

第三章 降维

3.1 为什么要降维

3.2 主成分分析(PCA)

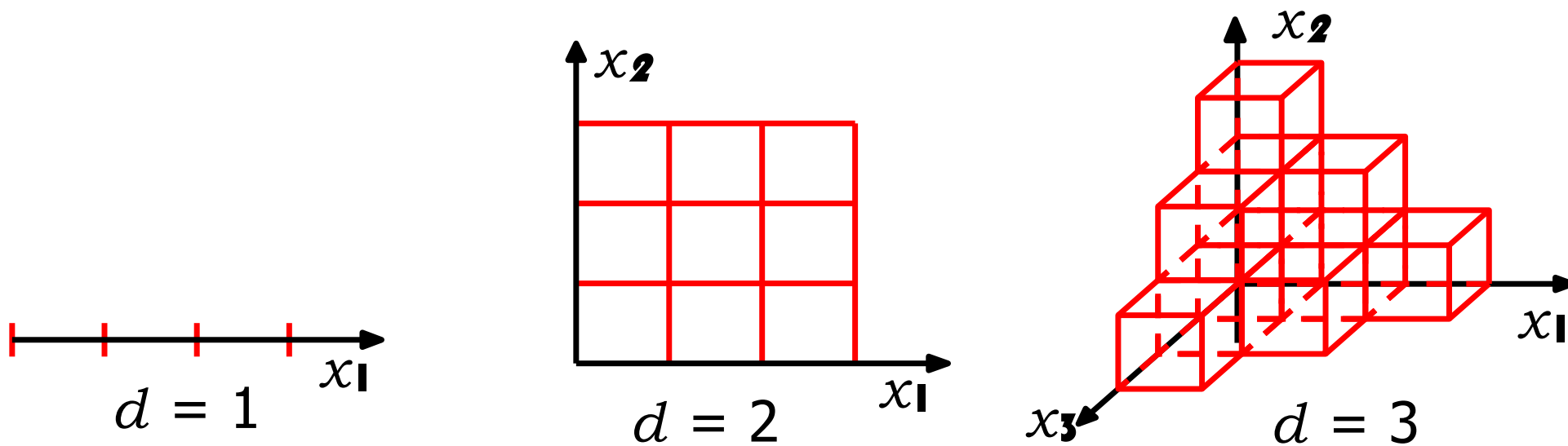
3.2.1 目标函数

3.2.2 算法

3.1 为什么要降维(dimension reduction)?

原始特征太多(高维空间)导致: 计算量和占用存储大、容易过拟合

“维度灾难”(curse of dimensionality): 样本空间大小与维度成指数增长。



数据稀疏问题导致过拟合 (overfitting)。

为什么要降维 (dimension reduction)?

目的：

- 1.数据压缩，减小存储大小，减小计算量；
- 2.更有利于后续处理（分类、聚类），防止过拟合；
- 3.方便可视化。

减少特征数目的方式：

- 特征选择：从给定特征集中选择部分最有用的特征
- 特征提取：基于当前数据集产生新的特征

降维又叫低维“嵌入”（embedding），指从原始的高维空间降到低维空间表示，使得映射后的低维表示（投影）尽量保留原数据的基本特性。

选择什么样的映射？

映射后得到的低维表示与原始高维表示之间的差别最小或保留信息最多：

- 尽量保留重要的/主要的维度，丢弃不重要的/次要的维度；
- 尽量保留样本之间的相似性关系：如在原始空间，样本1与样本2的距离很近，与样本3的距离很远。那么在低维空间里也应该保留以上关系。

主要方法：主成分分析(PCA)、多维尺度变化(MDS)等。

第三章 降维

3.1 为什么要降维

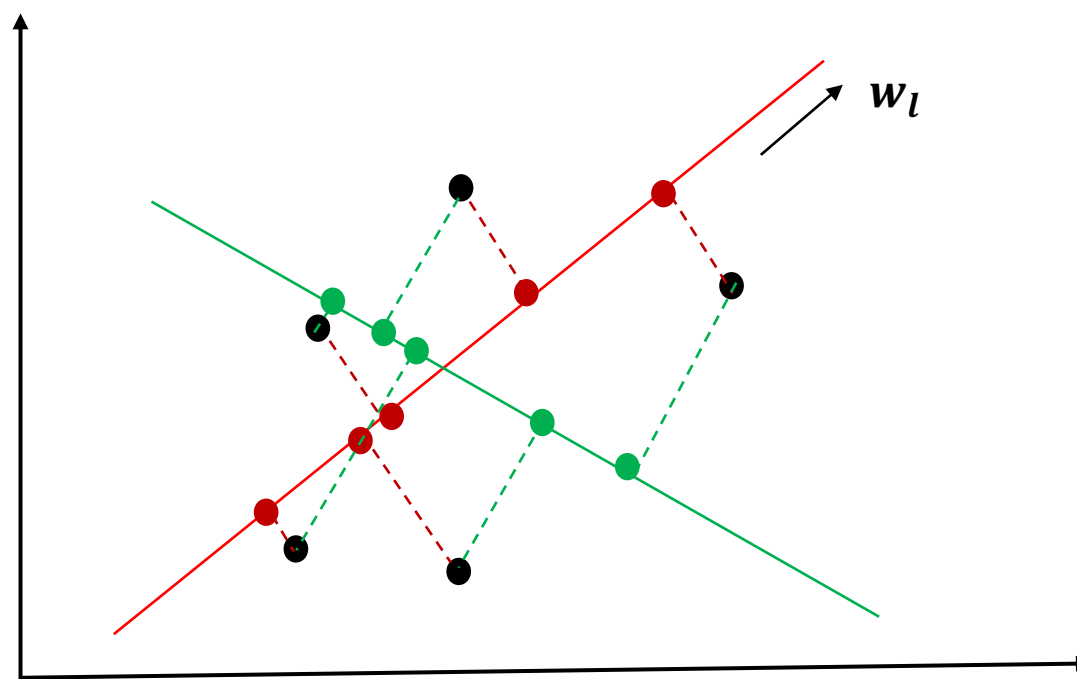
3.2 主成分分析(PCA)

3.2.1 目标函数

3.2.2 算法

3.2 主成分分析 (Principle Component Analysis)

目标：把数据从原始的高维空间映射到一个低维空间，使得映射后的样本之间具有最大可分性。



如何衡量可分性？

↓
方差

低维投影下方差的计算

假设数据集表示为对象-属性矩阵 $\mathbf{X}_{n \times d}$ ，通过矩阵 $\mathbf{W}_{d \times k}$ 把样本 \mathbf{x}_i 线性映射到 k 维向量 \mathbf{z}_i 。

$$\begin{matrix} & \mathbf{X}_{n \times d} & & \mathbf{W}_{d \times k} & & \mathbf{Z}_{n \times k} \\ \mathbf{x}_i & \begin{bmatrix} \vdots \\ \vdots \end{bmatrix} & \times & \begin{bmatrix} \vdots \\ \vdots \end{bmatrix} & = & \begin{bmatrix} \vdots \\ \vdots \end{bmatrix} \\ (1 \times d) & & & \mathbf{w}_l (d \times 1) & & \end{matrix}$$

$z_{il} = \mathbf{x}_i \mathbf{w}_l$

\mathbf{z}_i 就是 \mathbf{x}_i 在 k 维空间中的投影。若 $k < d$ ，则实现降维。

映射后在第 l 个维度的方差：

$$\frac{1}{n} \sum_{i=1}^n (z_{il} - \bar{z}_l)^2 = \frac{1}{n} \sum_{i=1}^n (\mathbf{x}_i \mathbf{w}_l - \bar{\mathbf{x}} \mathbf{w}_l)^2 = \mathbf{w}_l^T \mathbf{C} \mathbf{w}_l$$

其中， $\mathbf{C} = \frac{1}{n} \sum_{i=1}^n (\mathbf{x}_i - \bar{\mathbf{x}})^T (\mathbf{x}_i - \bar{\mathbf{x}})$ 为 \mathbf{X} 的协方差矩阵。

方差、协方差的无偏估计用系数 $\frac{1}{n-1}$ ，但是这个系数不影响 \mathbf{w}_l 的解。



推导: $\bar{z}_l = \bar{\mathbf{x}} \mathbf{w}_l$

$\bar{\mathbf{x}}$ 是 \mathbf{X} 按列求均值得到的横向量;

$\bar{z}_l = \frac{1}{n} \sum_{i=1}^n z_{il}$ 为 \mathbf{Z} 第 l 列的均值。

3.2.1 方差最大化

主成分分析可以建模为最大化方差，同时满足映射向量为单位向量的条件，即：

$$\begin{aligned} \max \mathbf{w}_l^T \mathbf{C} \mathbf{w}_l \\ \text{s.t. } \mathbf{w}_l^T \mathbf{w}_l = 1 \end{aligned}$$

以上问题是一个带约束条件的求极值问题，用拉格朗日乘子法

$$\begin{aligned} L(\mathbf{w}_l, \lambda) &= \mathbf{w}_l^T \mathbf{C} \mathbf{w}_l + \lambda(1 - \mathbf{w}_l^T \mathbf{w}_l) \\ \text{令 } \frac{\partial L(\mathbf{w}_l, \lambda)}{\partial \mathbf{w}_l} &= 2\mathbf{C} \mathbf{w}_l - 2\lambda \mathbf{w}_l = 0 \end{aligned}$$

由以上公式得到

$$\mathbf{C} \mathbf{w}_l = \lambda \mathbf{w}_l$$

即 \mathbf{w}_l 是矩阵 \mathbf{C} 的特征向量，对应特征值为 λ 。


PCA算法

输入：包含 n 个 d 维样本的数据集 $\mathbf{X}_{n \times d}$ 、输出维度大小 k 。

输出： $\mathbf{W} = [w_1, w_2, \dots, w_k]$, 其中 w_l 为第 l 个特征向量，对应第 l 大的特征值。

步骤：

1. 计算协方差矩阵 $\mathbf{C}_{d \times d}$ ；
2. 计算 \mathbf{C} 的 k 个最大特征值对应的特征向量；
3. 由得到的特征向量输出映射矩阵 \mathbf{W} ；



什么时候计算量大？
有没有其他方法得到 \mathbf{C} 的特征向量？

中心化后进行PCA

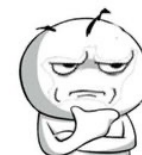
假设数据 $\mathbf{X}_{n \times d}$ 已经经过中心化处理(每列均值 $\bar{\mathbf{x}}=0$)，那么 $\bar{\mathbf{x}} \mathbf{w}_l = 0$

映射后在第 l 个维度的方差退化为：

$$\frac{1}{n} \sum_{i=1}^n (z_{il} - \bar{z}_l)^2 = \frac{1}{n} \sum_{i=1}^n (\mathbf{x}_i \mathbf{w}_l - \bar{\mathbf{x}} \mathbf{w}_l)^2$$



$$\frac{1}{n} \sum_{i=1}^n z_{il}^2 = \frac{1}{n} \sum_{i=1}^n (\mathbf{x}_i \mathbf{w}_l)^2 = \mathbf{w}_l^T \mathbf{C} \mathbf{w}_l$$



为什么进行PCA前一般要中心化？

其中， $\mathbf{C} = \frac{1}{n} \mathbf{X}^T \mathbf{X}$ 为中心化后的协方差矩阵。

用SVD求解中心化数据的协方差矩阵**C**的特征向量

不用计算**C**而直接对**X**进行奇异值分解(Singular Value Decomposition)。

SVD把矩阵**X**_{*n*×*d*}分解成两个正交矩阵与一个对角矩阵**Σ**：**X** = **UΣV**^{*T*}，
其中

$$\Sigma_{r \times r} = \begin{bmatrix} \sigma_1 & \cdots & 0 \\ \vdots & \ddots & \vdots \\ 0 & \cdots & \sigma_r \end{bmatrix}, \quad \sigma_i \text{ 为 } \mathbf{X} \text{ 的奇异值,}$$

左奇异向量**U**_{*n*×*r*} = [**u**₁, ... **u**_{*r*}] 为 (**XX**^{*T*})_{*n*×*n*}的特征向量，

右奇异向量**V**_{*d*×*r*} = [**v**₁, ... **v**_{*r*}] 为 (**X**^{*T*}**X**)_{*d*×*d*}的特征向量，即**C**的特征向量。

假设 $\sigma_1 > \sigma_2 \cdots > \sigma_r$ ，则 $\sigma_k = \sqrt{\lambda_k}$ ， σ_k 对应的向量**v**_{*k*}即为对应**C**第*k*大特征值 λ_k 的特征向量**w**_{*k*}，即**v**_{*k*} = **w**_{*k*}。

降维后的表示与重构

对 $\mathbf{X}_{n \times d}$ 进行PCA降维后得到正交变换矩阵 $\mathbf{W}_{d \times k}$ ，则样本在这个低维空间的表示为 $\mathbf{Z}_{n \times k} = \mathbf{XW}$ 。基于这个表示对原始数据进行重构，即映射回高维空间，得到 $\mathbf{X}'_{n \times d} = \mathbf{ZW}^T = \mathbf{XWW}^T$ 。

若对 \mathbf{X} 进行SVD分解后取前面最大的 k 个奇异值对应的 \mathbf{U} , $\mathbf{\Sigma}$, \mathbf{V} ，则

$\mathbf{X}'_{n \times d} = \mathbf{U}_{n \times k} \mathbf{\Sigma}_{k \times k} \mathbf{V}_{k \times d}^T$ ，对应 $\mathbf{W} = \mathbf{V}$ ， $\mathbf{Z}_{n \times k} = \mathbf{U}\mathbf{\Sigma}$ 或 \mathbf{XW} 。

重构误差为

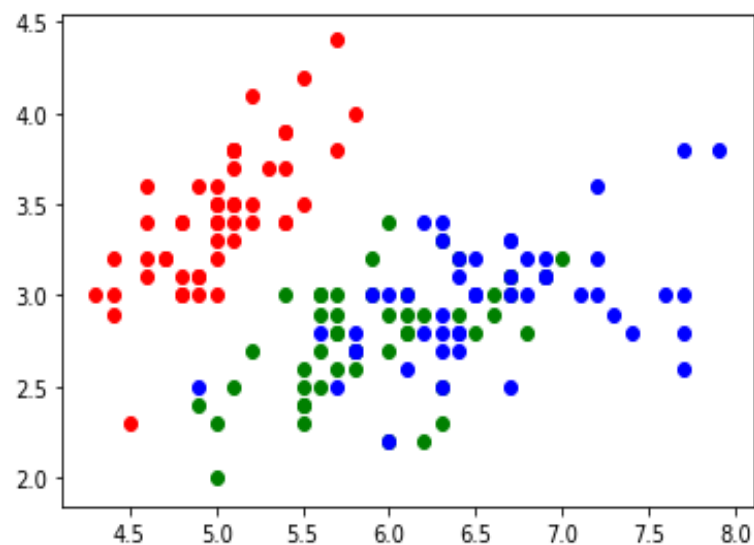
$$e = \|\mathbf{X} - \mathbf{X}'\|_2^2 = \sum_{i=1}^n \sum_{j=1}^n (x_{ij} - x'_{ij})^2$$

PCA方差最大化（零均值）
的目标函数等价于重构误差最小化：

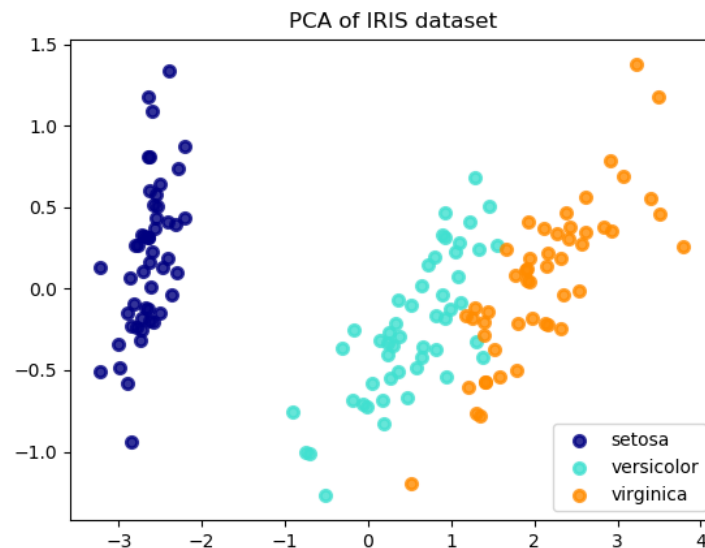
$$\begin{aligned} \min & \|\mathbf{X} - \mathbf{XWW}^T\|_2^2 \\ \text{s.t. } & \mathbf{W}^T \mathbf{W} = \mathbf{I} \end{aligned}$$

PCA降维后可视化

- 由PCA算法得到映射矩阵 \mathbf{W} 后，计算映射后的数据 $\mathbf{Z}_{n \times k} = \mathbf{X}\mathbf{W}$ 。



原始特征1和2



两个主成分

Python实现特征向量分解（以随机矩阵为例）

调用`scipy.linalg.eig`

```
import numpy as np
from scipy import linalg          #线性代数库
n, d=9, 6
x=np.random.rand(n, d)
covM=np.cov(x.T)                #np.cov() 输入矩阵中行为特征，所以转置
# eigvals 为特征值， eigvecs为对应特征向量
eigvals,eigvecs= linalg.eig(covM)
sort_index=np.argsort(-eigvals) #降序排列，该方法本身为升序，所以乘-1
k=2;
w=eigvecs[:,sort_index[0:k]]
#基于分解结果重构矩阵，计算重构误差
z=np.dot(x,w)
x1 = np.dot(z,w.T )
err=((x-x1)*(x-x1)).sum()
```

注意：此处`x`为numpy的array对象，`.sum()`表示对矩阵所有元素相加。但如果`x`为pandas的DataFrame对象，`x.sum()`只得到每列的和，`x.sum().sum()`才是所有元素的和。

Python实现SVD（以随机矩阵为例）

调用scipy.linalg.svd

```
import numpy as np
from scipy import linalg
n, d=9, 6
x=np.random.rand(n, d)
#s为r个奇异值（已从大到小排列）组成的向量（r=min(n, d)）
U, s, Vh = linalg.svd(x, full_matrices=False)
#转成对角矩阵
sigma = np.diag(s)
#基于分解结果重构矩阵，计算重构误差
k=4
U_k=U[:, :k]
sigma_k=sigma[:k, :k]
Vh_k=Vh[:k]
x2 = U_k.dot(sigma_k).dot(Vh_k)
err_k = ((x-x2)*(x-x2)).sum()
```

注意：SVD分解对一般矩阵都可以应用，但是只有对均值为零的矩阵的SVD分解得到的右奇异向量才等价于协方差矩阵的特征向量，即这个时候的SVD分解才等价于PCA。

推导: $\bar{X} = 0$ 时, 即数据已作零均值处理, $C = X^T X$

方差最大化 等价 最小化重构损失

$$\max \text{Tr}(W^T C W) \Leftrightarrow \min \|X - X W W^T\|^2$$

$$\text{s.t. } W^T W = I$$

$$\text{s.t. } W^T W = I$$

需要到以下公式: ① $\|A\|^2 = \text{Tr}(A^T A)$, A 为矩阵

$$\textcircled{2} \text{Tr}(AB) = \text{Tr}(BA) \quad \textcircled{3} (AB)^T = B^T A^T$$

根据公式①, 右边重构损失为: $\text{Tr}(X - X W W^T)^T (X - X W W^T)$

$$= \text{Tr}(X^T - W W^T X^T)(X - X W W^T)$$

$$\min - \text{Tr}(W^T C W)$$

$$= \underbrace{\text{Tr}(X^T X)}_{\text{常数}} - \underbrace{\text{Tr}(X^T X W W^T)}_a - \underbrace{\text{Tr}(W W^T X^T X)}_b + \underbrace{\text{Tr}(W W^T X^T X W W^T)}_c$$

$$\Leftrightarrow \max \text{Tr}(W^T C W)$$

证明结束

$$\underline{a=b=c} - \text{Tr}(W^T X^T X W) = -\text{Tr}(W^T C W)$$

$$\begin{aligned} & \text{Tr}(W W^T X^T X W W^T) \\ & \Downarrow \\ & \text{Tr}(W^T X^T X W W^T W) \\ & \Downarrow \\ & \text{Tr}(W^T X^T X W) \end{aligned}$$