



统计方法与机器学习

第四章：多重共线性 - 3

倪 蓓

DaSE@ECNU
(lni@dase.ecnu.edu.cn)



目录

① 主成分回归

主成分分析

主成分分析

主成分回归的定义

主成分回归的性质

选择主成分的个数

主成分回归

基本思想

主成分回归 = 主成分分析 + 回归分析

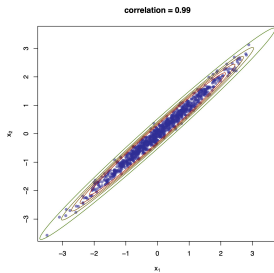
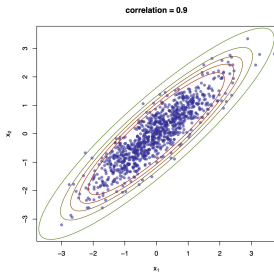
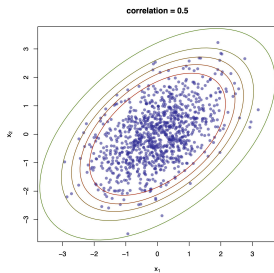
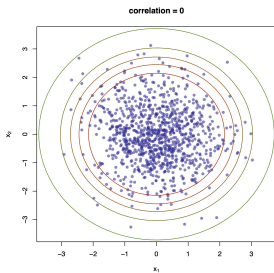
- 回归分析：研究因变量 y 与自变量 $\mathbf{x} = (x_1, x_2, \dots, x_p)'$ 之间的关系，即

$$y = \mathbf{x}'\boldsymbol{\beta} + \varepsilon.$$

- 主成分分析：用 k 个自变量的线性变换（主成分）代替原本的 p 个自变量 ($k < p$)；
- 关键问题：
 - 如何求主成分？
 - 如何利用主成分来估计回归参数 $\boldsymbol{\beta}$ ？
 - 主成分回归估计 $\hat{\boldsymbol{\beta}}_{\text{PC}}$ 与最小二乘估计 $\hat{\boldsymbol{\beta}}$ 有什么关系？

主成分分析

动机



主成分分析

动机

- 由于自变量个数太多，往往自变量之间存在着一定的相关性，因而使得所观测到的数据在一定程度上反映的信息有所重叠。
- 当自变量较多时，我们自然想到能用**较少**的综合变量代替原本多个自变量，而这几个综合变量有能够**尽可能多**地反映原本变量的信息，并且彼此之间互**不相关**。

主成分分析

基本想法

- 设 $\mathbf{x} = (x_1, x_2, \dots, x_p)'$ 是 p 维随机向量
 - 均值 $E(\mathbf{x}) = \boldsymbol{\mu}$;
 - 方差-协方差矩阵 $\text{Var}(\mathbf{x}) = \boldsymbol{\Sigma}$.
- 考虑 \mathbf{x} 的线性变换, 即

$$\begin{cases} z_1 = \mathbf{a}'_1 \mathbf{x} = a_{11}x_1 + a_{21}x_2 + \dots + a_{p1}x_p \\ z_2 = \mathbf{a}'_2 \mathbf{x} = a_{12}x_1 + a_{22}x_2 + \dots + a_{p2}x_p \\ \vdots \\ z_p = \mathbf{a}'_p \mathbf{x} = a_{1p}x_1 + a_{2p}x_2 + \dots + a_{pp}x_p \end{cases}$$

- 易知

$$\text{Var}(z_i) = \mathbf{a}_i' \boldsymbol{\Sigma} \mathbf{a}_i, \quad \text{Cov}(z_i, z_j) = \mathbf{a}_i' \boldsymbol{\Sigma} \mathbf{a}_j$$

主成分分析

基本想法

- 假如，我们希望用一个变量 z_1 来代替原来的 p 个变量 x_1, x_2, \dots, x_p ，那么这个“新”变量 z_1 需要满足
 - z_1 能够尽可能多地反映原来 p 个变量的信息。
- 问题：如何度量 p 维变量的“信息量”？
- 通常，信息越不确定，信息量越大；
- 在统计学中，常常采用“方差”度量数据的不确定性；
- 于是，选取一个新变量 z_1 。如果 $\text{Var}(z_1)$ 越大，那么 z_1 包含的信息量越多。

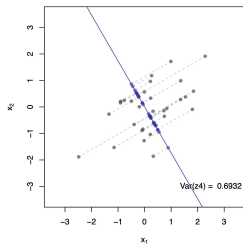
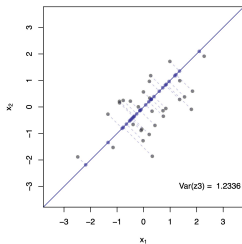
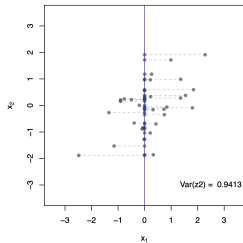
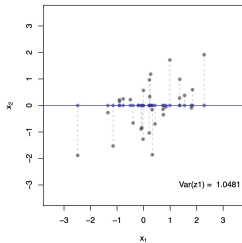
主成分分析

基本想法

- 假如，我们希望用一个变量 z_1 来代替原来的 p 个变量 x_1, x_2, \dots, x_p ，那么这个“新”变量 z_1 需要满足
 - z_1 能够尽可能多地反映原来 p 个变量的信息。
- 问题：如何度量 p 维变量的“信息量”？
- 通常，信息越不确定，信息量越大；
- 在统计学中，常常采用“方差”度量数据的不确定性；
- 于是，选取一个新变量 z_1 。如果 $\text{Var}(z_1)$ 越大，那么 z_1 包含的信息量越多。

主成分分析

动机



主成分分析

基本想法

- 令 $z = \mathbf{a}'\mathbf{x} = a_1x_1 + a_2x_2 + \cdots + a_px_p$ 。
- 我们想要选取的变量 $z_1 = \mathbf{a}'_1\mathbf{x}$ ，使得

$$z_1 = \arg \max \text{Var}(z) \quad \Leftrightarrow \quad \mathbf{a}_1 = \arg \max_{\mathbf{a}} \mathbf{a}'\Sigma\mathbf{a}.$$

- 如果 $\|\mathbf{a}\| > \|\tilde{\mathbf{a}}\|$ ，那么

$$\mathbf{a}'\Sigma\mathbf{a} > \tilde{\mathbf{a}}'\Sigma\tilde{\mathbf{a}}.$$

- 为了避免这个问题，我们需要对 \mathbf{a} 做出一些限制。
- 最常用的限制是： $\mathbf{a}'\mathbf{a} = 1$ 。
- 若存在满足以上约束 \mathbf{a}_1 ，使得 $\text{Var}(z_1)$ 达到最大，称 z_1 为**第一主成分**。

主成分分析

基本想法

- 如果第一主成分不足以代表原来的 p 个变量的绝大部分信息，我们需要进一步考虑 x 的第二个主成分 z_2 。
- 为了有效地代表原始变量的信息， z_1 已包含的信息并不希望体现在 z_2 中。
- 从统计学的角度来看，要求

$$\text{Cov}(z_2, z_1) = \mathbf{a}_2' \Sigma \mathbf{a}_1 = 0$$

- 同样，在两个约束 $\mathbf{a}_2' \mathbf{a}_2 = 1$ 和 $\mathbf{a}_2' \Sigma \mathbf{a}_1 = 0$ 下，当 $\text{Var}(z_2)$ 达到最大时，确定 \mathbf{a}_2 。
- 类似地，可以求第三主成分、第四主成分、...

主成分分析

定义 (主成分)

设

$$\mathbf{x} = (x_1, x_2, \dots, x_p)'$$

为 p 维随机向量且

$$\mathbf{a}_i = (a_{1i}, a_{2i}, \dots, a_{pi})'$$

是个 p 维常数向量。

如果 $z_i = \mathbf{a}_i' \mathbf{x}$ 是 \mathbf{x} 的线性组合, 且满足

- $\mathbf{a}_i' \mathbf{a}_i = 1$;
- 当 $i > 1$ 时, $\mathbf{a}_i' \Sigma \mathbf{a}_j = 0$;
- $\text{Var}(z_i) = \max_{\mathbf{a}' \mathbf{a} = 1, \mathbf{a}' \Sigma \mathbf{a}_j = 0, j=1, \dots, i-1} \text{Var}(\mathbf{a}' \mathbf{x})$ 。

那么, 称 z_i 为 \mathbf{x} 的第 i 个主成分。

主成分分析

主成分的求法：以第一主成分为例

- 设 p 维随机向量 \boldsymbol{x} 的均值 $E(\boldsymbol{x}) = 0$;
- \boldsymbol{x} 的方差-协方差矩阵 $\text{Var}(\boldsymbol{x}) = \boldsymbol{\Sigma}$.
- 问题：如何求第一主成分 $z_1 = \boldsymbol{a}'_1 \boldsymbol{x}$?
- 求 $\boldsymbol{a}_1 = (a_{11}, a_{21}, \dots, a_{p1})'$ 满足

$$\boldsymbol{a}_1 = \arg \max_{\boldsymbol{a}} \text{Var}(\boldsymbol{a}' \boldsymbol{x}) \quad \text{s.t.} \quad \boldsymbol{a}' \boldsymbol{a} = 1$$

主成分分析

主成分的求法：以第一主成分为例

- 采用拉格朗日乘子法，令

$$l(\mathbf{a}) = \text{Var}(\mathbf{a}'\mathbf{x}) - \lambda(\mathbf{a}'\mathbf{a} - 1) = \mathbf{a}'\Sigma\mathbf{a} - \lambda(\mathbf{a}'\mathbf{a} - 1)$$

- 于是有

$$\begin{cases} \frac{\partial l}{\partial \mathbf{a}} = 2(\Sigma - \lambda \mathbf{I})\mathbf{a} = 0 \\ \frac{\partial l}{\partial \lambda} = \mathbf{a}'\mathbf{a} - 1 = 0. \end{cases}$$

- 因为 $\mathbf{a} \neq 0$ ，所以， $|\Sigma - \lambda \mathbf{I}| = 0$ ，而且 $\Sigma\mathbf{a} = \lambda\mathbf{a}$ 。
- 于是，求解第一主成分的问题等价于求 Σ 的特征值和特征向量问题。

主成分分析

定理 3-4

设

$$\mathbf{x} = (x_1, \cdots, x_p)'$$

是 p 维随机向量, 且 $\text{Var}(\mathbf{x}) = \Sigma$ 且满足

- Σ 的特征值为 $\lambda_1 \geq \lambda_2 \geq \cdots \geq \lambda_p \geq 0$;
- $\mathbf{a}_1, \mathbf{a}_2, \cdots, \mathbf{a}_p$ 为相应的单位正交特征向量.

则 \mathbf{x} 的第 i 主成分为

$$z_i = \mathbf{a}_i' \mathbf{x}, i = 1, 2, \cdots, p$$

- 可证明

$$\text{Var}(\mathbf{z}) = \Lambda,$$

其中 $\mathbf{z} = (z_1, z_2, \cdots, z_p)'$ 且 $\Lambda = \text{diag}\{\lambda_1, \lambda_2, \cdots, \lambda_p\}$ 。

主成分分析

讨论

- 考虑以下情况：

$$\Sigma_1 = \begin{pmatrix} 100 & 5 \\ 5 & 1 \end{pmatrix} \quad \text{且} \quad \Sigma_2 = \begin{pmatrix} 1 & 0.5 \\ 0.5 & 1 \end{pmatrix}$$

- Σ_1 的特征值分别为 100.25, 0.75；相应的特征向量为 $(0.9987, 0.0503)'$ 和 $(-0.0503, 0.9987)'$ 。
- Σ_1 的特征值分别为 1.5, 0.5；相应的特征向量为 $(0.7071, 0.7071)'$ 和 $(-0.7071, 0.7071)'$ 。
- 根据方差-协方差矩阵 Σ 来求主成分时，会是优先考虑方差大的变量，有时方差大是根据变量的量纲造成的。
- 对于方差差异大的变量，利用方差-协方差矩阵来求得主成分，可能得到不合理的结果。

主成分分析

讨论

- 考虑以下情况：

$$\Sigma_1 = \begin{pmatrix} 100 & 5 \\ 5 & 1 \end{pmatrix} \quad \text{且} \quad \Sigma_2 = \begin{pmatrix} 1 & 0.5 \\ 0.5 & 1 \end{pmatrix}$$

- Σ_1 的特征值分别为 100.25, 0.75；相应的特征向量为 $(0.9987, 0.0503)'$ 和 $(-0.0503, 0.9987)'$ 。
- Σ_1 的特征值分别为 1.5, 0.5；相应的特征向量为 $(0.7071, 0.7071)'$ 和 $(-0.7071, 0.7071)'$ 。
- 根据方差-协方差矩阵 Σ 来求主成分时，会是优先考虑方差大的变量，有时方差大是根据变量的量纲造成的。
- 对于方差差异大的变量，利用方差-协方差矩阵来求得主成分，可能得到不合理的结果。

主成分分析

讨论

- 为了消除由于量纲不同所带来的不合理结果，我们可以采用“标准化”后的方差-协方差矩阵来计算主成分。
- 令 $\text{Corr}(\mathbf{x})$ 表示 \mathbf{x} 的相关阵。
- 假定由 $\text{Corr}(\mathbf{x})$ 所确定的主成分为 $\mathbf{z}^* = (z_1^*, \dots, z_p^*)'$ ，则 \mathbf{z}^* 的性质如下：
 - 主成分的协方差阵为

$$\text{Var}(\mathbf{z}^*) = \Lambda^* = \text{diag}\{\lambda_1^*, \lambda_2^*, \dots, \lambda_p^*\},$$

其中， $\lambda_1^* \geq \dots \geq \lambda_p^*$;

- 特征值满足

$$\sum_{i=1}^p \lambda_i^* = p$$

主成分分析

如何计算主成分？

- 样本数据（设计矩阵）为

$$\mathbf{X} = \begin{pmatrix} x_{11} & x_{12} & \cdots & x_{1p} \\ x_{21} & x_{22} & \cdots & x_{2p} \\ \vdots & \vdots & & \vdots \\ x_{n1} & x_{n2} & \cdots & x_{np} \end{pmatrix} = \begin{pmatrix} \mathbf{x}'_1 \\ \mathbf{x}'_2 \\ \vdots \\ \mathbf{x}'_n \end{pmatrix}.$$

主成分分析

如何计算主成分？

- 样本方差-协方差矩阵 S ，即

$$S = \frac{1}{n-1} \sum_{i=1}^n (\mathbf{x}_i - \bar{\mathbf{x}})(\mathbf{x}_i - \bar{\mathbf{x}})' \stackrel{\text{def}}{=} (s_{kl})_{p \times p}$$

其中，

$$\bar{\mathbf{x}} = \frac{1}{n} \sum_{i=1}^n \mathbf{x}_i \stackrel{\text{def}}{=} (\bar{x}_1, \dots, \bar{x}_p)'.$$

$$s_{kl} = \frac{1}{n-1} \sum_{i=1}^n (x_{ik} - \bar{x}_k)(x_{il} - \bar{x}_l).$$

主成分分析

如何计算主成分？

- 样本相关阵 \mathbf{R} 为

$$\mathbf{R} = (r_{kl})_{p \times p}$$

其中

$$r_{kl} = \frac{s_{kl}}{\sqrt{s_{kk}s_{ll}}}$$

主成分分析

基本定义

- 假定标准化后的矩阵为

$$\mathbf{X}^* = \begin{pmatrix} x_{11}^* & x_{12}^* & \cdots & x_{1p}^* \\ x_{21}^* & x_{22}^* & \cdots & x_{2p}^* \\ \vdots & \vdots & & \vdots \\ x_{n1}^* & x_{n2}^* & \cdots & x_{np}^* \end{pmatrix}$$

- x_{ij} 与 x_{ij}^* 之间的关系

$$x_{ij} = \frac{x_{ij} - \bar{x}_j}{\sqrt{l_{jj}}}$$

其中, $l_{jj} = \sum_{i=1}^n (x_{ij} - \bar{x}_j)^2 = (n-1)s_{jj}$

主成分分析

基本定义

- 考虑 $(\mathbf{X}^*)' \mathbf{X}^*$ 中的每一个元素

$$\begin{aligned}\sum_{i=1}^n (x_{ik}^*)(x_{il}^*) &= \sum_{i=1}^n \frac{x_{ik} - \bar{x}_k}{\sqrt{l_{kk}}} \times \frac{x_{il} - \bar{x}_l}{\sqrt{l_{ll}}} \\ &= \frac{1}{n-1} \sum_{i=1}^n \frac{(x_{ik} - \bar{x}_k)(x_{il} - \bar{x}_l)}{\sqrt{s_{kk}s_{ll}}} \\ &= r_{kl}\end{aligned}$$

- $(\mathbf{X}^*)' \mathbf{X}^*$ 等价于原始设计矩阵 \mathbf{X} 的样本相关阵 \mathbf{R} 。

主成分回归

动机

- 假定设计矩阵 X 已标准化。假设 $X'X$ 的
 - 特征值为 $\lambda_1 \geq \lambda_2 \geq \cdots \geq \lambda_p$;
 - 其相应单位正交化后的特征向量为 v_1, v_2, \cdots, v_p .
- 我们令
 - $\Lambda = \text{diag}\{\lambda_1, \lambda_2, \cdots, \lambda_p\}$;
 - $V' = (v_1, v_2, \cdots, v_p)$.
- 易知,
 - V' 是由相互正交的列向量组成, 且 $VV' = V'V = I$;
 - 特征值分解为 $X'X = V'\Lambda V$ 即

$$V(X'X)V' = (XV')'(XV') = \Lambda.$$

主成分回归

动机

- 令 $Z = XV'$, 那么

$$Z'Z = \Lambda.$$

- 在矩阵 $Z_{n \times p} = (z_1, z_2, \dots, z_p)$ 中第 j 个列向量 z_j 为 $n \times 1$ 向量, 且满足

$$z_j = Xv_j = (z_{1j}, z_{2j}, \dots, z_{nj})'$$

- 第 j 个主成分为

$$z_j = v_{1j}\mathbf{x}_1 + v_{2j}\mathbf{x}_2 + \dots + v_{pj}\mathbf{x}_p$$

主成分回归

动机

- 由于 X 是已经标准化, 即

$$\sum_{i=1}^n x_{ij} = 0, \quad \sum_{i=1}^n x_{ij}^2 = 1, \quad j = 1, 2, \dots, p.$$

- 而且 $Z = XV'$ 和 $Z'Z = \Lambda$, 我们有

$$\sum_{i=1}^n z_{ij} = 0, \quad \sum_{i=1}^n z_{ij}^2 = \lambda_j, \quad \sum_{i=1}^n z_{ij}z_{ik} = 0 \quad (j \neq k).$$

- 这表明了
 - 矩阵 Z 的各列之间正交;
 - 当 $\lambda_j \approx 0$ 时, $z_{1j}, z_{2j}, \dots, z_{nj}$ 均近似为 0;

主成分回归

主成分估计的定义

- 当 $|X'X| \approx 0$ 时, 存在一个 k , 使得 $\lambda_{k+1}, \dots, \lambda_p$ 均近似为 0。因此, z_{k+1}, \dots, z_p 近似为 0。
- 线性回归模型简化为

$$y = X\beta + \varepsilon = \mathbf{XV}'\mathbf{V}\beta + \varepsilon = \mathbf{Z}_{n \times p}\alpha_{p \times 1} + \varepsilon$$

- 我们将矩阵 Z 和向量 α 按以下方式拆分

$$\begin{aligned}\alpha &= (\alpha_1, \dots, \alpha_k, \alpha_{k+1}, \dots, \alpha_p)' = (\alpha'_1, \alpha'_2)', \\ Z &= (z_1, \dots, z_k, z_{k+1}, \dots, z_p) = (Z_1, Z_2).\end{aligned}$$

- 回归模型也可以写为

$$y = Z_1\alpha_1 + Z_2\alpha_2 + \varepsilon = Z_1\alpha_1 + \varepsilon$$

主成分回归

主成分估计的定义

- 线性回归模型

$$\mathbf{y} = \mathbf{Z}_1 \boldsymbol{\alpha}_1 + \boldsymbol{\varepsilon},$$

- $\boldsymbol{\alpha}_1$ 的最小二乘估计为

$$\hat{\boldsymbol{\alpha}}_1 = (\mathbf{Z}_1' \mathbf{Z}_1)^{-1} \mathbf{Z}_1' \mathbf{y} = \boldsymbol{\Lambda}_1^{-1} \mathbf{Z}_1' \mathbf{y}$$

其中, $\boldsymbol{\Lambda}_1 = \text{diag}\{\lambda_1, \dots, \lambda_k\}$, $\boldsymbol{\Lambda}_2 = \text{diag}\{\lambda_{k+1}, \dots, \lambda_p\}$.

- 那么, $\boldsymbol{\Lambda}$ 可以拆分为 $\boldsymbol{\Lambda}_1$ 和 $\boldsymbol{\Lambda}_2$, 即

$$\boldsymbol{\Lambda} = \begin{pmatrix} \boldsymbol{\Lambda}_1 & \mathbf{0} \\ \mathbf{0} & \boldsymbol{\Lambda}_2 \end{pmatrix}$$

- 相应地, 我们也可以将 \mathbf{V}' 做拆分, 即 $\mathbf{V}' = (\mathbf{V}_1', \mathbf{V}_2')$.

主成分回归

主成分估计的定义

- 根据

$$\beta = V' \alpha,$$

- 回归参数 β 的估计为

$$\hat{\beta}_{\text{PC}} = V' \begin{pmatrix} \hat{\alpha}_1 \\ 0 \end{pmatrix} = V_1' \hat{\alpha}_1 = V_1' \Lambda_1^{-1} Z_1' y$$

- 称 $\hat{\beta}_{\text{PC}}$ 为 β 的**主成分估计**。

主成分回归

性质 1

- 主成分估计与最小二乘估计有什么关系吗？
- 主成分估计可写为最小二乘估计的线性变换，即

$$\begin{aligned}\hat{\beta}_{\text{PC}} &= \mathbf{V}_1' \Lambda_1^{-1} \mathbf{Z}_1' \mathbf{y} \\&= \mathbf{V}_1' \Lambda_1^{-1} \mathbf{V}_1 \mathbf{X}' \mathbf{y} \\&= \mathbf{V}_1' \Lambda_1^{-1} \mathbf{V}_1 \mathbf{X}' \mathbf{X} (\mathbf{X}' \mathbf{X})^{-1} \mathbf{X}' \mathbf{y} \\&= \mathbf{V}_1' \Lambda_1^{-1} \mathbf{V}_1 \mathbf{X}' \mathbf{X} \hat{\beta} \\&= \mathbf{V}_1' \Lambda_1^{-1} \mathbf{V}_1 \mathbf{V}' \Lambda \mathbf{V} \hat{\beta} \\&= \mathbf{V}_1' \Lambda_1^{-1} \mathbf{V}_1 (\mathbf{V}_1', \mathbf{V}_2') \begin{pmatrix} \Lambda_1 & \mathbf{0} \\ \mathbf{0} & \Lambda_2 \end{pmatrix} \begin{pmatrix} \mathbf{V}_1 \\ \mathbf{V}_2 \end{pmatrix} \hat{\beta} \\&= \mathbf{V}_1' \mathbf{V}_1 \hat{\beta}.\end{aligned}$$

主成分回归

性质 1

- 主成分估计与最小二乘估计有什么关系吗？
- 主成分估计可写为最小二乘估计的线性变换，即

$$\begin{aligned}\hat{\beta}_{\text{PC}} &= \mathbf{V}_1' \Lambda_1^{-1} \mathbf{Z}_1' \mathbf{y} \\&= \mathbf{V}_1' \Lambda_1^{-1} \mathbf{V}_1 \mathbf{X}' \mathbf{y} \\&= \mathbf{V}_1' \Lambda_1^{-1} \mathbf{V}_1 \mathbf{X}' \mathbf{X} (\mathbf{X}' \mathbf{X})^{-1} \mathbf{X}' \mathbf{y} \\&= \mathbf{V}_1' \Lambda_1^{-1} \mathbf{V}_1 \mathbf{X}' \mathbf{X} \hat{\beta} \\&= \mathbf{V}_1' \Lambda_1^{-1} \mathbf{V}_1 \mathbf{V}' \Lambda \mathbf{V} \hat{\beta} \\&= \mathbf{V}_1' \Lambda_1^{-1} \mathbf{V}_1 (\mathbf{V}_1', \mathbf{V}_2') \begin{pmatrix} \Lambda_1 & \mathbf{0} \\ \mathbf{0} & \Lambda_2 \end{pmatrix} \begin{pmatrix} \mathbf{V}_1 \\ \mathbf{V}_2 \end{pmatrix} \hat{\beta} \\&= \mathbf{V}_1' \mathbf{V}_1 \hat{\beta}.\end{aligned}$$

主成分回归

性质 2

- 主成分估计是无偏估计吗?
- 主成分估计的期望为

$$\begin{aligned}E(\hat{\beta}_{\text{PC}}) &= E(\mathbf{V}_1' \mathbf{V}_1 \hat{\beta}) \\&= \mathbf{V}_1' \mathbf{V}_1 E(\hat{\beta}) \\&= \mathbf{V}_1' \mathbf{V}_1 \beta\end{aligned}$$

- 由于 $\mathbf{I}_p = \mathbf{V}'\mathbf{V} = \mathbf{V}_1'\mathbf{V}_1 + \mathbf{V}_2'\mathbf{V}_2$, 那么 $\mathbf{V}_1'\mathbf{V}_1 = \mathbf{I}_p - \mathbf{V}_2'\mathbf{V}_2$.
- 当 $k < p$ 时, $\mathbf{V}_1'\mathbf{V}_1\beta = (\mathbf{I} - \mathbf{V}_2'\mathbf{V}_2)\beta \neq \beta$ 。
- 此时, 主成分估计是有偏估计。

主成分回归

性质 3

- 虽然主成分估计是有偏估计，但是在均方误差的意义下，主成分估计是否优于最小二乘估计？

定理 3-5

当设计矩阵 X 存在多重共线性，选择合适的 k ，可使得

$$\text{MSE}(\hat{\beta}_{\text{PC}}) < \text{MSE}(\hat{\beta})$$

主成分回归

性质 3

证明：由于

$$\hat{\beta}_{\text{PC}} = \mathbf{V}' \begin{pmatrix} \hat{\alpha}_1 \\ \mathbf{0} \end{pmatrix}$$

我们有

$$\begin{aligned} \text{MSE}(\hat{\beta}_{\text{PC}}) &= E(\hat{\beta}_{\text{PC}} - \beta)'(\hat{\beta}_{\text{PC}} - \beta) \\ &= E \left(\mathbf{V}' \begin{pmatrix} \hat{\alpha}_1 \\ \mathbf{0} \end{pmatrix} - \mathbf{V}' \alpha \right)' \left(\mathbf{V}' \begin{pmatrix} \hat{\alpha}_1 \\ \mathbf{0} \end{pmatrix} - \mathbf{V}' \alpha \right) \\ &= E \left(\left(\begin{pmatrix} \hat{\alpha}_1 \\ \mathbf{0} \end{pmatrix} - \begin{pmatrix} \alpha_1 \\ \alpha_2 \end{pmatrix} \right)' \mathbf{V} \mathbf{V}' \left(\begin{pmatrix} \hat{\alpha}_1 \\ \mathbf{0} \end{pmatrix} - \begin{pmatrix} \alpha_1 \\ \alpha_2 \end{pmatrix} \right) \right) \\ &= E \left(\left(\begin{pmatrix} \hat{\alpha}_1 \\ \mathbf{0} \end{pmatrix} - \begin{pmatrix} \alpha_1 \\ \alpha_2 \end{pmatrix} \right)' \left(\begin{pmatrix} \hat{\alpha}_1 \\ \mathbf{0} \end{pmatrix} - \begin{pmatrix} \alpha_1 \\ \alpha_2 \end{pmatrix} \right) \right) \end{aligned}$$

主成分回归

性质 3

证明:

$$\begin{aligned}\text{MSE}(\hat{\beta}_{\text{PC}}) &= E(\hat{\alpha}_1 - \alpha_1)'(\hat{\alpha}_1 - \alpha_1) + \|\alpha_2\|^2 \\&= E(\epsilon' \mathbf{Z}_1 \Lambda_1^{-2} \mathbf{Z}_1' \epsilon) + \|\alpha_2\|^2 \\&= \sigma^2 \text{tr}(\Lambda_1^{-1}) + \|\alpha_2\|^2 \\&= \sigma^2 \sum_{j=1}^k \lambda_j^{-1} + \sum_{j=k+1}^p \alpha_j^2 \\&= \text{MSE}(\hat{\beta}) + \left(\sum_{j=k+1}^p \alpha_j^2 - \sigma^2 \sum_{j=k+1}^p \lambda_j^{-1} \right)\end{aligned}$$

主成分回归

性质 3

证明：我们已证明了

$$\text{MSE}(\hat{\beta}_{\text{PC}}) = \text{MSE}(\hat{\beta}) + \left(\sum_{j=k+1}^p \alpha_j^2 - \sigma^2 \sum_{j=k+1}^p \lambda_j^{-1} \right)$$

由于设计矩阵存在多重共线性，因此有一部分特征值 λ_j 非常接近于零。

不妨设后 $p - k$ 个特征值接近于零，则 $\sum_{j=k+1}^p \lambda_j^{-1}$ 将会很大，这导致了第二项为负。所以，

$$\text{MSE}(\hat{\beta}_{\text{PC}}) < \text{MSE}(\hat{\beta})$$

主成分回归

性质 4

- 主成分估计的长度 vs 最小二乘估计的长度
- 主成分估计的模的平方为

$$\begin{aligned}\|\hat{\beta}_{\text{PC}}\|^2 &= (\hat{\beta}_{\text{PC}})' \hat{\beta}_{\text{PC}} \\ &= (\mathbf{V}_1' \mathbf{V}_1 \hat{\beta})' (\mathbf{V}_1' \mathbf{V}_1 \hat{\beta}) \\ &= \hat{\beta}' \mathbf{V}_1' \mathbf{V}_1 \mathbf{V}_1' \mathbf{V}_1 \hat{\beta} \\ &= \hat{\beta}' \mathbf{V}_1' \mathbf{V}_1 \hat{\beta} \\ &\leq \hat{\beta}' \hat{\beta} \\ &= \|\hat{\beta}\|^2\end{aligned}$$

- 所以, $\|\hat{\beta}_{\text{PC}}\| \leq \|\hat{\beta}\|$ 主成分估计是压缩估计。

主成分回归

如何选择主成分的个数？

- **保留特征值比重大的主成分**

- 由于 $\sum_{j=1}^p \lambda_j = p$, 通常称 $\frac{\lambda_j}{p}$ 为第 j 个主成分 z_j 的贡献率。而 $\sum_{j=1}^k \frac{\lambda_j}{p}$ 为前 k 个主成分的累积贡献率。
- 具体方案：给定一个定值 $c_{pc}(0 < c_{pc} < 1)$, 如果存在 k , 使得

$$\sum_{j=1}^{k-1} \frac{\lambda_j}{p} < c_{pc}, \quad \sum_{j=1}^k \frac{\lambda_j}{p} \geq c_{pc}.$$

由此选取 k 。

- 通常, $c_{pc} = 70\%, 75\%$ 或 80% 。

主成分回归

如何选择主成分的个数？

- 删除特征值接近于零的主成分

- 具体方案：给定一个定值 c_0 ，如果

$$\lambda_k \geq c_0, \lambda_{k+1} < c_0$$

由此选取 k 。

- 均方误差确定 k

- 由于 $\sum_{j=1}^p \lambda_j^{-1}$ 与估计的均方误差有关，我们并不希望这个值太大，我们可以选取 k 满足

$$\sum_{j=1}^k \lambda_j^{-1} \leq 5k.$$