

降维

Dimension Reduction

① 背景

② 样本均值 & 样本方差

③ PCA

④ SVD 角度看 PCA 和 PCoA

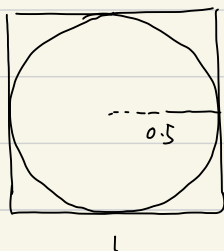
⑤ 概率角度看 PCA

背景 (Background)

① 维度灾难 (Curse of Dimensionality)

从几何角度解释.

1)



作假设维度为2, 有一个圆内嵌于正方形.

正方形边长为1,

$$\text{则 } V_{\text{正}} = 1^2 \quad \xrightarrow{\text{维度为3}} \quad V_{\text{立方体}} = 1^3$$

$$V_{\text{圆}} = \pi \cdot 0.5^2$$

$$V_{\text{球体}} = \frac{4}{3} \pi \cdot 0.5^3$$

$$\text{维度为 } D \quad \Rightarrow \quad V_{\text{超立方体}} = 1^D$$

$$V_{\text{超球体}} = K \cdot 0.5^D$$

由上面计算可知: 当维度不断增大, 则 $V_{\text{超球体}}$ 不断逼近于0.

又由于超球体内切于超立方体, 则可见超立方体的体积几乎分布于壳上, 即样本几乎分布于超立方体表面, 内部几乎没有样本.

2)



用同心圆举例. 如上述例子所示, 令 $r=1$

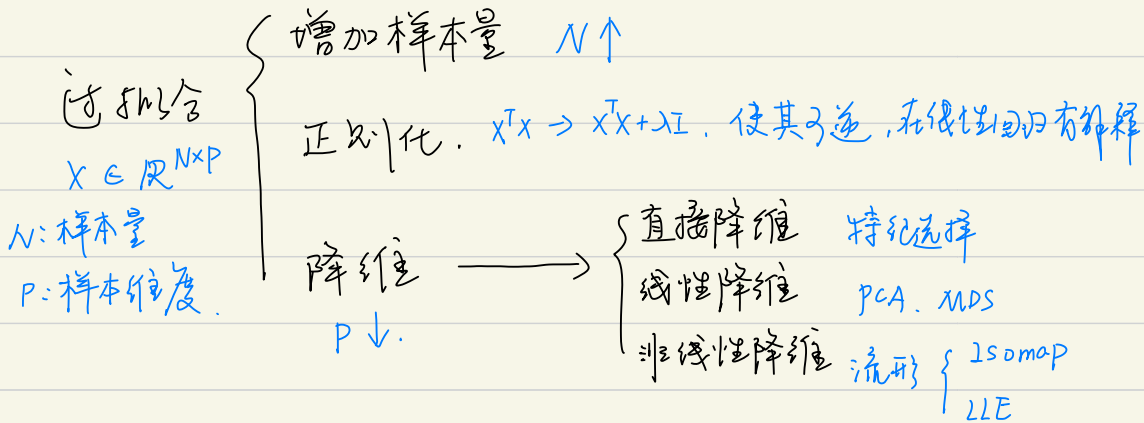
$$V_{\text{外}} = K, \quad V_{\text{壳}} = V_{\text{外}} - V_{\text{内}} = K - K(1-\varepsilon)^D$$

$$\text{则 } \frac{V_{\text{壳}}}{V_{\text{外}}} = 1 - (1-\varepsilon)^D, \quad \text{当 } D \rightarrow \infty \text{ 时, } (1-\varepsilon)^D \rightarrow 0.$$

$$\text{则 } \frac{V_{\text{壳}}}{V_{\text{外}}} = 1, \quad \text{即体积集中在球壳上.}$$

上述两个例子都证明了一件事. 数据在高维空间中分布的稀疏性. 这就是维度灾难.

(2)



样本均值 & 样本方差

(Sample Mean & Variance Matrix)

$$X = (x_1, x_2, \dots, x_N)^T, \quad x_i \in \mathbb{R}^p, \quad X \in \mathbb{R}^{N \times p}$$

$$\text{Sample Mean: } \bar{X} = \frac{1}{N} \sum_{i=1}^N x_i = \frac{1}{N} X^T \mathbf{1}_N$$

$$\text{Sample Covariance: } S = \frac{1}{N} \sum_{i=1}^N (x_i - \bar{X})(x_i - \bar{X})^T = \frac{1}{N} X^T H X$$

转换为矩阵表示. 令 $\mathbf{1}_N = \begin{pmatrix} 1 \\ \vdots \\ 1 \end{pmatrix} \in \mathbb{R}^{N \times 1}$, $I_N \in \mathbb{R}^{N \times N}$

$$\bar{X} = \frac{1}{N} \sum_{i=1}^N x_i$$

$$= \frac{1}{N} (x_1 \ x_2 \ \dots \ x_N) \begin{pmatrix} 1 \\ \vdots \\ 1 \end{pmatrix}$$

$$= \frac{1}{N} X^T \cdot \mathbf{1}_N$$

$$S = \frac{1}{N} \sum_{i=1}^N (x_i - \bar{x})(x_i - \bar{x})^T$$

$$= \frac{1}{N} \begin{pmatrix} x_1 - \bar{x} & x_2 - \bar{x} & \dots & x_N - \bar{x} \end{pmatrix} \begin{pmatrix} (x_1 - \bar{x})^T \\ (x_2 - \bar{x})^T \\ \vdots \\ (x_N - \bar{x})^T \end{pmatrix}$$

$$\text{对 } \begin{pmatrix} x_1 - \bar{x} & x_2 - \bar{x} & \dots & x_N - \bar{x} \end{pmatrix} = \begin{pmatrix} x_1 & x_2 & \dots & x_N \end{pmatrix} - \bar{x} \begin{pmatrix} 1 & 1 & \dots & 1 \end{pmatrix}$$

$$= X^T - \frac{1}{N} X^T \mathbf{1}_N \mathbf{1}_N^T$$

$$= X^T \left(I_N - \frac{1}{N} \mathbf{1}_N \mathbf{1}_N^T \right)$$

$H_N \rightarrow$ centering matrix

$$= X^T H$$

$$\text{则 } S = \frac{1}{N} X^T H \cdot (X^T H)^T = \frac{1}{N} X^T H H^T X$$

$$\text{对 } H = I_N - \frac{1}{N} \mathbf{1}_N \mathbf{1}_N^T$$

$$H^T = (I_N - \frac{1}{N} \mathbf{1}_N \mathbf{1}_N^T)^T = I_N - \frac{1}{N} \mathbf{1}_N \mathbf{1}_N^T = H$$

$$H^2 = H \cdot H = (I_N - \frac{1}{N} \mathbf{1}_N \mathbf{1}_N^T) (I_N - \frac{1}{N} \mathbf{1}_N \mathbf{1}_N^T)$$

$$= I_N - \frac{2}{N} \mathbf{1}_N \mathbf{1}_N^T + \underbrace{\left[\frac{1}{N} \mathbf{1}_N \mathbf{1}_N^T \mathbf{1}_N \mathbf{1}_N^T \right]}_{= \frac{1}{N} \mathbf{1}_N \mathbf{1}_N^T} \Rightarrow \frac{1}{N} \mathbf{1}_N \mathbf{1}_N^T$$

$$= I_N - \frac{1}{N} \mathbf{1}_N \mathbf{1}_N^T$$

$$= H$$

$$\text{则 } H^N = H$$

$$\Rightarrow S = \frac{1}{N} X^T H \cdot H^T X = \frac{1}{N} X^T H X$$

PCA

一个中心: 原始特征空间的重构.

相关 \rightarrow 无关

一所学校中每个学生都有学号和姓名的两个属性
在学校维度, 大家都一样, 方差为0, 不含信息.

↑

两个基本点: ① 最大投影方差

投影后的数据方差尽可能大, 代表所含信息量越多.

② 最小重构误差

降维后的数据重构回原始数据的代价最小.

① 最大投影方差 (maximum variance perspective).

假设投影的轴为 u_1 , 且 $u_1^T u_1 = 1$. (不考虑大小, 只考虑方向)

首先, 对数据中心化: $x_i - \bar{x}$

中心化后的数据在 u_1 上的投影为 $u_1^T (x_i - \bar{x})$

投影后的数据依然是中心化的, $\rightarrow u_1^T (x_i - \bar{x}) = u_1^T x_i - u_1^T \bar{x}$

可以看作先投影, 再中心化

$$\begin{aligned} \text{则} | \text{方差为} & \frac{1}{N} \sum_{i=1}^N [u_1^T (x_i - \bar{x})]^2 = \frac{1}{N} \sum_{i=1}^N u_1^T (x_i - \bar{x}) (x_i - \bar{x})^T u_1 \\ & = u_1^T \left[\frac{1}{N} \sum_{i=1}^N (x_i - \bar{x}) (x_i - \bar{x})^T \right] u_1 \\ & = u_1^T S u_1 \end{aligned}$$

$$\text{则} | \hat{u}_1 = \arg \max_{u_1} u_1^T S u_1$$

$$\text{s.t. } u_1^T u_1 = 1$$

用拉格朗日求子法, 得:

$$\mathcal{L}(u_1, \lambda) = u_1^T S u_1 + \lambda(1 - u_1^T u_1)$$

$$\frac{\partial \mathcal{L}(u_1, \lambda)}{\partial u_1} = 2S u_1 - 2\lambda u_1 = 0.$$

$$\text{则} \begin{cases} S u_1 = \lambda u_1 \end{cases}$$

即 u_1 为 S 的关于特征值 λ 的特征向量

$$\text{又} \because \hat{u}_1 = \underset{u_1}{\operatorname{argmax}} u_1^T S u_1$$

对 $S u_1 = \lambda u_1$ 两边左乘 u_1^T , 得 $u_1^T S u_1 = \lambda$.

$\therefore \lambda_1$ 为 S 的最大特征值

u_1 为 S 的关于最大特征值的特征向量.

上面这是降到一维, 如果要降到 m 维, 就需要取 m 根轴.

(u_1, u_2, \dots, u_m) 其中, 对应的特征值为 $\lambda_1 > \lambda_2 > \dots > \lambda_m$

也就是依次取前 m 个大的特征值

② 最小重构误差 (minimum error perspective)

PCA 其实是用 S 的所有特征值对应的, 相互正交的特征向量重构特征空间

$x_i \in \mathbb{R}^P$, 则用 (u_1, u_2, \dots, u_p) 重构特征空间. 假设 x_i 已经经过中心化

x_i 在 u_k 方向上的投影为: $x_i^T u_k$, 加上方向 $\Rightarrow (x_i^T u_k) u_k$.

则经过降维后, $\hat{x}_i = \sum_{k=1}^q (x_i^T u_k) u_k$ 其中, q 表示选取前 q 个特征向量, 指降到 q

原始数据: $x_i = \sum_{k=1}^P (x_i^T u_k) u_k$ 维, 但是这里 \hat{x}_i 还是 P 维向量.

则重构误差为: $x_i - \hat{x}_i$

则目标是最小化总体的误差: $J = \frac{1}{N} \sum_{i=1}^N \|x_i - \hat{x}_i\|^2 = \frac{1}{N} \sum_{i=1}^N \left\| \sum_{k=q+1}^P (x_i^T u_k) u_k \right\|^2$

由于 $\|au\|_2^2 = a^2 \|u\|_2^2$ 且 $u_k^T u_k = 1$

$$\begin{aligned} \text{则 } J &= \frac{1}{N} \sum_{i=1}^N \sum_{k=q+1}^P (x_i^T u_k)^2 \\ &= \frac{1}{N} \sum_{i=1}^N \sum_{k=q+1}^P u_k^T x_i x_i^T u_k \\ &= \sum_{k=q+1}^P u_k^T \left[\frac{1}{N} \sum_{i=1}^N x_i x_i^T \right] u_k. \end{aligned}$$

由于 x_i 已经过中心化.

$$\text{则 } J = \sum_{k=q+1}^P u_k^T S u_k.$$

$$\text{则 } \{u_k\} = \underset{\{u_k\}}{\operatorname{argmin}} \sum_{k=q+1}^P u_k^T S u_k$$

$$\text{s.t. } u_k^T u_k = 1$$

和最大投影方差求解方式一样 特别地, 由于 u_k 之间线性无关, 则每个 u_k 可以单独进行求解. 即 $u_k = \underset{u_k}{\operatorname{argmin}} u_k^T S u_k$. s.t. $u_k^T u_k = 1$.

所以, 最终结果是 $\{u_k\}$ 是 S 的前 $(P-q)$ 个小的特征值所对应的特征向量, 即降维是选前 q 个大小特征向量.

和①中的结论相同

SVD 角度看 PCA 和 PCoA

① 特征分解

对于方阵 A , $\Rightarrow Ax = \lambda x \Rightarrow AX = X\Lambda \Rightarrow A = X\Lambda X^{-1}$

则 $A = X\Lambda X^{-1}$ 称为 A 的特征分解.

特别地, 当 A 是对称矩阵时, 特征向量是正交的.

即 $XX^T = X^T X = I$, $\Rightarrow A = X\Lambda X^T$

② 奇异值分解

设 $A \in \mathbb{R}^{m \times n}$, $\text{rank}(A) = r$, 则存在 m 阶、 n 阶酉阵 U, V .

使 $A = U \begin{pmatrix} \Sigma & 0 \\ 0 & 0 \end{pmatrix} V^H$, 其中 $\Sigma = \text{diag}(\sigma_1, \sigma_2, \dots, \sigma_r)$

σ_i 为矩阵 A 的非零奇异值, $\underbrace{U^H U = I}_{\text{正交}}, V^H V = V V^H = I$.

③ PCA

X 为数据集.

则对 X 进行中心化, 并对中心化后的数据进行特征值分解.

$HX = U\Sigma V^T$ (为了写得方便, 这里使用约化后的奇异值分解, 并将转置转

则 $S = X^T HX = X^T H^T HX = V\Sigma U^T U\Sigma V^T = \underbrace{V\Sigma^2 V^T}_{\text{特征分解置写为转置}}$.

则 V 为 S 的特征向量组成的矩阵, Σ^2 为特征值组成的对角矩阵.

即对 HX 进行 PCA 后的 V 即为重构空间的一组基底.

假设一个基底为 u_k , 则重构后, x_i 在 u_k 方向上的坐标为 $x_i^T u_k$.

写成矩阵形式为 XV .

则中心化后的数据在新空间的坐标为 HXV

化简 $HXV = U\Sigma V^T V = U\Sigma$.

令 $T = HXX^T H^T = U\Sigma V^T V\Sigma U^T = U\Sigma^2 U^T$ (U 为 T 的特征向量组成的矩阵)

$$TU\Sigma = U\Sigma^2 U^T U\Sigma = U\Sigma^3 = (U\Sigma)\Sigma^2$$

可以发现, $U\Sigma$ 也为 T 的特征向量组成的矩阵, 可以看作是 U 的缩放.

那么, 意味着对 T 进行特征值分解能直接得到坐标.

综上, S 和 T 有相同的 eigen value

且 $S \rightarrow$ 特征分解 \rightarrow 得到主成份 (基质的方向)

$T \rightarrow$ 特征分解 \rightarrow 得到坐标.

那么 S 代表主成份分析 (principal component analysis, PCA)

T 代表主坐标分析 (principal coordinate analysis, PCoA)

概率角度PCA (probabilistic PCA)

P-PCA

$$x \in \mathbb{R}^p, z \in \mathbb{R}^q, q < p$$



observed data



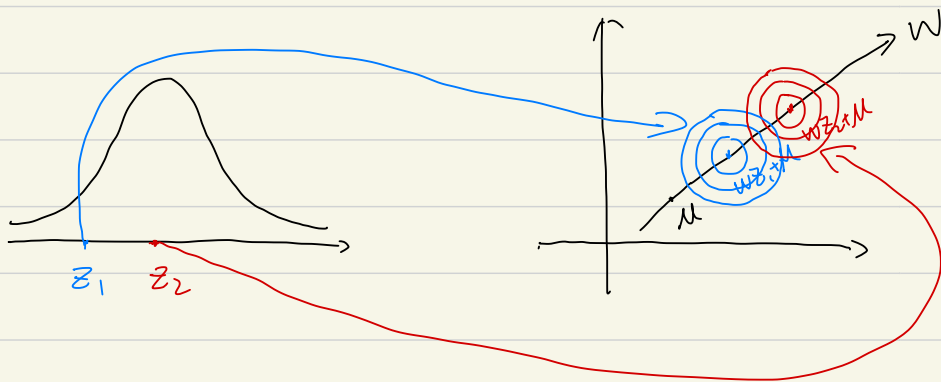
latent variable.

设

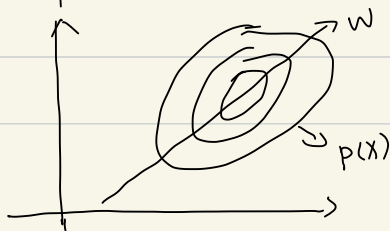
$$\begin{cases} z \sim \mathcal{N}(0_q, I_q) \\ x \sim Wz + \mu + \varepsilon \\ \varepsilon \sim \mathcal{N}(0_p, \sigma^2 I_p), \varepsilon \perp z. \end{cases}$$

\Rightarrow Linear Gaussian Model.

生成角度看 P-PCA. 取 $q=1, p=2$.



当 z 有足够多的采样数之后, 所有 z 的分布形成的分布如下所示



这就是 x 的分布

P-PCA 与 GMM 不同之处在于, GMM 中高斯分布的数量是确定的, 而 P-PCA 中高斯分布的数量是采样次数, 即 GMM 中, z 是离散, P-PCA 中, z 是连续的.

对于如何求解 P-PCA, $\begin{cases} \text{inference} \rightarrow p(z|x) \\ \text{learning} \rightarrow W, \mu, \sigma^2 \end{cases}$ (可以用最大似然, 也可以用 EM, 这里不讨论 learning).

本节主要关心 inference, 即如何求得 $p(z|x)$.

由已知的线性高斯模型, 可得:

$$E[x|z] = Wz + \mu + E[\varepsilon] = Wz + \mu$$

$$\text{Var}[x|z] = \text{Var}[\varepsilon] = \sigma^2 I.$$

$$\text{则 } x|z \sim N(Wz + \mu, \sigma^2 I).$$

$$\therefore E[x] = E[Wz + \mu] + E[\varepsilon] = \mu$$

$$\text{Var}[x] = \text{Var}[Wz] + \text{Var}[\varepsilon] = W W^T + \sigma^2 I = W W^T + \sigma^2 I$$

$$\text{则 } x \sim N(\mu, W W^T + \sigma^2 I)$$

为了求解 $z|x$, 可以使用数学基础部分求条件概率的方法.

$$\text{构造} \begin{pmatrix} x \\ z \end{pmatrix} \sim N \left(\begin{bmatrix} \mu \\ 0 \end{bmatrix}, \begin{bmatrix} WW^T + \sigma^2 I & \Delta \\ \Delta^T & I \end{bmatrix} \right)$$

$$\Delta = \text{Cov}(x, z)$$

$$= E[(x - \mu)(z - 0)^T]$$

$$= E[(x - \mu)z^T]$$

$$= E[(Wz + \varepsilon)z^T]$$

$$= E[Wzz^T] + \underbrace{E[\varepsilon z^T]}_{\varepsilon \perp z}$$

$$= W E[zz^T] + E[\varepsilon] E[z^T]$$

$$= W \cdot I$$

$$= W$$

$$\text{则} \begin{pmatrix} x \\ z \end{pmatrix} \sim N \left(\begin{bmatrix} \mu \\ 0 \end{bmatrix}, \begin{bmatrix} WW^T + \sigma^2 I & W \\ W^T & I \end{bmatrix} \right)$$

求得 x 看 $r|x_a$, z 看 $r|x_b$. 套用

$$\begin{cases} x_{b \cdot a} = x_b - \Sigma_{ba} \Sigma_{aa}^{-1} x_a \\ \mu_{b \cdot a} = \mu_b - \Sigma_{ba} \Sigma_{aa}^{-1} \mu_a \Rightarrow x_b | x_a \sim N(\mu_{b \cdot a} + \Sigma_{ba} \Sigma_{aa}^{-1} x_a, \Sigma_{bb \cdot a}) \\ \Sigma_{bb \cdot a} = \Sigma_{bb} - \Sigma_{ba} \Sigma_{aa}^{-1} \Sigma_{ab} \end{cases}$$

$$\text{得: } z|x \sim N(W^T(WW^T + \sigma^2 I)^{-1}(x - \mu), I - W^T(WW^T + \sigma^2 I)^{-1}W)$$