

# 降维

在特征向量所处的高维空间中，包含很多的冗余和噪音，我们希望通过降维的方式来寻找数据内部的特征，从而提升特征的表达能力，降低训练复杂度。

**常见的降维方法：**

- 主成分分析 (Principal Component Analysis)
- 线性判别分析 (Linear Discriminant Analysis)
- 等距映射 (Isometry)
- 局部先行嵌入 (Locally Linear Embedding)
- 拉普拉斯特征映射 (Laplacian Eigenmaps)

## PCA主成分分析

### 1. 什么是PCA？什么是主成分？

- PCA是一种线性、非监督、全局的降维算法
- PCA旨在找到数据的主成分，通过主成分来表征原始数据
- 比如在三维空间中的数据点分布在一个过原点的平面上，那么我们可以通过旋转坐标轴，使得xy轴与平面重合，就可以只通过两个维度来表示这些点，旋转后的坐标轴包含的信息就是主成分
- 主成分是特征向量变化最大的方向 (the direction of largest variation)，也可以理解为数据在这个方向上的投影分布得更加分散（方差更大），所以PCA其实是找最佳投影方向
- 主成分可以使数据的信噪比最大化，因为信号的方差越大，噪声的方差越小

## 2. 如何通过最大方差理论找到主成分？PCA如何求解？

- 先将数据进行去中心化：  $x_1, x_2, \dots, x_n = v_1 - \mu, v_2 - \mu, \dots, v_n - \mu$
- 向量内积是第一个向量投影到第二个向量上的长度，所以向量  $x_i$  在  $\omega$  上的投影坐标为  $x_i^T \omega$
- PCA的目标是找到一个投影方向  $\omega$  使得投影后的点  $x_i^T \omega$  的方差最大
- 因为去中心化后投影点的均值为0，所以投影点的方差为：

$$\begin{aligned} D(x) &= \frac{1}{n} \sum_{i=1}^n (x_i^T \omega)^2 = \frac{1}{n} \sum_{i=1}^n (x_i^T \omega)^T (x_i^T \omega) \\ &= \frac{1}{n} \sum_{i=1}^n \omega^T x_i x_i^T \omega \\ &= \omega^T \left( \frac{1}{n} \sum_{i=1}^n x_i x_i^T \right) \omega \end{aligned}$$

- 括号中的部分其实是数据的协方差矩阵（covariance matrix）
- 因此PCA的优化问题变成了找到方向  $\omega$  使得：

$$\begin{cases} \max \{ \omega^T \Sigma \omega \}, \\ s.t. \quad \omega^T \omega = 1. \end{cases}$$

- 限制：  $\omega$  是归一化后的单位方向向量，不然max函数只会让  $\omega$  越来越大
- 对  $\omega$  求导后可以知道，  $x$  投影后的方差就是协方差矩阵的特征值（eigine value），所以求投影点的最大方差，就是求  $x$  协方差矩阵的的最大特征值
- 次佳的投影方向则位于最佳投影方向的正交空间（orthogonal）
- 我们将  $x$  协方差矩阵的特征值从大到小排序，取前d大的特征值对应的特征向量，即为我们的主成分  $w_1, w_2, \dots, w_d$
- 最后通过向量内积，将n维的样本映射到d维中：

$$x'_i = \begin{bmatrix} \omega_1^T x_i \\ \omega_2^T x_i \\ \vdots \\ \omega_d^T x_i \end{bmatrix}.$$

- 降维后的信息占比为：

$$\eta = \sqrt{\frac{\sum_{i=1}^d \lambda_i^2}{\sum_{i=1}^n \lambda_i^2}}.$$

• PCA是线性降维法，可以通过核方法来应对非线性可投影的数据

### 3. 如何从回归的角度（最小平方误差和理论）定义PCA的目标并求解？

- PCA可以理解为在高维空间中，找到一个d维超平面，使得样本点到这个超平面的距离平方和最小，点 $x_k$ 到超平面 $D$ 的距离为：

$$\text{distance}(\mathbf{x}_k, D) = \|\mathbf{x}_k - \widetilde{\mathbf{x}}_k\|_2,$$

- 最佳投影超平面由d个正交基构成 $W = \{w_1, w_2, \dots, w_d\}$ ，点 $x_k$ 在超平面 $D$ 上的投影是：

$$\widetilde{\mathbf{x}}_k = \sum_{i=1}^d (\omega_i^T \mathbf{x}_k) \omega_i,$$

- 从这个角度出发，PCA的优化目标是：

$$\begin{cases} \arg \min_{\omega_1, \dots, \omega_d} \sum_{k=1}^n \|\mathbf{x}_k - \widetilde{\mathbf{x}}_k\|_2^2, \\ s.t. \quad \omega_i^T \omega_j = \delta_{ij} = \begin{cases} 1, i = j; \\ 0, i \neq j. \end{cases} \end{cases}$$

- 可以简化为求矩阵的迹（对角线元素之和）：

$$\begin{cases} \arg \max_W \text{tr}(W^T X X^T W), \\ s.t. \quad W^T W = I. \end{cases}$$

## LDA线性判别分析

PCA没有考虑数据的标签，只是把原数据映射到方差较大的方向上，映射后不同类别的数据会完全混合在一起，很难被区分开。

1. 对于有类别标签的数据，如何设计目标函数使得降维的过程中不损失类别信息？

- 找到一个投影方向 $w$ ，使得投影后的样本尽可能按照原始分类分开
- 例如二分类任务，我们希望最大化投影点的均值差：

$$D(C_1, C_2) = \|\widetilde{\mu}_1 - \widetilde{\mu}_2\|_2^2$$

$$\mu_1 = \frac{1}{N_1} \sum_{x \in C_1} x, \quad \mu_2 = \frac{1}{N_2} \sum_{x \in C_2} x$$

$$\widetilde{\mu}_1 = \omega^T \mu_1, \quad \widetilde{\mu}_2 = \omega^T \mu_2,$$

- 因此，LDA的优化目标是：

$$\begin{cases} \max_{\omega} \|\omega^T (\mu_1 - \mu_2)\|_2^2, \\ s.t. \quad \omega^T \omega = 1. \end{cases}$$

- 但只是最大化两类点投影后的中心距离，投影后的点会有重叠的部分，所以我们还需要最小化两类点投影后的各自的内部方差，即：最大化类间距离和最小化类内距离

- 所以LDA的目标函数是类间距离和类内距离的比值：

$$\max_{\omega} J(\omega) = \frac{\|\omega^T (\mu_1 - \mu_2)\|_2^2}{D_1 + D_2},$$

- $D_1$ 和 $D_2$ 是两类点投影后的内部方差：

$$D_1 = \sum_{x \in C_1} (\omega^T x - \omega^T \mu_1)^2 =$$

$$\sum_{x \in C_1} \omega^T (x - \mu_1)(x - \mu_1)^T \omega,$$

- 定义类内散度矩阵 $S_B$ 和类内散度矩阵 $S_w$ ，则目标函数可以被简化为：

$$J(\omega) = \frac{\omega^T (\mu_1 - \mu_2)(\mu_1 - \mu_2)^T \omega}{\sum_{x \in C_i} \omega^T (x - \mu_i)(x - \mu_i)^T \omega}.$$

$$J(\omega) = \frac{\omega^T S_B \omega}{\omega^T S_w \omega}$$

- 求解LDA，只需要针对 $w$ 求偏导，并令导数等于0，得：

$$S_w^{-1} S_B \omega = \lambda \omega$$

- $\lambda$ 是目标函数 $J(w)$ s

- 所以求解LDA只需要求出样本的均值差和类内方差：

$$\omega = S_w^{-1} (\mu_1 - \mu_2)$$

- LDA比PCA更适合降维有类别信息的数据，但是LDA对数据做了很强的假设：
  - 每个类的数据都是高斯分布
  - 各个类的协方差相等
- 线性模型对噪声的鲁棒性很好，但是表达能力有限，可以引入核方法进行扩展

## 2. 多类别LDA的求解过程：

- 多类别情况下的LDA，类间散度是全局散度和类内散度的差，即每个类别的中心经过投影后和全局中心的距离
- 多类别LDA的求解步骤：
  - 计算每个类别的均值 $\mu_j$ 和全局均值 $\mu$
  - 计算类内散度矩阵 $S_w$ ，全局散度矩阵 $S_t$ ：

$$S_w = \sum_{x \in C_i} (x - \mu_i)(x - \mu_i)^T$$

$$S_t = \sum_{i=1}^n (x_i - \mu)(x_i - \mu)^T$$

- 计算类间散度矩阵 $S_b$ ：

$$\begin{aligned} S_b &= S_t - S_w \\ &= \sum_{i=1}^n (x_i - \mu)(x_i - \mu)^T - \sum_{x \in C_i} (x - \mu_i)(x - \mu_i)^T \\ &= \sum_{j=1}^N \left( \sum_{x \in C_j} (x - \mu)(x - \mu)^T - \sum_{x \in C_j} (x - \mu_j)(x - \mu_j)^T \right) \\ &= \sum_{j=1}^N m_j (\mu_j - \mu)(\mu_j - \mu)^T, \end{aligned}$$

- 定义目标函数：

$$J(W) = \frac{\text{tr}(W^T S_b W)}{\text{tr}(W^T S_w W)},$$

- 对矩阵 $S_w^{-1} S_b$ 进行特征值分解，将特征值从大到小排序
- 取特征值前d大的特征值对应的特征向量 $w_1, w_2, \dots, w_d$ ，将n维样本映射到d维空间：

$$x'_i = \begin{bmatrix} \omega_1^T x_i \\ \omega_2^T x_i \\ \vdots \\ \omega_d^T x_i \end{bmatrix}.$$

### 3. LDA和PCA从原理和应用上有什么异同?

- 无监督任务使用PCA（无监督算法），有监督任务使用LDA（监督算法）
- 从求解过程看，相似度很高，但是原理有所区别：
  - PCA旨在最大化投影点的方差（假设方差越大，信息量越多），用主成分来表示原始数据可以去除冗余的维度
  - LDA旨在最小化类内方差和最大化类间方差，利用数据标签找到数据中最具判别性的维度，使得投影后的数据更容易被区分
- 应用场景区别，例如语音识别：
  - 可以使用PCA过滤掉固定频率（方差较小）的背景噪音
  - 但是目标是识别人声的话，就需要通过LDA来考虑标签，使降维后的数据有区分性
- PCA和LDA在人脸识别的应用：
  - 基于PCA的方法叫特征脸（Eigenface）方法，对人脸特征的协方差矩阵做特征分解，较大的特征值对应的特征向量具有与人脸相似的形状