

10.1 引言

一. PCA的主要功能

在信息损失最小的前提下，对高维空间进行降维处理。

数据类型：样本点×变量（定量变量）

$$\begin{array}{c} e_1 \\ e_2 \\ \vdots \\ e_n \end{array} \begin{array}{cccc} x_1 & x_2 & \text{L} & x_p \\ \left[\begin{array}{cccc} x_{11} & x_{12} & \text{L} & x_{1p} \\ x_{21} & x_{22} & \text{L} & x_{2p} \\ \text{M} & \text{M} & \text{L} & \text{M} \\ x_{n1} & x_{n2} & \text{L} & x_{np} \end{array} \right]_{n \times p} \end{array} \xrightarrow{\text{PCA}} \begin{array}{c} e_1 \\ e_2 \\ \vdots \\ e_n \end{array} \begin{array}{ccc} y_1 & \text{L} & y_m \\ \left[\begin{array}{ccc} y_{11} & \text{L} & y_{1m} \\ y_{21} & \text{L} & y_{2m} \\ \text{M} & \text{ML} & \text{M} \\ y_{n1} & \text{L} & y_{nm} \end{array} \right]_{n \times m} \end{array} \quad m \ll p$$

10.3 数据的原则化处理

(一) “中心化”处理—平移变换

$$\tilde{x}_{ij} = x_{ij} - \bar{x}_{.j}, \quad i=1,2,\dots,n, \quad j=1,2,\dots,p$$

性质：不变化样本点集中点与点的相互位置；

$$\begin{array}{c} \left[\begin{array}{ccc} x_{11} & \text{L} & x_{1p} \\ \text{M} & & \text{M} \\ x_{n1} & \text{L} & x_{np} \end{array} \right]_{n \times p} \xrightarrow{\text{中心化}} \left[\begin{array}{ccc} \tilde{x}_{11} & \text{L} & \tilde{x}_{1p} \\ \text{M} & & \text{M} \\ \tilde{x}_{n1} & \text{L} & \tilde{x}_{np} \end{array} \right]_{n \times p} \\ \bar{x}_1 \quad \text{L} \quad \bar{x}_p \end{array}$$

(二) 原则化处理：中心化——压缩

$$y_{ij} = \frac{x_{ij} - \bar{x}_j}{S_j}$$

性质：

➤ $g^* = 0$ (均值为0)

➤ $s_j^* = 1, \quad j=1,2,\dots,p$ (方差等于1)

➤ $r_{jk}^* = \frac{S_{jk}^*}{S_j^* S_k^*} = S_{jk}^*$

对于原则化数据表:

$$\tilde{x}_{ij} = \frac{x_{ij} - \bar{x}_j}{s_j}$$

(1) 变量方差均等于1

(2) 有关系数矩阵 = 协方差矩阵

$$\mathbf{R} = \begin{bmatrix} 1 & r_{12} & \text{L} & r_{1p} \\ r_{21} & 1 & \text{L} & r_{2p} \\ \text{M} & \text{M} & & \text{M} \\ r_{p1} & r_{p2} & \text{L} & 1 \end{bmatrix} = \begin{bmatrix} s_1^2 & s_{12} & \text{L} & s_{1p} \\ s_{21} & s_2^2 & \text{L} & s_{2p} \\ \text{M} & \text{M} & & \text{M} \\ s_{p1} & s_{p2} & \text{L} & s_p^2 \end{bmatrix} = \mathbf{V}$$

10.4 PCA的算法

一. PCA对数据系统做“最佳简化”的含意

PCA可在确保信息损失的前提下，经线性变换和舍弃一小部分信息，以少数线性无关的新综合变量取代原始采用的多维有关变量。

$$\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_p \xrightarrow{\text{PCA: 平移} \oplus \text{旋转}} \mathbf{y}_1, \mathbf{y}_2, \dots, \mathbf{y}_m \quad (m < p)$$

称： $L = L(\mathbf{v}_1, \mathbf{v}_2, \dots, \mathbf{v}_m)$ 为“主超平面”；

称： $L_i = L(\mathbf{v}_1, \mathbf{v}_2, \dots, \mathbf{v}_i)$ 为“主平面”；

二. PCA算法中的几种要素

输入—输出:

$$\begin{array}{c} x_1 \quad x_2 \quad \text{L} \quad x_p \\ e_1 \left[\begin{array}{cccc} x_{11} & x_{12} & \text{L} & x_{1p} \\ e_2 & x_{21} & x_{22} & \text{L} & x_{2p} \\ \text{M} & \text{M} & \text{M} & \text{L} & \text{M} \\ e_n & x_{n1} & x_{n2} & \text{L} & x_{np} \end{array} \right]_{n \times p} \xrightarrow{\text{平移} \oplus \text{旋转}} \begin{array}{c} y_1 \quad \text{L} \quad y_m \\ e_1 \left[\begin{array}{ccc} y_{11} & \text{L} & y_{1m} \\ e_2 & y_{21} & \text{L} & y_{2m} \\ \text{M} & \text{ML} & \text{M} \\ e_n & y_{n1} & \text{L} & y_{nm} \end{array} \right]_{n \times m} \end{array} \quad m < p$$

(1) 平移变换: 把原点移到重心:

$$\tilde{x}_{ij} = x_{ij} - \bar{x}_j, \quad i = 1, 2, \text{L}, n,$$

$$j = 1, 2, \text{L}, p$$

(2) 旋转变换，得到“**主轴**”： $u_1, u_2, \dots, u_p \in R^p$

其中， u_1 相应数据变异最大的方向， u_2 与 u_1 垂直，相应于数据变异第二大方向，...

所以 u_1, \dots, u_p 是原则正交的，即： $u_j' u_k = \begin{cases} 1 & j \neq k \\ 0 & j = k \end{cases}$

(3) 求样本点 e_i 在 u_h 轴上的投影坐标

$$y_h(i) = e_i' \cdot u_h$$

全部样本点在 u_h 上的投影构成“**第 h 主成份 y_h** ”：

$$y_h = (y_h(1), y_h(2), \dots, y_h(n))' \in R^n$$

在主成份中， $\text{Var}(y_1) \rightarrow \max$

而 $y_2 \perp y_1$ ，且 $\text{Var}(y_2)$ 是次大的……

(4) 在 \mathbf{u}_h 主轴上, \mathbf{e}_i 的投影坐标是 $y_h(i)$

$$y_h(i) = \mathbf{e}'_i \cdot \mathbf{u}_h, \quad i = 1, 2, \dots, n$$

第 h 主成份为:

$$\mathbf{y}_h = (y_h(1), y_h(2), \dots, y_h(n))'$$

$$= \begin{pmatrix} \mathbf{e}'_1 \cdot \mathbf{u}_h \\ \mathbf{M} \\ \mathbf{e}'_n \cdot \mathbf{u}_h \end{pmatrix} = \begin{pmatrix} \mathbf{e}'_1 \\ \mathbf{M} \\ \mathbf{e}'_n \end{pmatrix} \cdot \mathbf{u}_h = \mathbf{X}\mathbf{u}_h$$

$$= (\mathbf{x}_1, \dots, \mathbf{x}_p) \begin{pmatrix} u_h(1) \\ \mathbf{M} \\ u_h(p) \end{pmatrix} = \sum_{j=1}^p u_h(j) \mathbf{x}_j$$

y_h 是原变量 x_1, \dots, x_p 的线性组合, 组合系数为 $u_h(1), \dots, u_h(p)$

PAC算法推导:

不妨设变量 $\mathbf{X}_1, \mathbf{X}_2, \dots, \mathbf{X}_p$ 都是中心化的, $\mathbf{Y}_h \in L(\mathbf{X}_1, \mathbf{X}_2, \dots, \mathbf{X}_p)$

求第主成份 $\mathbf{Y}_h = \sum_{j=1}^p u_{hj} \mathbf{X}_j$ 经过旋转变换得到的Y, 是X的线性组合

$$D(\mathbf{Y}_h) = \frac{1}{n} \sum_{i=1}^n (y_{ih} - 0)^2 = \frac{1}{n} \sum_{i=1}^n y_{ih}^2 = \frac{1}{n} \|\mathbf{Y}_h\|^2 = \frac{1}{n} \langle \mathbf{Y}_h, \mathbf{Y}_h \rangle$$

$$= \frac{1}{n} \left\langle \left(u_{h1} \mathbf{X}_1 + u_{h2} \mathbf{X}_2 + \dots + u_{hp} \mathbf{X}_p \right), \left(u_{h1} \mathbf{X}_1 + u_{h2} \mathbf{X}_2 + \dots + u_{hp} \mathbf{X}_p \right) \right\rangle$$

$$= \frac{1}{n} \begin{pmatrix} u_{h1} & u_{h2} & \dots & u_{hp} \end{pmatrix} \begin{pmatrix} \langle \mathbf{X}_1, \mathbf{X}_1 \rangle & \langle \mathbf{X}_1, \mathbf{X}_2 \rangle & \dots & \langle \mathbf{X}_1, \mathbf{X}_p \rangle \\ \langle \mathbf{X}_2, \mathbf{X}_1 \rangle & \langle \mathbf{X}_2, \mathbf{X}_2 \rangle & \dots & \langle \mathbf{X}_2, \mathbf{X}_p \rangle \\ \dots & \dots & \dots & \dots \\ \langle \mathbf{X}_p, \mathbf{X}_1 \rangle & \langle \mathbf{X}_p, \mathbf{X}_2 \rangle & \dots & \langle \mathbf{X}_p, \mathbf{X}_p \rangle \end{pmatrix} \begin{pmatrix} u_{h1} \\ u_{h2} \\ \dots \\ u_{hp} \end{pmatrix}$$

$$= \mathbf{u}'_h \mathbf{V} \mathbf{u}_h$$

$$\max \sum_{h=1}^m \mathbf{u}'_h \mathbf{V} \mathbf{u}_h$$

$$\text{s.t.} \begin{cases} \|\mathbf{u}_h\| = 1 \\ \mathbf{u}'_h \mathbf{u}_l = 0 \quad h \neq l \\ \mathbf{u}'_1 \mathbf{V} \mathbf{u}_1 \geq \mathbf{u}'_2 \mathbf{V} \mathbf{u}_2 \geq L \geq \mathbf{u}'_m \mathbf{V} \mathbf{u}_m \\ k = 1, 2, L, p, l \neq k \end{cases}$$

所以 $\mathbf{u}_1, \mathbf{u}_2, L, \mathbf{u}_p$ 是矩阵 \mathbf{V} 的特征向量,

相应的特征值是 $\lambda_1 \geq \lambda_2 \geq L \geq \lambda_p$

三. PCA的计算措施 (一般情况下)

(1) 数据的原则化

$$x_{ij}^* = \frac{x_{ij} - \bar{x}_j}{s_j}$$

为以便起见, 仍记 $(x_{..}^*) = (x_{..}) = \mathbf{X}_{\dots}$ 。

(2) 计算原则化数据表 $\mathbf{X}_{n \times p}$ 的协方差矩阵 \mathbf{V} 。

(3) 求 \mathbf{V} 的前 m 个特征值 $\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_m > 0$,

以及相应的特征向量: $\mathbf{u}_1, \mathbf{u}_2, \dots, \mathbf{u}_m$ (主轴

)

它们是原则正交的: $\mathbf{u}'_h \mathbf{u}_l = \begin{cases} 1 & h=l \\ 0 & h \neq l \end{cases}$

(4) 在 \mathbf{u}_h 主轴上, \mathbf{e}_i 的投影坐标是 $y_h(i)$

$$y_h(i) = \mathbf{e}'_i \cdot \mathbf{u}_h, \quad i = 1, 2, \dots, n$$

第 h 主成份为:

$$\mathbf{y}_h = (y_h(1), y_h(2), \dots, y_h(n))'$$

$$= \begin{pmatrix} \mathbf{e}'_1 \cdot \mathbf{u}_h \\ \mathbf{M} \\ \mathbf{e}'_n \cdot \mathbf{u}_h \end{pmatrix} = \begin{pmatrix} \mathbf{e}'_1 \\ \mathbf{M} \\ \mathbf{e}'_n \end{pmatrix} \cdot \mathbf{u}_h = \mathbf{X} \mathbf{u}_h$$

$$= (\mathbf{x}_1, \dots, \mathbf{x}_p) \begin{pmatrix} u_h(1) \\ \mathbf{M} \\ u_h(p) \end{pmatrix} = \sum_{j=1}^p u_h(j) \mathbf{x}_j$$

y_h 是原变量 x_1, \dots, x_p 的线性组合, 组合系数为 $u_h(1), \dots, u_h(p)$

四、主成份的统计特征

第 h 主成份 $\mathbf{v}_h = (v_h(1), v_h(2), \dots, v_h(n))' \in \mathbf{R}^n$

$$y_h(i) = \mathbf{e}_i' \cdot \mathbf{u}_h, \quad i = 1, 2, \dots, n$$

① y_h 的均值为0。

$$\frac{1}{n} \sum_{i=1}^n y_h(i) = 0$$

② y_h 的方差等于 λ_h 。

$$\text{Var}(\mathbf{y}_h) = \frac{1}{n} \sum_{i=1}^n (y_h(i) - 0)^2 = \lambda_h$$

③ y_j 与 y_k 的协方差等于0:

$$\text{Cov}(\mathbf{y}_h, \mathbf{y}_l) = 0, \quad \forall h \neq l$$

总结：PCA算法的输入与输出

$$\mathbf{X}_{n \times p} = \begin{bmatrix} & & \mathbf{M} & & \\ \mathbf{L} & & x_{ij} & & \mathbf{L} \\ & & \mathbf{M} & & \end{bmatrix}_{n \times p}$$

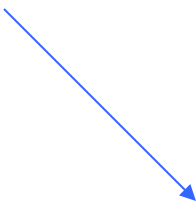
n 个样本点， p 个变量

① $\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_m$

$\text{Var}(\mathbf{y}_1), \text{Var}(\mathbf{y}_2), \dots, \text{Var}(\mathbf{y}_m)$

② $\mathbf{u}_1, \mathbf{u}_2, \dots, \mathbf{u}_m \in \mathbf{R}^p$ (主轴)

③ $\mathbf{y}_1, \mathbf{y}_2, \dots, \mathbf{y}_m \in \mathbf{R}^n$ (主成份)


$$\mathbf{Y}_{n \times m} = \begin{bmatrix} & & \mathbf{M} & & \\ \mathbf{L} & & y_{ij} & & \mathbf{L} \\ & & \mathbf{M} & & \end{bmatrix}_{n \times m}$$

n 个样本点， m 个变量

总结： 经过主成份分析，

$$\begin{bmatrix} x_{11} & x_{12} & \text{L} & x_{1p} \\ x_{21} & x_{22} & \text{L} & x_{2p} \\ \text{M} & \text{M} & \text{L} & \text{M} \\ x_{n1} & x_{n2} & \text{L} & x_{np} \end{bmatrix}_{n \times p} \xrightarrow{\text{PCA}} \begin{bmatrix} y_{11} & \text{L} & y_{1m} \\ y_{21} & \text{L} & y_{2m} \\ \text{ML} & & \text{M} \\ y_{n1} & \text{L} & y_{nm} \end{bmatrix}_{n \times m} \quad m < p$$

均值 $g = (\bar{x}_1, \bar{x}_2, \text{L}, \bar{x}_p) \rightarrow g = (0, \text{L}, 0)$

方差 $s_1^2, s_2^2, \text{L}, s_p^2 \rightarrow \lambda_1 \geq \lambda_2 \geq \text{L} \geq \lambda_m$

$\text{Cov}(x_j, x_k) \neq 0 \ (j \neq k) \rightarrow \text{Cov}(y_j, y_k) = 0 \ (j \neq k)$

10.5 PCA的辅助分析技术

一. 怎样选用精度合适的主超平面

1. m 维主超平面的精度测量

主成份分析前, $X_{n \times p}$ 数据中的全部变异信息:

$$\sum_{j=1}^p \text{Var}(x_j) = \sum_{j=1}^p s_j^2 \stackrel{p}{=}$$

原则化

主成份分析后保存的数据变差:

$$\text{Var}(\mathbf{y}_1) = \lambda_1, \text{Var}(\mathbf{y}_2) = \lambda_2, \dots, \text{Var}(\mathbf{y}_m) = \lambda_m$$

$$\sum_{h=1}^m \text{Var}(\mathbf{y}_h) = \sum_{h=1}^m \lambda_h$$

形象地看： $[\mathbf{x}_1, \mathbf{x}_2, \mathbf{L}, \mathbf{x}_p] \xrightarrow{PCA} [\mathbf{y}_1, \mathbf{y}_2, \mathbf{L}, \mathbf{y}_m]$

方差： $s_1^2, s_2^2, \mathbf{L}, s_p^2 \quad \lambda_1, \lambda_2, \mathbf{L}, \lambda_m$

注意： $\sum_{h=1}^p \lambda_h = \sum_{h=1}^p \text{Var}(\mathbf{y}_h) = \sum_{j=1}^p s_j^2$

所以，定义“**合计贡献率**”：

$$Q_m = \frac{\sum_{h=1}^m \lambda_h}{\sum_{j=1}^p s_j^2} = \frac{1}{p} \sum_{h=1}^m \lambda_h$$

原则化

2.、怎样选用合适精度的 u_1, \dots, u_m 。

根据合计贡献率能够拟定所要选用的成份的个数。

例.管理期刊评价

Total Variance Explained

Component	Initial Eigenvalues			Extraction Sums of Loadings	
	Total	% of Variance	Cumulative %	Total	% of Variance
1	1.532	38.293	38.293	1.532	38.293
2	1.026	25.644	63.937	1.026	25.644
3	.886	22.150	86.087		
4	.557	13.913	100.000		

(2) 若希望 Q_m 在80%左右，应选用3个主成份。

某些科技问题的合计贡献率要求在90%以上。但对复杂的社会科学、行为科学或经济学中的数据，能到达60%也能够考虑。

二. 主成份的命名

主成份 y_1, \dots, y_m 是原变量 x_1, \dots, x_p 的线性组合。
原变量 x_1, \dots, x_p 都有明确的物理含意。

问题： y_1, \dots, y_m 的物理含意是什么？

1. 作用：指出影响系统构造的主要原因和主要特征。

例 ①：分析各阶层人员生活状态

发展中国家： y_1 ——食品， y_2 ——穿着

发达国家： y_1 ——住宅， y_2 ——旅游

以此能够划分不同社会阶层的生活档次。

（在这个方向，人们的生活水平差距最大）

例②：中国城市经济分析：

1984： y_1 —综合水平， y_2 ——工农业投入国家。

1988： y_1 —综合水平， y_2 ——外贸，科技。

中国改革开放以来，因为开放程度不同，使中国各地域经济水平差距逐渐拉大。所以，加大开放力度，发展高科技产业是城市发展的主要工作方面。

2. 措施：专业知识 + 数学手段

数学手段：研究 y_h 与 x_1, \dots, x_p 的有关关系。

对于原则化数据能够证明：

$$r(v, x) = \sqrt{\lambda} \cdot u(i)$$

所以：

$$r(\mathbf{y}_1, \mathbf{x}_1) = \sqrt{\lambda_1} u_1(1)$$

$$r(\mathbf{y}_1, \mathbf{x}_2) = \sqrt{\lambda_1} u_1(2)$$

L L

$$r(\mathbf{y}_1, \mathbf{x}_p) = \sqrt{\lambda_1} u_1(p)$$

第一种主轴：

$$\mathbf{u}_1 = (u_1(1), u_1(2), \dots, u_1(p))'$$

由此可见，仅差一种常量倍 $\sqrt{\lambda_1}$ ：

$(u_1(1), u_1(2), \dots, u_1(p))$ 是 \mathbf{y}_1 与 $\mathbf{x}_1, \dots, \mathbf{x}_p$ 的有关系数。

所以，能够经过观察 $\mathbf{u}_1 = (u_1(1), u_1(2), \dots, u_1(p))'$ 来拟定 \mathbf{y}_1 的含意。

例.管理期刊分类评估

$$r(\mathbf{y}_1, \mathbf{x}_j) = \sqrt{\lambda_1} u_1(j)$$

$$r(\mathbf{y}_2, \mathbf{x}_i) = \sqrt{\lambda_2} u_2(j)$$

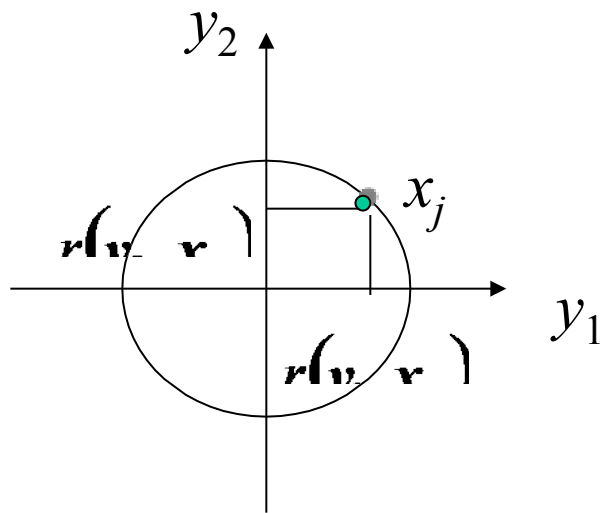
Component Matrix^a

	Component	
	1	2
BYCISHU	.784	-3.5E-02
ZAIWENL	.102	.948
YZQIKAN	.802	-.249
NSFC	.513	.254

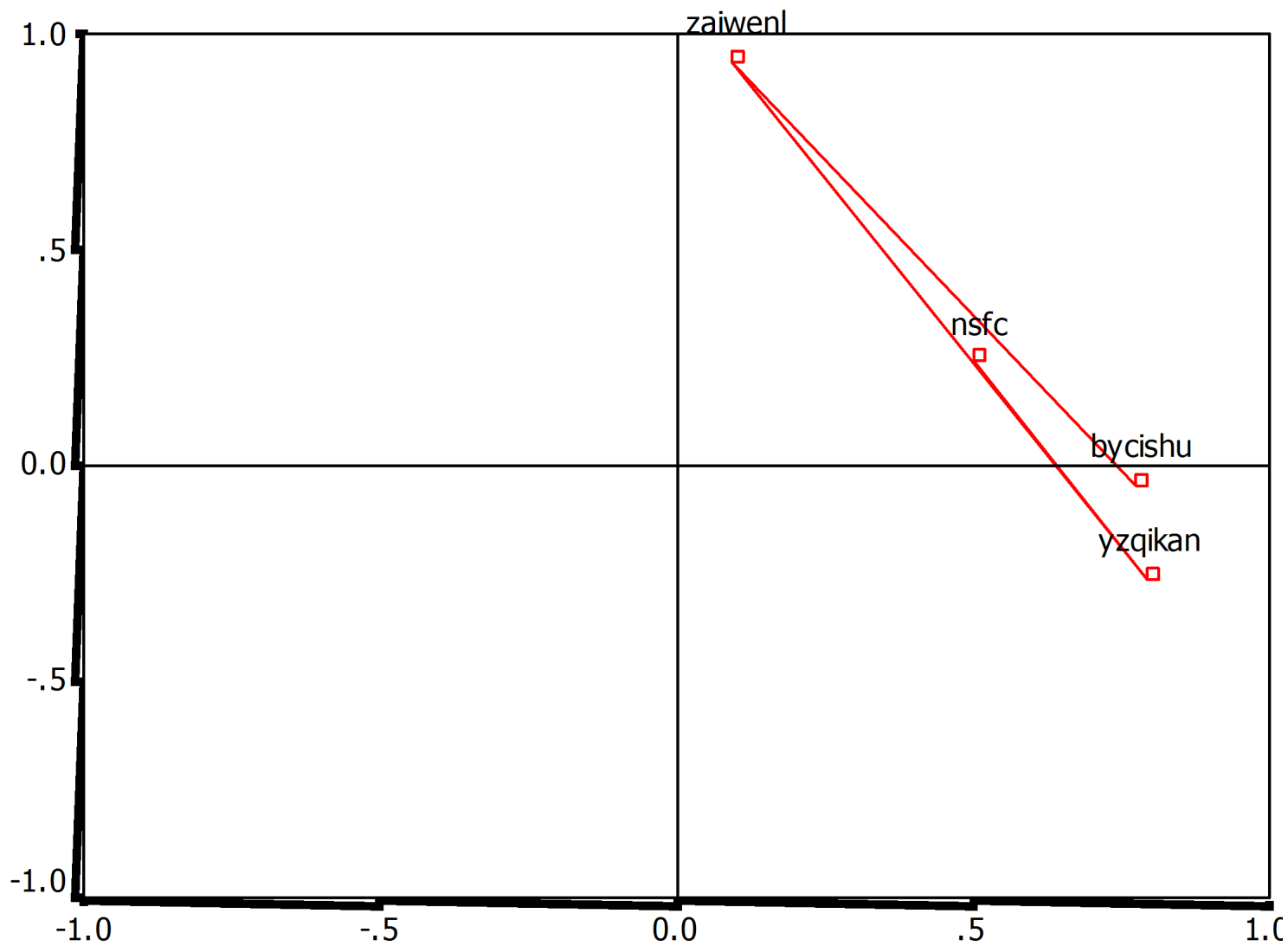
Extraction Method: Principal Component Ana

(2) 有关圆图 (Component Plot)

若 $m = 2$ $\begin{cases} \text{横坐标 } r(y_1, x_j) \\ \text{纵坐标 } r(y_2, x_j) \end{cases}$



Component Plot



三. 判断“特异点” (e_k)

“特异点”：

在PCA中，若有 e_k 远离数据分布的平均水平，能够用“点对主成份方差的贡献”来测量。

如：
$$\text{Var}(\mathbf{y}_1) = \lambda_1 = \frac{1}{n} \sum_{i=1}^n y_1^2(i)$$

则定义“ e_i 对 $\text{Var}(\mathbf{y}_1)$ 的贡献”为：

$$\frac{\frac{1}{n} y_1^2(i)}{\lambda_1} = \text{CTR}_1(i)$$

若 $\text{CTR}_1(\text{BJ}) = 1/3$ ，说明有1/3的 $\text{Var}(\mathbf{y}_1)$ 是由BJ提供的
所以，BJ是特异点。

一般地，定义“ e_i 对 $\text{Var}(y_h)$ 的贡献” $\text{CTR}_h(i)$:

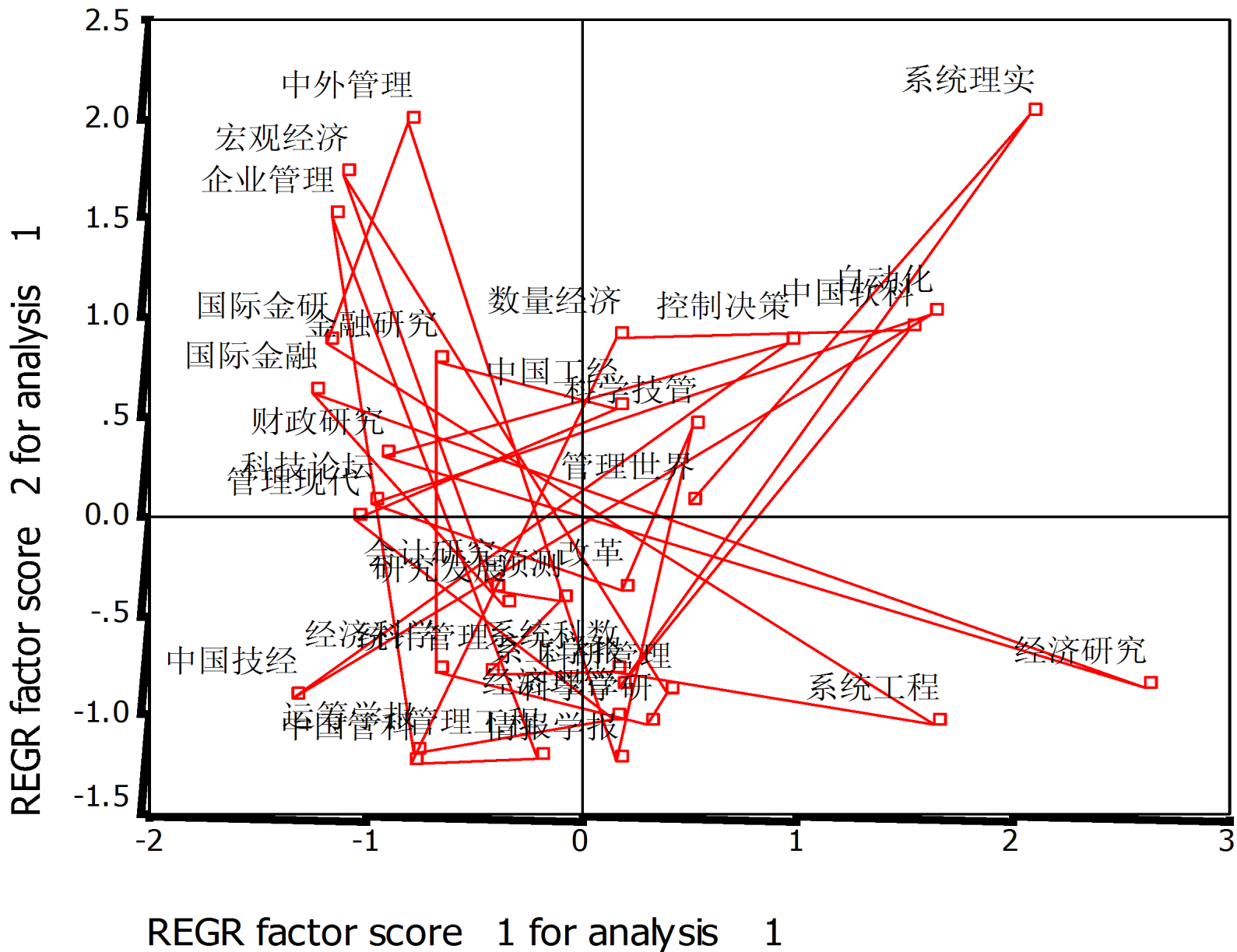
$$\text{CTR}_h(i) = \frac{y_h^2(i)}{n \cdot \lambda_h}$$

$\text{CTR}(i)$ 过大原因:

- (1) 数据本身的特异性 (BJ, SH, GZ, SZ, TJ)
- (2) 数据统计上的错误

处理措施: 除去这些特异点，能够提升分析精度，图示也愈加清楚。

四. 主平面图



10.6 利用主成份分析构造评估函数

PCA将一种高维变量系统有效的降至 1 维

例1: Kendall [英] 评估英国各地域农业生产水平。48个郡，10种农作物：小麦 (x_1)、大麦 (x_2)、燕麦 (x_3)、土豆 (x_4)、菜豆 (x_5)、马铃薯 (x_6)、萝卜 (x_7)、饲料甜菜 (x_8)、临时牧场干草 (x_9)、永久牧场干草 (x_{10})。（精度：47.6%）

$$Y_1 = 0.39 x_1 + 0.37 x_2 + 0.39 x_3 + 0.27 x_4 + 0.22 x_5 \\ + 0.30 x_6 + 0.32 x_7 + 0.26 x_8 + 0.24 x_9 + 0.34 x_{10}$$

第一主成份 y_1 与 x_1, \dots, x_{10} 均正有关。所以 y_1 称为——“**水平因子**”，可用于评估排序。即：某个样本点在 y_1 上取值很大时，它在 x_1, \dots, x_{10} 取值都会很大。

以上内容仅为本文档的试下载部分，为可阅读页数的一半内容。如要下载或阅读全文，请访问：<https://d.book118.com/358022076045006132>