

# 多元统计

陈崇双

西南交通大学数学学院统计系

[ccsmars@swjtu.edu.cn](mailto:ccsmars@swjtu.edu.cn)

2018-2019学年

## 1 主成分分析

- 问题背景
- 变量降维的可行性
- 基本原理
- 主成分的性质
- 由相关阵提取主成分
- 一些细节

# 第一节：问题背景

为了全面准确地反映事物的特征和变化规律，往往考虑与其有关的多个指标（或变量）。

# 问题背景



中国可持续发展指标体系，来源《中国经济网》，[http://www.ce.cn/xwzx/gnsz/gdxw/201801/08/t20180108\\_27611518.shtml](http://www.ce.cn/xwzx/gnsz/gdxw/201801/08/t20180108_27611518.shtml)

# 问题背景

一级指标	评价维度	表征指标
创新驱动	科技投入	研究与实验发展经费投入强度 (%)
	知识生产	国际科技论文被引次数
	科技价值	每万人口发明专利拥有量 (件)
集约高效	价值创造	人均GDP (万元)
	产业质量	工业增加值率 (%)
	空间集约	亩均增加值 (万元)
平衡普惠	产业结构	知识密集型服务业增加值占GDP的比例 (%)
	城乡均衡	城乡人均可支配收入比值
	分配公平	劳动收入占GDP比重 (%)
绿色生态	能源消耗	单位GDP能耗 (吨标准煤/万元)
	水消耗	单位GDP水耗 (立方米/万元)
	碳排放	单位GDP碳排放 (千克/万元)

高质量发展指标体系, 来源《凤凰网》,

<http://wemedia.ifeng.com/88805720/wemedia.shtml>

# 问题背景

- 优点：描述详尽，刻画细腻。
- 缺点：增加问题的复杂性，信息重叠，主次不清，难以获得直观清晰的把握。

# 问题背景

- 优点：描述详尽，刻画细腻。
- 缺点：增加问题的复杂性，信息重叠，主次不清，难以获得直观清晰的把握。

能否用较少的几项指标来代替原来的指标，并且能较多地反映原来指标的信息？

Hotelling于1933年首先提出主成分分析(Principal Component Analysis, PCA)。



Hotelling于1933年首先提出主成分分析(Principal Component Analysis, PCA)。

PCA并不是比较各指标的重要性，将不太重要的指标简单去掉，而是通过全面分析各项指标所携带的信息，从中提取出一些潜在的综合性指标(主成分)。用综合性指标替代原来较多的指标。

Hotelling于1933年首先提出主成分分析(Principal Component Analysis, PCA)。

PCA并不是比较各指标的重要性，将不太重要的指标简单去掉，而是通过全面分析各项指标所携带的信息，从中提取出一些潜在的综合性指标(主成分)。用综合性指标替代原来较多的指标。

例如，上海证券综合指数，简称上证综指，是上海证券交易所编制，以上海证券交易所挂牌上市的全部股票为计算范围，以发行量为权数综合，反映了上海证券交易所的总体走势。(来源搜狗百科)

截止至2018.06

快速检索

标准检索

专业检索

作者发文检索

科研基金检索

句子检索

来源期刊检索

### 1. 输入检索控制条件: ▲

☒ 期刊年期: 从  年到  年 指定期: 
☐ 更新时间:

来源期刊:    来源类别:

支持基金:

☒ ☐ 作者:   作者单位:

### 2. 输入内容检索条件:

☒ ☐ 篇名     并且包含     精确

☐ ☐ 或者      并且包含     精确

☐ 仅限优先出版论文 ☐ 中英文扩展检索

### 3. 您可以按如下文献分组排序方式选择文献: (分组只对前4万条记录分组, 排序只在800万条记录以内有效)

文献分组浏览: 学科类别 期刊名称 研究资助基金 研究层次 文献作者 作者单位 中文关键词 发表年度 不分组

“紫色”刊名为“中国知网”独家出版刊物

文献排序浏览: 发表时间 相关性 被引频次 下载频次

每页记录数: 10  50

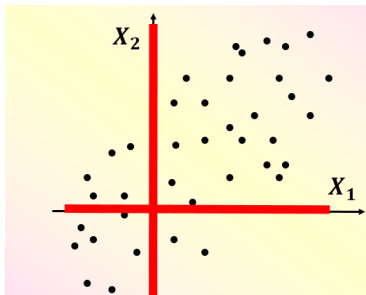
找到 1,621 条结果 浏览 1/82
          后页

计算机软件及应用(309)	轻工业手工业(192)
自动化技术(178)	化学(167)
环境科学与资源利用(107)	电力工业(74)
工业通用技术及设备(60)	电信技术(57)
物理学(54)	地质学(48)
农作物(44)	植物保护(40)
农业基础科学(39)	宏观经济管理(38)
建筑科学与工程(32)	金属学及金属工艺(32)
矿业工程(29)	公路与水路运输(28)
生物学(28)	有机化工(25)
中药学(25)	机械工业(25)
互联网技术(25)	一般化学工业(24)
地球物理学(22)	汽车工业(21)
自然地理学和测绘学(19)	

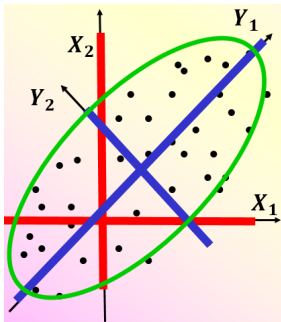
## 第二节：变量降维的可行性

### 例1

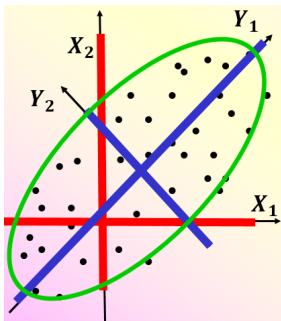
设有一批样品，每个样品涉及两个变量 $X_1$ 和 $X_2$ 。在二维平面中，样本点散布情况如图所示。



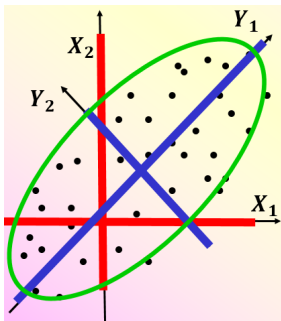
样本点沿 $X_1$ 和 $X_2$ 轴都有一定离散性。



- (1) 若以椭圆长短轴为坐标轴，相当于**平移旋转变换**。不论选用何种坐标系，同样的样品所反映的**信息量**肯定相同。

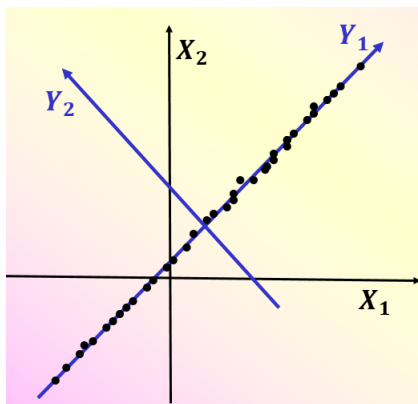


- (1) 若以椭圆长短轴为坐标轴，相当于**平移旋转变换**。不论选用何种坐标系，同样的样品所反映的**信息量**肯定相同。
- (2) 数据越离散，则分布越广泛，代表性越好，故可用**方差**表示信息量。一个变量 $Y_1$ 集中反映 $X_1$ 和 $X_2$ 的大部分信息，忽略 $Y_2$ 也无损大局。

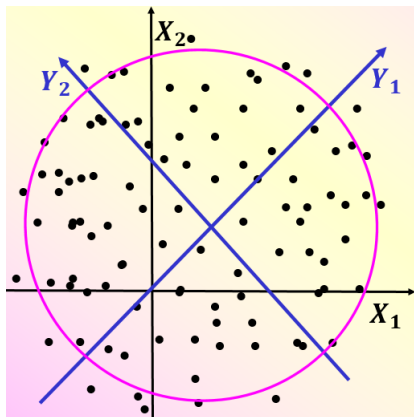


- (1) 若以椭圆长短轴为坐标轴，相当于**平移旋转变换**。不论选用何种坐标系，同样的样品所反映的**信息量**肯定相同。
- (2) 数据越离散，则分布越广泛，代表性越好，故可用**方差**表示信息量。一个变量 $Y_1$ 集中反映 $X_1$ 和 $X_2$ 的大部分信息，忽略 $Y_2$ 也无损大局。
- (3)  $Y_1$ 和 $Y_2$ 正交，避免了信息重叠。





用 $Y_1$ 来反映 $X_1$ 和 $X_2$ 共同的信息，几乎没有损失！



数据杂乱分散，一个变量难以反映 $X_1$ 和 $X_2$ 共同的信息！

# 变量降维的可行性

主成分与原始指标(变量)之间:

- (1) 经坐标轴的平移和旋转变换后, 样品的信息可在新坐标系中选用少数几个轴进行压缩。即, 主成分的数目少于原始指标的数目。

# 变量降维的可行性

主成分与原始指标(变量)之间:

- (1) 经坐标轴的平移和旋转变换后, 样品的信息可在新坐标系中选用少数几个轴进行压缩。即, 主成分的数目少于原始指标的数目。
- (2) 坐标轴进行平移旋转变换, 则样品在新坐标系( $Y_1, Y_2$ )下的坐标, 为原坐标系( $X_1, X_2$ )下坐标的线性组合。即, 每一个主成分都是各原始指标的线性组合。

# 变量降维的可行性

主成分与原始指标(变量)之间:

- (1) 经坐标轴的平移和旋转变换后, 样品的信息可在新坐标系中选用少数几个轴进行压缩。即, 主成分的数目少于原始指标的数目。
- (2) 坐标轴进行平移旋转变换, 则样品在新坐标系( $Y_1, Y_2$ )下的坐标, 为原坐标系( $X_1, X_2$ )下坐标的线性组合。即, 每一个主成分都是各原始指标的线性组合。
- (3) 少数几个主成分, 一方面保留了原始指标的绝大多数信息, 另一方面各主成分不重叠地反映某个方面的综合信息, 即不相关。

概括起来：

主成分分析方法：全面分析各项指标携带的信息，从中提取较少几项综合性指标（主成分），且互不相关，最大限度地保留指标所反映的信息，进而用这较少的几项综合性指标来刻画个体。

# 第三节：基本原理

主要内容：数学模型，求解

# 基本原理

设有一个 $p$ 维总体 $(X_1, X_2, \dots, X_p)$ ，从中抽取 $n$ 个样本，观测值向量分别为 $\mathbf{x}_{(i)}, i = 1, 2, \dots, n$ ，其中 $\mathbf{x}_{(i)} = (x_{i1}, x_{i2}, \dots, x_{ip})^\top, i = 1, 2, \dots, n$ 。希望通过这 $p$ 项可观测指标 $X_1, X_2, \dots, X_p$ 提取出 $m$ (远小于 $p$ )项综合性指标 $Y_1, Y_2, \dots, Y_m$ 。



# 基本原理

设有一个 $p$ 维总体 $(X_1, X_2, \dots, X_p)$ ，从中抽取 $n$ 个样本，观测值向量分别为 $\mathbf{x}_{(i)}, i = 1, 2, \dots, n$ ，其中 $\mathbf{x}_{(i)} = (x_{i1}, x_{i2}, \dots, x_{ip})^T, i = 1, 2, \dots, n$ 。希望通过这 $p$ 项可观测指标 $X_1, X_2, \dots, X_p$ 提取出 $m$ (远小于 $p$ )项综合性指标 $Y_1, Y_2, \dots, Y_m$ 。

PCA的处理办法：

- ① 将每个综合性指标都看成是各可观测指标的线性组合；
- ② 采用方差来度量一个随机变量所包含的信息量。
- ③ 有效性VS.充分性。

# 基本原理

建立如下数学模型：

$$\begin{cases} Y_1 = l_{11}X_1 + l_{12}X_2 + \cdots + l_{1p}X_p \triangleq \boldsymbol{l}_1^\top \boldsymbol{X} \\ Y_2 = l_{21}X_1 + l_{22}X_2 + \cdots + l_{2p}X_p \triangleq \boldsymbol{l}_2^\top \boldsymbol{X} \\ \dots \\ Y_m = l_{m1}X_1 + l_{m2}X_2 + \cdots + l_{mp}X_p \triangleq \boldsymbol{l}_m^\top \boldsymbol{X} \end{cases}$$

# 基本原理

建立如下数学模型：

$$\begin{cases} Y_1 = l_{11}X_1 + l_{12}X_2 + \cdots + l_{1p}X_p \triangleq \mathbf{l}_1^\top \mathbf{X} \\ Y_2 = l_{21}X_1 + l_{22}X_2 + \cdots + l_{2p}X_p \triangleq \mathbf{l}_2^\top \mathbf{X} \\ \cdots \\ Y_m = l_{m1}X_1 + l_{m2}X_2 + \cdots + l_{mp}X_p \triangleq \mathbf{l}_m^\top \mathbf{X} \end{cases}$$

其中 $\mathbf{l}_i = (l_{i1}, l_{i2}, \cdots, l_{ip})^\top$ ,  $i = 1, 2, \dots, m$ 是常向量；

$\mathbf{X} = (X_1, X_2, \cdots, X_p)^\top$ 的均值向量为 $\boldsymbol{\mu}$ ，协方差阵为 $\boldsymbol{\Sigma}$ 。

通过确定 $m$ 个常数向量 $\mathbf{l}_1, \mathbf{l}_2, \cdots, \mathbf{l}_m$ ，使得 $Y_i = \mathbf{l}_i^\top \mathbf{X}$ 的方差尽可能大，且各 $Y_i$ 之间互不相关。

# 基本原理

注1:  $Y_i$ 与 $Y_j$ 不相关, 则 $Cov(Y_i, Y_j) = Cov(l_i X, l_j X) = l_i^\top \Sigma l_j = 0$ 。

# 基本原理

注1:  $Y_i$ 与 $Y_j$ 不相关, 则 $Cov(Y_i, Y_j) = Cov(\mathbf{l}_i^T \mathbf{X}, \mathbf{l}_j^T \mathbf{X}) = \mathbf{l}_i^T \boldsymbol{\Sigma} \mathbf{l}_j = 0$ 。

注2: 由于 $D(Y_i) = D(\mathbf{l}_i^T \mathbf{X}) = \mathbf{l}_i^T \boldsymbol{\Sigma} \mathbf{l}_i$ , 则

$$D((c\mathbf{l}_i)^T \mathbf{X}) = (c\mathbf{l}_i)^T \boldsymbol{\Sigma} (c\mathbf{l}_i) = c^2 \mathbf{l}_i^T \boldsymbol{\Sigma} \mathbf{l}_i = c^2 D(Y_i), \forall c \in R$$

# 基本原理

注1:  $Y_i$ 与 $Y_j$ 不相关, 则 $Cov(Y_i, Y_j) = Cov(\mathbf{l}_i^T \mathbf{X}, \mathbf{l}_j^T \mathbf{X}) = \mathbf{l}_i^T \Sigma \mathbf{l}_j = 0$ 。

注2: 由于 $D(Y_i) = D(\mathbf{l}_i^T \mathbf{X}) = \mathbf{l}_i^T \Sigma \mathbf{l}_i$ , 则

$$D((c\mathbf{l}_i)^T \mathbf{X}) = (c\mathbf{l}_i)^T \Sigma (c\mathbf{l}_i) = c^2 \mathbf{l}_i^T \Sigma \mathbf{l}_i = c^2 D(Y_i), \forall c \in R$$

即有 $D(Y_i) \rightarrow \infty$  ( $\|\mathbf{l}_i\| = \sqrt{\mathbf{l}_i^T \mathbf{l}_i} \rightarrow \infty$ ), 从而无意义。

# 基本原理

因此，线性变换需满足如下约束：

- ①  $\mathbf{l}_i^\top \mathbf{l}_i = l_{i1}^2 + l_{i2}^2 + \cdots + l_{ip}^2 = 1, \quad i = 1, 2, \cdots, p;$
- ②  $\mathbf{l}_i^\top \Sigma \mathbf{l}_j = 0, \quad i, j = 1, 2, \cdots, m, \quad \text{且 } i \neq j.$

# 基本原理

因此，线性变换需满足如下约束：

①  $\mathbf{l}_i^\top \mathbf{l}_i = l_{i1}^2 + l_{i2}^2 + \cdots + l_{ip}^2 = 1, \quad i = 1, 2, \cdots, p;$

②  $\mathbf{l}_i^\top \Sigma \mathbf{l}_j = 0, \quad i, j = 1, 2, \cdots, m, \quad \text{且 } i \neq j。$

也即是说，在  $X_1, X_2, \cdots, X_p$  满足约束(1)的所有线性组合中，

①  $Y_1$  是方差最大者；

②  $Y_2$  是与  $Y_1$  不相关中的方差最大者；

③ ...

④  $Y_m$  是与  $Y_1, Y_2, \cdots, Y_{m-1}$  都不相关中的方差最大者。



# 基本原理

因此，线性变换需满足如下约束：

- ①  $\mathbf{l}_i^\top \mathbf{l}_i = l_{i1}^2 + l_{i2}^2 + \cdots + l_{ip}^2 = 1, \quad i = 1, 2, \cdots, p;$
- ②  $\mathbf{l}_i^\top \boldsymbol{\Sigma} \mathbf{l}_j = 0, \quad i, j = 1, 2, \cdots, m, \quad \text{且 } i \neq j。$

也即是说，在 $X_1, X_2, \cdots, X_p$ 满足约束(1)的所有线性组合中，

- ①  $Y_1$ 是方差最大者；
- ②  $Y_2$ 是与 $Y_1$ 不相关中的方差最大者；
- ③ ...
- ④  $Y_m$ 是与 $Y_1, Y_2, \cdots, Y_{m-1}$ 都不相关中的方差最大者。

综合变量 $Y_1, Y_2, \cdots, Y_m$ 分别称为原始变量的第一、第二、……、第 $m$ 个主成分。

## 定理1

设 $p$ 阶实对称矩阵 $A$ ，将其特征值按大小顺序排列为 $\lambda_1 \geq \lambda_2 \geq \cdots \geq \lambda_p$ ；各特征值对应的标准正交特征向量分别为 $\xi_1, \xi_2, \cdots, \xi_p$ 。则 $x \in \mathbb{R}^p$ 有

$$\max_{x \neq 0} \frac{x^T A x}{x^T x} = \lambda_1, \quad \min_{x \neq 0} \frac{x^T A x}{x^T x} = \lambda_p$$

## 定理1

设 $p$ 阶实对称矩阵 $A$ ，将其特征值按大小顺序排列为 $\lambda_1 \geq \lambda_2 \geq \cdots \geq \lambda_p$ ；各特征值对应的标准正交特征向量分别为 $\xi_1, \xi_2, \cdots, \xi_p$ 。则 $x \in \mathbb{R}^p$ 有

$$\max_{x \neq 0} \frac{x^T A x}{x^T x} = \lambda_1, \quad \min_{x \neq 0} \frac{x^T A x}{x^T x} = \lambda_p$$

且

$$\operatorname{argmax}_{x \neq 0} \frac{x^T A x}{x^T x} = \xi_1, \quad \operatorname{argmin}_{x \neq 0} \frac{x^T A x}{x^T x} = \xi_p$$

# 基本原理

Proof.

矩阵 $\mathbf{A}$ 与单位阵 $\mathbf{I}$ 有谱分解 $\mathbf{A} = \sum_{i=1}^p \lambda_i \boldsymbol{\xi}_i \boldsymbol{\xi}_i^T$ ,  $\mathbf{I} = \sum_{i=1}^p \boldsymbol{\xi}_i \boldsymbol{\xi}_i^T$ 。

# 基本原理

## Proof.

矩阵 $\mathbf{A}$ 与单位阵 $\mathbf{I}$ 有谱分解 $\mathbf{A} = \sum_{i=1}^p \lambda_i \boldsymbol{\xi}_i \boldsymbol{\xi}_i^T$ ,  $\mathbf{I} = \sum_{i=1}^p \boldsymbol{\xi}_i \boldsymbol{\xi}_i^T$ 。

对 $\mathbf{x} \in \mathbb{R}^p$ 有 $\mathbf{x} = \sum_{i=1}^p a_i \boldsymbol{\xi}_i$ , 进而

# 基本原理

## Proof.

矩阵 $\mathbf{A}$ 与单位阵 $\mathbf{I}$ 有谱分解 $\mathbf{A} = \sum_{i=1}^p \lambda_i \boldsymbol{\xi}_i \boldsymbol{\xi}_i^\top$ ,  $\mathbf{I} = \sum_{i=1}^p \boldsymbol{\xi}_i \boldsymbol{\xi}_i^\top$ 。

对 $\mathbf{x} \in \mathbb{R}^p$ 有 $\mathbf{x} = \sum_{i=1}^p a_i \boldsymbol{\xi}_i$ , 进而

$$\begin{aligned} \max_{\mathbf{x} \neq \mathbf{0}} \frac{\mathbf{x}^\top \mathbf{A} \mathbf{x}}{\mathbf{x}^\top \mathbf{x}} &= \max_{\mathbf{x} \neq \mathbf{0}} \frac{\sum_{i=1}^p \lambda_i a_i^2}{\sum_{i=1}^p a_i^2} \leq \frac{\lambda_1 \sum_{i=1}^p a_i^2}{\sum_{i=1}^p a_i^2} = \lambda_1 \\ \min_{\mathbf{x} \neq \mathbf{0}} \frac{\mathbf{x}^\top \mathbf{A} \mathbf{x}}{\mathbf{x}^\top \mathbf{x}} &= \min_{\mathbf{x} \neq \mathbf{0}} \frac{\sum_{i=1}^p \lambda_i a_i^2}{\sum_{i=1}^p a_i^2} \geq \frac{\lambda_p \sum_{i=1}^p a_i^2}{\sum_{i=1}^p a_i^2} = \lambda_p \end{aligned}$$



# 基本原理

类似地有如下结论:

## 命题1

设 $p$ 阶实对称矩阵 $A$ , 将其特征值按大小顺序排列为 $\lambda_1 \geq \lambda_2 \geq \cdots \geq \lambda_p$ ; 各特征值对应的标准正交特征向量分别为 $\xi_1, \xi_2, \cdots, \xi_p$ 。则 $x \in \mathbb{R}^p$ 有

$$\max_{x \neq 0; x^T \xi_i = 0, i=1, 2, \cdots, k} \frac{x^T A x}{x^T x} = \lambda_{k+1}$$

$$\min_{x \neq 0; x^T \xi_i = 0, i=1, 2, \cdots, k} \frac{x^T A x}{x^T x} = \lambda_p$$

# 基本原理

## 命题2

设随机向量  $\mathbf{X} = (X_1, X_2, \dots, X_p)^T$  的协方差矩阵为  $\Sigma$ ,  $\Sigma$  的特征值从大到小排列为  $\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_p$ , 各特征值对应的标准正交特征向量分别为  $\xi_1, \xi_2, \dots, \xi_p$ , 则第  $i$  个主成分为

$$Y_i = \xi_{i1}X_1 + \xi_{i2}X_2 + \dots + \xi_{ip}X_p = \xi_i^T \mathbf{X}, i = 1, 2, \dots, m$$

且

$$\text{Cov}(Y_i, Y_j) = \xi_i^T \Sigma \xi_j = \begin{cases} \lambda_i, & i = j; \\ 0, & i \neq j. \end{cases}$$



# 基本原理

## Proof.

由定理1可知, 对 $\forall \mathbf{x} \in R^p$ , 有

$$D(Y_1) = \max_{\mathbf{x} \neq \mathbf{0}} \frac{\mathbf{x}^\top \Sigma \mathbf{x}}{\mathbf{x}^\top \mathbf{x}}$$

# 基本原理

## Proof.

由定理1可知, 对 $\forall \mathbf{x} \in R^p$ , 有

$$D(Y_1) = \max_{\mathbf{x} \neq \mathbf{0}} \frac{\mathbf{x}^T \Sigma \mathbf{x}}{\mathbf{x}^T \mathbf{x}} = \frac{\boldsymbol{\xi}_1^T \Sigma \boldsymbol{\xi}_1}{\boldsymbol{\xi}_1^T \boldsymbol{\xi}_1}$$

# 基本原理

## Proof.

由定理1可知, 对 $\forall \mathbf{x} \in R^p$ , 有

$$D(Y_1) = \max_{\mathbf{x} \neq \mathbf{0}} \frac{\mathbf{x}^\top \Sigma \mathbf{x}}{\mathbf{x}^\top \mathbf{x}} = \frac{\boldsymbol{\xi}_1^\top \Sigma \boldsymbol{\xi}_1}{\boldsymbol{\xi}_1^\top \boldsymbol{\xi}_1} = \lambda_1$$

# 基本原理

## Proof.

由定理1可知, 对 $\forall \mathbf{x} \in R^p$ , 有

$$D(Y_1) = \max_{\mathbf{x} \neq \mathbf{0}} \frac{\mathbf{x}^\top \Sigma \mathbf{x}}{\mathbf{x}^\top \mathbf{x}} = \frac{\boldsymbol{\xi}_1^\top \Sigma \boldsymbol{\xi}_1}{\boldsymbol{\xi}_1^\top \boldsymbol{\xi}_1} = \lambda_1$$

类似地可得:

$$D(Y_{k+1}) = \max_{\mathbf{x} \neq \mathbf{0}; \mathbf{x}^\top \boldsymbol{\xi}_i = 0, i=1, 2, \dots, k} \frac{\mathbf{x}^\top \mathbf{A} \mathbf{x}}{\mathbf{x}^\top \mathbf{x}}$$

# 基本原理

## Proof.

由定理1可知, 对 $\forall \mathbf{x} \in R^p$ , 有

$$D(Y_1) = \max_{\mathbf{x} \neq \mathbf{0}} \frac{\mathbf{x}^\top \Sigma \mathbf{x}}{\mathbf{x}^\top \mathbf{x}} = \frac{\boldsymbol{\xi}_1^\top \Sigma \boldsymbol{\xi}_1}{\boldsymbol{\xi}_1^\top \boldsymbol{\xi}_1} = \lambda_1$$

类似地可得:

$$D(Y_{k+1}) = \max_{\mathbf{x} \neq \mathbf{0}; \mathbf{x}^\top \boldsymbol{\xi}_i = 0, i=1, 2, \dots, k} \frac{\mathbf{x}^\top A \mathbf{x}}{\mathbf{x}^\top \mathbf{x}} = \frac{\boldsymbol{\xi}_{k+1}^\top \Sigma \boldsymbol{\xi}_{k+1}}{\boldsymbol{\xi}_{k+1}^\top \boldsymbol{\xi}_{k+1}}$$

# 基本原理

## Proof.

由定理1可知, 对 $\forall \mathbf{x} \in R^p$ , 有

$$D(Y_1) = \max_{\mathbf{x} \neq \mathbf{0}} \frac{\mathbf{x}^\top \Sigma \mathbf{x}}{\mathbf{x}^\top \mathbf{x}} = \frac{\boldsymbol{\xi}_1^\top \Sigma \boldsymbol{\xi}_1}{\boldsymbol{\xi}_1^\top \boldsymbol{\xi}_1} = \lambda_1$$

类似地可得:

$$D(Y_{k+1}) = \max_{\mathbf{x} \neq \mathbf{0}; \mathbf{x}^\top \boldsymbol{\xi}_i = 0, i=1, 2, \dots, k} \frac{\mathbf{x}^\top A \mathbf{x}}{\mathbf{x}^\top \mathbf{x}} = \frac{\boldsymbol{\xi}_{k+1}^\top \Sigma \boldsymbol{\xi}_{k+1}}{\boldsymbol{\xi}_{k+1}^\top \boldsymbol{\xi}_{k+1}} = \lambda_{k+1}$$

# 基本原理

## Proof.

由定理1可知, 对 $\forall \mathbf{x} \in R^p$ , 有

$$D(Y_1) = \max_{\mathbf{x} \neq \mathbf{0}} \frac{\mathbf{x}^\top \Sigma \mathbf{x}}{\mathbf{x}^\top \mathbf{x}} = \frac{\xi_1^\top \Sigma \xi_1}{\xi_1^\top \xi_1} = \lambda_1$$

类似地可得:

$$D(Y_{k+1}) = \max_{\mathbf{x} \neq \mathbf{0}; \mathbf{x}^\top \xi_i = 0, i=1, 2, \dots, k} \frac{\mathbf{x}^\top A \mathbf{x}}{\mathbf{x}^\top \mathbf{x}} = \frac{\xi_{k+1}^\top \Sigma \xi_{k+1}}{\xi_{k+1}^\top \xi_{k+1}} = \lambda_{k+1}$$

$$\text{Cov}(Y_i, Y_j) = \xi_i^\top \Sigma \xi_j = \xi_i^\top \left( \sum_{k=1}^p \lambda_k \xi_k \xi_k^\top \right) \xi_j$$

# 基本原理

## Proof.

由定理1可知, 对 $\forall \mathbf{x} \in R^p$ , 有

$$D(Y_1) = \max_{\mathbf{x} \neq \mathbf{0}} \frac{\mathbf{x}^\top \Sigma \mathbf{x}}{\mathbf{x}^\top \mathbf{x}} = \frac{\xi_1^\top \Sigma \xi_1}{\xi_1^\top \xi_1} = \lambda_1$$

类似地可得:

$$D(Y_{k+1}) = \max_{\mathbf{x} \neq \mathbf{0}; \mathbf{x}^\top \xi_i = 0, i=1, 2, \dots, k} \frac{\mathbf{x}^\top A \mathbf{x}}{\mathbf{x}^\top \mathbf{x}} = \frac{\xi_{k+1}^\top \Sigma \xi_{k+1}}{\xi_{k+1}^\top \xi_{k+1}} = \lambda_{k+1}$$

$$\text{Cov}(Y_i, Y_j) = \xi_i^\top \Sigma \xi_j = \xi_i^\top \left( \sum_{k=1}^p \lambda_k \xi_k \xi_k^\top \right) \xi_j = \sum_{k=1}^p \lambda_k \xi_i^\top \xi_k \xi_k^\top \xi_j$$





**总结：**主成分的系数向量即为协差阵的标准正交特征向量，重根按重数计算。

**总结：**主成分的系数向量即为协差阵的标准正交特征向量，重根按重数计算。

无论协差阵 $\Sigma$ 的各特征根是否相等，对应的标准化特征向量 $\xi_1, \xi_2, \dots, \xi_p$ 总是存在，进一步施以正交化（如施密特正交化法）。这样，求主成分的问题就转化为求特征根与特征向量。

## 第四节：主成分的性质

性质1：主成分向量 $\mathbf{Y} = (Y_1, Y_2, \dots, Y_p)^\top$ 的协差阵为对角阵 $\text{diag}(\lambda_1, \lambda_2, \dots, \lambda_p) \triangleq \mathbf{\Lambda}$ 。

## 第四节：主成分的性质

性质1：主成分向量 $\mathbf{Y} = (Y_1, Y_2, \dots, Y_p)^\top$ 的协差阵为对角阵 $\text{diag}(\lambda_1, \lambda_2, \dots, \lambda_p) \triangleq \mathbf{\Lambda}$ 。

性质2：记 $\mathbf{\Sigma} = (\sigma_{ij})_{p \times p}$ ，有 $\sum_{i=1}^p \lambda_i = \sum_{i=1}^p \sigma_{ii}$ 。

## 第四节：主成分的性质

性质1：主成分向量 $\mathbf{Y} = (Y_1, Y_2, \dots, Y_p)^\top$ 的协差阵为对角阵 $\text{diag}(\lambda_1, \lambda_2, \dots, \lambda_p) \triangleq \mathbf{\Lambda}$ 。

性质2：记 $\mathbf{\Sigma} = (\sigma_{ij})_{p \times p}$ ，有 $\sum_{i=1}^p \lambda_i = \sum_{i=1}^p \sigma_{ii}$ 。

注1：关于证明。记 $\mathbf{Q} = (\boldsymbol{\xi}_1, \boldsymbol{\xi}_2, \dots, \boldsymbol{\xi}_p)$ ，则 $\mathbf{\Sigma} = \mathbf{Q}\mathbf{\Lambda}\mathbf{Q}^\top$ ，且有

## 第四节：主成分的性质

性质1：主成分向量 $\mathbf{Y} = (Y_1, Y_2, \dots, Y_p)^\top$ 的协差阵为对角阵 $\text{diag}(\lambda_1, \lambda_2, \dots, \lambda_p) \triangleq \mathbf{\Lambda}$ 。

性质2：记 $\mathbf{\Sigma} = (\sigma_{ij})_{p \times p}$ ，有 $\sum_{i=1}^p \lambda_i = \sum_{i=1}^p \sigma_{ii}$ 。

注1：关于证明。记 $\mathbf{Q} = (\boldsymbol{\xi}_1, \boldsymbol{\xi}_2, \dots, \boldsymbol{\xi}_p)$ ，则 $\mathbf{\Sigma} = \mathbf{Q}\mathbf{\Lambda}\mathbf{Q}^\top$ ，且有

$$\sum_{i=1}^p \sigma_{ii} = \text{tr}(\mathbf{\Sigma})$$

## 第四节：主成分的性质

性质1：主成分向量 $\mathbf{Y} = (Y_1, Y_2, \dots, Y_p)^\top$ 的协差阵为对角阵 $\text{diag}(\lambda_1, \lambda_2, \dots, \lambda_p) \triangleq \mathbf{\Lambda}$ 。

性质2：记 $\mathbf{\Sigma} = (\sigma_{ij})_{p \times p}$ ，有 $\sum_{i=1}^p \lambda_i = \sum_{i=1}^p \sigma_{ii}$ 。

注1：关于证明。记 $\mathbf{Q} = (\boldsymbol{\xi}_1, \boldsymbol{\xi}_2, \dots, \boldsymbol{\xi}_p)$ ，则 $\mathbf{\Sigma} = \mathbf{Q}\mathbf{\Lambda}\mathbf{Q}^\top$ ，且有

$$\sum_{i=1}^p \sigma_{ii} = \text{tr}(\mathbf{\Sigma}) = \text{tr}(\mathbf{Q}\mathbf{\Lambda}\mathbf{Q}^\top)$$

## 第四节：主成分的性质

性质1：主成分向量 $\mathbf{Y} = (Y_1, Y_2, \dots, Y_p)^\top$ 的协差阵为对角阵 $\text{diag}(\lambda_1, \lambda_2, \dots, \lambda_p) \triangleq \mathbf{\Lambda}$ 。

性质2：记 $\mathbf{\Sigma} = (\sigma_{ij})_{p \times p}$ ，有 $\sum_{i=1}^p \lambda_i = \sum_{i=1}^p \sigma_{ii}$ 。

注1：关于证明。记 $\mathbf{Q} = (\boldsymbol{\xi}_1, \boldsymbol{\xi}_2, \dots, \boldsymbol{\xi}_p)$ ，则 $\mathbf{\Sigma} = \mathbf{Q}\mathbf{\Lambda}\mathbf{Q}^\top$ ，且有

$$\sum_{i=1}^p \sigma_{ii} = \text{tr}(\mathbf{\Sigma}) = \text{tr}(\mathbf{Q}\mathbf{\Lambda}\mathbf{Q}^\top) = \text{tr}(\mathbf{\Lambda}\mathbf{Q}\mathbf{Q}^\top)$$



## 第四节：主成分的性质

性质1：主成分向量 $\mathbf{Y} = (Y_1, Y_2, \dots, Y_p)^\top$ 的协差阵为对角阵 $\text{diag}(\lambda_1, \lambda_2, \dots, \lambda_p) \triangleq \mathbf{\Lambda}$ 。

性质2：记 $\mathbf{\Sigma} = (\sigma_{ij})_{p \times p}$ ，有 $\sum_{i=1}^p \lambda_i = \sum_{i=1}^p \sigma_{ii}$ 。

注1：关于证明。记 $\mathbf{Q} = (\boldsymbol{\xi}_1, \boldsymbol{\xi}_2, \dots, \boldsymbol{\xi}_p)$ ，则 $\mathbf{\Sigma} = \mathbf{Q}\mathbf{\Lambda}\mathbf{Q}^\top$ ，且有

$$\sum_{i=1}^p \sigma_{ii} = \text{tr}(\mathbf{\Sigma}) = \text{tr}(\mathbf{Q}\mathbf{\Lambda}\mathbf{Q}^\top) = \text{tr}(\mathbf{\Lambda}\mathbf{Q}\mathbf{Q}^\top) = \text{tr}(\mathbf{\Lambda})$$

## 第四节：主成分的性质

性质1：主成分向量 $\mathbf{Y} = (Y_1, Y_2, \dots, Y_p)^\top$ 的协差阵为对角阵 $\text{diag}(\lambda_1, \lambda_2, \dots, \lambda_p) \triangleq \mathbf{\Lambda}$ 。

性质2：记 $\mathbf{\Sigma} = (\sigma_{ij})_{p \times p}$ ，有 $\sum_{i=1}^p \lambda_i = \sum_{i=1}^p \sigma_{ii}$ 。

注1：关于证明。记 $\mathbf{Q} = (\boldsymbol{\xi}_1, \boldsymbol{\xi}_2, \dots, \boldsymbol{\xi}_p)$ ，则 $\mathbf{\Sigma} = \mathbf{Q}\mathbf{\Lambda}\mathbf{Q}^\top$ ，且有

$$\sum_{i=1}^p \sigma_{ii} = \text{tr}(\mathbf{\Sigma}) = \text{tr}(\mathbf{Q}\mathbf{\Lambda}\mathbf{Q}^\top) = \text{tr}(\mathbf{\Lambda}\mathbf{Q}\mathbf{Q}^\top) = \text{tr}(\mathbf{\Lambda}) = \sum_{i=1}^p \lambda_i$$

## 第四节：主成分的性质

性质1：主成分向量 $\mathbf{Y} = (Y_1, Y_2, \dots, Y_p)^\top$ 的协差阵为对角阵 $\text{diag}(\lambda_1, \lambda_2, \dots, \lambda_p) \triangleq \mathbf{\Lambda}$ 。

性质2：记 $\mathbf{\Sigma} = (\sigma_{ij})_{p \times p}$ ，有 $\sum_{i=1}^p \lambda_i = \sum_{i=1}^p \sigma_{ii}$ 。

注1：关于证明。记 $\mathbf{Q} = (\boldsymbol{\xi}_1, \boldsymbol{\xi}_2, \dots, \boldsymbol{\xi}_p)$ ，则 $\mathbf{\Sigma} = \mathbf{Q}\mathbf{\Lambda}\mathbf{Q}^\top$ ，且有

$$\sum_{i=1}^p \sigma_{ii} = \text{tr}(\mathbf{\Sigma}) = \text{tr}(\mathbf{Q}\mathbf{\Lambda}\mathbf{Q}^\top) = \text{tr}(\mathbf{\Lambda}\mathbf{Q}\mathbf{Q}^\top) = \text{tr}(\mathbf{\Lambda}) = \sum_{i=1}^p \lambda_i$$

注2：意义说明。 $\mathbf{\Sigma}$ 的主对角线元素 $\sigma_{11}, \sigma_{22}, \dots, \sigma_{pp}$ 分别是 $\mathbf{X}$ 各分量 $X_1, X_2, \dots, X_p$ 的方差。

## 第四节：主成分的性质

性质1：主成分向量 $\mathbf{Y} = (Y_1, Y_2, \dots, Y_p)^\top$ 的协差阵为对角阵 $\text{diag}(\lambda_1, \lambda_2, \dots, \lambda_p) \triangleq \mathbf{\Lambda}$ 。

性质2：记 $\mathbf{\Sigma} = (\sigma_{ij})_{p \times p}$ ，有 $\sum_{i=1}^p \lambda_i = \sum_{i=1}^p \sigma_{ii}$ 。

注1：关于证明。记 $\mathbf{Q} = (\boldsymbol{\xi}_1, \boldsymbol{\xi}_2, \dots, \boldsymbol{\xi}_p)$ ，则 $\mathbf{\Sigma} = \mathbf{Q}\mathbf{\Lambda}\mathbf{Q}^\top$ ，且有

$$\sum_{i=1}^p \sigma_{ii} = \text{tr}(\mathbf{\Sigma}) = \text{tr}(\mathbf{Q}\mathbf{\Lambda}\mathbf{Q}^\top) = \text{tr}(\mathbf{\Lambda}\mathbf{Q}\mathbf{Q}^\top) = \text{tr}(\mathbf{\Lambda}) = \sum_{i=1}^p \lambda_i$$

注2：意义说明。 $\mathbf{\Sigma}$ 的主对角线元素 $\sigma_{11}, \sigma_{22}, \dots, \sigma_{pp}$ 分别是 $\mathbf{X}$ 各分量 $X_1, X_2, \dots, X_p$ 的方差。即 $\sum_{i=1}^p \sigma_{ii}$ 刻画了原 $p$ 个变量 $X_1, X_2, \dots, X_p$ 所携带的信息总量。

## 第四节：主成分的性质

性质1：主成分向量 $\mathbf{Y} = (Y_1, Y_2, \dots, Y_p)^\top$ 的协差阵为对角阵 $\text{diag}(\lambda_1, \lambda_2, \dots, \lambda_p) \triangleq \mathbf{\Lambda}$ 。

性质2：记 $\mathbf{\Sigma} = (\sigma_{ij})_{p \times p}$ ，有 $\sum_{i=1}^p \lambda_i = \sum_{i=1}^p \sigma_{ii}$ 。

注1：关于证明。记 $\mathbf{Q} = (\boldsymbol{\xi}_1, \boldsymbol{\xi}_2, \dots, \boldsymbol{\xi}_p)$ ，则 $\mathbf{\Sigma} = \mathbf{Q}\mathbf{\Lambda}\mathbf{Q}^\top$ ，且有

$$\sum_{i=1}^p \sigma_{ii} = \text{tr}(\mathbf{\Sigma}) = \text{tr}(\mathbf{Q}\mathbf{\Lambda}\mathbf{Q}^\top) = \text{tr}(\mathbf{\Lambda}\mathbf{Q}\mathbf{Q}^\top) = \text{tr}(\mathbf{\Lambda}) = \sum_{i=1}^p \lambda_i$$

注2：意义说明。 $\mathbf{\Sigma}$ 的主对角线元素 $\sigma_{11}, \sigma_{22}, \dots, \sigma_{pp}$ 分别是 $\mathbf{X}$ 各分量 $X_1, X_2, \dots, X_p$ 的方差。即 $\sum_{i=1}^p \sigma_{ii}$ 刻画了原 $p$ 个变量 $X_1, X_2, \dots, X_p$ 所携带的信息总量。而 $\sum_{i=1}^p \lambda_i$ 为 $p$ 个主成分 $Y_1, Y_2, \dots, Y_p$ 所携带的信息量之和。即主成分数等于变量数时，无信息损失。

# 主成分的性质

## 定义1

第 $k$ 个主成分 $Y_k$ 的信息量在全部信息量中所占的比例,  $\lambda_k / \sum_{i=1}^p \lambda_i$ 称为第 $k$ 个主成分的贡献率。前 $k$ 个主成分的贡献率之和  $\sum_{i=1}^k \lambda_i / \sum_{i=1}^p \lambda_i$ 称为前 $k$ 个主成分 $Y_1, Y_2, \dots, Y_k$ 的累积贡献率。

# 主成分的性质

## 定义1

第 $k$ 个主成分 $Y_k$ 的信息量在全部信息量中所占的比例,  $\lambda_k / \sum_{i=1}^p \lambda_i$ 称为第 $k$ 个主成分的贡献率。前 $k$ 个主成分的贡献率之和  $\sum_{i=1}^k \lambda_i / \sum_{i=1}^p \lambda_i$ 称为前 $k$ 个主成分 $Y_1, Y_2, \dots, Y_k$ 的累积贡献率。

注1: 主成分的贡献率越大, 可认为其综合能力越强。

# 主成分的性质

## 定义1

第 $k$ 个主成分 $Y_k$ 的信息量在全部信息量中所占的比例,  $\lambda_k / \sum_{i=1}^p \lambda_i$ 称为第 $k$ 个主成分的贡献率。前 $k$ 个主成分的贡献率之和  $\sum_{i=1}^k \lambda_i / \sum_{i=1}^p \lambda_i$ 称为前 $k$ 个主成分 $Y_1, Y_2, \dots, Y_k$ 的累积贡献率。

注1: 主成分的贡献率越大, 可认为其综合能力越强。

注2: 若前 $k$ 个主成分的累积贡献率大于85%, 可认为前 $k$ 个主成分已综合了原始变量的大部分信息, 可不再提取新的主成分了。



# 主成分的性质

## 定义2

主成分 $Y_k$ 与原始变量 $X_i$ 的相关系数 $\rho(Y_k, X_i) = \frac{\sqrt{\lambda_k}}{\sqrt{\sigma_{ii}}} l_{ki}, k, i = 1, 2, \dots, p$ , 称之为 $Y_k$ 关于 $X_i$ 的因子负荷量。

# 主成分的性质

## 定义2

主成分 $Y_k$ 与原始变量 $X_i$ 的相关系数 $\rho(Y_k, X_i) = \frac{\sqrt{\lambda_k}}{\sqrt{\sigma_{ii}}} l_{ki}$ ,  $k, i = 1, 2, \dots, p$ , 称之为 $Y_k$ 关于 $X_i$ 的因子负荷量。

注1: 关于证明。记 $\mathbf{e}_i = (0, \dots, 0, 1, 0, \dots, 0)^\top$ , 则

# 主成分的性质

## 定义2

主成分 $Y_k$ 与原始变量 $X_i$ 的相关系数 $\rho(Y_k, X_i) = \frac{\sqrt{\lambda_k}}{\sqrt{\sigma_{ii}}} l_{ki}$ ,  $k, i = 1, 2, \dots, p$ , 称之为 $Y_k$ 关于 $X_i$ 的因子负荷量。

注1: 关于证明。记 $\mathbf{e}_i = (0, \dots, 0, 1, 0, \dots, 0)^\top$ , 则

$$\begin{aligned}\rho(Y_k, X_i) &= \frac{\text{Cov}(Y_k, X_i)}{\sqrt{\text{Var}(Y_k)}\sqrt{\text{Var}(X_i)}} = \frac{\text{Cov}(\mathbf{l}_k^\top \mathbf{X}, \mathbf{e}_i^\top \mathbf{X})}{\sqrt{\lambda_k}\sqrt{\sigma_{ii}}} \\ &= \frac{\mathbf{l}_k^\top \Sigma \mathbf{e}_i}{\sqrt{\lambda_k}\sqrt{\sigma_{ii}}} = \frac{\lambda_k \mathbf{l}_k^\top \mathbf{e}_i}{\sqrt{\lambda_k}\sqrt{\sigma_{ii}}} = \frac{\sqrt{\lambda_k}}{\sqrt{\sigma_{ii}}} l_{ki}\end{aligned}$$

# 主成分的性质

注2: 意义说明。 $\rho(Y_k, X_i)$ 的绝对值越大, 表示 $Y_k$ 所反映的信息与 $X_i$ 越密切。如果 $\rho(Y_k, X_i)$ 符号为正, 说明主成分 $Y_k$ 与 $X_i$ 正相关, 二者变化趋势相同; 否则为负相关, 变化趋势相反。

# 主成分的性质

注2: 意义说明。 $\rho(Y_k, X_i)$ 的绝对值越大, 表示 $Y_k$ 所反映的信息与 $X_i$ 越密切。如果 $\rho(Y_k, X_i)$ 符号为正, 说明主成分 $Y_k$ 与 $X_i$ 正相关, 二者变化趋势相同; 否则为负相关, 变化趋势相反。

注3: 主成分 $Y_k$ 与原变量 $X_i$ 的线性组合系数 $l_{ki}$ , 表达了原变量 $X_i$ 每增减一个单位, 主成分 $Y_k$ 相应的增减量。因此可以通过观察各组合系数 $l_{k1}, l_{k2}, \dots, l_{kp}$ 的符号、大小, 并结合原变量 $X_i$ 的实际含义, 对主成分 $Y_k$ 的综合含义做出解释并命名。

# 主成分的性质

性质3:  $\sum_{k=1}^p \rho^2(Y_k, X_i) = 1$ 。

# 主成分的性质

性质3:  $\sum_{k=1}^p \rho^2(Y_k, X_i) = 1$ 。

注1: 关于证明。由前面分析可知, 第 $i$ 个主成分的系数向量 $\mathbf{l}_i$ , 为第 $i$ 大特征根所对应的单位正交特征向量 $\xi_i$ ,  $i = 1, 2, \dots, p$ 。一方面,

$$\mathbf{l}_i^T \mathbf{\Lambda} \mathbf{l}_i = (l_{i1}, l_{i2}, \dots, l_{ip}) \text{diag}(\lambda_1, \lambda_2, \dots, \lambda_p) (l_{i1}, l_{i2}, \dots, l_{ip})^T = \sum_{k=1}^p \lambda_k l_{ki}^2$$

# 主成分的性质

性质3:  $\sum_{k=1}^p \rho^2(Y_k, X_i) = 1$ 。

注1: 关于证明。由前面分析可知, 第 $i$ 个主成分的系数向量 $\mathbf{l}_i$ , 为第 $i$ 大特征根所对应的单位正交特征向量 $\xi_i$ ,  $i = 1, 2, \dots, p$ 。一方面,

$$\mathbf{l}_i^T \mathbf{\Lambda} \mathbf{l}_i = (l_{i1}, l_{i2}, \dots, l_{ip}) \text{diag}(\lambda_1, \lambda_2, \dots, \lambda_p) (l_{i1}, l_{i2}, \dots, l_{ip})^T = \sum_{k=1}^p \lambda_k l_{ki}^2$$

另一方面, 若记 $\mathbf{Q} = (\xi_1, \xi_2, \dots, \xi_p)$ ,  $\mathbf{e}_i = (0, \dots, 0, 1, 0, \dots, 0)^T$ , 则 $\mathbf{l}_i = \mathbf{Q} \mathbf{e}_i$ 并且

$$\mathbf{l}_i^T \mathbf{\Lambda} \mathbf{l}_i = \mathbf{e}_i^T \mathbf{Q}^T \mathbf{\Lambda} \mathbf{Q} \mathbf{e}_i = \mathbf{e}_i^T \mathbf{\Sigma} \mathbf{e}_i = \sigma_{ii}$$



# 主成分的性质

性质3:  $\sum_{k=1}^p \rho^2(Y_k, X_i) = 1$ 。

注1: 关于证明。由前面分析可知, 第*i*个主成分的系数向量 $\mathbf{l}_i$ , 为第*i*大特征根所对应的单位正交特征向量 $\xi_i$ ,  $i = 1, 2, \dots, p$ 。一方面,

$$\mathbf{l}_i^T \mathbf{\Lambda} \mathbf{l}_i = (l_{i1}, l_{i2}, \dots, l_{ip}) \text{diag}(\lambda_1, \lambda_2, \dots, \lambda_p) (l_{i1}, l_{i2}, \dots, l_{ip})^T = \sum_{k=1}^p \lambda_k l_{ki}^2$$

另一方面, 若记 $\mathbf{Q} = (\xi_1, \xi_2, \dots, \xi_p)$ ,  $\mathbf{e}_i = (0, \dots, 0, 1, 0, \dots, 0)^T$ , 则 $\mathbf{l}_i = \mathbf{Q} \mathbf{e}_i$ 并且

$$\mathbf{l}_i^T \mathbf{\Lambda} \mathbf{l}_i = \mathbf{e}_i^T \mathbf{Q}^T \mathbf{\Lambda} \mathbf{Q} \mathbf{e}_i = \mathbf{e}_i^T \mathbf{\Sigma} \mathbf{e}_i = \sigma_{ii}$$

綜上有,  $\sum_{k=1}^p \rho^2(Y_k, X_i) = \frac{1}{\sigma_{ii}} \sum_{k=1}^p \lambda_k l_{ki}^2 = 1$ 。得证。

# 主成分的性质

性质3:  $\sum_{k=1}^p \rho^2(Y_k, X_i) = 1$ 。

注1: 关于证明。由前面分析可知, 第*i*个主成分的系数向量 $\mathbf{l}_i$ , 为第*i*大特征根所对应的单位正交特征向量 $\xi_i$ ,  $i = 1, 2, \dots, p$ 。一方面,

$$\mathbf{l}_i^\top \mathbf{\Lambda} \mathbf{l}_i = (l_{i1}, l_{i2}, \dots, l_{ip}) \text{diag}(\lambda_1, \lambda_2, \dots, \lambda_p) (l_{i1}, l_{i2}, \dots, l_{ip})^\top = \sum_{k=1}^p \lambda_k l_{ki}^2$$

另一方面, 若记 $\mathbf{Q} = (\xi_1, \xi_2, \dots, \xi_p)$ ,  $\mathbf{e}_i = (0, \dots, 0, 1, 0, \dots, 0)^\top$ , 则 $\mathbf{l}_i = \mathbf{Q} \mathbf{e}_i$ 并且

$$\mathbf{l}_i^\top \mathbf{\Lambda} \mathbf{l}_i = \mathbf{e}_i^\top \mathbf{Q}^\top \mathbf{\Lambda} \mathbf{Q} \mathbf{e}_i = \mathbf{e}_i^\top \mathbf{\Sigma} \mathbf{e}_i = \sigma_{ii}$$

綜上有,  $\sum_{k=1}^p \rho^2(Y_k, X_i) = \frac{1}{\sigma_{ii}} \sum_{k=1}^p \lambda_k l_{ki}^2 = 1$ 。得证。

# 主成分的性质

注2：意义说明。既然主成分 $\mathbf{Y}$ 是原始变量 $\mathbf{X}$ 的线性组合，因此 $X_i$ 也可以表示成 $Y_1, Y_2, \dots, Y_p$ 的线性组合。

# 主成分的性质

注2: 意义说明。既然主成分 $\mathbf{Y}$ 是原始变量 $\mathbf{X}$ 的线性组合, 因此 $X_i$ 也可以表示成 $Y_1, Y_2, \dots, Y_p$ 的线性组合。从线性回归的角度来看,  $X_i$ 完全由 $Y_1, Y_2, \dots, Y_p$ 解释而没有随机扰动, 从而拟合优度 $R^2 = 1$ 。也即有,  $X_i$ 和 $Y_1, Y_2, \dots, Y_p$ 的全相关系数的平方等于1。

# 主成分的性质

注2：意义说明。既然主成分 $\mathbf{Y}$ 是原始变量 $\mathbf{X}$ 的线性组合，因此 $X_i$ 也可以表示成 $Y_1, Y_2, \dots, Y_p$ 的线性组合。从线性回归的角度来看， $X_i$ 完全由 $Y_1, Y_2, \dots, Y_p$ 解释而没有随机扰动，从而拟合优度 $R^2 = 1$ 。也即有， $X_i$ 和 $Y_1, Y_2, \dots, Y_p$ 的全相关系数的平方等于1。同时 $Y_1, Y_2, \dots, Y_p$ 之间互不相关，所以 $X_i$ 与 $Y_1, Y_2, \dots, Y_p$ 的全相关系数的平方，也就是 $X_i$ 与每个解释变量的相关系数的平方和，从而 $\sum_{i=1}^p \rho^2(Y_k, X_i) = 1$ 。

# 主成分的性质

注2: 意义说明。既然主成分 $\mathbf{Y}$ 是原始变量 $\mathbf{X}$ 的线性组合, 因此 $X_i$ 也可以表示成 $Y_1, Y_2, \dots, Y_p$ 的线性组合。从线性回归的角度来看,  $X_i$ 完全由 $Y_1, Y_2, \dots, Y_p$ 解释而没有随机扰动, 从而拟合优度 $R^2 = 1$ 。也即有,  $X_i$ 和 $Y_1, Y_2, \dots, Y_p$ 的全相关系数的平方等于1。同时 $Y_1, Y_2, \dots, Y_p$ 之间互不相关, 所以 $X_i$ 与 $Y_1, Y_2, \dots, Y_p$ 的全相关系数的平方, 也就是 $X_i$ 与每个解释变量的相关系数的平方和, 从而 $\sum_{i=1}^p \rho^2(Y_k, X_i) = 1$ 。

# 主成分的性质

## 定义3

$X_i$ 与前 $m$ 个主成分 $Y_1, Y_2, \dots, Y_m$ 的全相关系数的平方和, 称为 $Y_1, Y_2, \dots, Y_m$ 对原始变量 $X_i$ 的**方差贡献率**, 即

$$v_i = \frac{1}{\sigma_{ii}} \sum_{k=1}^m \lambda_k l_{ki}^2, i = 1, 2, \dots, p$$

# 主成分的性质

## 定义3

$X_i$ 与前 $m$ 个主成分 $Y_1, Y_2, \dots, Y_m$ 的全相关系数的平方和, 称为 $Y_1, Y_2, \dots, Y_m$ 对原始变量 $X_i$ 的**方差贡献率**, 即

$$v_i = \frac{1}{\sigma_{ii}} \sum_{k=1}^m \lambda_k l_{ki}^2, i = 1, 2, \dots, p$$

注1:  $v_i \leq 1$ , 当 $m = p$ 时,  $v_i = 1$ 。



# 主成分的性质

## 定义3

$X_i$ 与前 $m$ 个主成分 $Y_1, Y_2, \dots, Y_m$ 的全相关系数的平方和, 称为 $Y_1, Y_2, \dots, Y_m$ 对原始变量 $X_i$ 的**方差贡献率**, 即

$$v_i = \frac{1}{\sigma_{ii}} \sum_{k=1}^m \lambda_k l_{ki}^2, i = 1, 2, \dots, p$$

注1:  $v_i \leq 1$ , 当 $m = p$ 时,  $v_i = 1$ 。

注2: 方差贡献率刻画了, 提取的主成分(前 $m$ 个)表达了原始变量的信息, 也即是解释原始变量的能力。

## 第五节：由相关阵提取主成分

实际问题中，原 $p$ 个变量 $X_1, X_2, \dots, X_p$ 可能有不同的量纲，因而可能由于单位选择不合理导致方差相差很悬殊。

## 第五节：由相关阵提取主成分

实际问题中，原 $p$ 个变量 $X_1, X_2, \dots, X_p$ 可能有不同的量纲，因而可能由于单位选择不合理导致方差相差很悬殊。此时通过协差阵提取主成分，会过分关照方差大的变量，而漠视方差小的变量，从而给出不合理的主成分分析结果。

## 第五节：由相关阵提取主成分

实际问题中，原 $p$ 个变量 $X_1, X_2, \dots, X_p$ 可能有不同的量纲，因而可能由于单位选择不合理导致方差相差很悬殊。此时通过协差阵提取主成分，会过分关照方差大的变量，而漠视方差小的变量，从而给出不合理的主成分分析结果。

为了避免这种不合理结果，往往先将各变量标准化，而标准化变量的协差阵正好是原变量 $X_1, X_2, \dots, X_p$ 的相关阵，由此提取主成分。

# 由相关阵提取主成分

令  $Z_i = \frac{X_i - \mu_i}{\sqrt{\sigma_{ii}}}, i = 1, 2, \dots, p$ , 其中  $\mu_i$  与  $\sigma_{ii}$  分别表示变量  $X_i$  的期望与方差, 则  $E(Z_i) = 0, D(Z_i) = 1$ 。

# 由相关阵提取主成分

令  $Z_i = \frac{X_i - \mu_i}{\sqrt{\sigma_{ii}}}, i = 1, 2, \dots, p$ , 其中  $\mu_i$  与  $\sigma_{ii}$  分别表示变量  $X_i$  的期望与方差, 则  $E(Z_i) = 0, D(Z_i) = 1$ 。并令  $\Sigma^* = \text{diag}(\sigma_{11}, \sigma_{22}, \dots, \sigma_{pp})$ , 则  $\mathbf{Z} = (Z_1, Z_2, \dots, Z_p)^\top = (\Sigma^*)^{-1/2}(\mathbf{X} - \boldsymbol{\mu})$ ,

# 由相关阵提取主成分

令  $Z_i = \frac{X_i - \mu_i}{\sqrt{\sigma_{ii}}}$ ,  $i = 1, 2, \dots, p$ , 其中  $\mu_i$  与  $\sigma_{ii}$  分别表示变量  $X_i$  的期望与方差, 则  $E(Z_i) = 0$ ,  $D(Z_i) = 1$ 。并令  $\Sigma^* = \text{diag}(\sigma_{11}, \sigma_{22}, \dots, \sigma_{pp})$ , 则  $\mathbf{Z} = (Z_1, Z_2, \dots, Z_p)^\top = (\Sigma^*)^{-1/2}(\mathbf{X} - \boldsymbol{\mu})$ , 且均值向量  $E(\mathbf{Z}) = \mathbf{0}$ , 协差阵

$$\text{Cov}(\mathbf{Z}) = (\Sigma^*)^{-1/2} \Sigma (\Sigma^*)^{-1/2}$$

# 由相关阵提取主成分

令  $Z_i = \frac{X_i - \mu_i}{\sqrt{\sigma_{ii}}}$ ,  $i = 1, 2, \dots, p$ , 其中  $\mu_i$  与  $\sigma_{ii}$  分别表示变量  $X_i$  的期望与方差, 则  $E(Z_i) = 0$ ,  $D(Z_i) = 1$ 。并令  $\Sigma^* = \text{diag}(\sigma_{11}, \sigma_{22}, \dots, \sigma_{pp})$ , 则  $\mathbf{Z} = (Z_1, Z_2, \dots, Z_p)^\top = (\Sigma^*)^{-1/2}(\mathbf{X} - \boldsymbol{\mu})$ , 且均值向量  $E(\mathbf{Z}) = \mathbf{0}$ , 协差阵

$$\text{Cov}(\mathbf{Z}) = (\Sigma^*)^{-1/2} \Sigma (\Sigma^*)^{-1/2} = \begin{pmatrix} 1 & \rho_{12} & \cdots & \rho_{1p} \\ \rho_{12} & 1 & \cdots & \rho_{2p} \\ \vdots & \vdots & & \vdots \\ \rho_{1p} & \rho_{2p} & \cdots & 1 \end{pmatrix} \triangleq \mathbf{R}$$



# 由相关阵提取主成分

由相关阵求主成分的过程与主成分个数的确定准则，与由协差阵情形相同。仍用 $\lambda_i$ ,  $\xi_i$ 分别表示相关阵 $\mathbf{R}$ 的特征值与对应的标准正交特征向量。

# 由相关阵提取主成分

由相关阵求主成分的过程与主成分个数的确定准则，与由协差阵情形相同。仍用 $\lambda_i$ ,  $\boldsymbol{\xi}_i$ 分别表示相关阵 $\mathbf{R}$ 的特征值与对应的标准正交特征向量。则主成分与原始变量的关系式为：

$$Y_i = \boldsymbol{\xi}_i^T \mathbf{Z} = \boldsymbol{\xi}_i^T (\boldsymbol{\Sigma}^*)^{-1/2} (\mathbf{X} - \boldsymbol{\mu}), i = 1, 2, \dots, p$$

# 由相关阵提取主成分

性质1:  $\sum_{i=1}^p D(Y_i) = \sum_{i=1}^p D(Z_i) = p$ 。

# 由相关阵提取主成分

性质1:  $\sum_{i=1}^p D(Y_i) = \sum_{i=1}^p D(Z_i) = p$ 。

注1: 关于证明: 由于  $Y_i = \xi_i^T \mathbf{Z}$ , 则  $D(Y_i) = \xi_i^T D(Z_i) \xi_i = \xi_i^T \xi_i = 1$ 。

# 由相关阵提取主成分

性质1:  $\sum_{i=1}^p D(Y_i) = \sum_{i=1}^p D(Z_i) = p$ 。

注1: 关于证明: 由于  $Y_i = \xi_i^T \mathbf{Z}$ , 则  $D(Y_i) = \xi_i^T D(Z_i) \xi_i = \xi_i^T \xi_i = 1$ 。

注2: 关于直观意思。经过标准化后, 只要主成分数等于变量数时, 信息就没有损失。

# 由相关阵提取主成分

第 $k$ 个主成分的方差占总方差的比例，即第 $k$ 个主成分的方差贡献率为 $\lambda_k/p$ ，前 $k$ 个主成分的累积方差贡献率为 $\sum_{i=1}^k \lambda_i/p$ 。公共因子 $Y_k$ 关于 $Z_i$ 的因子负荷量为 $\rho(Y_k, Z_i) = \sqrt{\lambda_k} l_{ki}$ ,  $k, i = 1, 2, \dots, p$

## 第六节：一些细节

主要内容：协差阵或相关阵选取，协差阵或相关阵估计，主分分的解释与命名，主成分分析的应用

# 协方差阵或相关阵选取

- 基于协方差阵或相关阵都可求得主成分，二者一般来说有差别，甚至差别很大。



# 协差阵或相关阵选取

- 基于协差阵或相关阵都可求得主成分，二者一般来说有差别，甚至差别很大。
- 度量单位不同或指标取值范围差异非常大时，由相关阵求主成分。

# 协差阵或相关阵选取

- 基于协差阵或相关阵都可求得主成分，二者一般来说有差别，甚至差别很大。
- 度量单位不同或指标取值范围差异非常大时，由相关阵求主成分。
- 指标为同度量或取值范围为同量级时，根据协差阵求解主成分。

# 协差阵或相关阵选取

- 基于协差阵或相关阵都可求得主成分，二者一般来说有差别，甚至差别很大。
- 度量单位不同或指标取值范围差异非常大时，由相关阵求主成分。
- 指标为同度量或取值范围为同量级时，根据协差阵求解主成分。
- 标准化抹杀了原始数据的一部分重要信息，使得标准化后各变量对主成分的作用趋于相等。

# 协方差阵或相关阵估计

以上分析的前提是，随机向量 $\mathbf{X} = (X_1, X_2, \dots, X_p)^\top$ 的协方差阵 $\mathbf{\Sigma}$ 或相关阵 $\mathbf{R}$ 已知。如果二者未知，那么需要先对其进行估计。

# 协方差阵或相关阵估计

以上分析的前提是，随机向量 $\mathbf{X} = (X_1, X_2, \dots, X_p)^\top$ 的协方差阵 $\Sigma$ 或相关阵 $\mathbf{R}$ 已知。如果二者未知，那么需要先对其进行估计。

设 $\mathbf{x}_{(i)} = (x_{i1}, x_{i2}, \dots, x_{ip})^\top$ 是来自总体 $\mathbf{X} = (X_1, X_2, \dots, X_p)^\top$ 的 $n$ 个个体的， $i = 1, 2, \dots, n$ 。则均值向量的估计为

$$\bar{\mathbf{x}} = \frac{1}{n} \sum_{i=1}^n \mathbf{x}_{(i)} = (\bar{x}_1, \bar{x}_2, \dots, \bar{x}_p)^\top$$

# 协方差阵或相关阵估计

以上分析的前提是，随机向量 $\mathbf{X} = (X_1, X_2, \dots, X_p)^\top$ 的协方差阵 $\Sigma$ 或相关阵 $\mathbf{R}$ 已知。如果二者未知，那么需要先对其进行估计。

设 $\mathbf{x}_{(i)} = (x_{i1}, x_{i2}, \dots, x_{ip})^\top$ 是来自总体 $\mathbf{X} = (X_1, X_2, \dots, X_p)^\top$ 的 $n$ 个个体的， $i = 1, 2, \dots, n$ 。则均值向量的估计为

$$\bar{\mathbf{x}} = \frac{1}{n} \sum_{i=1}^n \mathbf{x}_{(i)} = (\bar{x}_1, \bar{x}_2, \dots, \bar{x}_p)^\top$$

协方差阵估计为

$$\mathbf{S} = \frac{1}{n-1} \sum_{i=1}^n (\mathbf{x}_i - \bar{\mathbf{x}})(\mathbf{x}_i - \bar{\mathbf{x}})^\top$$

# 协方差阵或相关阵估计

以上分析的前提是, 随机向量 $\mathbf{X} = (X_1, X_2, \dots, X_p)^\top$ 的协方差阵 $\Sigma$ 或相关阵 $\mathbf{R}$ 已知。如果二者未知, 那么需要先对其进行估计。

设 $\mathbf{x}_{(i)} = (x_{i1}, x_{i2}, \dots, x_{ip})^\top$ 是来自总体 $\mathbf{X} = (X_1, X_2, \dots, X_p)^\top$ 的 $n$ 个个体的,  $i = 1, 2, \dots, n$ 。则均值向量的估计为

$$\bar{\mathbf{x}} = \frac{1}{n} \sum_{i=1}^n \mathbf{x}_{(i)} = (\bar{x}_1, \bar{x}_2, \dots, \bar{x}_p)^\top$$

协方差阵估计为

$$\mathbf{S} = \frac{1}{n-1} \sum_{i=1}^n (\mathbf{x}_i - \bar{\mathbf{x}})(\mathbf{x}_i - \bar{\mathbf{x}})^\top$$

相关阵估计为 $\mathbf{R} = (r_{ij})_{p \times p}$ , 其中 $r_{ij} = \frac{S_{ij}}{\sqrt{S_{ii}S_{jj}}}$ 。

# 主成分的解释与命名

提取出的每个主成分 $Y_k$ 都是原 $p$ 个变量 $X_1, X_2, \dots, X_p$ 的线性组合

$$Y_k = \mathbf{l}_k^\top \mathbf{X} = l_{k1}X_1 + l_{k2}X_2 + \dots + l_{kp}X_p$$

主成分 $Y_k$ 可看成是对原变量 $X_1, X_2, \dots, X_p$ 中某一类信息的综合。



# 主分的解释与命名

提取出的每个主成分 $Y_k$ 都是原 $p$ 个变量 $X_1, X_2, \dots, X_p$ 的线性组合

$$Y_k = \mathbf{l}_k^\top \mathbf{X} = l_{k1}X_1 + l_{k2}X_2 + \dots + l_{kp}X_p$$

主成分 $Y_k$ 可看成是对原变量 $X_1, X_2, \dots, X_p$ 中某一类信息的综合。原变量 $X_1, X_2, \dots, X_p$ 都有明确的实际含义，那么线性组合后得到的新变量 $Y_k$ 的含义又是什么？这就是主分的解释与命名。

# 主成分的解释与命名

提取出的每个主成分 $Y_k$ 都是原 $p$ 个变量 $X_1, X_2, \dots, X_p$ 的线性组合

$$Y_k = \mathbf{l}_k^T \mathbf{X} = l_{k1}X_1 + l_{k2}X_2 + \dots + l_{kp}X_p$$

主成分 $Y_k$ 可看成是对原变量 $X_1, X_2, \dots, X_p$ 中某一类信息的综合。原变量 $X_1, X_2, \dots, X_p$ 都有明确的实际含义，那么线性组合后得到的新变量 $Y_k$ 的含义又是什么？这就是主成分的解释与命名。

主成分的解释与命名，除了需结合与问题有关的专业知识之外，数学手段是分析主成分与原变量的相关性方向和程度。

# 主成分的解释与命名

## 例2

考察某地区上市股票石化板块的五种股票。记 $X_1, X_2, X_3$ 分别表示三家化工企业的股票回升率,  $X_4, X_5$ 分别表示两家石油公司的股票回升率。这五项指标虽较详尽地刻画了石化板块在一周内的股票涨跌, 但进行长年累月的趋势分析就会感到指标多了。为此进行主成分分析。

# 主分的解释与命名

## 例2

考察某地区上市股票石化板块的五种股票。记 $X_1, X_2, X_3$ 分别表示三家化工企业的股票回升率,  $X_4, X_5$ 分别表示两家石油公司的股票回升率。这五项指标虽较详尽地刻画了石化板块在一周内的股票涨跌, 但进行长年累月的趋势分析就会感到指标多了。为此进行主成分分析。

根据过去一年累计50个交易周这五种股票的周回升率数据, 估计出协差阵 $\mathbf{S}$ , 然后求其特征根与特征向量, 仅列出前两个

$$\lambda_1 = 2.854, \quad \mathbf{l}_1^T = (0.464, 0.457, 0.470, 0.421, 0.421)$$

$$\lambda_2 = 0.809, \quad \mathbf{l}_2^T = (0.240, 0.509, 0.260, -0.526, -0.582)$$

且 $(\lambda_1 + \lambda_2) / \sum_{i=1}^5 \lambda_i > 0.85$ 。

# 主分的解释与命名

由此得到两个主成分

$$Y_1 = 0.464X_1 + 0.457X_2 + 0.470X_3 + 0.421X_4 + 0.421X_5$$

$$Y_2 = 0.240X_1 + 0.509X_2 + 0.260X_3 - 0.526X_4 - 0.582X_5$$

# 主分的解释与命名

由此得到两个主成分

$$Y_1 = 0.464X_1 + 0.457X_2 + 0.470X_3 + 0.421X_4 + 0.421X_5$$

$$Y_2 = 0.240X_1 + 0.509X_2 + 0.260X_3 - 0.526X_4 - 0.582X_5$$

原来通过五项指标( $X_1, X_2, X_3, X_4, X_5$ )刻画石化板块的股市行情, 现就只需要两项综合性指标( $Y_1, Y_2$ )了。

# 主分的解释与命名

由此得到两个主成分

$$Y_1 = 0.464X_1 + 0.457X_2 + 0.470X_3 + 0.421X_4 + 0.421X_5$$

$$Y_2 = 0.240X_1 + 0.509X_2 + 0.260X_3 - 0.526X_4 - 0.582X_5$$

原来通过五项指标( $X_1, X_2, X_3, X_4, X_5$ )刻画石化板块的股市行情, 现就只需要两项综合性指标( $Y_1, Y_2$ )了。如已知五支股票的回升率为(3%, 2.5%, 4%, -1%, 1.5%), 则对应的综合性指标为(4.687%, 2.686%)。

# 主成分的解释与命名

由此得到两个主成分

$$Y_1 = 0.464X_1 + 0.457X_2 + 0.470X_3 + 0.421X_4 + 0.421X_5$$

$$Y_2 = 0.240X_1 + 0.509X_2 + 0.260X_3 - 0.526X_4 - 0.582X_5$$

原来通过五项指标( $X_1, X_2, X_3, X_4, X_5$ )刻画石化板块的股市行情, 现就只需要两项综合性指标( $Y_1, Y_2$ )了。如已知五支股票的回升率为(3%, 2.5%, 4%, -1%, 1.5%), 则对应的综合性指标为(4.687%, 2.686%)。

$Y_1$ 的系数皆为正且大小差异不大。可认为 $Y_1$ 是石化股票的一种综合指数。该指数越大, 表示石化板块的股票涨势越强。



# 主分的解释与命名

由此得到两个主成分

$$Y_1 = 0.464X_1 + 0.457X_2 + 0.470X_3 + 0.421X_4 + 0.421X_5$$

$$Y_2 = 0.240X_1 + 0.509X_2 + 0.260X_3 - 0.526X_4 - 0.582X_5$$

原来通过五项指标( $X_1, X_2, X_3, X_4, X_5$ )刻画石化板块的股市行情, 现就只需要两项综合性指标( $Y_1, Y_2$ )了。如已知五支股票的回升率为(3%, 2.5%, 4%, -1%, 1.5%), 则对应的综合性指标为(4.687%, 2.686%)。

$Y_1$ 的系数皆为正且大小差异不大。可认为 $Y_1$ 是石化股票的一种综合指数。该指数越大, 表示石化板块的股票涨势越强。

$Y_2$ 的系数向量的前三个分量(对应化工企业)为正, 后两个分量(对应石油公司)为负。可认为 $Y_2$ 是石化板块中的行业对比指数, 该指数越大说明加工相对于生产更受投资者追捧,

# 主分的解释与命名

由此得到两个主成分

$$Y_1 = 0.464X_1 + 0.457X_2 + 0.470X_3 + 0.421X_4 + 0.421X_5$$

$$Y_2 = 0.240X_1 + 0.509X_2 + 0.260X_3 - 0.526X_4 - 0.582X_5$$

原来通过五项指标( $X_1, X_2, X_3, X_4, X_5$ )刻画石化板块的股市行情, 现就只需要两项综合性指标( $Y_1, Y_2$ )了。如已知五支股票的回升率为(3%, 2.5%, 4%, -1%, 1.5%), 则对应的综合性指标为(4.687%, 2.686%)。

$Y_1$ 的系数皆为正且大小差异不大。可认为 $Y_1$ 是石化股票的一种综合指数。该指数越大, 表示石化板块的股票涨势越强。

$Y_2$ 的系数向量的前三个分量(对应化工企业)为正, 后两个分量(对应石油公司)为负。可认为 $Y_2$ 是石化板块中的行业对比指数, 该指数越大说明加工相对于生产更受投资者追捧, 且 $X_2$ 的系数明显大于 $X_1, X_3$ 的系数, 故 $X_2$ 所对应的企业在加工行业中起着龙头作用。

# 主成分分析的应用

- ①降维：用少数的主成分替换较多的原始变量；
- ②主元回归：用主成分作为自变量进行回归。

# 本章作业

自选一个关注的问题，作为因变量；收集可能的影响因素（多个）作为自变量。对于这些影响因素合理地提取主成分，并将因变量关于这些主成分进行回归，并对比将原始变量直接进行线性回归的效果，尝试解释其中的原因。

要求：限5人组队；需提供数据出处，原始数据；若采用spss等不编程，需注明操作步骤并解释计算结果，严禁直接粘贴软件的图表；若采用matlab, R等编程，需提供代码；提交纸质版，封面为组内成员的姓名和学号。