



上海交通大学
SHANGHAI JIAO TONG UNIVERSITY



第二讲

主成份分析

Principal Component Analysis (PCA)

彭志科

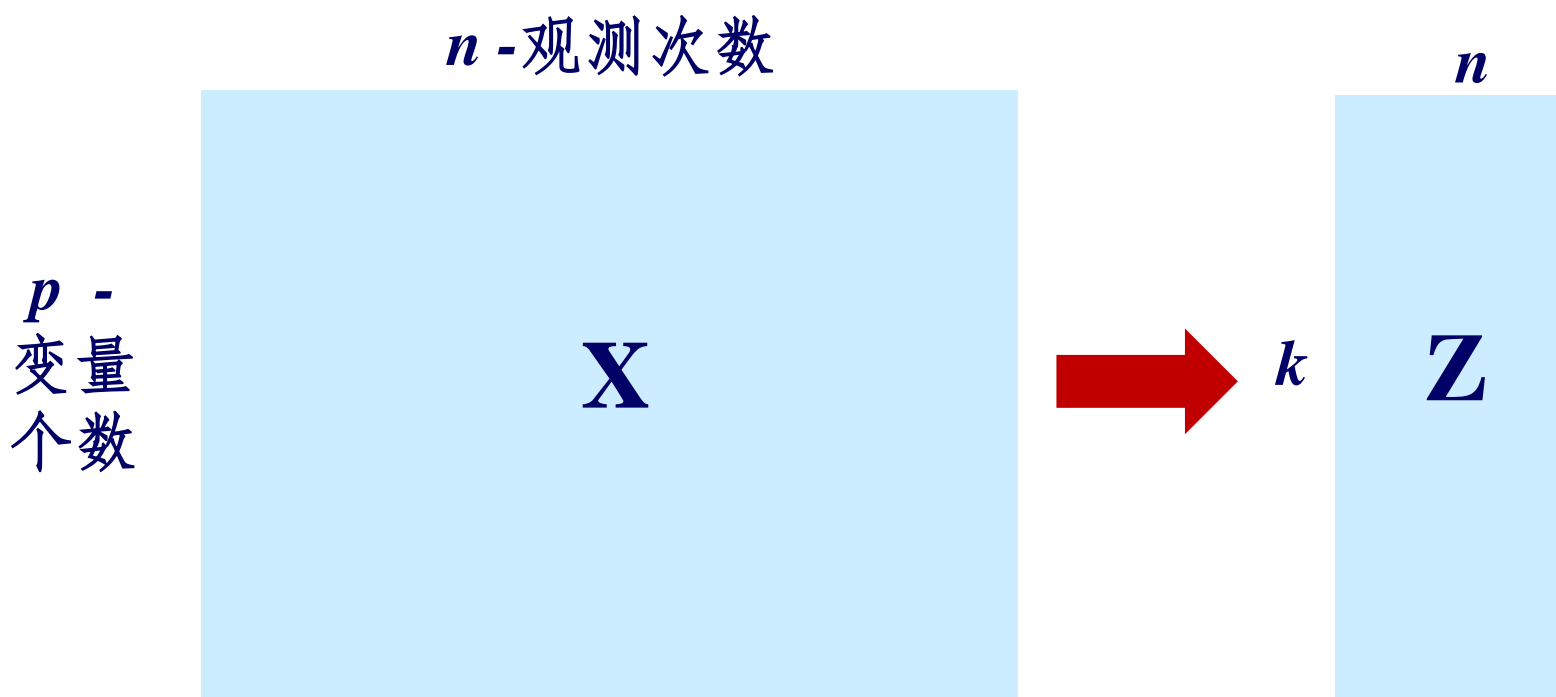
Email: z.peng@sjtu.edu.cn

上海交通大学

机械系统与振动国家重点实验室

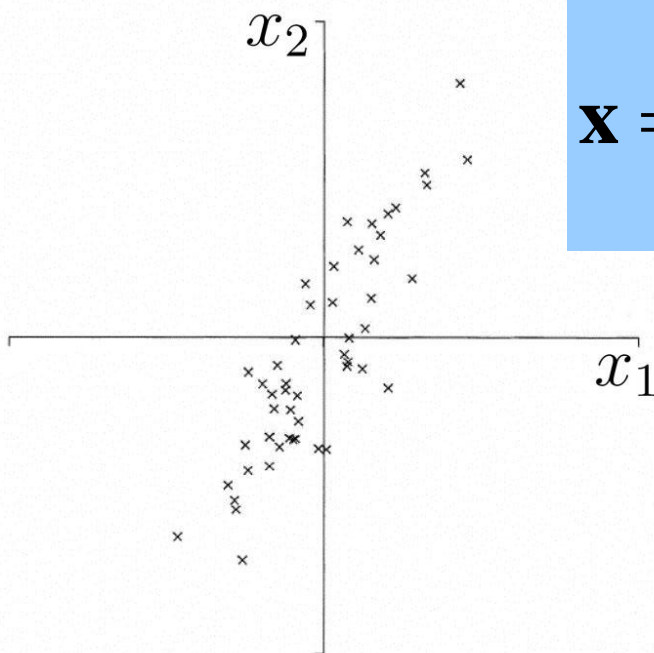
基本概念

- PCA的基本目的是对高维观测数据进行线性变换以构造出一组低维数据，在尽量不损失观测信息的同时，达到数据压缩的目的。



基本概念

示例



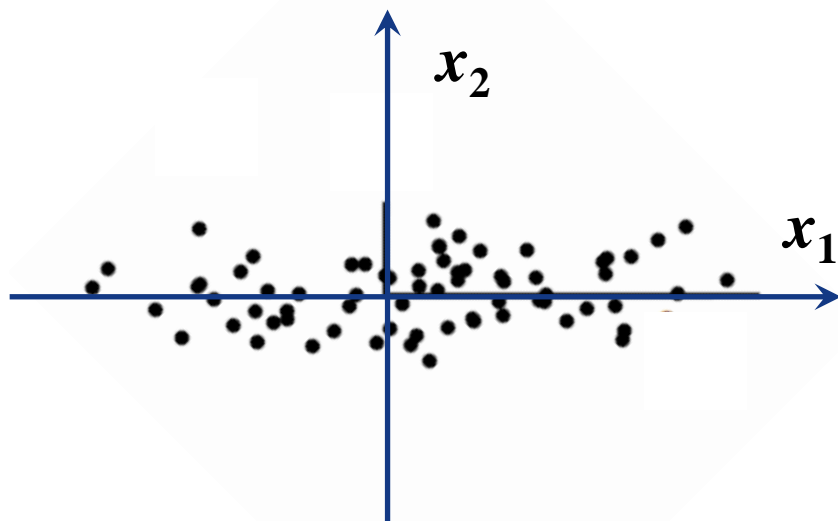
$$\mathbf{X} = \begin{pmatrix} \mathbf{x}_1 \\ \mathbf{x}_2 \end{pmatrix} = \begin{pmatrix} x_{11} & \cdots & x_{1n} \\ x_{21} & \cdots & x_{2n} \end{pmatrix}$$

2D空间的 n 次观测

问题：能否从 x_1 和 x_2 的两个量的 n 次观测数据构造出一个变量，它可以最大的保留原始观测信息？

基本概念

信息量



预处理: $\bar{x} = 0$

$$\text{var}(x) = \frac{1}{n} \sum_{k=1}^n x_k^2$$

$$\text{var}(x) = \frac{1}{n} \langle x_k, x_k \rangle$$

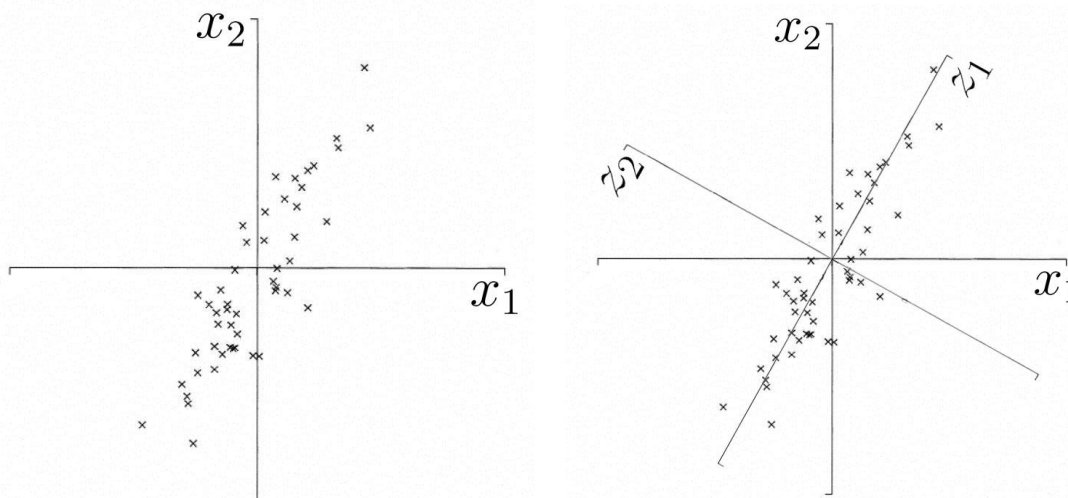
问题：两个观测量 x_1 和 x_2 ，谁携带的信息大？

答案： x_1 ，因为 $\text{var}(x_1) > \text{var}(x_2)$

另外， $SNR = \frac{\text{var}(\text{signal})}{\text{var}(\text{noise})}$ ， $SNR(x_1) > SNR(x_2)$

基本概念

续问题：能否从 x_1 和 x_2 的两个量的 n 次观测数据构造出一个变量，它可以最大的保留原始观测信息？



答案：
$$z_1 = c_{11}x_1 + c_{12}x_2; (c_{11}, c_{12}) = \arg \max_{(\lambda_1, \lambda_2)} (\text{var}(\lambda_1 x_1 + \lambda_2 x_2))$$

z_1 ：第1个主成份； z_2 ：第2个主成份，垂直于 z_1

主成份：是一系列对观测数据的最小二乘线性拟合，且相互正交。

基本概念

高维情况

$$\begin{pmatrix} \mathbf{x}_1 \\ \vdots \\ \mathbf{x}_p \end{pmatrix} = \begin{pmatrix} x_{11} & \cdots & x_{1n} \\ \vdots & \mathbf{X} & \vdots \\ x_{p1} & \cdots & x_{pn} \end{pmatrix} \quad \begin{pmatrix} \mathbf{c}_1 \\ \vdots \\ \mathbf{c}_p \end{pmatrix} = \begin{pmatrix} c_{11} & \cdots & c_{1p} \\ \vdots & \mathbf{C} & \vdots \\ c_{p1} & \cdots & c_{pp} \end{pmatrix} \quad \begin{pmatrix} \mathbf{z}_1 \\ \vdots \\ \mathbf{z}_p \end{pmatrix} = \begin{pmatrix} z_{11} & \cdots & z_{1n} \\ \vdots & \mathbf{Z} & \vdots \\ z_{p1} & \cdots & z_{pn} \end{pmatrix}$$

$$\mathbf{Z} = \mathbf{C}\mathbf{X}$$

第1个主成份: $\mathbf{z}_1 = \mathbf{c}_1\mathbf{X}$

其中, $(c_{11} \cdots c_{1p})$ 的选择, 是要使 $\text{var}(\mathbf{z}_1)$ 最大。

并且 $\langle \mathbf{c}_1, \mathbf{c}_1 \rangle = 1$

基本概念

高维情况

第 k 个主成份: $\mathbf{z}_k = \mathbf{c}_k \mathbf{X}$

其中 $(\mathbf{c}_{k1} \cdots \mathbf{c}_{kp})$ 的选择, 是要使在满足下列条件的情况下 $\text{var}(\mathbf{z}_k)$ 最大。

1) $\text{cov}(\mathbf{z}_k, \mathbf{z}_l) = 0, \text{ for } k > l > 1$

$$\rightarrow \frac{1}{n} \langle \mathbf{z}_k, \mathbf{z}_l \rangle = 0$$

2) $\langle \mathbf{c}_k, \mathbf{c}_k \rangle = 1$

$$\langle \mathbf{c}_k, \mathbf{c}_l \rangle = \delta(k - l)$$

$$\{\mathbf{c}_i, i = 1, \dots, p\}$$

单位正交基

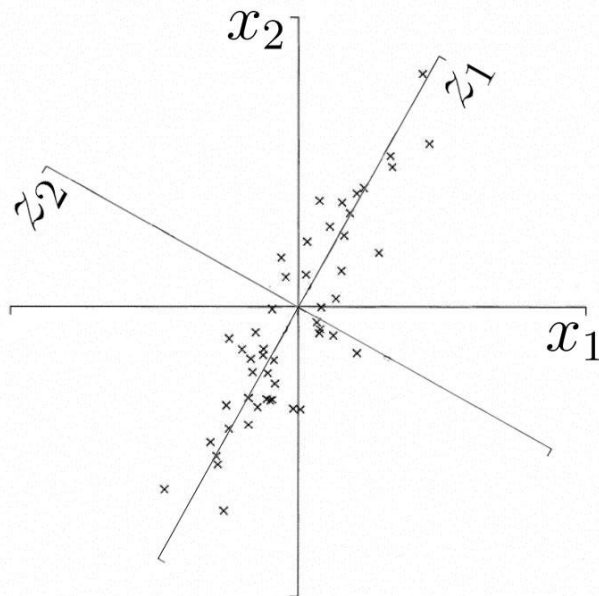
基本概念

$$\mathbf{Z} = \mathbf{C}\mathbf{X}$$

$$\{\mathbf{c}_i, i = 1, \dots, p\}$$

单位正交基

PCA的实质是将高维观测数据投影到一组新的正交坐标基，使得观测数据的信息主要集中在少数几个坐标上。



算法 - I

如何确定变换矩阵 \mathbf{C}

$$\mathbf{Z} = \mathbf{C}\mathbf{X}$$

$$\begin{aligned}\text{cov}(\mathbf{Z}) &= \frac{1}{n} \mathbf{Z}\mathbf{Z}^T = \frac{1}{n} (\mathbf{C}\mathbf{X})(\mathbf{C}\mathbf{X})^T \\ &= \frac{1}{n} \mathbf{C}\mathbf{X}\mathbf{X}^T \mathbf{C}^T \\ &= \mathbf{C} \text{cov}(\mathbf{X}) \mathbf{C}^T\end{aligned}$$

因为 $\text{cov}(z_k, z_l) = 0, \text{ for } k > l > 1$

$\text{cov}(\mathbf{Z})$ 必须为对角阵

算法 - I

选择变换矩阵 \mathbf{C} 使得 $\text{cov}(\mathbf{Z})$ 为对角阵

$$\text{cov}(\mathbf{Z}) = \mathbf{C} \text{cov}(\mathbf{X}) \mathbf{C}^T$$

因为 $\text{cov}(\mathbf{X}) = \mathbf{E} \mathbf{D} \mathbf{E}^T$, \mathbf{D} 为对角阵, \mathbf{E} 为特征向量

选择

$$\mathbf{C} = \mathbf{E}^T; \mathbf{E} = \mathbf{C}^T$$

$$\mathbf{E}^{-1} = \mathbf{E}^T$$

$$\mathbf{E} \mathbf{E}^T = \mathbf{I}$$

$$\begin{aligned} \text{cov}(\mathbf{Z}) &= \mathbf{C} (\mathbf{E} \mathbf{D} \mathbf{E}^T) \mathbf{C}^T \\ &= (\mathbf{C} \mathbf{C}^T) \mathbf{D} (\mathbf{C} \mathbf{C}^T) \\ &= \mathbf{D} \end{aligned}$$

Matlab Code

```
function [z,C,V] = pca(x)
[M,N] = size(x);
mx = mean(x,2);
x = x - repmat(mx,1,N);
co = 1 / (N-1) * x * x';
[C, V] = eig(co);
V = diag(V);
[temp, iInd] = sort(V);
V = V(iInd);
C = C(:,iInd);
z = C'*x;
```

算法 - II

方阵特征值分解

任意方阵 \mathbf{X}

$$\mathbf{X}\mathbf{v} = \lambda\mathbf{v}$$

λ -特征值; \mathbf{v} -特征向量

矩阵特征值分解

$$\mathbf{X} = \mathbf{Q}\mathbf{\Lambda}\mathbf{Q}^{-1}$$

$\mathbf{\Lambda}$ - 对角阵; \mathbf{Q} -正交阵

奇异值分解 (SVD)

任意矩阵 \mathbf{X}

$$\mathbf{X}_{p \times n} = \mathbf{U}_{p \times p} \mathbf{\Lambda}_{p \times n} \mathbf{V}_{n \times n}^T$$

\mathbf{U}, \mathbf{V} - 正交阵

$\mathbf{\Lambda}$ - 对角元非零 (奇异值)

算法 - II

SVD与PCA的关系

$$\mathbf{X}_{p \times n} = \mathbf{U}_{p \times p} \mathbf{\Lambda}_{p \times n} \mathbf{V}_{n \times n}^T$$



$$\mathbf{Z} = \mathbf{U}_{p \times p}^T \mathbf{X}_{p \times n} = \mathbf{\Lambda}_{p \times n} \mathbf{V}_{n \times n}^T$$



$$\begin{aligned} \text{cov}(\mathbf{Z}) &= \mathbf{\Lambda}_{p \times n} \mathbf{V}_{n \times n}^T \mathbf{V}_{n \times n} \mathbf{\Lambda}_{n \times p}^T \\ &= \mathbf{D}_{p \times p} \end{aligned}$$

$$d_{ii} = \sigma_{ii}^2$$

$$\mathbf{Z} = \mathbf{U}\mathbf{X}$$

Matlab Code

```
function [z,C,V] = pca_svd(x)
```

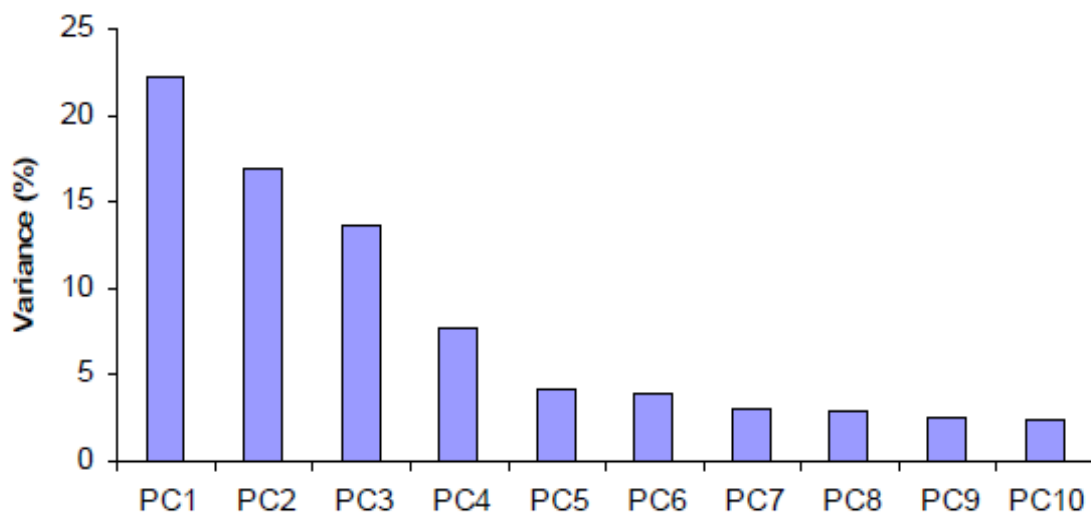
```
[M,N] = size(x);  
mx = mean(x,2);  
x = x - repmat(mx,1,N);  
Y = x' / sqrt(N-1);  
[u, S, C] = svd(Y);  
S = diag(S);  
V = S .* S;  
z = C' * x;
```

算法 - II

PC个数的确定

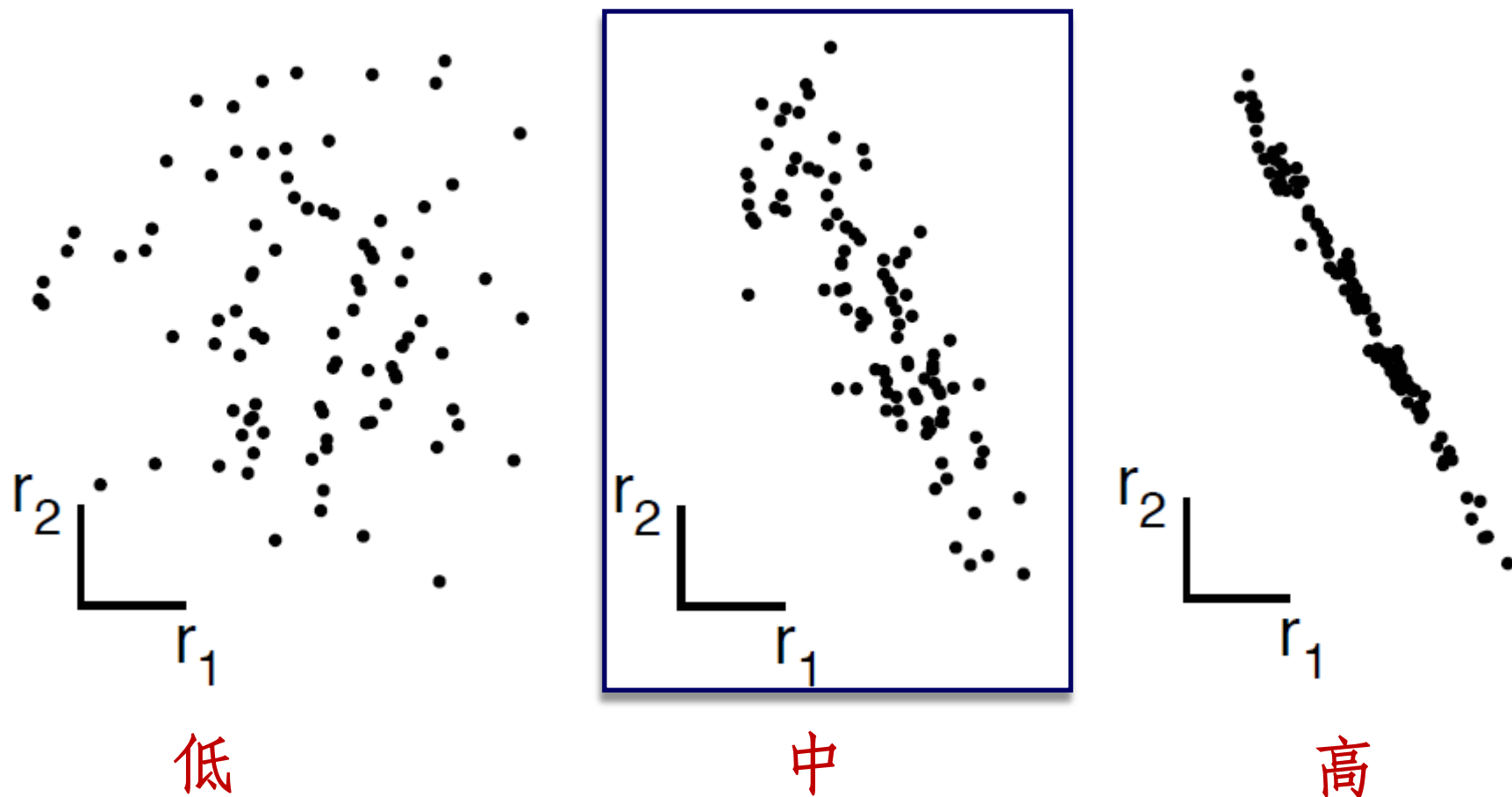
$$\mathbf{Z} = \mathbf{U}_{p \times p}^T \mathbf{X}_{p \times n} = \mathbf{\Lambda}_{p \times n} \mathbf{V}_{n \times n}^T$$

$$\begin{aligned} \mathbf{X}_{p \times n} &\approx \mathbf{U}_{p \times r} \mathbf{\Lambda}_{r \times r} \mathbf{V}_{r \times n}^T \\ &\approx \mathbf{U}_{p \times r} \mathbf{Z}_{r \times n} \end{aligned}$$



讨论

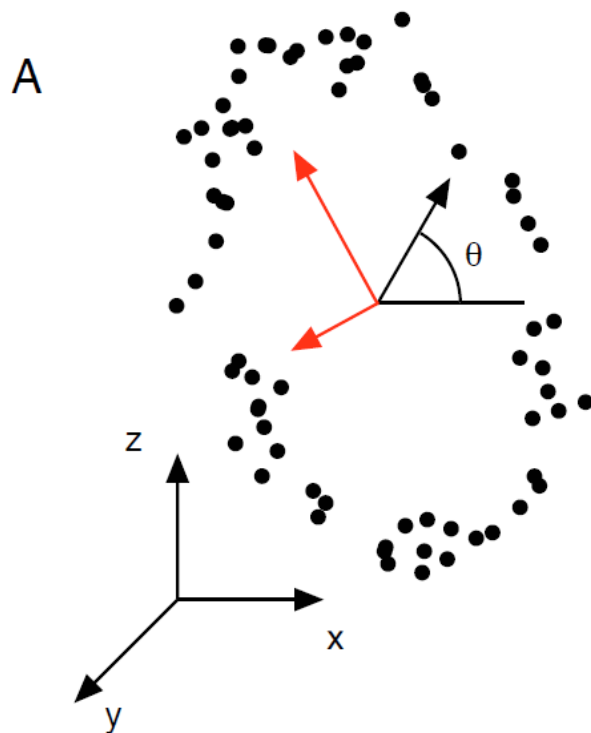
信息冗余度



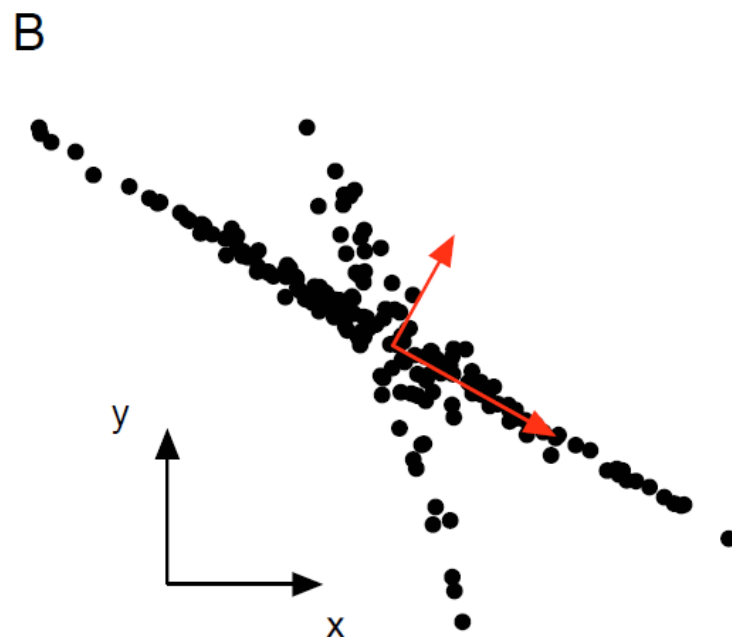
讨论

PCA优点： 非参数化 / 变换矩阵正交

PCA缺点



Kernel PCA



ICA

讨论

模态分析与PCA之间的关系

$$\mathbf{M}\ddot{\mathbf{x}} + \mathbf{K}\mathbf{x} = \mathbf{F}$$

\mathbf{M} , \mathbf{K} 都为对称阵, 存在同一正交阵 \mathbf{P} , 使得

$$\mathbf{P}^T \mathbf{M} \mathbf{P} = \begin{pmatrix} \hat{m}_1 & 0 & 0 \\ 0 & \ddots & 0 \\ 0 & 0 & \hat{m}_n \end{pmatrix} \quad \mathbf{P}^T \mathbf{K} \mathbf{P} = \begin{pmatrix} \hat{k}_1 & 0 & 0 \\ 0 & \ddots & 0 \\ 0 & 0 & \hat{k}_n \end{pmatrix}$$

讨论

模态分析与PCA之间的关系

$$\mathbf{M}\ddot{\mathbf{x}} + \mathbf{K}\mathbf{x} = \mathbf{F}$$

令 $\mathbf{y} = \mathbf{P}^T \mathbf{x} \rightarrow \mathbf{x} = \mathbf{P}\mathbf{y}$



$$\mathbf{M}\mathbf{P}\ddot{\mathbf{y}} + \mathbf{K}\mathbf{P}\mathbf{y} = \mathbf{F}$$



$$\mathbf{P}^T \mathbf{M}\mathbf{P}\ddot{\mathbf{y}} + \mathbf{P}^T \mathbf{K}\mathbf{P}\mathbf{y} = \mathbf{P}^T \mathbf{F}$$



$$\begin{pmatrix} \hat{m}_1 & 0 & 0 \\ 0 & \ddots & 0 \\ 0 & 0 & \hat{m}_n \end{pmatrix} \ddot{\mathbf{y}} + \begin{pmatrix} \hat{k}_1 & 0 & 0 \\ 0 & \ddots & 0 \\ 0 & 0 & \hat{k}_n \end{pmatrix} \mathbf{y} = \begin{pmatrix} \hat{f}_1 \\ \vdots \\ \hat{f}_n \end{pmatrix}$$

讨论

模态分析与PCA之间的关系

$$\begin{pmatrix} \hat{m}_1 & 0 & 0 \\ 0 & \ddots & 0 \\ 0 & 0 & \hat{m}_n \end{pmatrix} \ddot{\mathbf{y}} + \begin{pmatrix} \hat{k}_1 & 0 & 0 \\ 0 & \ddots & 0 \\ 0 & 0 & \hat{k}_n \end{pmatrix} \mathbf{y} = \begin{pmatrix} \hat{f}_1 \\ \vdots \\ \hat{f}_n \end{pmatrix}$$

假设

$$\|\mathbf{y}_1\| > \|\mathbf{y}_2\| \cdots > \|\mathbf{y}_n\|$$

$$\mathbf{x}_{n \times 1} = \mathbf{P}_{n \times n} \mathbf{y}_{n \times 1}$$



$$\mathbf{x}_{n \times 1} \approx \hat{\mathbf{P}}_{n \times k} \mathbf{y}_{k \times 1}$$

$k \ll n$ 有限元

$$\mathbf{X}_{n \times R} \approx \hat{\mathbf{P}}_{n \times k} \mathbf{Y}_{k \times R}$$

R – 采样点数

谢谢聆听
欢迎交流