

主成分分析与数据压缩

特征提取与K-L变换

- 特征提取：用映射（或变换）的方法把原始特征变换为较少的新特征
- K-L (Karhunen-Loeve)变换：最优正交线性变换，相应的特征提取方法被称为PCA方法
- PCA (Principle Component Analysis)方法：进行特征降维变换，不能完全地表示原有的对象，能量总会有损失。希望找到一种能量最为集中的的变换方法使损失最小。

K-L变换

➤离散K-L变换：对向量 \mathbf{x} 用确定的完备正交归一向量系 \mathbf{u}_j 展开

$$\mathbf{x} = \sum_{j=1}^{\infty} y_j \mathbf{u}_j$$

$$\mathbf{u}_i^T \mathbf{u}_j = \delta_{ij}$$

$$\mathbf{x} \rightarrow \mathbf{y} \qquad y_j = \mathbf{u}_j^T \mathbf{x}$$

K-L变换

➤用有限项估计 \mathbf{x} :

$$\hat{\mathbf{x}} = \sum_{j=1}^d y_j \mathbf{u}_j \quad y_j = \mathbf{u}_j^T \mathbf{x}$$

➤该估计的均方误差:

$$\varepsilon = E \left[(\mathbf{x} - \hat{\mathbf{x}})^T (\mathbf{x} - \hat{\mathbf{x}}) \right]$$

$$\varepsilon = E \left[\sum_{j=d+1}^{\infty} y_j^2 \right] = E \left[\sum_{j=d+1}^{\infty} \mathbf{u}_j^T \mathbf{x} \mathbf{x}^T \mathbf{u}_j \right]$$

$$\mathbf{R} = \left[r_{ij} = E(x_i x_j) \right] = E \left[\mathbf{x} \mathbf{x}^T \right]$$

$$\varepsilon = \sum_{j=d+1}^{\infty} \mathbf{u}_j^T E \left[\mathbf{x} \mathbf{x}^T \right] \mathbf{u}_j = \sum_{j=d+1}^{\infty} \mathbf{u}_j^T \mathbf{R} \mathbf{u}_j$$

K-L变换：求解最小均方误差正交基

➤用Lagrange乘子法：

$$\text{if } \mathbf{R}\mathbf{u}_j = \lambda_j \mathbf{u}_j \text{ then } \varepsilon = \sum_{j=d+1}^{\infty} \mathbf{u}_j^T \mathbf{R} \mathbf{u}_j \text{ 取得极值}$$

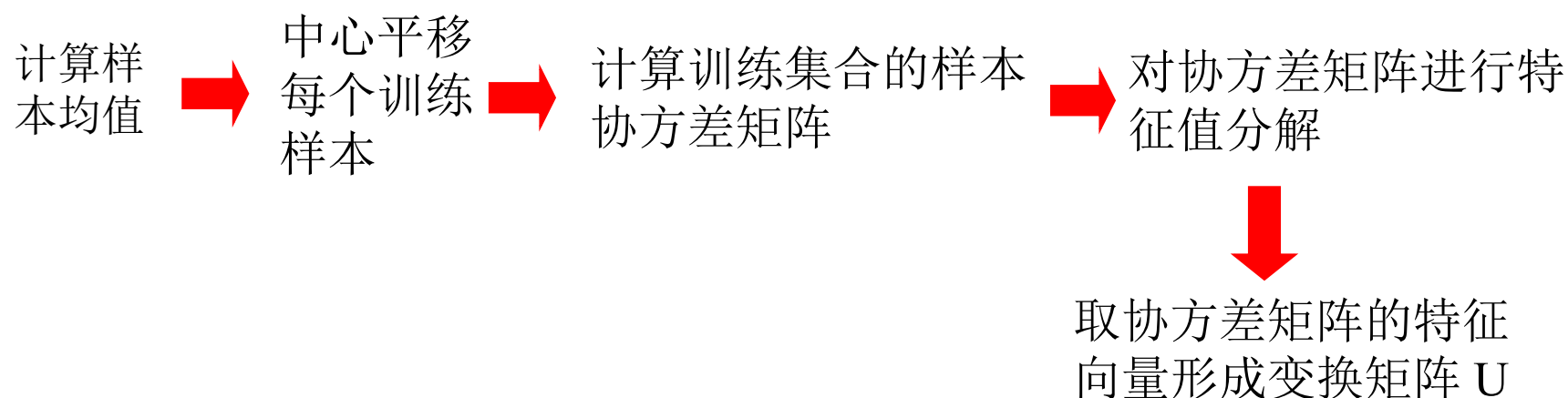
➤结论：以相关矩阵 \mathbf{R} 的 d 个特征向量为基向量来展开 \mathbf{x} 时，其均方误差为：

$$\varepsilon = \sum_{j=d+1}^{\infty} \lambda_j$$

➤K-L变换：当取矩阵 \mathbf{R} 的 d 个最大特征值对应的本征向量来展开 \mathbf{x} 时，其截断均方误差最小。这 d 个本征向量组成的正交坐标系称作 \mathbf{x} 所在的 D 维空间的 d 维K-L变换坐标系， \mathbf{x} 在K-L坐标系上的展开系数向量 \mathbf{y} 称作 \mathbf{x} 的K-L变换。

PCA：计算方法

➤ 计算过程为：



➤ 变换 $\mathbf{y} = \mathbf{U}^T (\mathbf{x} - \bar{\mathbf{x}})$

➤ 重构 $\mathbf{x} = \bar{\mathbf{x}} + \mathbf{U} \mathbf{y}$

K-L变换性质

➤ \mathbf{y} 的相关矩阵是对角矩阵:

$$\begin{aligned} E[y_i y_j] &= E[\mathbf{u}_i^T \mathbf{x} \mathbf{x}^T \mathbf{u}_j] = \mathbf{u}_i^T E[\mathbf{x} \mathbf{x}^T] \mathbf{u}_j \\ &= \mathbf{u}_i^T \mathbf{R} \mathbf{u}_j = \mathbf{u}_i^T \lambda_j \mathbf{u}_j = \lambda_i \delta_{ij} \end{aligned}$$

$$\begin{aligned} E[\mathbf{y} \mathbf{y}^T] &= E[U^T \mathbf{x} \mathbf{x}^T U] \\ &= U^T \mathbf{R} U = \Lambda \end{aligned}$$

\mathbf{x} 的协方差矩阵的特征值与 \mathbf{y} 的相等

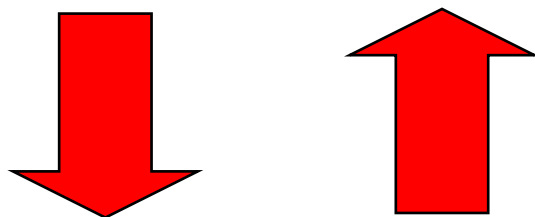
K-L变换性质

➤ K-L坐标系把矩阵 \mathbf{R} 对角化，即通过K-L变换消除原有向量 \mathbf{x} 的各分量间的相关性，从而有可能去掉那些带有较少信息的分量以达到降低特征维数的目的。

$$\mathbf{\Lambda} = \begin{bmatrix} \lambda_1 & & & 0 \\ & \lambda_2 & & \\ & & \ddots & \\ 0 & & & \lambda_d \end{bmatrix}$$

PCA：用于降维和重构

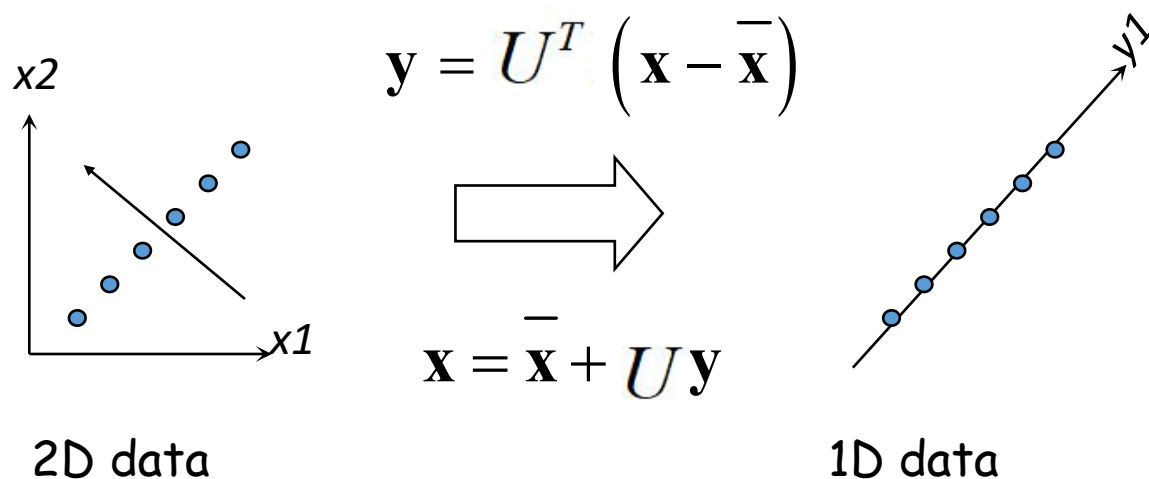
原始高维数据



压缩后低维数据

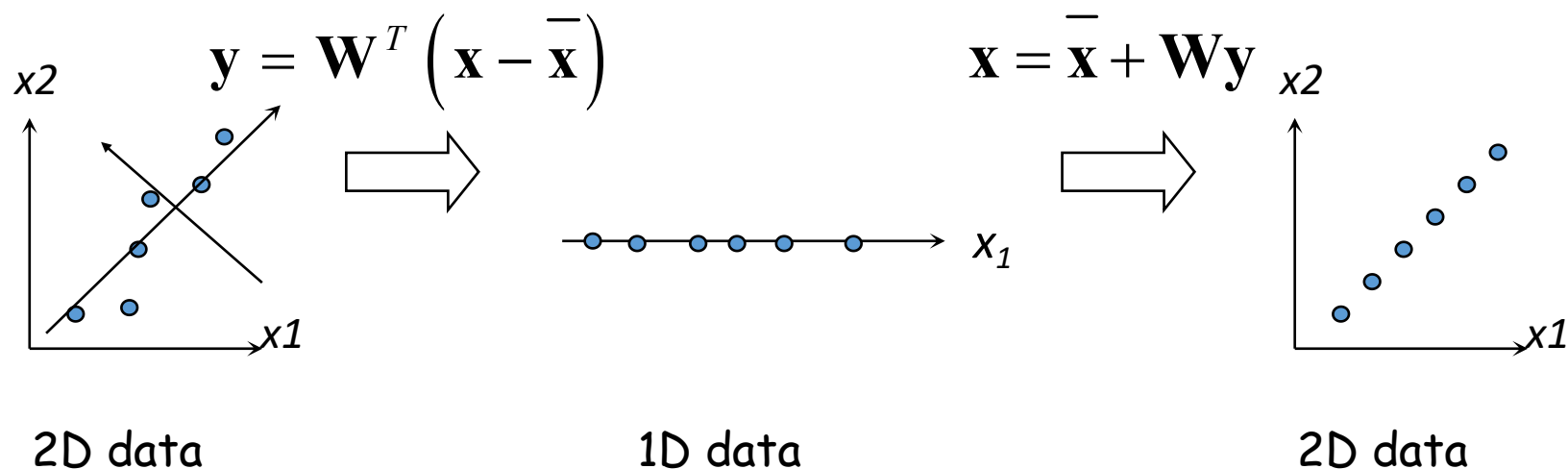
- 在变换后的特征空间中，每个特征向量对应的特征值的大小代表该特征向量所描述的方向上的总体方差的大小。
- 从 W 中去掉那些对应较小特征值的特征向量，意味着在信息丢失最小的意义上降维！

数据降维：理想情况



➤原始数据空间中，其中一维数据的方差为0，没有信息，可以完全去掉，而没有任何损失！

数据降维：非理想情况



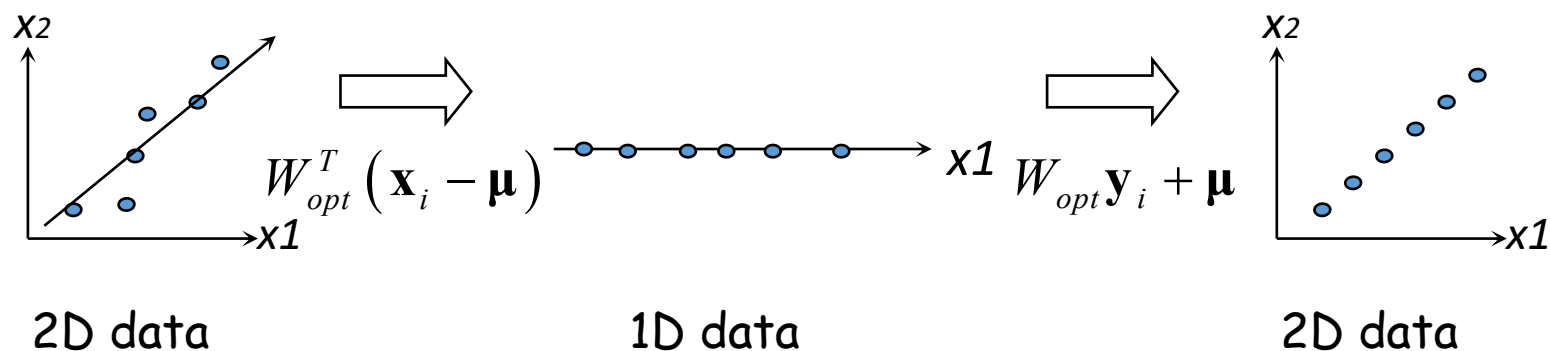
➤原始数据空间中，其中一维数据的方差比较小，包含少量信息，去掉后有少量损失！

PCA降维: 数据损失分析

➤ 投影后数据部分丢失，但是可以证明：在只使用前 d 个特征向量的情况下， \mathbf{x}_i 与其逆PCA重构 \mathbf{x}'_i 之间的均方误差为：

$$\sum_{j=1}^n \lambda_j - \sum_{j=1}^d \lambda_j = \sum_{j=d+1}^n \lambda_j$$

➤ 最小均方误差意义下的最佳变换



PCA小结

■ 一种多元统计分析方法

- 变换后各维数据之间的相关性最小
- 最小均方误差意义下的最佳变换

$$\mathbf{y} = \mathbf{U}^T (\mathbf{x} - \bar{\mathbf{x}})$$

$$\mathbf{x} = \bar{\mathbf{x}} + \mathbf{U} \mathbf{y}$$

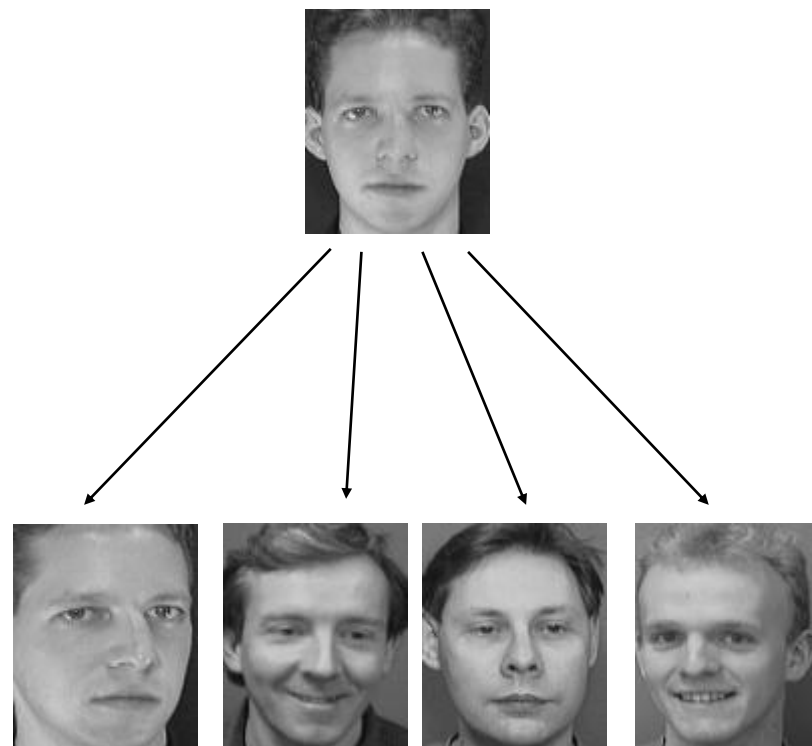
基于PCA的人脸识别

■建模

- 2D 灰度矩阵，按行向量化为1D向量，所有图像均表示为这样的向量
- 对其进行K-L变换，求得协方差矩阵的特征值和特征向量，特征向量构成的空间就是投影空间

■识别

- 计算输入图像的向量的投影系数与已知人脸库中所有向量投影系数的相似度，排序即可给出识别结果



Eigenface人脸识别方法

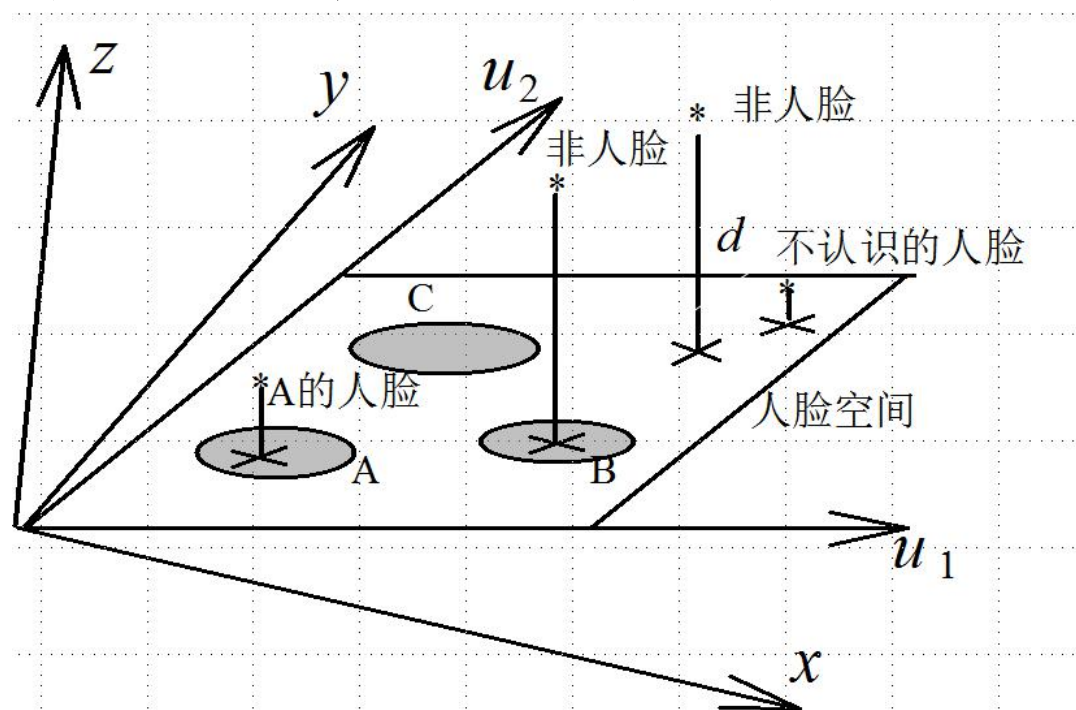
➤ \mathbf{x} 为输入图像

$$\mathbf{y} = U^T (\mathbf{x} - \bar{\mathbf{x}})$$

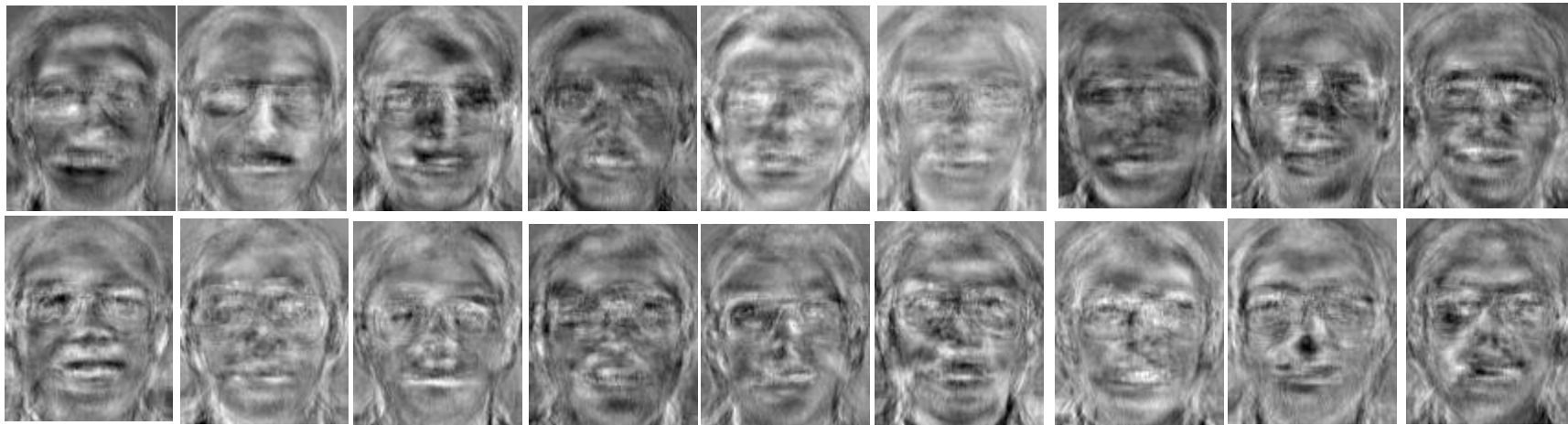
➤ \mathbf{y} 作为提取的特征

$$\mathbf{x} = \bar{\mathbf{x}} + U \mathbf{y}$$

➤ 可以采用欧式距离，也可以采用Cosine进行识别



特征脸



➤ Eigenface

➤ 可视化的“特征脸”