) 有

$$\begin{aligned} \text{cov}(\mathbf{Z}_i, \mathbf{Z}_i) &= \text{cov}\left(\sum_{k=1}^m a_{ik}\mathbf{F}_k + \boldsymbol{\theta}_i, \sum_{k=1}^m a_{ik}\mathbf{F}_k + \boldsymbol{\theta}_i\right) = \\ &= \text{cov}[(a_{i1}, a_{i2}, \dots, a_{im})\mathbf{F}, (a_{i1}, a_{i2}, \dots, a_{im})\mathbf{F}] + \text{cov}(\boldsymbol{\theta}_i, \boldsymbol{\theta}_i) = \\ &= (a_{i1}, a_{i2}, \dots, a_{im})\text{cov}(\mathbf{F}, \mathbf{F})(a_{i1}, a_{i2}, \dots, a_{im})^\top + \sigma_i^2 = \\ &= a_{i1}^2 + a_{i2}^2 + \dots + a_{im}^2 + \sigma_i^2 = 1。 \end{aligned}$$

3) 根据条件 iv) 和式 (6) ~ (7) 得

$$r_{ij} = \frac{\text{cov}(\mathbf{Z}_i, \mathbf{F}_j)}{\sqrt{\text{cov}(\mathbf{Z}_i, \mathbf{Z}_i)}\sqrt{\text{cov}(\mathbf{F}_j, \mathbf{F}_j)}} = a_{ij}。 \quad (8)$$

4) 由式 (5) 求 \mathbf{Z} 的协方差矩阵, 可得

$$\text{cov}(\mathbf{Z}, \mathbf{Z}) = \mathbf{A}\text{cov}(\mathbf{F}, \mathbf{F})\mathbf{A}^\top + \text{cov}(\boldsymbol{\theta}, \boldsymbol{\theta})。 \quad (9)$$

根据等式

$$\text{cov}(\mathbf{Z}, \mathbf{Z}) = \frac{1}{p-1} \mathbf{Z}\mathbf{Z}^\top = \mathbf{R}, \quad (7)$$

$$\text{cov}(\mathbf{F}, \mathbf{F}) = \text{diag}(1, 1, \dots, 1),$$

$$\text{cov}(\boldsymbol{\theta}, \boldsymbol{\theta}) = \text{diag}(\sigma_1^2, \sigma_2^2, \dots, \sigma_n^2),$$

将式 (9) 化为

$$\mathbf{R} = \mathbf{A}\mathbf{A}^\top + \text{diag}(\sigma_1^2, \sigma_2^2, \dots, \sigma_n^2)。 \quad (10)$$

推论 1 在模型 $\mathbf{Z} = \mathbf{A}\mathbf{F} + \boldsymbol{\theta}$ 中, 对于任意 m 阶正交矩阵 \mathbf{T} , 若 $\mathbf{A} = \mathbf{A}^*\mathbf{T}$, 则有

$$\mathbf{R} = \mathbf{A}^*\mathbf{A}^{*\top} + \text{diag}(\sigma_1^2, \sigma_2^2, \dots, \sigma_n^2)。$$

定理 2 在式 (10) 中, 令

$$\mathbf{R}^* = \mathbf{R} - \text{diag}(\sigma_1^2, \sigma_2^2, \dots, \sigma_n^2),$$

$$h_j = \mathbf{A}_j^\top \mathbf{A}_j (j=1, 2, \dots, m),$$

若满足条件 $h_1 \geq h_2 \geq \dots \geq h_m$, 则有如下结论:

1) 矩阵 \mathbf{R}^* 的特征值 $\lambda_i = h_i (i=1, 2, \dots, m)$;

2) 相应特征值 λ_i 的单位特征向量

$$\mathbf{t}_i = \frac{1}{\sqrt{\lambda_i}} \mathbf{A}_i (i=1, 2, \dots, m);$$

3) 载荷矩阵 $\mathbf{A} = (\sqrt{\lambda_1}\mathbf{t}_1, \sqrt{\lambda_2}\mathbf{t}_2, \dots, \sqrt{\lambda_m}\mathbf{t}_m)。$

证明 令 $\mathbf{R}^* = (r_{ij}^*)_{n \times n}$, 由式 (10) 得

$$r_{ij}^* = \sum_{k=1}^m a_{ik}a_{jk} = \begin{cases} r_{ij}, & i \neq j, \\ r_{ij} - \sigma_i^2, & i = j, \end{cases} i, j=1, 2, \dots, n。 \quad (11)$$

首先考虑在满足条件 $r_{ij}^* = \sum_{k=1}^m a_{ik}a_{jk}$ 时, 使 $h_1 = \mathbf{A}_1^\top \mathbf{A}_1$ 取得最大值, 求向量 $\mathbf{A}_1 = (a_{11}, a_{21}, \dots, a_{n1})^\top$ 的情形。

根据拉格朗日乘数法，构造目标函数

$$L(a_{11}, a_{21}, \dots, a_{n1}) = A_1^T A_1 + \sum_{i=1}^n \sum_{j=1}^n u_{ij} \left(r_{ij}^* - \sum_{k=1}^m a_{ik} a_{jk} \right). \quad (12)$$

由于 $R^* = (r_{ij}^*)_{n \times n}$ 是对称矩阵，在式 (12) 中有

$u_{ik} = u_{ki}$ ，分别对 $a_{i1} (i=1, 2, \dots, n)$ 求偏导数得

$$\frac{\partial L}{\partial a_{i1}} = 2a_{i1} - 2 \sum_{j=1}^n u_{ij} a_{j1} = 0 \quad (i=1, 2, \dots, n), \quad (13)$$

$$\frac{\partial L}{\partial a_{ij}} = -2 \sum_{k=1}^n u_{ik} a_{kj} = 0 \quad (j \neq 1; i=1, 2, \dots, n). \quad (14)$$

将式 (13) 和 (14) 改写为

$$\sum_{k=1}^n u_{ik} a_{kj} - \tau_{1j} a_{i1} = 0 \quad (i=1, 2, \dots, n; j=1, 2, \dots, m), \quad (15)$$

$$\text{式中 } \tau_{1j} = \begin{cases} 1, & j=1, \\ 0, & j \neq 1. \end{cases}$$

用 a_{i1} 乘式 (15) 并对 i 求和，得

$$\sum_{k=1}^n \left(\sum_{i=1}^n u_{ik} a_{i1} \right) a_{kj} - \tau_{1j} \sum_{i=1}^n a_{i1}^2 = 0. \quad (16)$$

结合式 (13)，将式 (16) 化为

$$\sum_{k=1}^n a_{k1} a_{kj} - \tau_{1j} h_1 = 0 \quad (j=1, 2, \dots, m). \quad (17)$$

用 a_{ij} 乘式 (17) 并对 j 求和，得

$$\sum_{k=1}^n a_{k1} \left(\sum_{j=1}^m a_{ij} a_{kj} \right) - \sum_{j=1}^m \tau_{1j} h_1 a_{ij} = 0. \quad (18)$$

结合式 (11)，由式 (18) 得

$$\begin{pmatrix} r_{11}^* & r_{12}^* & \dots & r_{1n}^* \\ r_{21}^* & r_{22}^* & \dots & r_{2n}^* \\ \vdots & \vdots & \ddots & \vdots \\ r_{n1}^* & r_{n2}^* & \dots & r_{nn}^* \end{pmatrix} \begin{bmatrix} a_{11} \\ a_{21} \\ \vdots \\ a_{n1} \end{bmatrix} = h_1 a_{i1} \quad (i=1, 2, \dots, n). \quad (19)$$

将式 (19) 的 n 个等式用矩阵表示即为

$$(R^* - h_1 I) A_1 = 0. \quad (20)$$

式 (20) 表明， h_1 是矩阵 R^* 的最大特征值， A_1 是相应于 h_1 的特征向量，若设 R^* 的最大特征值为 λ_1 ，相应的单位特征向量为 t_1 ，则可得定理 2 中的结论：

$$\lambda_1 = h_1, t_1 = \frac{1}{\sqrt{\lambda_1}} A_1, \text{ 或改写为 } A_1 = \sqrt{\lambda_1} t_1.$$

由 R^* 的谱分解式

$$R^* = A_1 A_1^T + \sum_{i=2}^m A_i A_i^T = \lambda_1 t_1 t_1^T + \sum_{i=2}^m \lambda_i t_i t_i^T,$$

同理可得矩阵 R^* 由大到小的其余 $m-1$ 个特征值 $\lambda_2 \geq \lambda_3 \geq \dots \geq \lambda_m$ ，其中 $\lambda_j = h_j (j=2, 3, \dots, m)$ ，以及相应的

单位特征向量 $t_j = \frac{1}{\sqrt{\lambda_j}} A_j$ ，或改写为 $A_j = \sqrt{\lambda_j} t_j$

($j=2, 3, \dots, m$)。

定理 2 证毕。

由于样本相似矩阵 R 是已知的，特殊因子 θ_i 的方差 σ_i^2 通常是一个待估计的量，等式

$$R^* = R - \text{diag}(\sigma_1^2, \sigma_2^2, \dots, \sigma_n^2)$$

中的矩阵 R^* 为未知。因此定理 2 中由条件 $R^* = AA^T$ 求载荷矩阵 A 不可行。在实际应用中通常将 R^* 作近似处理，将由条件 $R^* = AA^T$ 求载荷矩阵 A 的问题转化为

在条件 $R \approx AA^T$ 下，由协方差矩阵 $\frac{1}{p-1} ZZ^T = R$ 求载

荷矩阵 A 的近似矩阵 A^* ，再利用回归的思想估计特殊因子的方差。因为样本离差矩阵 ZZ^T 是一个 n 阶的

对称矩阵，而 $Z^T Z$ 是一个 p 阶对称矩阵，样本数 n 通常远大于样本所具有的属性个数 p ，因此求 n 阶标准

矩阵 ZZ^T 的特征值和特征向量的计算量较求 p 阶矩阵

$Z^T Z$ 的特征值和特征向量复杂得多。为此，由矩阵

$\frac{1}{p-1} Z^T Z$ 的特征值及特征向量得到 $\frac{1}{p-1} ZZ^T$ 的特征值

及特征向量，即可减少求载荷矩阵 A 的近似矩阵 A^*

的计算量。

推论 2 设矩阵 $Z^T Z$ 的 m 个非零特征值为

$\lambda_i (i=1, 2, \dots, m)$ ，其排列的顺序由大到小，相应于第

i 个特征值 λ_i 的单位特征向量为 η_i 。对于样本相似矩

阵 $R = \frac{1}{p-1} ZZ^T$ ，则有如下结论：

1) $u_i = \frac{\lambda_i}{p-1} (i=1, 2, \dots, m)$ 是矩阵 R 的特征值；

2) 矩阵 R 与矩阵 ZZ^T 具有相同的单位特征向量

ξ_i ，且 $\xi_i = \frac{1}{\sqrt{\lambda_i}} Z \eta_i (i=1, 2, \dots, m)$ ；

3) 载荷矩阵

$$A^* = \left(\frac{1}{\sqrt{p-1}} Z \eta_1, \frac{1}{\sqrt{p-1}} Z \eta_2, \dots, \frac{1}{\sqrt{p-1}} Z \eta_m \right).$$

由引理 3 和定理 2 易证，故略去证明。

3 实例

例 1 10 名学生的数学与语文考试成绩见表 1。

表 1 学生成绩

Table 1 Student's score

样本 序号	数学 y_{i1}	语文 y_{i2}	平均 \bar{y}_i	样本 序号	数学 y_{i1}	语文 y_{i2}	平均 \bar{y}_i
1	77	64	70.5	6	66	52	59.0
2	67	65	66.0	7	77	72	74.5
3	80	74	77.0	8	83	61	72.0
4	74	84	79.0	9	86	41	63.5
5	78	62	70.0	10	65	84	74.5

以这 10 名学生作为样本观测点, 共 10 个样本。第 i 个样本用 Y_i 表示, 它是由数学成绩 y_{i1} 与语文成绩 y_{i2} 构成的二维数组, 记为 $Y_i = (y_{i1}, y_{i2})^T (i=1, 2, \dots, 10)$ 。试

$$Z = \begin{bmatrix} \frac{\sqrt{2}}{2} & \frac{\sqrt{2}}{2} & \frac{\sqrt{2}}{2} & -\frac{\sqrt{2}}{2} & \frac{\sqrt{2}}{2} & \frac{\sqrt{2}}{2} & \frac{\sqrt{2}}{2} & \frac{\sqrt{2}}{2} & \frac{\sqrt{2}}{2} & -\frac{\sqrt{2}}{2} \\ -\frac{\sqrt{2}}{2} & -\frac{\sqrt{2}}{2} & -\frac{\sqrt{2}}{2} & \frac{\sqrt{2}}{2} & -\frac{\sqrt{2}}{2} & -\frac{\sqrt{2}}{2} & -\frac{\sqrt{2}}{2} & -\frac{\sqrt{2}}{2} & -\frac{\sqrt{2}}{2} & \frac{\sqrt{2}}{2} \end{bmatrix}^T,$$

从而得 $Z^T Z = \begin{bmatrix} 5 & -5 \\ -5 & 5 \end{bmatrix}$ 。

矩阵 $Z^T Z$ 有一个非零特征值 $\lambda=10$, 对应的单位特征向量 $\eta = \begin{bmatrix} -\frac{\sqrt{2}}{2} & \frac{\sqrt{2}}{2} \end{bmatrix}^T$; 根据矩阵 $Z^T Z$ 与矩阵 ZZ^T 特征值及特征向量之间的关系, 得矩阵 ZZ^T 的单位特征向量

$$\xi = \frac{1}{\sqrt{\lambda}} Z \eta = \frac{1}{\sqrt{10}} [-1 \ -1 \ -1 \ 1 \ -1 \ -1 \ -1 \ -1 \ -1 \ 1]^T.$$

于是得载荷矩阵

$$A = [-1 \ -1 \ -1 \ 1 \ -1 \ -1 \ -1 \ -1 \ -1 \ 1]^T,$$

提取的公共因子只有一个, 即 F_1 , 样本 $Z = (Z_1, Z_2, \dots, Z_{10})^T$ 关于公共因子 F_1 和特殊因子变量 θ 的数学模型表示为 $Z = AF_1 + \theta$ 。

根据定理 1 可知, 第 i 个样本与公共因子 F_1 的相关系数 $r_{i1} = a_{i1} (i=1, 2, \dots, 10)$ 见表 2。

表2 相关系数分布

Table 2 The distribution of correlation coefficient

样本 Z_i	Z_1	Z_2	Z_3	Z_4	Z_5	Z_6	Z_7	Z_8	Z_9	Z_{10}
相关系数 r_{i1}	-1	-1	-1	1	-1	-1	-1	-1	-1	1

表 2 表明, Z_4 和 Z_{10} 相关于公共因子 F_1 的正方向; $Z_1, Z_2, Z_3, Z_5, Z_6, Z_7, Z_8, Z_9$ 相关于公共因子 F_1 的反方向。因此, 可将这 10 名学生分为 2 类, 第一类由 4 号与 10 号学生组成; 第二类为余下的 8 名学生组成。公共因子 F_1 的正方向表明学生的语文成绩优于数学成绩, 反方向表明学生的语文成绩劣于数学成绩。

4 结语

对于 p 种属性的 n 个样本, 样本相似矩阵 $R = \frac{1}{p-1} ZZ^T$ 是一个 n 阶的高阶对称矩阵, 矩阵 $Z^T Z$

用因子分析法, 按样本与因子相似的程度将这 10 名学生分类, 且作出合理的解释。

解 由 $z_{ij} = \frac{y_{ij} - \bar{y}_i}{\sqrt{v_{ij}}}$ 将原始数据标准化得矩阵

的阶是一个比 R 矩阵阶数低得多的 p 阶对称矩阵。由推论 2 获得求载荷矩阵 A^* 的方法, 其计算的复杂度远低于通过 R 矩阵获得载荷矩阵 A^* 的方法。

参考文献:

- [1] 何晓群. 多元统计分析[M]. 北京: 中国人民大学出版社, 2012: 142-144.
He Xiaoqun. Multivariate Statistical Analysis[M]. Beijing: China Renmin University Press, 2012: 142-144.
- [2] 虞欣, 郑肇葆. 基于 Q 型因子分析的训练样本的选择[J]. 测绘学报, 2007, 36(1): 67-71.
Yu Xin, Zheng Zhaobao. Selection of Training Samples Based on Q-Factor Analysis[J]. Acta Geodaetica et Cartographica Sinica, 2007, 36(1): 67-71.
- [3] 殷瑞飞, 朱建平. 关于利用因子分析方法对变量分类的探讨[J]. 统计与决策, 2005(2): 20-21.
Yin Ruifei, Zhu Jianping. Using the Factor Analysis Method for the Classification Variables[J]. Statistics and Decision, 2005(2): 20-21.
- [4] 张秋瑾. 主成分分析法在多变量变动分析中的应用[J]. 数学的实践与认识, 2012, 42(17): 29-33.
Zhang Qiujin. The Application of Principal Component Analysis Method in Multivariate Analysis of Changes[J]. Mathematics in Practice and Theory, 2012, 42(17): 29-33.
- [5] Ramsay J Q. Functional Components of Variation in Handwriting[J]. Journal of the American Statistic Association, 2000, 95(449): 9-15.
- [6] 周勇, 朱砾. 线性代数[M]. 上海: 复旦大学出版社, 2012: 129-131.
Zhou Yong, Zhu Li. Linear Algebra[M]. Shanghai: Fudan University Press, 2012: 129-131.

(责任编辑: 邓光辉)

