

## PCA 误差分析

数据的变换、降维：

$m$  个特征方向

$$\begin{aligned} \mathbf{y} &= [y_1, y_2, \dots, y_m]^T \\ &= [\mathbf{w}_1^T \mathbf{x}, \mathbf{w}_2^T \mathbf{x}, \dots, \mathbf{w}_m^T \mathbf{x}]^T, \\ &= \mathbf{W}^T \mathbf{x} \end{aligned} \quad \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_m \end{bmatrix} = \begin{bmatrix} \mathbf{w}_1^T \\ \mathbf{w}_2^T \\ \vdots \\ \mathbf{w}_m^T \end{bmatrix} \mathbf{x}$$

选取  $l$  个方向：

$$\begin{aligned} \mathbf{y} &= [y_1, y_2, \dots, y_l]^T \\ &= [\mathbf{w}_1^T \mathbf{x}, \mathbf{w}_2^T \mathbf{x}, \dots, \mathbf{w}_l^T \mathbf{x}]^T, \\ &= \mathbf{W}^T \mathbf{x} \end{aligned} \quad \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_l \end{bmatrix} = \begin{bmatrix} \mathbf{w}_1^T \\ \mathbf{w}_2^T \\ \vdots \\ \mathbf{w}_l^T \end{bmatrix} \mathbf{x}, \quad l < m$$

数据的反变换、重构：

无损截断重构（舍弃特征值为 0 的方向）

$$\begin{aligned} \mathbf{x} &= \mathbf{W} \mathbf{y} \\ &= \sum_{j=1}^m y_j \mathbf{w}_j, \quad \mathbf{W}^T \mathbf{W} = \mathbf{I} \end{aligned}$$

有损截断重构（舍弃部分特征值大于 0 的方向）

$$\hat{\mathbf{x}} = \sum_{j=1}^l y_j \mathbf{w}_j = [\mathbf{w}_1, \mathbf{w}_2, \dots, \mathbf{w}_l] \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_l \end{bmatrix}, \quad l < m$$

## 误差分析

**证明 1：误差向量和重构空间是正交的**

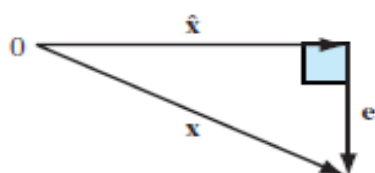
误差向量 = 无损重构 - 有损重构

$$\mathbf{e} = \mathbf{x} - \hat{\mathbf{x}} = \sum_{i=l+1}^m a_i \mathbf{w}_i, \quad \text{代入下式}$$

$$\begin{aligned} \mathbf{e}^T \hat{\mathbf{x}} &= \sum_{i=l+1}^m a_i \mathbf{w}_i^T \sum_{j=1}^l a_j \mathbf{w}_j \\ &= \sum_{i=l+1}^m \sum_{j=1}^l a_i a_j \mathbf{w}_i^T \mathbf{w}_j \\ &= 0 \end{aligned}$$

**证明 2：重构误差等于截断特征值之和**

误差向量的模长？



分别算一下无损重构、有损重构、误差的 2 范数，

$$\|\mathbf{e}\|^2 = \|\mathbf{x}\|^2 - \|\hat{\mathbf{x}}\|^2 = \sum_{j=l+1}^m a_j^2 = \sum_{j=l+1}^m \mathbf{w}_j^T \mathbf{x} \mathbf{x}^T \mathbf{w}_j, \quad ,$$

$$\|\hat{\mathbf{x}}\|^2 = \sum_{j=1}^l a_j^2, \quad \text{其中, } a_j = \mathbf{x}^T \mathbf{w}_j$$

$$\|\mathbf{e}\|^2 = \|\mathbf{x}\|^2 - \|\hat{\mathbf{x}}\|^2 = \sum_{j=l+1}^m a_j^2 = \sum_{j=l+1}^m \mathbf{w}_j^T \mathbf{x} \mathbf{x}^T \mathbf{w}_j$$

对于样本集（ $N$  个样本）的平均误差：

$$\begin{aligned}
 \varepsilon &= \frac{1}{N} \sum_{n=1}^N \|\mathbf{e}_n\|^2 \\
 &= \frac{1}{N} \sum_{j=l+1}^m \mathbf{w}_j^T \left( \sum_{n=1}^N \mathbf{x} \mathbf{x}^T \right) \mathbf{w}_j \\
 &= \sum_{j=l+1}^m \mathbf{w}_j^T \left( \frac{1}{N} \sum_{n=1}^N \mathbf{x} \mathbf{x}^T \right) \mathbf{w}_j \\
 &= \sum_{j=l+1}^m \mathbf{w}_j^T \mathbf{S} \mathbf{w}_j \\
 &= \sum_{j=l+1}^m \sigma_j^2 = \sum_{j=l+1}^m \lambda_j
 \end{aligned}$$

因此样本集的方差、重构后的方差，以及重构误差表示为：

$$\begin{aligned}
 \sum_{j=1}^m \sigma_j^2 &= \sum_{j=1}^m \lambda_j \\
 \sum_{j=1}^l \sigma_j^2 &= \sum_{j=1}^l \lambda_j \\
 \sum_{j=l+1}^m \sigma_j^2 &= \sum_{j=l+1}^m \lambda_j
 \end{aligned}$$