

第5章

降维与度量学习

Dimensionality Reduction and Metric Learning

赫 然

rhe@nlpr.ia.ac.cn

<https://rhe-web.github.io/>

智能感知与计算研究中心 (CRIPAC)

中科院自动化研究所 模式识别国家重点实验室

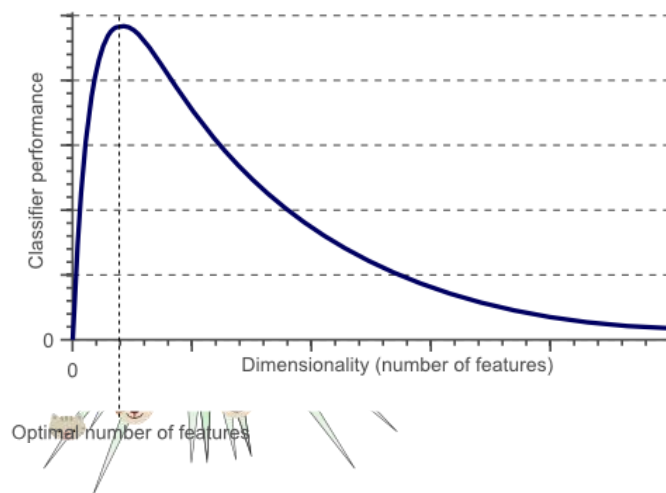
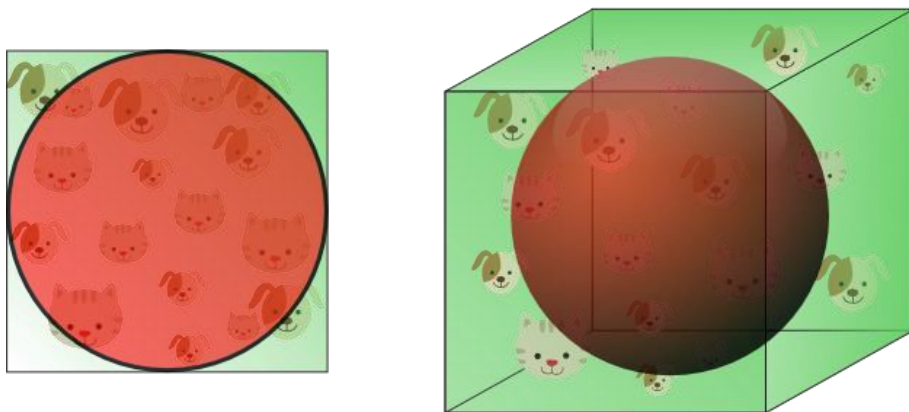
内容提要

- 引言
- 主成分分析
- 多维缩放
- 流形学习方法
- 距离度量学习

引言

● 维度灾难 (Curse of Dimensionality)

— 维数灾难最早由理查德·贝尔曼 (Richard E. Bellman) 在考虑优化问题时提出的，它用来描述当空间维度增加时分析和组织高维空间中的数据会遇到各种问题。



引言

● 维度灾难 (Curse of Dimensionality)

- 随着维数的增加，**计算量**呈指数倍增长。
- 随着维数的增加，具有**相同距离**的两个样本其**相似程度**可以相差很远。
- 当维度增加时，**空间的体积**增加得很快，可用**数据变得稀疏**。
- 稀疏性对于任何要求“**具有统计学意义的方法**”而言都是一个问题。但是，为了获得在统计学上正确并且有可靠的结果，用来支撑这一结果所需要的数据量通常**随着维数的增加而呈指数级增长**。

引言

● 维度缩减

— 缓解维数灾难的一个重要途径是**降维**，即通过某种数学变换将原始高维特征空间变换至某个低维“子空间”。在该子空间中，样本密度大幅度提高，距离计算也变得更加容易。

— 为什么能降维：

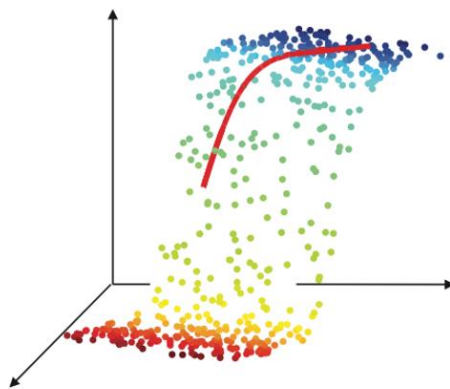
- 在很多时候，人们观测或收集到的数据虽然是高维的，但与学习任务密切相关的特征通常位于某个**低维分布上**，即高维空间中的一个低维“**嵌入**” (embedding)。

引言

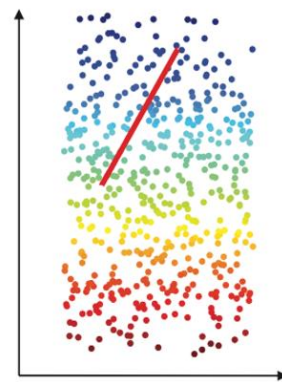
● 维度缩减

— 缓解维数灾难的一个重要途径是**降维**，即通过某种数学变换将原始高维特征空间变换至某个低维“子空间”。在该子空间中，样本密度大幅度提高，距离计算也变得更为容易。

— **为什么能降维：**



(a) 三维空间中观察到的样本点



(b) 二维空间中的曲面

内容提要

- 引言
- 主成分分析
- 多维缩放
- 流形学习方法
- 距离度量学习

主成分分析(Principle Component Analysis)

● 线性降维法

- 对高维空间中的样本 x 进行线性变换：

$$y = W^T x \quad x \in R^m, W \in R^{m \times d}, y \in R^d, d < m$$

- 变换矩阵 $W = [w_1, w_2, \dots, w_d]$ 可视为 m 维空间中由 d 个基向量组成的矩阵。

- $y = W^T x$ 可视为样本 x 与 d 个基向量分别做内积运算而得到，即 x 在新坐标系下的坐标。显然，新空间中的特征是原空间中特征的线性组合。

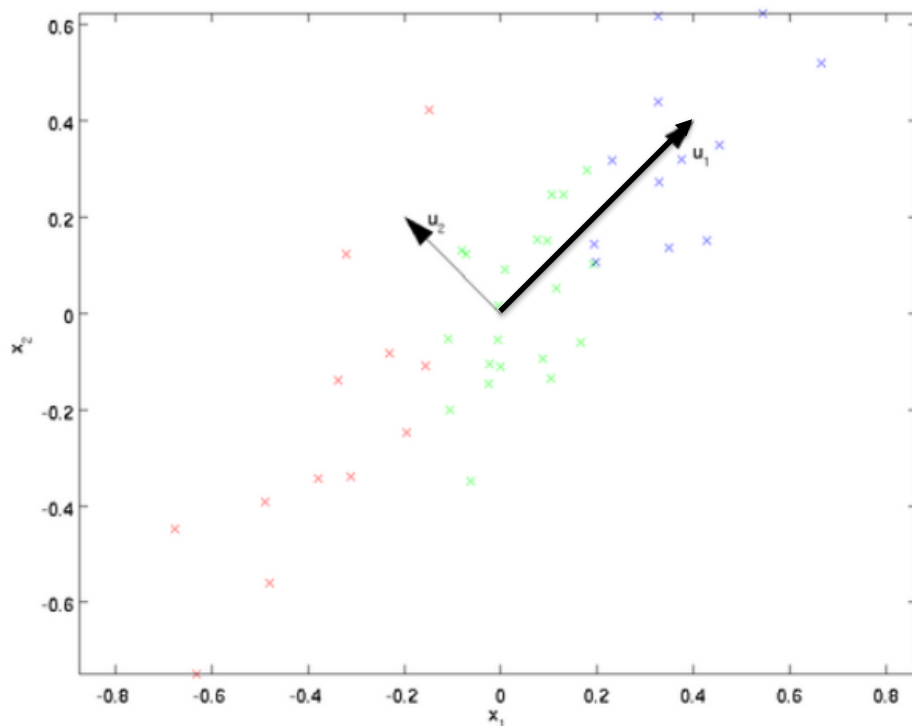
- 不同方法的不同之处在于对低维子空间的性质有不同的要求，即对 W 施加不同的约束。

主成分分析

- PCA (Principal Component Analysis) 基本思想

- 如何仅用一个超平面从整体上对所有样本

$\{x_1, x_2, \dots, x_n\} \in R^m$ 进行恰当表示？



PCA是一种统计方法。通过正交变换将一组可能存在相关性的变量转换为一组线性不相关的变量，转换后的这组变量叫主成分。

主成分分析

- PCA (Principal Component Analysis) 基本思想

- 我们如何仅用一个超平面从整体上对所有样本

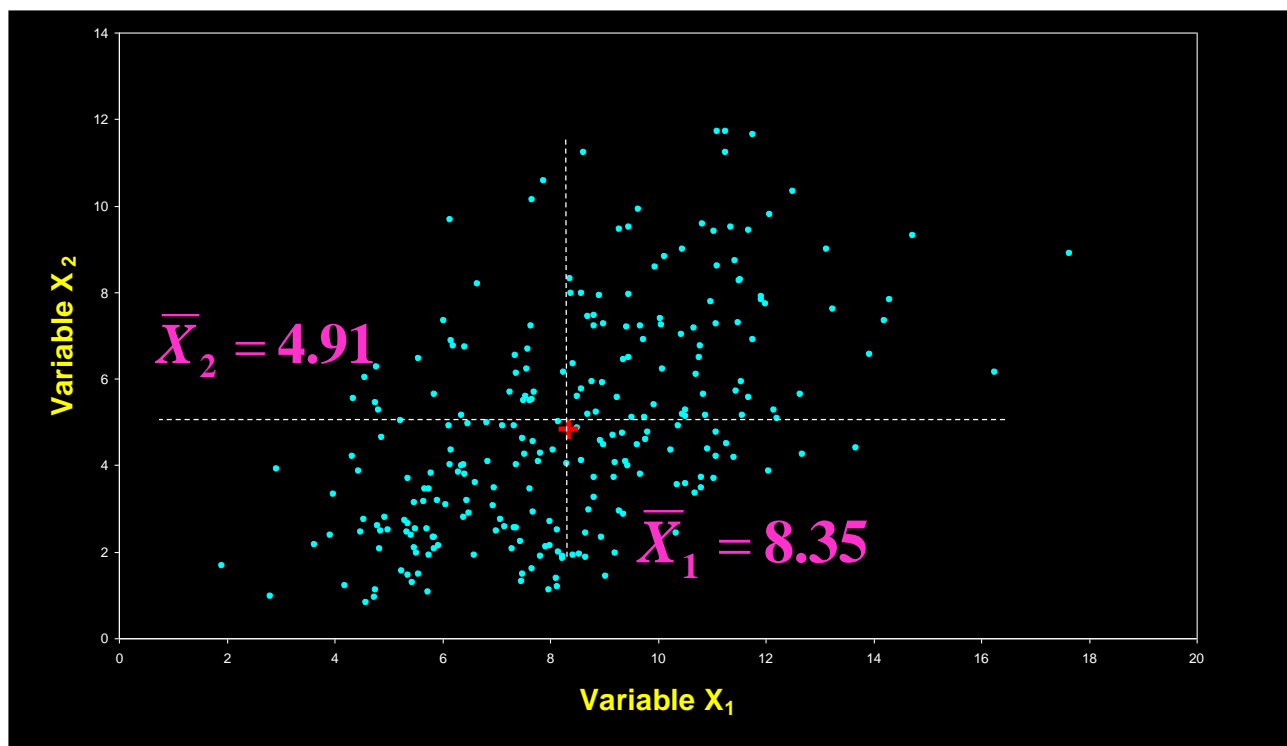
- $\{x_1, x_2, \dots, x_n\} \in R^m$ 进行恰当表示？

- 通常有如下两种思路：

- **可重构性**：样本到这个超平面的距离都足够近；**基于最小投影距离**。
 - **可区分性**：样本点在这个超平面上的投影能够尽可能地分开。**基于最大投影方差**。

主成分分析1- Minimum Error Formulation

- m 维空间的完备正交基组 (complete orthonormal set), 即新坐标系由 $W = [w_1, w_2, \dots, w_m]$, 且



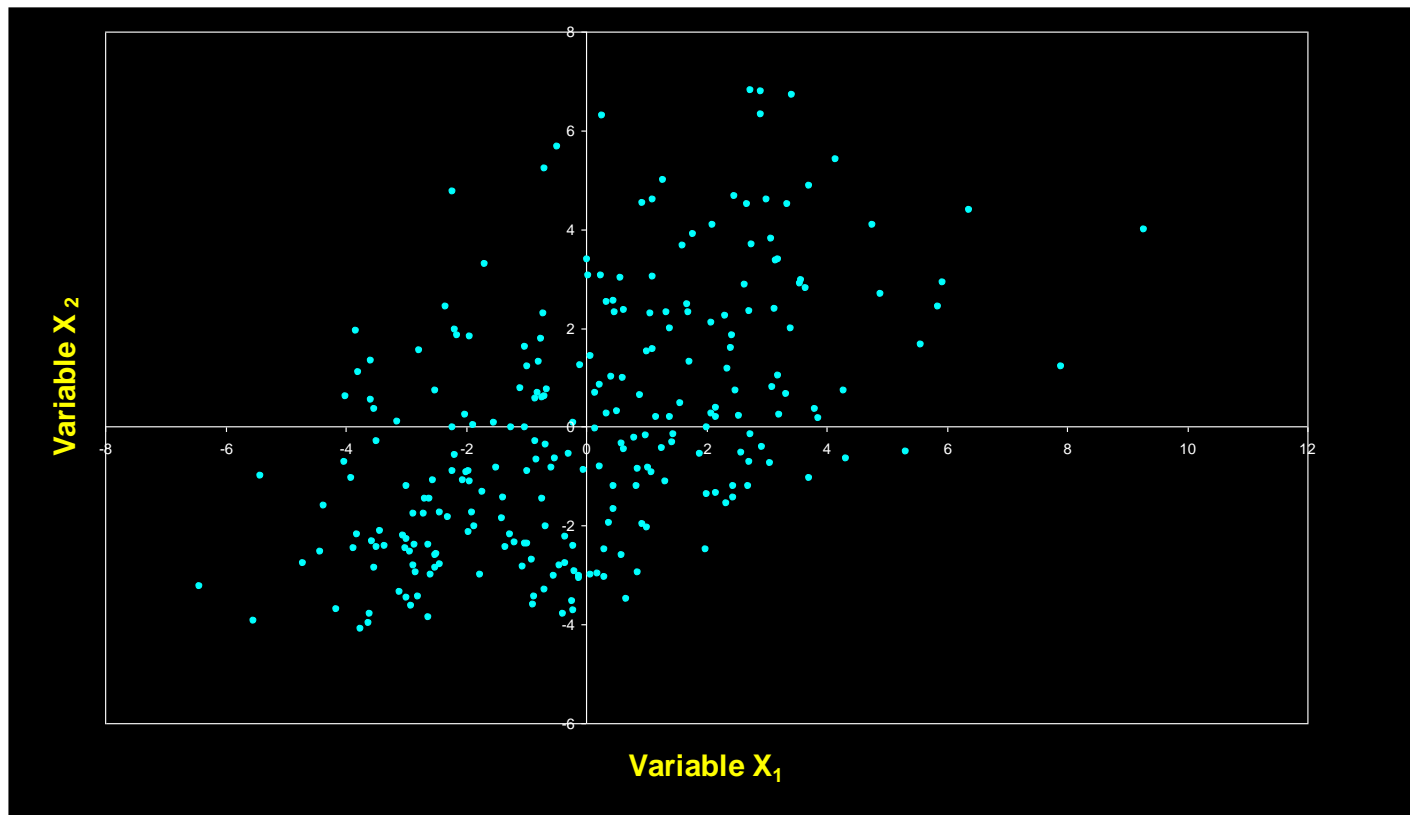
$$V_1 = 6.67$$

$$V_2 = 6.24$$

$$C_{1,2} = 3.42$$

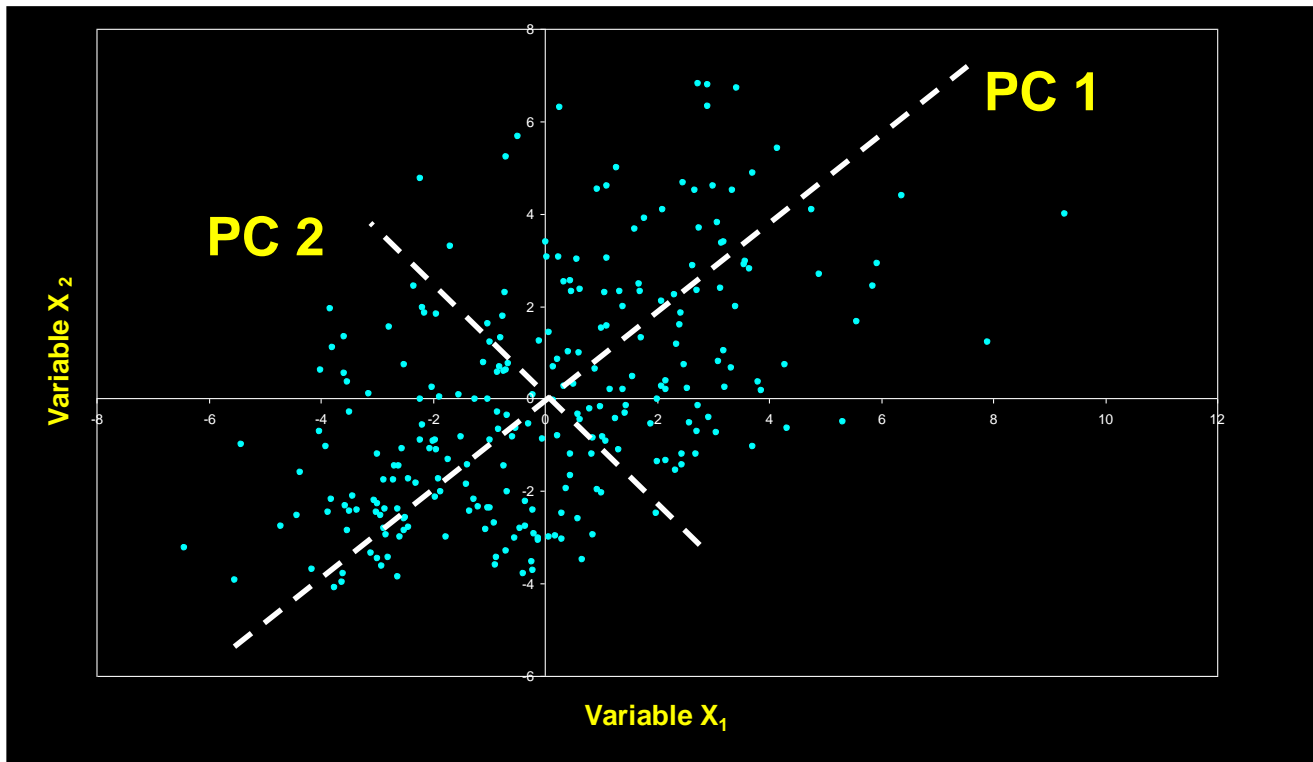
主成分分析1- Minimum Error Formulation

- m 维空间的完备正交基组(complete orthonormal set), 即新坐标系由 $W = [w_1, w_2, \dots, w_m]$, 且



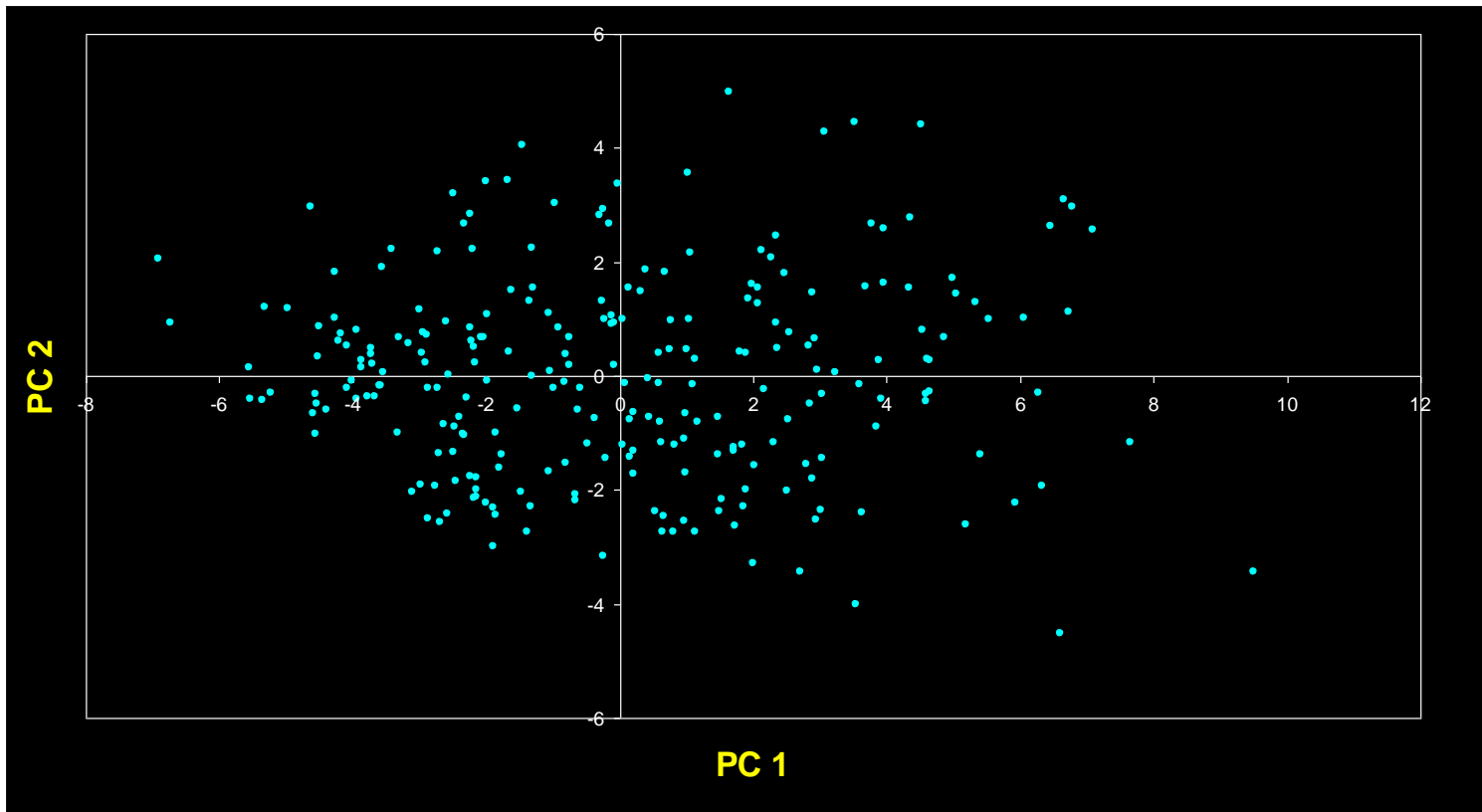
主成分分析1- Minimum Error Formulation

- m 维空间的完备正交基组(complete orthonormal set), 即新坐标系由 $W = [w_1, w_2, \dots, w_m]$, 且



主成分分析1- Minimum Error Formulation

- m 维空间的完备正交基组(complete orthonormal set), 即新坐标系由 $W = [w_1, w_2, \dots, w_m]$, 且



主成分分析1- Minimum Error Formulation

- 核心思想是通过旋转坐标系以最小化误差(projection error)
- m 维空间的完备正交基组(complete orthonormal set), 即新坐标系由 $W = [w_1, w_2, \dots, w_m]$, 且

$$w_k^T w_j = \delta_{kj} = \begin{cases} 1, & k = j \\ 0, & k \neq j \end{cases}$$



$$x_i = \sum_{j=1}^m y_{ij} w_j$$

主成分分析1- Minimum Error Formulation

$$x_i = \sum_{k=1}^m y_{ik} w_k$$



$$w_j^T x_i = \sum_{k=1}^m w_j^T y_{ik} w_k \Rightarrow w_j^T x_i = w_j^T y_{ij} w_j \Rightarrow y_{ij} = x_i^T w_j$$



$$x_i = \sum_{j=1}^m (x_i^T w_j) w_j$$

主成分分析1- Minimum Error Formulation

- 根据我们对已有数据的理解，我们知道有的冗余信息较多 (high redundant information)，有的信息冗余信息较少 (low redundant information)。我们希望达到的目的是将两者分开，因此认为将投影后的估计值进行线性拆分

$$\tilde{x}_i = \sum_{j=1}^{\underline{d}} z_{ij} w_j + \sum_{j=\underline{d}+1}^m b_j w_j$$

- 其中 b_j 只跟维度相关

主成分分析1- Minimum Error Formulation

- 目标函数 (objective function, loss function)

$$\begin{aligned} J &= \frac{1}{n} \sum_{i=1}^n \|x_i - \tilde{x}_i\|_2^2 \\ &= \frac{1}{n} \sum_{i=1}^n \left\| x_i - \left[\sum_{j=1}^d z_{ij} w_j - \sum_{j=d+1}^m b_j w_j \right] \right\|_2^2 \end{aligned}$$

- 对 b_j 求导并令其为零, 得到

$$b_j = \bar{x}^T w_j, \quad j = d+1, d+2, \dots, m \quad \bar{x} = \frac{1}{n} \sum_{i=1}^n x_i$$

主成分分析1- Minimum Error Formulation

- 对 z_{ij} 求导得到,

$$\begin{aligned}\frac{\partial J}{\partial z_{ij}} &= \frac{\partial J}{\partial \tilde{x}_i} \frac{\partial \tilde{x}_i}{\partial z_{ij}} = \frac{\partial J}{\partial \tilde{x}_i} \frac{\partial}{\partial z_{ij}} \left(\sum_{k=1}^d z_{ik} w_k \right) & \frac{\partial J}{\partial \tilde{x}_i} &= -\frac{2}{n} (x_i - \tilde{x}_i)^T \\ &= -\frac{2}{n} (x_i - \tilde{x}_i)^T w_j \\ &= -\frac{2}{n} \left(x_i - \sum_{k=1}^d z_{ik} w_k \right)^T w_j = -\frac{2}{n} (x_i^T w_j - z_{ij})\end{aligned}$$

- 令其为零, 我们得到

$$z_{ij} = x_i^T w_j, \quad j = 1, 2, \dots, d$$

主成分分析1- Minimum Error Formulation

- 将 z_{ij} 和 b_j 表达式分别带入 $x_i - \tilde{x}_i$ ，得到，

$$\begin{aligned}x_i - \tilde{x}_i &= \sum_{j=1}^m (x_i^T w_j) w_j - \left(\sum_{j=1}^d z_{ij} w_j + \sum_{j=d+1}^m b_j w_j \right) \\&= \sum_{j=1}^m (x_i^T w_j) w_j - \sum_{j=1}^d (x_i^T w_j) w_j - \sum_{j=d+1}^m (\bar{x}^T w_j) w_j \\&= \sum_{j=d+1}^m ((x_i - \bar{x})^T w_j) w_j\end{aligned}$$

- 带入到目标函数 J 中

主成分分析1- Minimum Error Formulation

- 目标函数 J 可重写为之关于 w_j 的形式,

$$\begin{aligned} J &= \frac{1}{n} \sum_{i=1}^n \left\{ \left[\sum_{j=d+1}^m ((x_i - \bar{x})^T w_j) w_j \right]^T \left[\sum_{j=d+1}^m ((x_i - \bar{x})^T w_j) w_j \right] \right\} \\ &= \frac{1}{n} \sum_{i=1}^n \left\{ \sum_{j=d+1}^m (x_i^T w_j - \bar{x}^T w_j)^2 w_j^T w_j \right\} \\ &= \frac{1}{n} \sum_{i=1}^n \sum_{j=d+1}^m (x_i^T w_j - \bar{x}^T w_j)^2 \\ &= \sum_{j=d+1}^m w_j^T S w_j \end{aligned}$$
$$S = \frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})(x_i - \bar{x})^T$$

主成分分析1- Minimum Error Formulation

- 目标函数 J 可重写为关于 w_j 的形式,

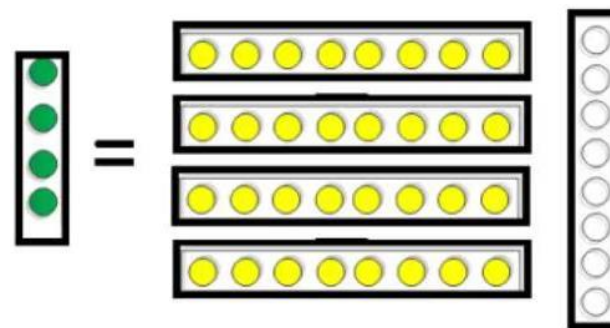
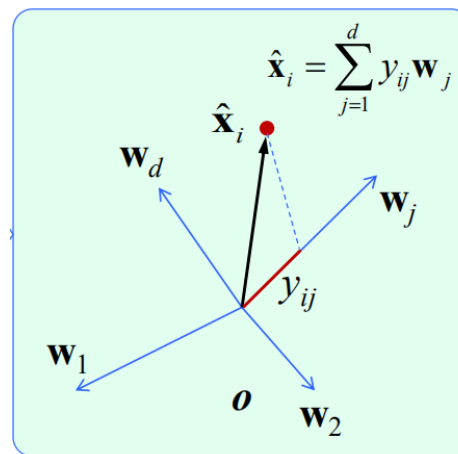
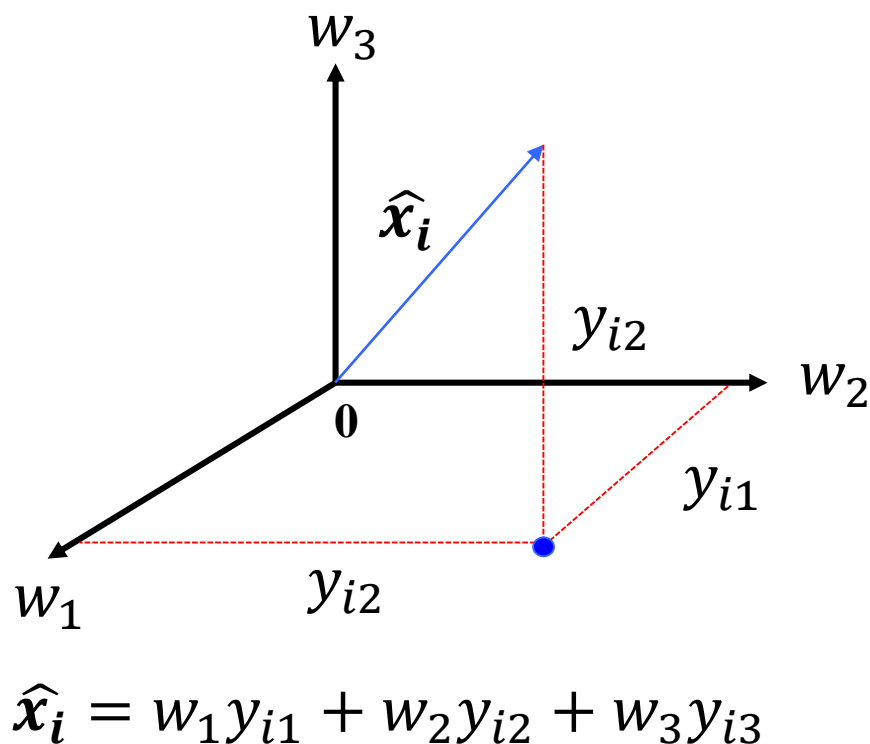
$$\begin{aligned}\min J &= \min \frac{1}{n} \sum_{i=1}^n \|x_i - \tilde{x}_i\|_2^2 \\ &= \min_{w_j} \sum_{j=d+1}^m w_j^T S w_j = \min_W \text{tr}(W^T S W)\end{aligned}$$

- w_j 为协方差矩阵 S 的最小特征值对应的特征向量

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

主成分分析2- 重构角度

— 向量空间中一个矢量的表示方法（零均值）



$$y = W^T x$$

主成分分析2- 重构角度

— 由 W 定义新坐标系：假定投影变换是**正交变换**，即新坐标系由 $W = [w_1, w_2, \dots, w_d]$ 来表示 ($d < m$)， w_i 的模等于1， w_i 与 w_j **两两正交**。

— 设样本点 x_i 在新坐标系下的坐标为： $y_i = [y_{i1}, y_{i2}, \dots, y_{id}]^T \in R^d$

— 在正交坐标系下，对样本点 x_i ，有新坐标：

$$y_{ij} = w_j^T x_i, w_j \in R^m \quad j = 1, 2, \dots, d$$

— 在新坐标系下，可得 x_i 的新表示：

$$\hat{x}_i = \sum_{j=1}^d y_{ij} w_j \quad i = 1, 2, \dots, n$$

主成分分析

$$\sum_{i=1}^n \|x_i - \hat{x}_i\|_2^2 = \sum_{i=1}^n \|x_i - \sum_{j=1}^d y_{ij} w_j\|_2^2 = \sum_{i=1}^n \|x_i - W y_i\|_2^2$$

$$x \in R^m$$

$$W \in R^{m \times d}, y_i \in R^d$$

重构误差

主成分分析



$$\begin{aligned}\sum_{i=1}^n \|x_i - \hat{x}_i\|_2^2 &= \sum_{i=1}^n \left\| x_i - \sum_{j=1}^d y_{ij} w_j \right\|_2^2 = \sum_{i=1}^n \|x_i - W y_i\|_2^2 \\ &= \sum_{i=1}^n ((W y_i)^T W y_i - 2x_i^T W y_i + x_i^T x_i)\end{aligned}$$

主成分分析



$$\begin{aligned}\sum_{i=1}^n \|x_i - \hat{x}_i\|_2^2 &= \sum_{i=1}^n \|x_i - \sum_{j=1}^d y_{ij} w_j\|_2^2 = \sum_{i=1}^n \|x_i - W y_i\|_2^2 \\ &= \sum_{i=1}^n ((W y_i)^T W y_i - 2x_i^T W y_i + x_i^T x_i) \\ \because W^T W &= I \quad = \sum_{i=1}^n (y_i^T y_i - 2y_i^T y_i + x_i^T x_i)\end{aligned}$$

主成分分析



$$\sum_{i=1}^n \|x_i - \hat{x}_i\|_2^2 = \sum_{i=1}^n \|x_i - \sum_{j=1}^d y_{ij} w_j\|_2^2 = \sum_{i=1}^n \|x_i - W y_i\|_2^2$$

$$= \sum_{i=1}^n ((W y_i)^T W y_i - 2x_i^T W y_i + x_i^T x_i)$$

$$\because W^T W = I$$

$$= \sum_{i=1}^n y_i^T y_i - 2y_i^T y_i + x_i^T x_i$$

$$\because y_i = W^T x_i$$
$$= - \sum_{i=1}^n y_i^T y_i + \text{常量} = - \sum_{i=1}^n (W^T x_i)^T (W^T x_i) + const$$

主成分分析



$$\sum_{i=1}^n \|x_i - \hat{x}_i\|_2^2 = \sum_{i=1}^n \|x_i - \sum_{j=1}^d y_{ij} w_j\|_2^2 = \sum_{i=1}^n \|x_i - W y_i\|_2^2$$

$$= \sum_{i=1}^n ((W y_i)^T W y_i - 2x_i^T W y_i + x_i^T x_i)$$

$$\because W^T W = I$$

$$= \sum_{i=1}^n (y_i^T y_i - 2y_i^T x_i + x_i^T x_i)$$

$$\because y_i = W^T x_i$$

$$= - \sum_{i=1}^n y_i^T x_i + \text{const} = - \sum_{i=1}^n (W^T x_i)^T (W^T x_i) + \text{const}$$

$$= -\text{tr} \left(\sum_{i=1}^n (W^T x_i)^T (W^T x_i) \right) + \text{const}$$

主成分分析



$$\sum_{i=1}^n \|x_i - \hat{x}_i\|_2^2 = \sum_{i=1}^n \|x_i - \sum_{j=1}^d y_{ij} w_j\|_2^2 = \sum_{i=1}^n \|x_i - W y_i\|_2^2$$

$$= \sum_{i=1}^n ((W y_i)^T W y_i - 2x_i^T W y_i + x_i^T x_i)$$

$$\because W^T W = I$$

$$= \sum_{i=1}^n (y_i^T y_i - 2y_i^T x_i + x_i^T x_i)$$

$$\because y_i = W^T x_i$$

$$= - \sum_{i=1}^n y_i^T x_i + \text{const} = - \sum_{i=1}^n (W^T x_i)^T (W^T x_i) + \text{const}$$

$$= -\text{tr} \left(\sum_{i=1}^n (W^T x_i)^T (W^T x_i) \right) + \text{const}$$

$$\because \text{tr}(AB) = \text{tr}(BA)$$

$$\because \text{tr}(A) + \text{tr}(B) = \text{tr}(A + B) \quad = -\text{tr} \left(W^T \sum_{i=1}^n x_i x_i^T W \right) + \text{const}$$

主成分分析2- 重构角度

— 进一步，假定数据已经零均值化， $\sum_{i=1}^n x_i = 0$ 即

$$X = [x_1, x_2, \dots, x_n] \in R^{m \times n}$$

$$\begin{aligned} \sum_{i=1}^n \|x_i - \hat{x}_i\|_2^2 &= -tr \left(W^T \sum_{i=1}^n x_i x_i^T W \right) + const \\ &= -tr \left(W^T X X^T W \right) + const \end{aligned}$$

主成分分析2- 重构角度

— 进一步，假定数据已经零均值化， $\sum_{i=1}^n x_i = 0$ 即

$$X = [x_1, x_2, \dots, x_n] \in R^{m \times n}$$

$$\begin{aligned} \sum_{i=1}^n \|x_i - \hat{x}_i\|_2^2 &= -tr \left(W^T \sum_{i=1}^n x_i x_i^T W \right) + const \\ &= -tr \left(W^T X X^T W \right) + const \end{aligned}$$

于是，获得主成分分析的最优化模型：

主成分分析2- 重构角度

— 进一步，假定数据已经零均值化， $\sum_{i=1}^n x_i = 0$ 即

$$X = [x_1, x_2, \dots, x_n] \in R^{m \times n}$$

$$\begin{aligned} \sum_{i=1}^n \|x_i - \hat{x}_i\|_2^2 &= -tr \left(W^T \sum_{i=1}^n x_i x_i^T W \right) + const \\ &= -tr \left(W^T X X^T W \right) + const \end{aligned}$$

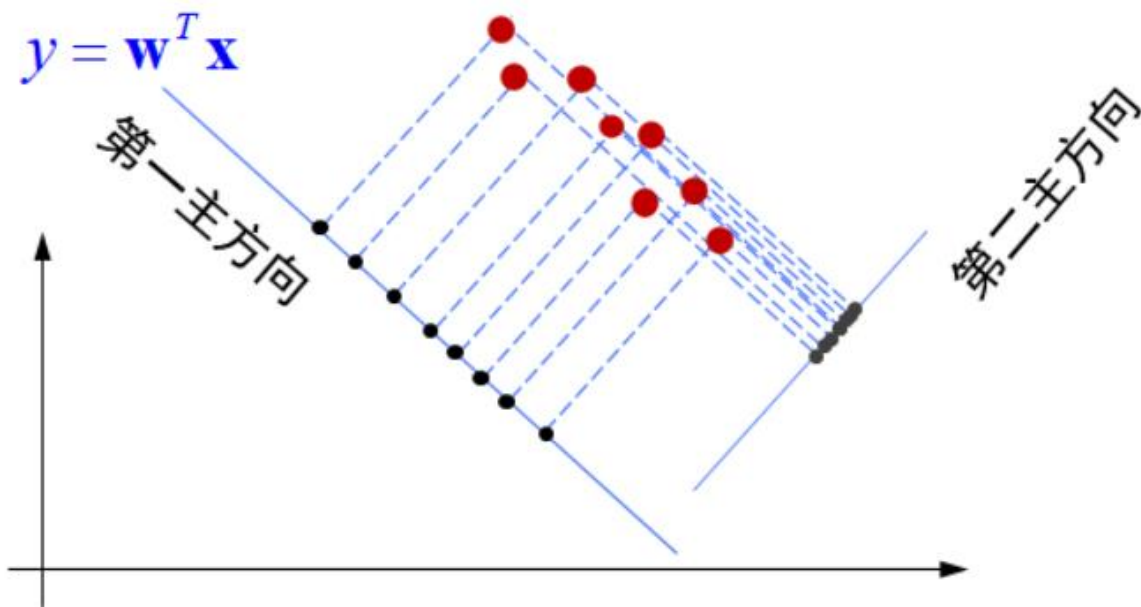
于是，获得主成分分析的最优化模型：

$$\max_{W \in R^{m \times d}} tr(W^T X X^T W) + const, \text{ s.t. } W^T W = I$$

主成分分析3- 可区分角度

— 使所有样本点的投影尽可能地分开，则需**最大化投影点的方差**，投影后获得的样本点为：

$$y_i = \mathbf{w}^T \mathbf{x}_i \in R^d, i = 1, 2, \dots, n$$



主成分分析3- 可区分角度

— 使所有样本点的投影尽可能地分开，则需**最大化投影点的方差**，投影后获得的样本点为：

$$y_i = w^T x_i \in R^d, i = 1, 2, \dots, n$$

— 由于数据点是**零均值化**的，则：

$$\sum_{i=1}^n y_i = w^T \sum_{i=1}^n x_i = 0$$

主成分分析3- 可区分角度

- 使所有样本点的投影尽可能地分开，则需**最大化投影点的方差**，投影后获得的样本点为：

$$y_i = w^T x_i \in R^d, i = 1, 2, \dots, n$$

- 由于数据点是**零均值化**的，则：

$$\sum_{i=1}^n y_i = w^T \sum_{i=1}^n x_i = 0$$

- 因此投影后的样本点的**（协）方差**为：

$$\sum_{i=1}^n w^T x_i x_i^T w = w^T X X^T w$$

主成分分析3- 可区分角度

要使数据具有最大可分性，就应该使数据尽量分散开来，因此应该使其**方差最大**。

$$\max_{W \in \mathbb{R}^{m \times d}} W^T X X^T W \quad s.t. W^T W = I$$

主成分分析

- PCA求解

—采用拉格朗日乘子法，经过简单矩阵运算，我们有：

$$XX^T w = \lambda w$$

$$J(W) = w^T XX^T w + \lambda(1 - w^T w)$$

w 求导 $XX^T w - \lambda w = 0$

$$XX^T w = \lambda w$$

主成分分析

- PCA求解

- 采用拉格朗日乘子法，经过简单矩阵运算，我们有：

$$XX^T w = \lambda w$$

则主成分分析的解为：

主成分分析

● PCA求解

—采用拉格朗日乘子法，经过简单矩阵运算，我们有：

$$XX^T w = \lambda w$$

则主成分分析的解为：

对协方差矩阵 XX^T 进行特征值分解，并对特征值进行排序 $\lambda_1 \geq \lambda_2 \geq \lambda_3 \geq \dots \geq \lambda_d$ 取前 d 个特征值对应的特征向量构成变换矩阵 W 。

主成分分析

● PCA求解

—采用拉格朗日乘子法，经过简单矩阵运算，我们有：

$$XX^T w = \lambda w$$

则主成分分析的解为：

对协方差矩阵 XX^T 进行特征值分解，并对特征值进行排序 $\lambda_1 \geq \lambda_2 \geq \lambda_3 \geq \dots \geq \lambda_d$ 取前 d 个特征值对应的特征向量构成变换矩阵 W 。

λ_1 对应的特征向量称为第一主成分，其他依此类推。

主成分分析

● PCA求解

PCA算法步骤

输入样本集 $D = \{x_1, x_2, \dots, x_m\}$, 维数 d'

1 对所有样本进行中心化: $x_i \leftarrow x_i - \frac{1}{m} \sum_{i=1}^m x_i$

2 计算样本的协方差矩阵 XX^T

3 对协方差矩阵 XX^T 做特征分解

4 取最大的 d' 个特征值对应的特征向量 $w_1, w_2, \dots, w_{d'}$

输出投影矩阵 $W = (w_1, w_2, \dots, w_{d'})$ 。

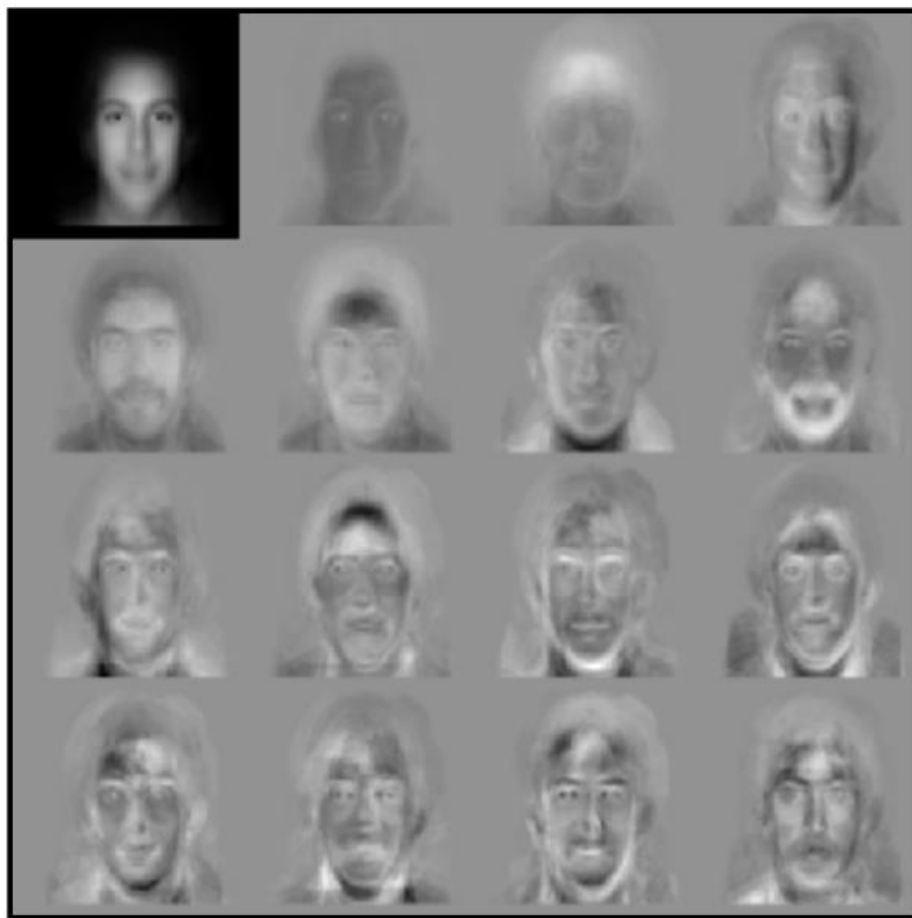
主成分分析

● 讨论

- 降低至多少维：
- 采用交叉验证，结合**最近邻分类器**来选择合适的维度 d 。
- 公式法：
$$\frac{\sum_{i=1}^d \lambda_i}{\sum_{i=1}^m \lambda_i} \geq t \quad (\text{比如, } t=95\%)$$
- 舍弃 $m-d$ 个特征值对应的特征向量导致了维数缩减。
 - 舍弃这些信息之后能使样本的采样密度增大，这正是降低维的重要动机。
 - 另外，当数据受到噪声影响时，**最小的特征值所对应的特征向量往往与噪声有关**，将它们舍弃可在一定程度上起到去噪的效果。

主成分分析应用

- 讨论-物理意义 (Meanface, EigenFaces)



Eigenfaces
from 7562
images:

**top left image
is linear
combination
of rest.**

**Sirovich & Kirby (1987)
Turk & Pentland (1991)**

主成分分析- 概率PCA

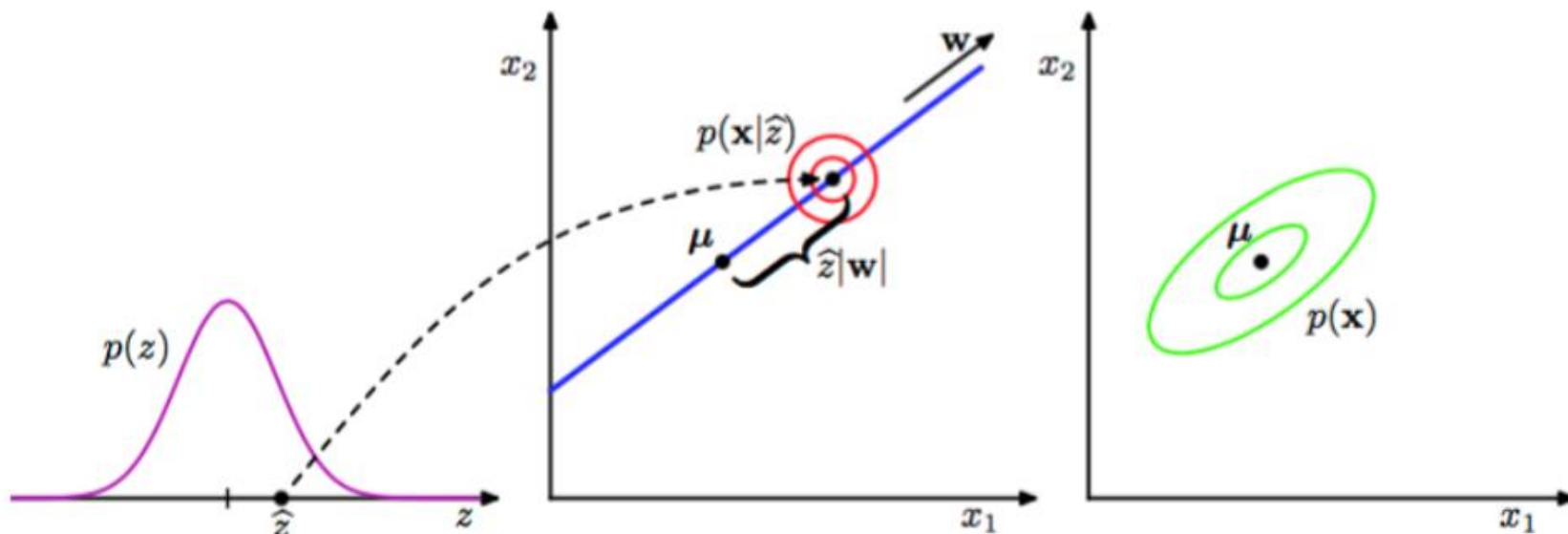
—构建一个概率模型，对于一个预测数据 x ，可以认为它是由一个隐变量 z 生成并且 z 及其条件概率 $p(x|z)$ 均服从高斯分布：

$$\begin{aligned}p(z) &= N(z|0, I) \\p(x|z) &= N(x|Wz + \mu, \sigma^2 I) \\x &= Wz + \mu + \epsilon\end{aligned}$$

其中， x 的均值是 z 的一个一般线性函数，由 $m \times d$ 的矩阵 W （列张成数据空间的线性子空间）和 m 维向量 μ 控制。而 σ^2 控制着条件概率分布的方差， ϵ 为 m 维0均值高斯分布噪声，协方差为 $\sigma^2 I$ 。

主成分分析- 概率PCA

构建一个概率模型，对于一个预测数据 \mathbf{x} ，可以认为它是由一个隐变量 z 生成并且 z 及其条件概率 $p(\mathbf{x}|z)$ 均服从高斯分布，生成过程如下：



主成分分析- 概率PCA

线性高斯模型，还是高斯分布

$$p(x) = \int p(x|z)p(z)dz \stackrel{\downarrow}{=} N(x|\mu, C)$$

$$C = WW^T + \sigma^2 I$$

$$\begin{aligned} E(x) &= E(Wz + \mu + \epsilon) = \mu \\ cov(x) &= E([X - E(x)]^2) = E[(Wz + \epsilon)(Wz + \epsilon)^T] \\ cov(x) &= E(Wzz^T W^T) + E(\epsilon\epsilon^T) = WW^T + \sigma^2 I \end{aligned}$$

z, ϵ 相互独立的随机变量, $p(z) = N(z|0, I)$, $p(\epsilon) = N(\epsilon|0, \sigma^2 I)$

主成分分析

- 求解

- 对应的线性高斯模型，它的边缘分布还是高斯模型：

$$p(x) = N(x|\mu, C)$$

C 是 $D \times D$ 协方差矩阵, $D = \{x_1, x_2, \dots, x_n\}$

$$C = WW^T + \sigma^2 I$$

- 极大似然估计函数为：

$$\ln p(X|W, u, \sigma^2) = -\frac{n}{2} \{D \ln(2\pi) + \ln|C| + \text{tr}(C^{-1}S)\}$$
$$S = \frac{1}{n} \sum (x_i - u)(x_i - u)^T$$

主成分分析

- 求解

- 解为: $W = U_M(L_M - \sigma^2 I)^{\frac{1}{2}}R$ L_M 是包含特征值的对角矩阵

- 特殊情况: $R = I$

$$W = U_M(L_M - \sigma^2 I)^{\frac{1}{2}}$$

W 中包含特征值 L_M
起到缩放到标准高斯的作用

- 考虑特征值的标准PCA的映射关系

$$x = Wz = U_M L_M^{\frac{1}{2}} z$$

- 伸缩矩阵: $L_0 = (L_M - \sigma^2 I)^{\frac{1}{2}}$

主成分分析

● 求解

– **W 的列**决定了标准PCA的主子空间, 对于 σ^2 , 对应的最大化似然解为:

$$\sigma^2 = \frac{1}{m-d} \sum_{j=d+1}^m \lambda_j$$

此时 σ^2 是与丢弃维度相关的平均方差, 如果 u 是主子空间 (principal subspace) 中的一个向量:

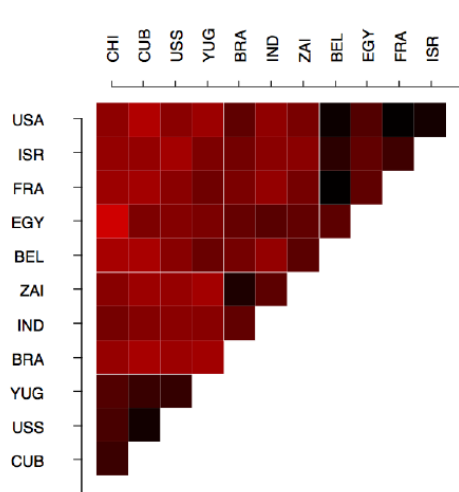
$$\begin{aligned} u^T C u &= u^T (W W^T + \sigma^2 I) u = (u^T U_M (L_M - \sigma^2 I) U_M^T u) + \sigma^2 \\ &= (u^T \sum_{j=1}^d \lambda_j u_j u_j^T u - \sigma^2) + \sigma^2 \\ &= (\lambda - \sigma^2) + \sigma^2 = \lambda \end{aligned}$$

内容提要

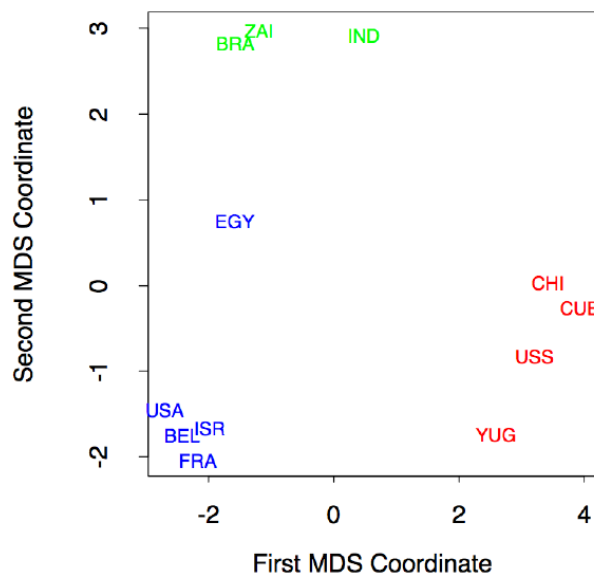
- 引言
- 主成分分析
- 多维缩放
- 流形学习方法
- 距离度量学习

多维缩放(Multiple Dimensional Scaling, MDS)

- 假定 m 维空间的 n 个样本 $\{x_1, x_2, \dots, x_n\} \in R^m$ 在原始空间的距离矩阵 $D \in R^{n \times n}$, 其第 i 行 j 列的元素为样本 x_i 到 x_j 的距离。

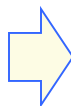


Reordered Dissimilarity Matrix



First MDS Coordinate

已知距离矩阵



低维的MDS坐标

多维缩放(Multiple Dimensional Scaling, MDS)

- 假定 m 维空间的 n 个样本 $\{x_1, x_2, \dots, x_n\} \in R^m$ 在原始空间的**距离矩阵** $D \in R^{n \times n}$, 其第 i 行 j 列的元素为样本 x_i 到 x_j 的距离。
- **目标**: 获得这 n 个样本 $Z = [z_1, z_2, \dots, z_n] \in R^{d \times n}$ 在 $d(d < m)$ 维空间中的表示。
- **准则**: 假定降维后的样本仍**保持两两之间的距离**。

MDS不是保留数据的最大可分性, 而是**更加关注与高维数据内部的特征**。MDS算法集中于保留高维空间中的“**相似度**”信息, 即 $B \doteq Z^T Z$ 。

多维缩放(Multiple Dimensional Scaling, MDS)

- 假定 m 维空间的 n 个样本 $\{x_1, x_2, \dots, x_n\} \in R^m$ 在原始空间的距离矩阵 $D \in R^{n \times n}$, 其第 i 行 j 列的元素为样本 x_i 到 x_j 的距离。
- **目标**: 获得这 n 个样本 $Z = [z_1, z_2, \dots, z_n] \in R^{d \times n}$ 在 $d(d < m)$ 维空间中的表示。
- **准则**: 假定降维后的样本仍保持两两之间的距离。

$$\text{dist}_{ij}^2 = \|z_i\|_2^2 + \|z_j\|_2^2 - 2z_i^T z_j$$

$$\text{令 } b_{ij} = z_i^T z_j$$

$$\text{dist}_{ij}^2 = b_{ii} + b_{jj} - 2b_{ij}$$

$$b_{ij} = 1/2(b_{ii} + b_{jj} - \text{dist}_{ij}^2)$$

$$D \in R^{n \times n}$$

多维缩放(Multiple Dimensional Scaling, MDS)

- 如果给定一个距离矩阵 D （该矩阵中实际上保存着原始样本数据点之间的距离, 是一个实对称矩阵），我们能不能得到矩阵 B ，进而得到 Z ？
- 实际上如果不加上其它条件约束的话，这样是做不到的，因为整体平移和旋转样点是不改变样点间距离的，可以得到无数个满足要求的矩阵 B 。
- 不失一般性，令降维后的样本是**零均值化（零中心化）**的，即 $\sum_{i=1}^n z_i = 0 \in R^d$
- 矩阵 B 的行和列之和均为零，即 $\sum_{i=1}^n b_{ij} = \sum_{j=1}^n b_{ij} = 0$

$$\sum_{i=1}^n b_{ij} = (z_1^T z_j + \dots + z_n^T z_j) = (z_1^T + \dots + z_n^T) z_j = 0^T z_j = 0$$

多维缩放(Multiple Dimensional Scaling, MDS)

— 通过计算可以求得：

$$dist_{ij}^2 = \|z_i\|_2^2 + \|z_j\|_2^2 - 2z_i^T z_j$$

$$\sum_{i=1}^n b_{ij} = \sum_{j=1}^n b_{ij} = 0$$

$$\sum_{i=1}^n dist_{ij}^2 = \sum_{i=1}^n (b_{ii} + b_{jj} - 2b_{ij}) = tr(B) + nb_{jj}$$

$$\sum_{j=1}^n dist_{ij}^2 = \sum_{j=1}^n (b_{ii} + b_{jj} - 2b_{ij}) = tr(B) + nb_{ii}$$

$$\sum_{i=1}^n \sum_{j=1}^n dist_{ij}^2 = 2ntr(B)$$

其中 $tr(.)$ 表示矩阵的迹(trace), $tr(B) = \sum \|z_i\|_2^2$

同时 $tr(B)$ 可计算,

$$tr(B) = \frac{1}{2n} \sum_{i=1}^n \sum_{j=1}^n dist_{ij}^2$$

多维缩放(Multiple Dimensional Scaling, MDS)

— 引入如下新符号，并记为：

$$dist_{i.}^2 = \frac{1}{n} \sum_{j=1}^n dist_{ij}^2$$

$$\sum_{j=1}^n dist_{ij}^2 = tr(B) + nb_{ii}$$

$$dist_{.j}^2 = \frac{1}{n} \sum_{i=1}^n dist_{ij}^2$$

$$\sum_{i=1}^n dist_{ij}^2 = tr(B) + nb_{jj}$$

$$dist_{..}^2 = \frac{1}{n^2} \sum_{i=1}^n \sum_{j=1}^n dist_{ij}^2$$

$$\sum_{i=1}^n \sum_{j=1}^n dist_{ij}^2 = 2ntr(B)$$

可以算出，

$$tr(B) = \frac{n}{2} dist_{..}^2, \quad b_{ii} = dist_{i.}^2 - \frac{1}{2} dist_{..}^2, \quad b_{jj} = dist_{.j}^2 - \frac{1}{2} dist_{..}^2$$

多维缩放(Multiple Dimensional Scaling, MDS)

- 可计算出 $b_{ij} = \frac{1}{2}(dist_{i.}^2 + dist_{.j}^2 - dist_{..}^2 - dist_{ij}^2)$
- 在获得矩阵 B 之后, 则可对矩阵 B 进行特征值分解。
- 注意到: $B = Z^T Z \in R^{n \times n}$, 且 B 是对称矩阵, 于是有:

$$B = U \Lambda U^T \quad Z = \Lambda_d^{\frac{1}{2}} U_d^T \in R^{d \times n}$$

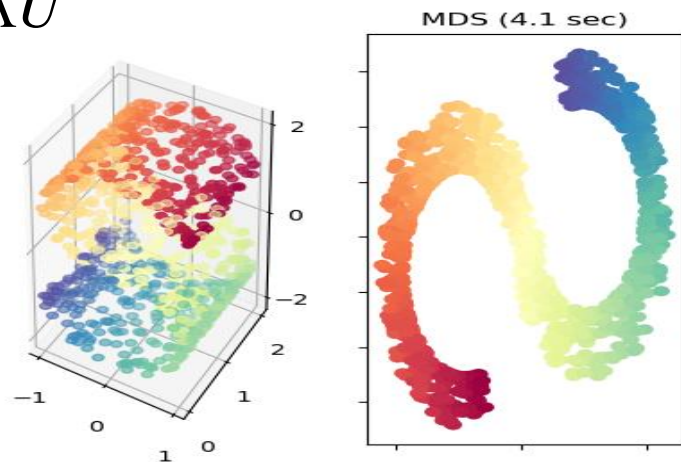
$$\Lambda_d^{\frac{1}{2}} = \text{diag}(\sqrt{\lambda_1}, \sqrt{\lambda_2}, \dots, \sqrt{\lambda_d}) \in R^{d \times d} \quad U_d = [u_1, u_2, \dots, u_d] \in R^{n \times d}$$

其中, $\Lambda_d^{\frac{1}{2}}$ 表示由矩阵 B 的前 d 个最大的特征值开根号后对应的对角矩阵; U_d 由前 d 个最大的特征值对应的特征向量组成。

多维缩放(Multiple Dimensional Scaling, MDS)

MDS算法步骤

- 1 给定距离矩阵 D ;
- 2 计算 $dist^2_i, dist^2_j, dist^2_{..}$;
- 3 计算矩阵 B : $b_{ij} = -0.5(dist^2_{ij} - dist^2_i - dist^2_j + dist^2_{..})$;
- 4 对矩阵 B 进行特征值分解: $B = U\Lambda U^T$
- 5 输出: $Z = \Lambda_d^{\frac{1}{2}} U_d^T \in R^{d \times n}$



MDS举例 - 飞行距离可视化

- 18个城市的飞行距离Airline distances (km)

TABLE 13.2. *Airline distances (km) between 18 cities. Source: Atlas of the World, Revised 6th Edition, National Geographic Society, 1995, p. 131.*

	Beijing	Cape Town	Hong Kong	Honolulu	London	Melbourne
Cape Town	12947					
Hong Kong	1972	11867				
Honolulu	8171	18562	8945			
London	8160	9635	9646	11653		
Melbourne	9093	10338	7392	8862	16902	
Mexico	12478	13703	14155	6098	8947	13557
Montreal	10490	12744	12462	7915	5240	16730
Moscow	5809	10101	7158	11342	2506	14418
New Delhi	3788	9284	3770	11930	6724	10192
New York	11012	12551	12984	7996	5586	16671
Paris	8236	9307	9650	11988	341	16793
Rio de Janeiro	17325	6075	17710	13343	9254	13227
Rome	8144	8417	9300	12936	1434	15987
San Francisco	9524	16487	11121	3857	8640	12644
Singapore	4465	9671	2575	10824	10860	6050
Stockholm	6725	10334	8243	11059	1436	15593
Tokyo	2104	14737	2893	6208	9585	8159

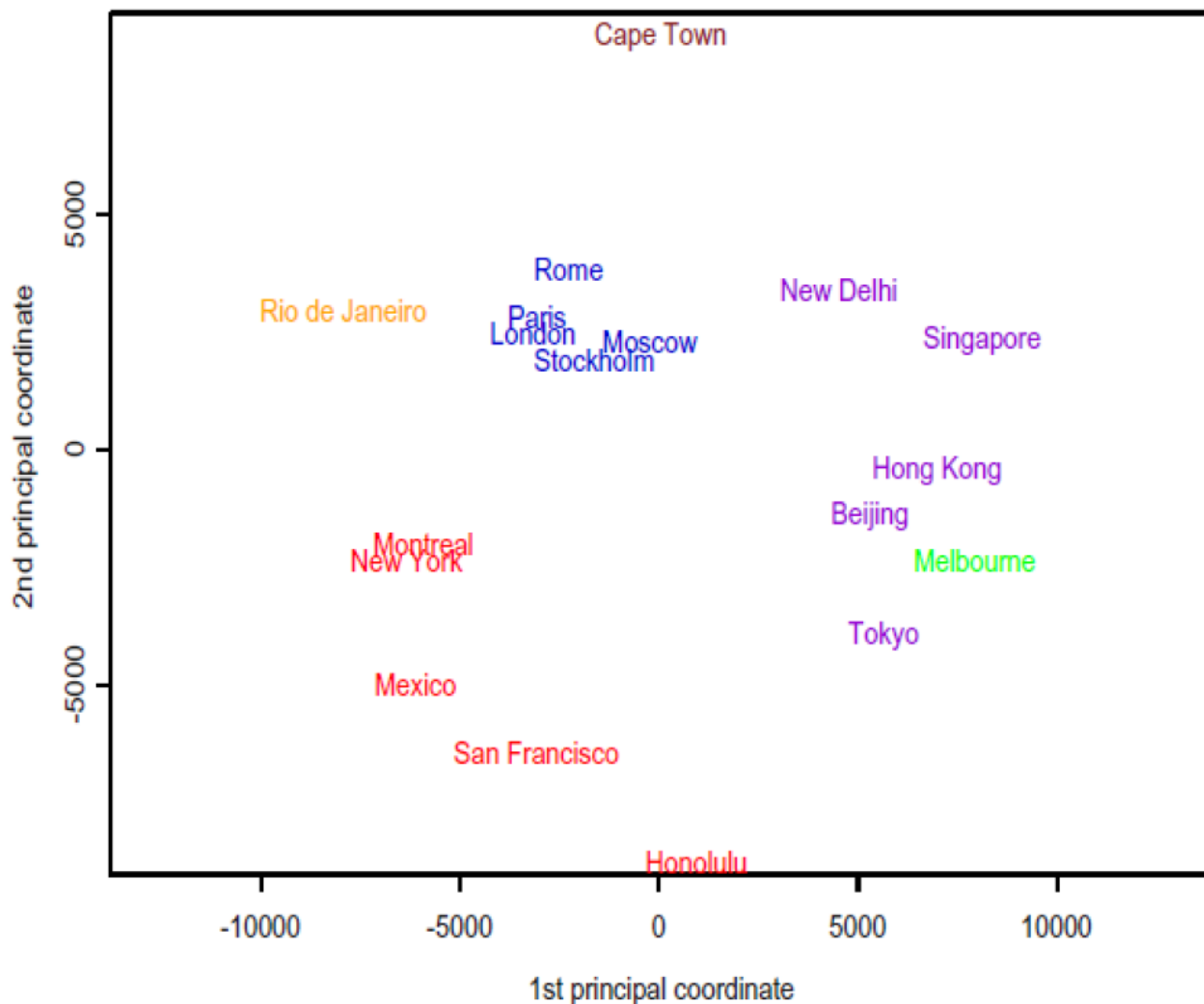
MDS举例 - 飞行距离可视化

TABLE 13.6. *Eigenvalues of B and the eigenvectors corresponding to the first three largest eigenvalues (in red) for the airline distances example.*

	Eigenvalues	Eigenvectors		
1	471582511	0.245	-0.072	0.183
2	316824787	0.003	0.502	-0.347
3	253943687	0.323	-0.017	0.103
4	-98466163	0.044	-0.487	-0.080
5	-74912121	-0.145	0.144	0.205
6	-47505097	0.366	-0.128	-0.569
7	31736348	-0.281	-0.275	-0.174
8	-7508328	-0.272	-0.115	0.094
9	4338497	-0.010	0.134	0.202
10	1747583	0.209	0.195	0.110
11	-1498641	-0.292	-0.117	0.061
12	145113	-0.141	0.163	0.196
13	-102966	-0.364	0.172	-0.473
14	60477	-0.104	0.220	0.163
15	-6334	-0.140	-0.356	-0.009
16	-1362	0.375	0.139	-0.054
17	100	-0.074	0.112	0.215
18	0	0.260	-0.214	0.173

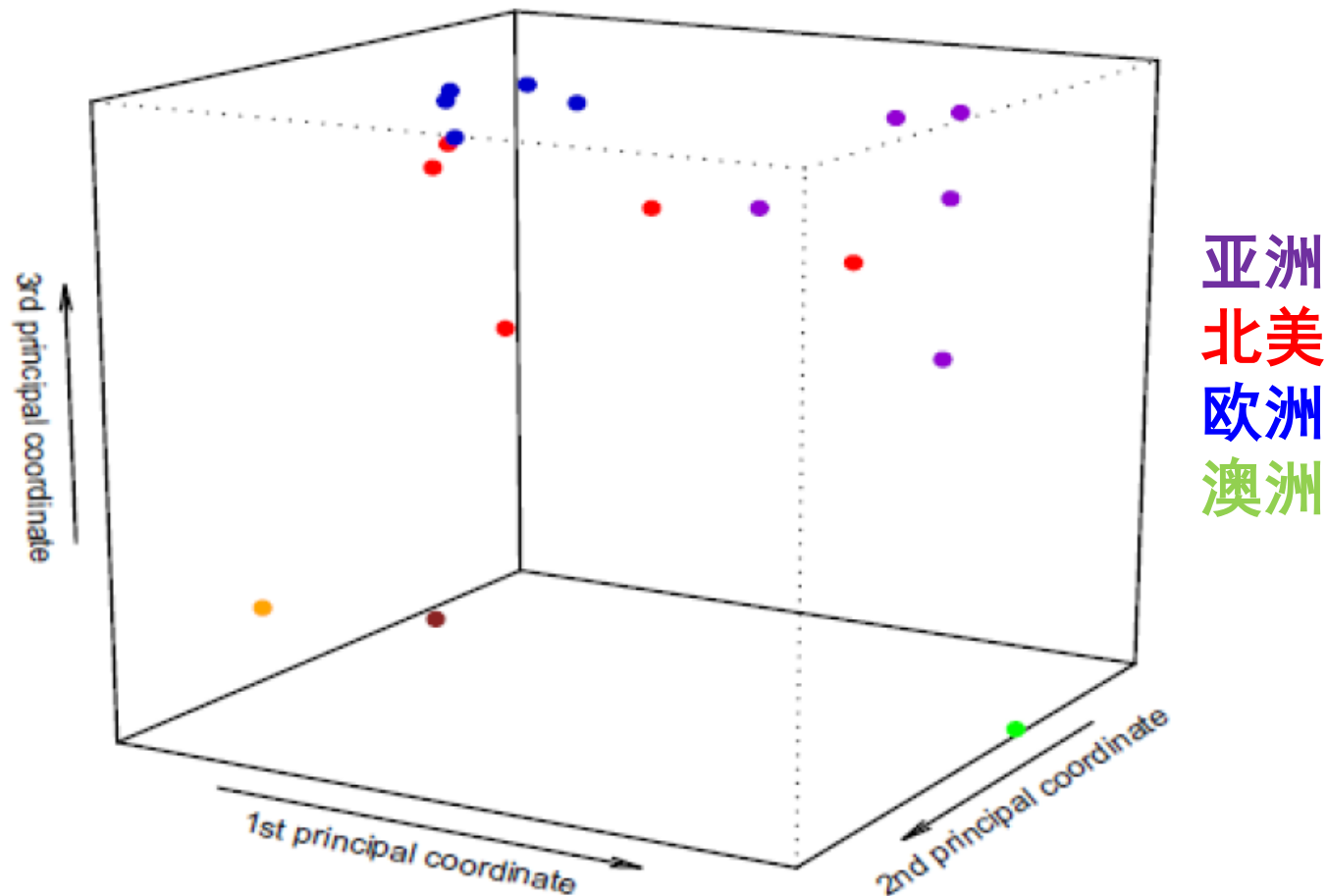
- Airline distance is non-Euclidean
- Take the first 3 largest eigenvalues (inspection of scree plot)

MDS举例 - 飞行距离可视化



二维表示（2个特征向量）

MDS举例



三维表示（3个特征向量）

内容提要

- 引言
- 主成分分析
- 多维缩放
- 流形学习方法
- 距离度量学习

流形学习-流形定义

- **定义**：流形上的每一个点的开邻域，与欧氏空间的开集同胚。
- **几何**：流形是一块一块欧氏空间拼装而成的弯曲空间。
- **直观**：流形是欧氏空间的一种推广，是在低维空间来表达高维空间所难以表达的空间结构。
- **在数学上**：流形用于描述一个几何形体，它在局部具有欧氏空间的性质。即可以应用欧氏距离来描述局部区域，但在全局欧氏距离不成立。

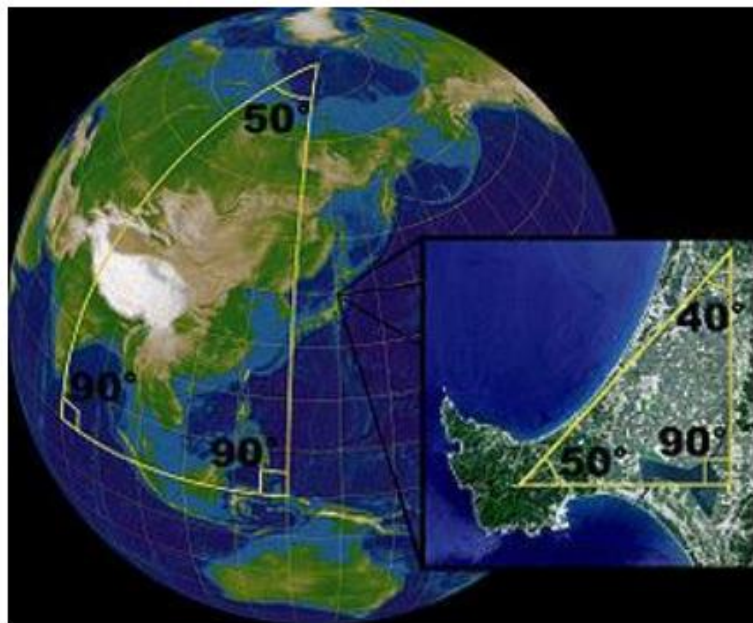
The **surface of a sphere** is a two-dimensional manifold as it can be represented by a collection of two-dimensional maps.

流形学习-流形定义

- **定义**：流形上的每一个点的开邻域，与欧氏空间的开集同胚。
- **几何**：流形是一块一块欧氏空间拼装而成的弯曲空间。
- **直观**：流形是欧氏空间的一种推广，是在低维空间来表达高维空间所难以表达的空间结构。

— **在数学上**：流形空间的性质。即可局部欧氏距离不成立

The **surface** can be represented

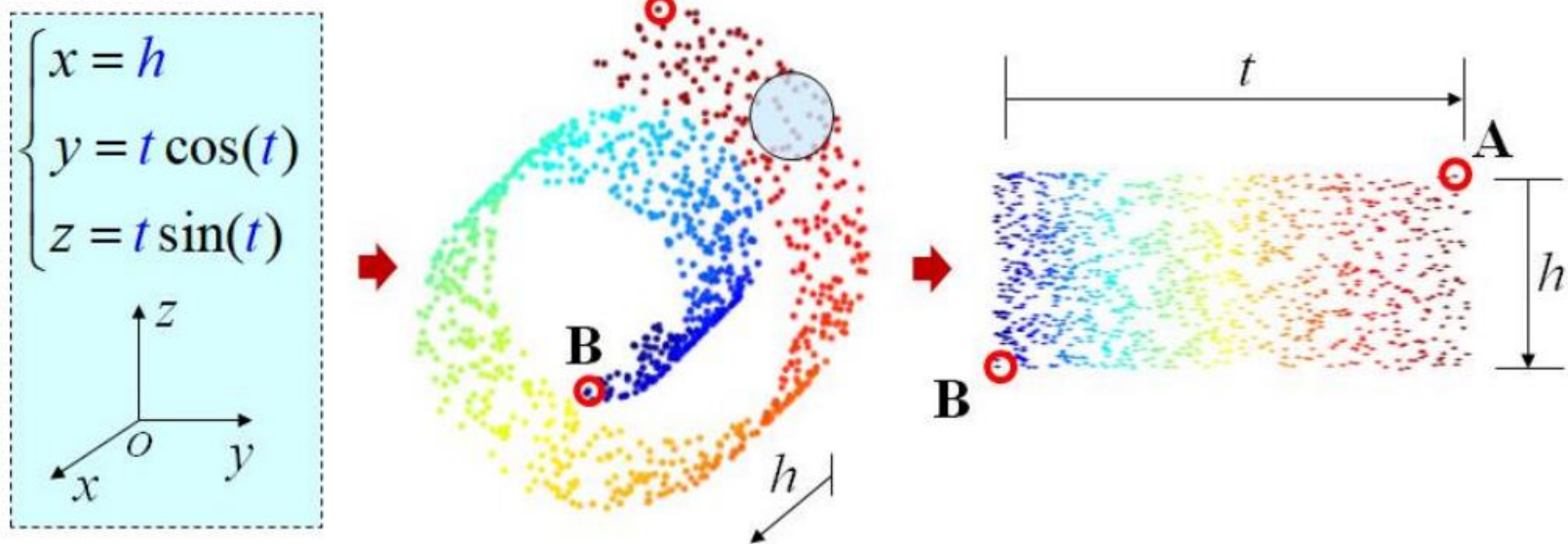


在局部具有欧氏空间的性质，但在全局

ifold as it is not Euclidean maps.

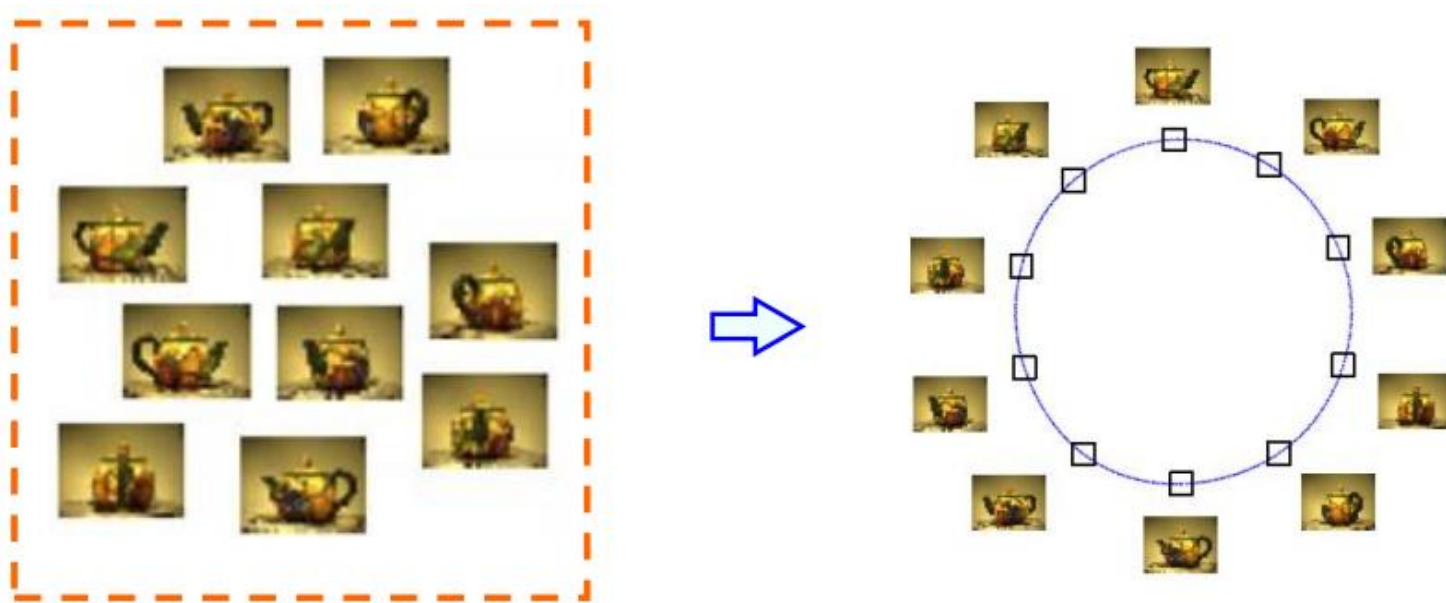
流形学习-流形示例

– **Swiss roll surface:** Swiss roll surface is a **2D** manifold.

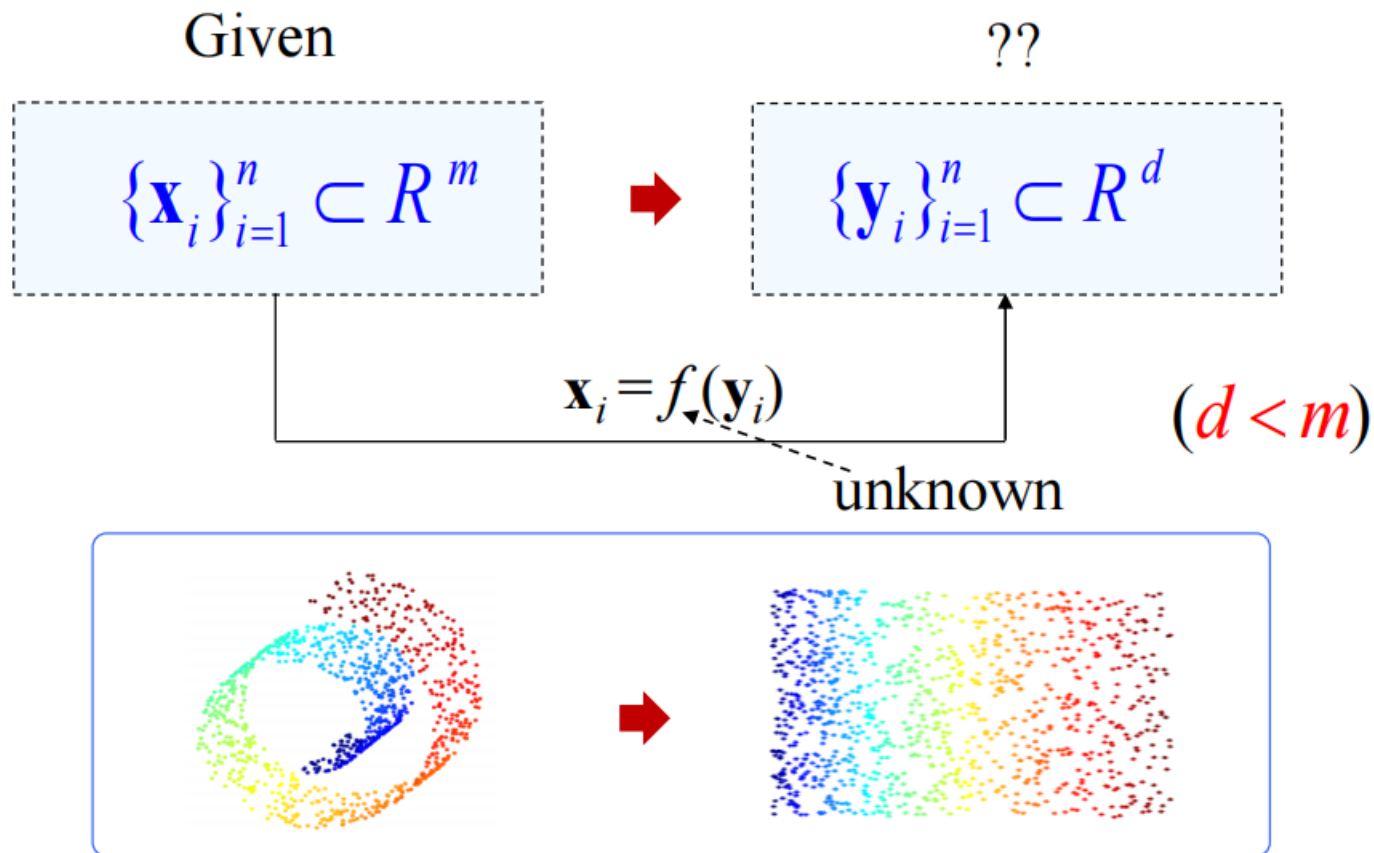


流形学习-流形示例

— **一个直观的例子**：400张360度全角度拍摄的图片会排列成一个圆。



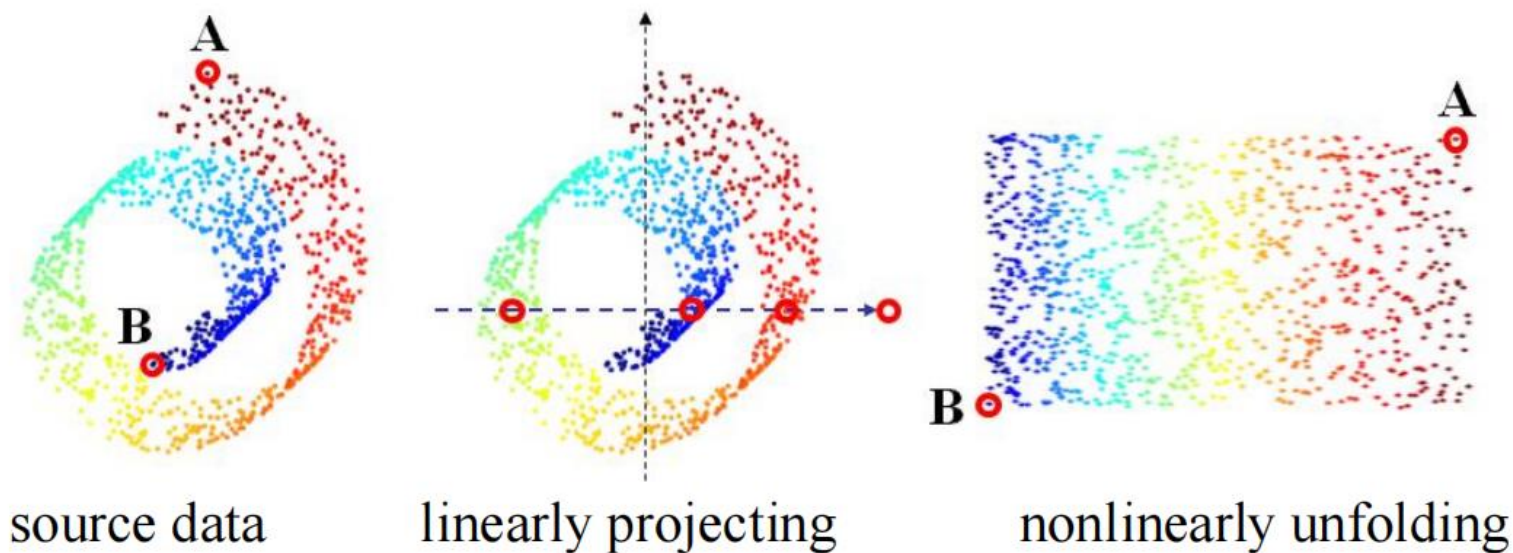
流形学习-数学描述



流形学习

● 非线性维数缩减

— **Problem formulation** in machine learning: in view of dimensionality reduction



通过线性投影将高维数据降到低维将难以展开非线性结构！

流形学习

● 一些假定

— 流形光滑 (smooth manifold) :

$$(f : C \subset R^d \rightarrow R^m)$$

— 密集采样 (densely sampling)

— 不可自相交 (no self-intersections)



● 经典算法

— LLE, Isomap, Laplacian Eigenmap, HLLE, MVU, LTSA, LSE, etc.

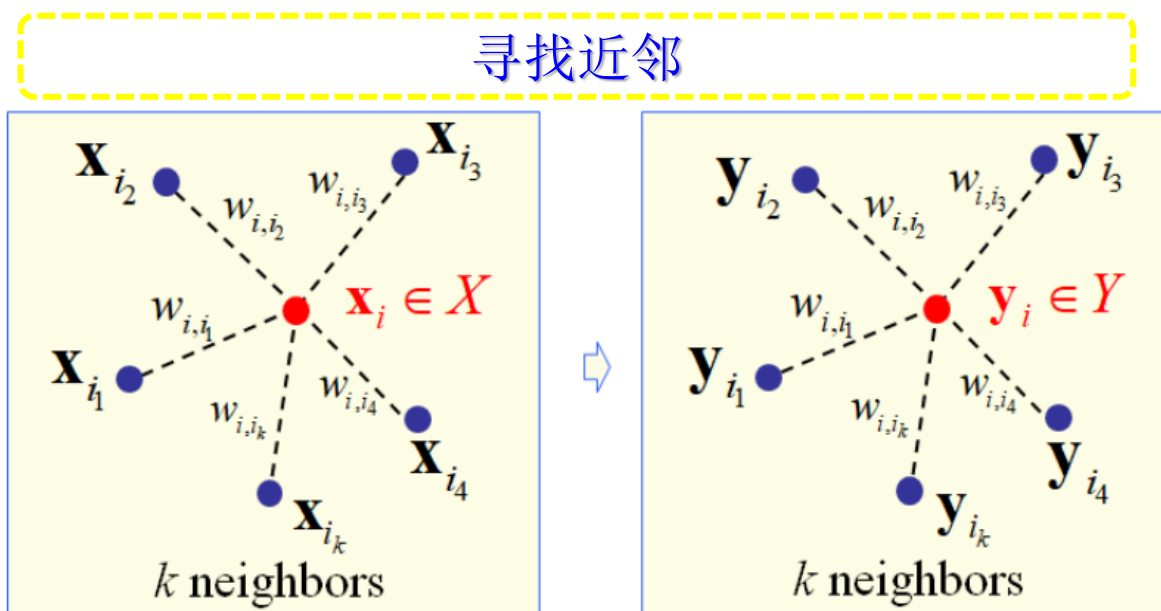
流形学习-LLE(Locally linear embedding)

— 基本思想：给定数据集，通过最近邻等方式构造一个数据图(data graph)。然后在每一个局部区域，高维空间中的样本线性重构关系在低维空间中均得以保持。

Sam Roweis & Lawrence Saul. *Nonlinear dimensionality reduction by locally linear embedding*. **Science**, v.290 [no.5500](#), Dec.22, 2000. pp.2323--2326.

流形学习-LLE(Locally linear embedding)

- 基本思想：给定数据集，通过最近邻等方式构造一个数据图(data graph)。然后在每一个局部区域，高维空间中的样本线性重构关系在低维空间中均得以保持。
- 基本步骤：寻找近邻，线性重构，低维嵌入



$$x_i \approx \sum_{j=1}^k w_{i,i_j} x_{i_j}$$

$$y_i \approx \sum_{j=1}^k w_{i,i_j} y_{i_j}$$

流形学习-LLE(Locally linear embedding)

— 最优线性表示系数

$$\min_{w_i} \sum_{i=1}^m \|x_i - \sum_{j=1}^k w_{i,j} x_{i_j}\|_2^2, \sum_{j=1}^k w_{i,j} = 1$$

通过拉格朗日乘子法，可得如下有关线性表示系数的解：

$$w_i = \frac{Z_i^{-1} e}{e^T Z_i^{-1} e}$$

流形学习-LLE(Locally linear embedding)

— 最优线性表示系数

$$\min_{w_i} \sum_{i=1}^m \|x_i - \sum_{j=1}^k w_{i,j} x_{i_j}\|_2^2, \sum_{j=1}^k w_{i,j} = 1$$

通过拉格朗日乘子法，可得如下有关线性表示系数的解：

$$w_i = \frac{Z_i^{-1} e}{e^T Z_i^{-1} e}$$

$$w_i = [w_{i,i_1}, w_{i,i_2}, \dots, w_{i,i_k}]^T \in R^k$$

$$X_i = [x_{i_1}, x_{i_2}, \dots, x_{i_k}]^T \in R^{m \times k}$$

$$N_i = [x_{i_1}, x_{i_2}, \dots, x_{i_k}]^T \in R^{m \times k}$$

$$e = [1, 1, \dots, 1]^T \in R^k$$

$$Z_i = (X_i - N_i)^T (X_i - N_i)$$

流形学习-LLE(Locally linear embedding)

— 推导过程

$$\begin{aligned} J(\mathbf{w}) &= \min_{\mathbf{w}_i} \sum_{i=1}^m \left\| \mathbf{x}_i - \sum_{j=1}^k w_{i,j} \mathbf{x}_{i_j} \right\|_2^2 \\ &= \sum_{i=1}^m \left\| \sum_{j=1}^k w_{i,j} \mathbf{x}_i - \sum_{j=1}^k w_{i,j} \mathbf{x}_{i_j} \right\|_2^2 = \sum_{i=1}^m \left\| \sum_{j=1}^k w_{i,j} (\mathbf{x}_i - \mathbf{x}_{i_j}) \right\|_2^2 \\ &= \sum_{i=1}^m \left\| (\mathbf{X}_i - \mathbf{N}_i) \mathbf{w}_i \right\|_2^2 \\ &= \sum_{i=1}^m \mathbf{w}_i^T \underbrace{(\mathbf{X}_i - \mathbf{N}_i)^T (\mathbf{X}_i - \mathbf{N}_i)}_{\mathbf{Z}_i} \mathbf{w}_i \end{aligned}$$

流形学习-LLE(Locally linear embedding)

— 推导过程

$$J(w) = \sum_{i=1}^m w_i^T Z_i w_i$$

对于限制项，可以写为：

$$\sum_{j=1}^k w_{i,i_j} = w_i^T e = 1$$

流形学习-LLE(Locally linear embedding)

— 拉格朗日乘子法

$$L(w) = \sum_{i=1}^m w_i^T Z_i w_i + \lambda(w_i^T e - 1)$$

对 w_i 求导可得:

$$2Z_i w_i + \lambda e = 0$$

$$w_i = \frac{1}{2} \lambda Z_i^{-1} e$$

对 w_i 归一化:

$$w_i = \frac{Z_i^{-1} e}{e^T Z_i^{-1} e}$$

流形学习-LLE(Locally linear embedding)

- 全局嵌入：利用在原始空间中获得的局部线性重构关系，在低维空间中重构对应的样本点：

$$y_i \approx \sum_{j=1}^k w_{i,i_j} y_{i_j} \quad i = 1, 2, \dots, n$$

- 考虑所有新样本点的重构误差，得到全局嵌入的目标函数：

$$\sum_{i=1}^n \|y_i - \sum_{j=1}^k w_{i,i_j} y_{i_j}\|_2^2 = \text{tr}(Y(I - W)(I - W)^T Y^T)$$

$$Y = [y_1, y_2, \dots, y_n]^T \in R^{d \times n}$$

W 为权重矩阵，其第 i 行记录对应样本点 x_i 的 k 个权重，只有在对应的邻居位置 i_1, i_2, \dots, i_k 处才有值，其余全为零。

流形学习-LLE(Locally linear embedding)

$$\sum_{i=1}^n \|y_i - \sum_{j=1}^k w_{i,j} y_{i_j}\|_2^2$$

i近邻点 $W_{ij} = w_{i,j}$ 否则 $W_{ij} = 0$

$$\sum_{j=1}^n W_{i,j} y_{i_j} = \sum_{j=1}^k w_{i,j} y_{i_j} = Y W_i$$

$$\sum_{i=1}^n \|Y(I_i - W_i)\|_2^2$$

$$= \text{tr}(Y(I - W)(I - W)^T Y^T)$$

$$\because A = [a_1, a_2, \dots, a_n], \therefore \sum_i (a_i)^2 = \sum_i a_i^T a_i = \text{tr}(A A^T)$$

流形学习-LLE(Locally linear embedding)

- LLE (Locally linear embedding)

- 全局嵌入学习模型:

$$\min_Y \text{tr} (Y(I - W)(I - W)^T Y^T) \text{ s.t. } YY^T = I$$

- 求解: 通过求矩阵 $(I - W)(I - W)^T$ 的特征值分解来得到。

- 取出该矩阵最小的 $k + 1$ 个特征值对应的特征向量;
 - 丢弃特征值零对应的分量全相等的特征向量;
 - 即采用第2至第 $d + 1$ 个最小的特征值对应的特征向量组成样本的新的坐标。

[1] S. Roweis and L. Saul, “Nonlinear dimensionality reduction by locally linear embedding,” Science, vol. 290, pp. 2323–2326, 2000.

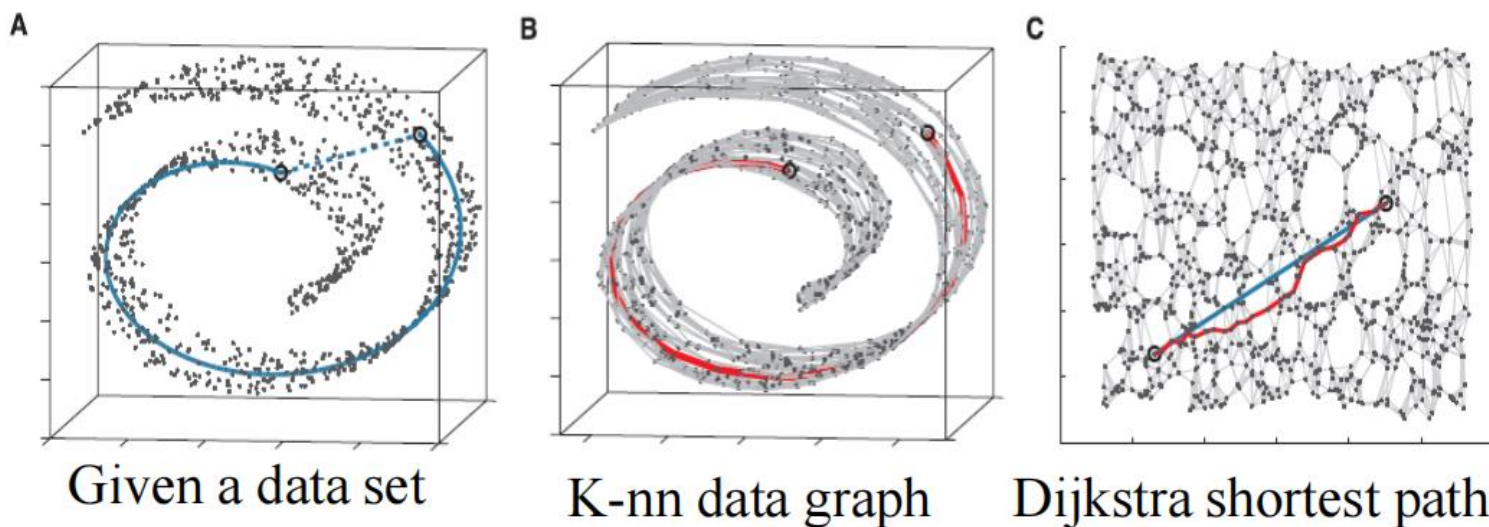
流形学习

LLE算法步骤

- 1 **Input**: 给出数据 $X = [x_1, x_2, \dots, x_n] \in R^{m \times n}$, 近邻参数 k 以及低维空间 d ;
- 2 for $i = 1, 2, \dots, n$
- 3 确定的 k 个近邻;
- 4 对 x_i 进行**线性最优表示**, 获取近邻重构权重;
- 5 end for
- 6 构造权重矩阵 W ;
- 7 求解: $\min_Y \text{tr}(Y(I - W)(I - W)^T Y^T) \text{ s.t. } YY^T = I$
- 8 采用第2至第 $d + 1$ 个**最小的特征值对应的特征向量**组成新坐标;
- 9 **输出结果**。

流形学习-Isomap(isometric feature mapping)

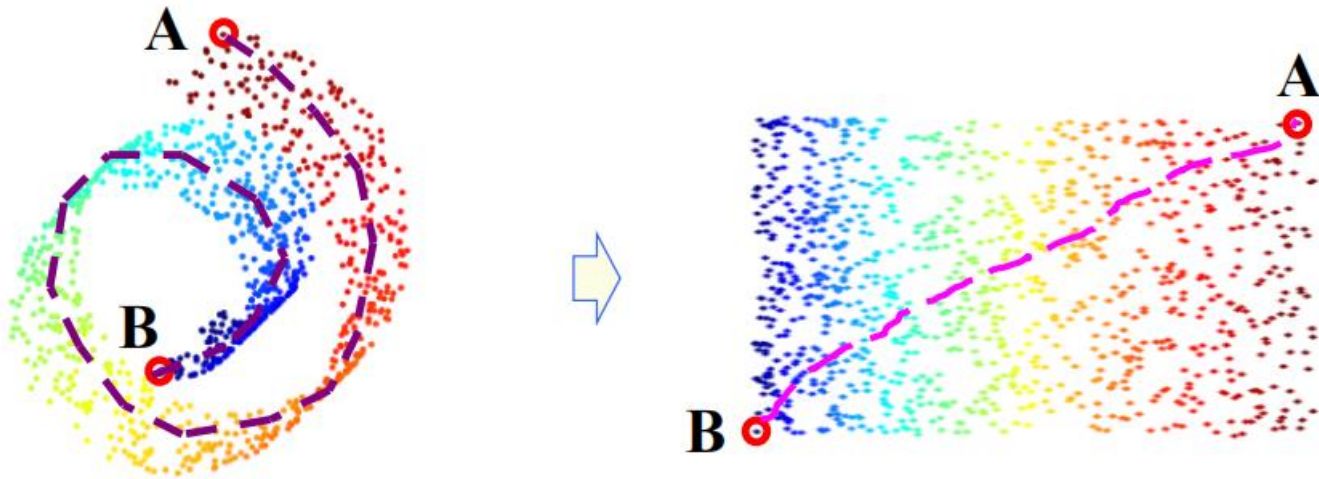
— 基本思想：给定数据集，通过最近邻等方式构造一个**数据图 (data graph)**。然后，计算任意两个点之间的**最短路径**（即**测地距离**）。对于所有的任意两个点对，期望在低维空间中保持其测地距离。



几乎所有的流形学习方法都需要首先构建一个关于数据的图

流形学习-Isomap(isometric feature mapping)

— 基本思想：给定数据集，通过最近邻等方式构造一个数据图(data graph)。然后，计算任意两个点之间的最短路径（即测地距离）。对于所有的任意两个点对，期望在低维空间中保持其测地距离。



流形学习-Isomap(isometric feature mapping)

Isomap算法步骤

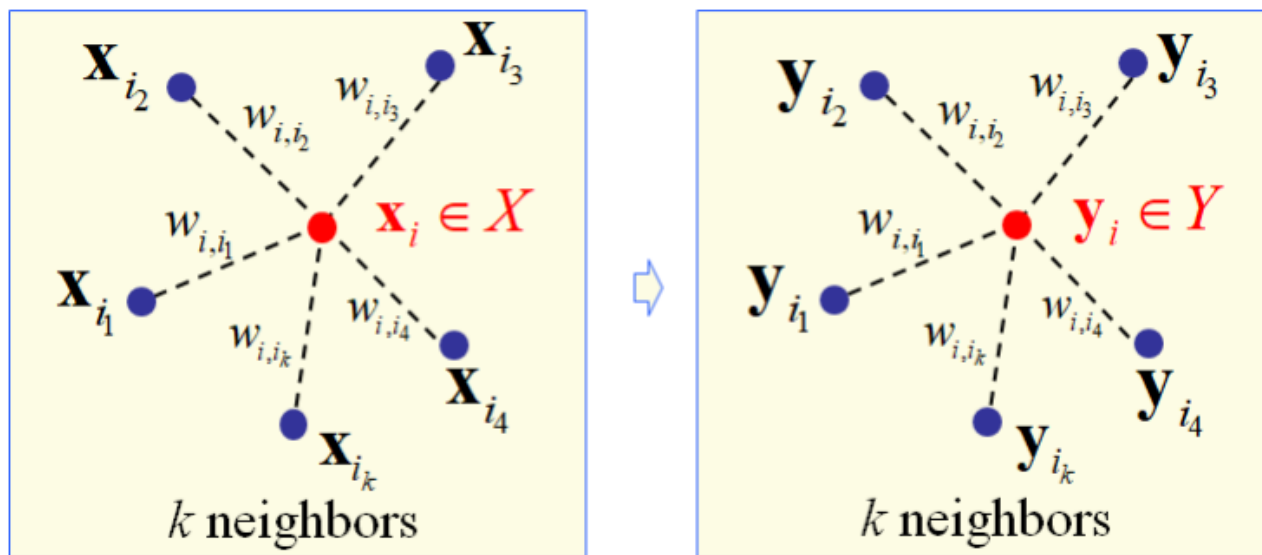
- 1 **Input**: 给出数据 $X = [x_1, x_2, \dots, x_n] \in R^{m \times n}$, 近邻参数 k 以及低维空间 d ;
- 2 for $i = 1, 2, \dots, n$
- 3 确定 x_i 的 k 个近邻;
- 4 x_i 与 k 个近邻点之间的距离设定为**欧氏距离**, 与非近邻点的距离设置为**无穷大**;
- 5 end for
- 6 调用**最短路径法**计算任意两样本点 x_i 与 x_j 之间的距离 d_{ij} 。由此可构造距离矩阵。
- 7 调用**MDS算法** (见MDS算法) ;
- 8 **MDS算法**的计算结果作为低维嵌入结果

流形学习-LE(Laplacian Eigenmapping)

— 基本思想：给定数据集，通过最近邻等方式构造一个数据图 (data graph)。然后，在每一个局部区域，计算点与点之间的亲合度（相似度），期望点对亲合度在低维空间中也得到保持。

流形学习-LE(Laplacian Eigenmapping)

- 基本思想：给定数据集，通过最近邻等方式构造一个数据图(data graph)。然后，在每一个局部区域，计算点与点之间的亲合度（相似度），期望点对亲合度在低维空间中也得到保持。



流形学习-LE(Laplacian Eigenmapping)

— 如何计算点对亲合度？

流形学习-LE(Laplacian Eigenmapping)

— 如何计算点对亲合度？

$$w_{i,i_j} = \exp\left(-\frac{\|x_i - x_{i_j}\|_2^2}{2\sigma^2}\right)$$

流形学习-LE(Laplacian Eigenmapping)

— 如何计算点对亲合度？

$$w_{i,i_j} = \exp\left(-\frac{\|x_i - x_{i_j}\|_2^2}{2\sigma^2}\right)$$

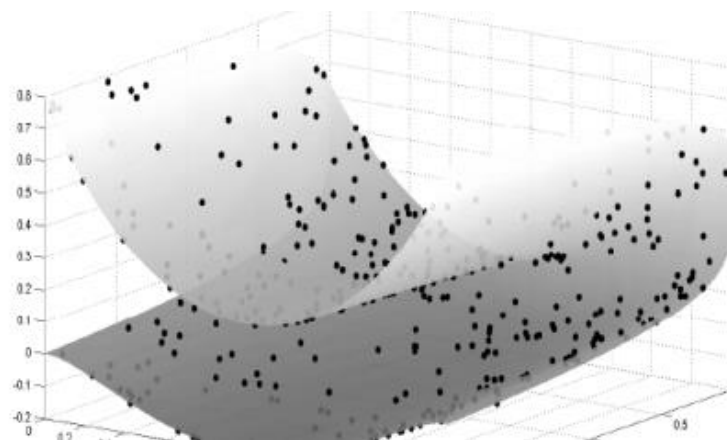
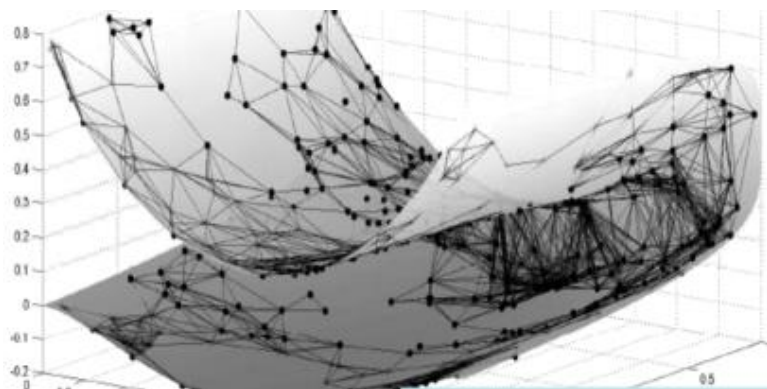
— 图构造 $G(V, E)$ （第六章）

流形学习-LE(Laplacian Eigenmapping)

— 如何计算点对亲合度？

$$w_{i,i_j} = \exp\left(-\frac{\|x_i - x_{i_j}\|_2^2}{2\sigma^2}\right)$$

— 图构造 $G(V, E)$ (第六章)



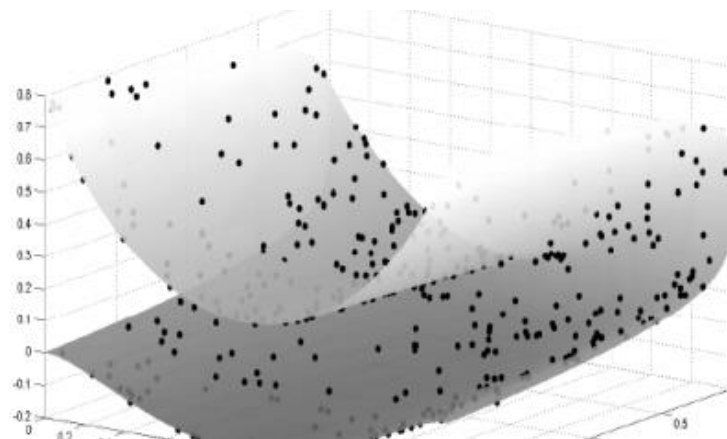
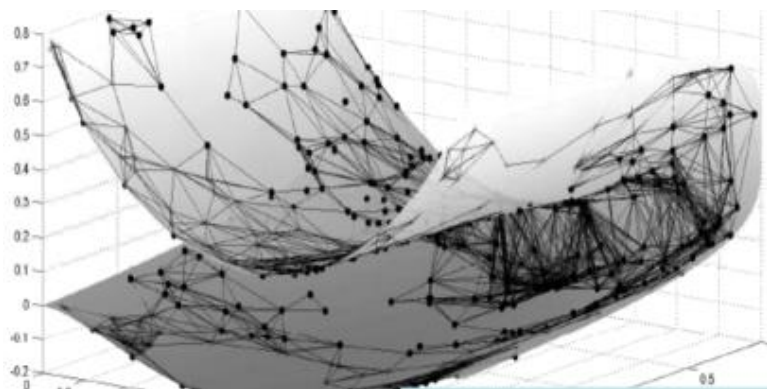
流形学习

— 如何计算点对亲合度？

$$w_{i,i_j} = \exp\left(-\frac{\|x_i - x_{i_j}\|_2^2}{2\sigma^2}\right)$$

— 图构造 $G(V, E)$ (第六章)

$$\begin{pmatrix} 0 & w_{12} & w_{12} & \bullet & \bullet & \bullet & w_{1n} \\ w_{21} & 0 & w_{22} & \bullet & \bullet & \bullet & w_{2n} \\ w_{31} & w_{32} & 0 & \bullet & \bullet & \bullet & w_{3n} \\ \bullet & \bullet & \bullet & \bullet & \bullet & \bullet & \bullet \\ \bullet & \bullet & \bullet & \bullet & \bullet & \bullet & \bullet \\ \bullet & \bullet & \bullet & \bullet & \bullet & \bullet & \bullet \\ w_{n1} & w_{n2} & w_{n4} & \bullet & \bullet & \bullet & 0 \end{pmatrix}$$



流形学习

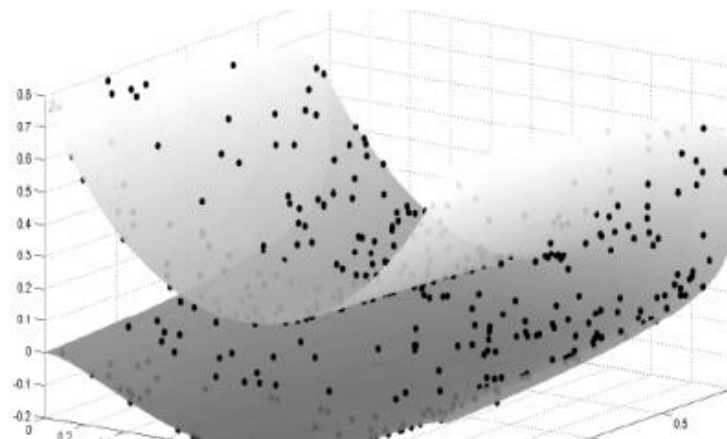
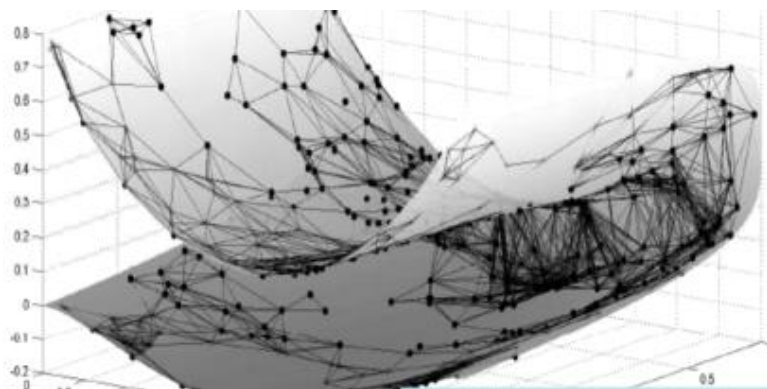
每行 k 个元素为零

— 如何计算点对亲合度？

$$w_{i,i_j} = \exp\left(-\frac{\|x_i - x_{i_j}\|_2^2}{2\sigma^2}\right)$$

— 图构造 $G(V, E)$ (第六章)

$$\begin{pmatrix} 0 & w_{12} & w_{12} & \bullet & \bullet & \bullet & w_{1n} \\ w_{21} & 0 & w_{22} & \bullet & \bullet & \bullet & w_{2n} \\ w_{31} & w_{32} & 0 & \bullet & \bullet & \bullet & w_{3n} \\ \bullet & \bullet & \bullet & \bullet & \bullet & \bullet & \bullet \\ \bullet & \bullet & \bullet & \bullet & \bullet & \bullet & \bullet \\ \bullet & \bullet & \bullet & \bullet & \bullet & \bullet & \bullet \\ w_{n1} & w_{n2} & w_{n4} & \bullet & \bullet & \bullet & 0 \end{pmatrix}$$



流形学习-LE(Laplacian Eigenmapping)

— 如何在低维空间保持亲合度？

构造如下目标函数：

流形学习-LE(Laplacian Eigenmapping)

— 如何在低维空间保持亲合度？

构造如下目标函数：

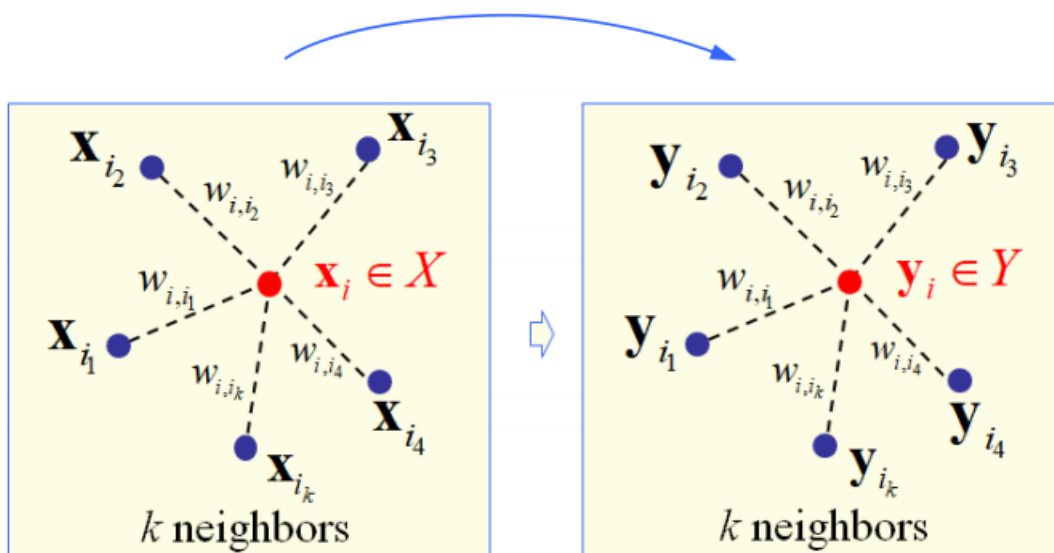
$$E(Y) = \sum_{i,j} w_{i,i_j} \|y_i - y_{i_j}\|_2^2$$

流形学习-LE(Laplacian Eigenmapping)

— 如何在低维空间保持亲合度？

构造如下目标函数：

$$E(Y) = \sum_{i,j} w_{i,j} \|y_i - y_{i_j}\|_2^2$$



流形学习-LE(Laplacian Eigenmapping)

— 考虑目标函数

- 对任意向量 $f = [f_1, f_2, \dots, f_n]^T$ 有如下结论:

令

$$D = \begin{pmatrix} d_1 & & & \\ & d_2 & & \\ & & \ddots & \\ & & & d_n \end{pmatrix}$$

$$d_i = \sum_{j=1}^k w_{i,i_j}$$

$$i = 1, 2, \dots, n$$

流形学习-LE(Laplacian Eigenmapping)

— 考虑目标函数

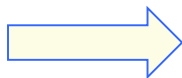
- 对任意向量 $\mathbf{f} = [f_1, f_2, \dots, f_n]^T$ 有如下结论：

亲和度矩阵

$$D = \begin{pmatrix} d_1 & & & \\ & d_2 & & \\ & & \ddots & \\ & & & d_n \end{pmatrix}$$

$$d_i = \sum_{j=1}^k w_{i,j}$$

$$i = 1, 2, \dots, n$$



$$\mathbf{f}^T (D - \boxed{W}) \mathbf{f} = \sum_{i=1}^n d_i f_i^2 - \sum_{i,j=1}^n f_i f_j w_{ij}$$

$$= \frac{1}{2} \left(\sum_{i=1}^n d_i f_i^2 - 2 \sum_{i,j=1}^n f_i f_j w_{ij} + \sum_{j=1}^n d_j f_j^2 \right)$$

$$= \frac{1}{2} \left(\sum_{i,j=1}^n w_{ij} (f_i - f_j)^2 \right)$$

$$\geq 0$$

流形学习-LE(Laplacian Eigenmapping)

— 考虑目标函数

$$Y = [y_1, y_2, \dots, y_n] = \begin{pmatrix} y_{11} & y_{12} & \bullet & \bullet & \bullet & y_{1n} \\ y_{21} & y_{22} & \bullet & \bullet & \bullet & y_{2n} \\ \bullet & \bullet & \bullet & \bullet & \bullet & \bullet \\ \bullet & \bullet & \bullet & \bullet & \bullet & \bullet \\ \bullet & \bullet & \bullet & \bullet & \bullet & \bullet \\ y_{d1} & y_{d2} & \bullet & \bullet & \bullet & y_{dn} \end{pmatrix} = \begin{pmatrix} f_1^T \\ f_2^T \\ \bullet \\ \bullet \\ \bullet \\ f_n^T \end{pmatrix} \in R^{d \times n}$$

其中, $f_i = [y_{i1}, y_{i2}, \dots, y_{in}]^T \in R, i = 1, 2, \dots, d$

$$E(Y) = \sum_{i,j} w_{i,i_j} \|y_i - y_{i_j}\|_2^2 = \sum_{i,j} w_{i,i_j} [(\mathbf{y}_{1i} - \mathbf{y}_{1i_j})^2 + \dots + (\mathbf{y}_{di} - \mathbf{y}_{di_j})^2] =$$
$$2(f_1^T (D - W)f_1 + f_2^T (D - W)f_2 + \dots + f_d^T (D - W)f_d) = \mathbf{2} * \mathbf{tr}(Y(D - W)Y^T)$$

流形学习-LE(Laplacian Eigenmapping)

— 学习模型

$$\min E(Y) = \text{tr}(Y(D - W)Y^T), \text{ s.t. } YY^T = \mathbf{1}$$

— 令 $L = D - W$, L 有一个特征值为0, 对应的特征向量全为1:

$$Le = (D - W)e = (D - W)(1, 1, \dots, 1)^T =$$



$$Le = 0e$$

$$\begin{pmatrix} d_1 - \sum_{j=1}^k w_{1,1j} \\ d_2 - \sum_{j=1}^k w_{2,2j} \\ \vdots \\ \vdots \\ \vdots \\ d_n - \sum_{j=1}^k w_{n,nj} \end{pmatrix} = \begin{pmatrix} 0 \\ 0 \\ \vdots \\ \vdots \\ \vdots \\ 0 \end{pmatrix} \in R^n$$

流形学习

LE算法步骤

- 1 Input: 给出数据 $X = [x_1, x_2, \dots, x_n] \in R^{m \times n}$, 近邻参数 k 以及低维空间 d ;
 - 2 确定 x_i 的 k 个近邻, 确定亲和度矩阵 W , 计算度矩阵 D ;
 - 3 求解模型 $\min E(Y), \text{s.t. } YY^T = I$;
 - 4 采用第2至第 $d + 1$ 个最小的特征值对应的特征向量组成低维嵌入 Y ;
 - 5 输出 $Y \in R^{d \times n}$
-

M. Belkin and P. Niyogi, "Laplacian eigenmaps for dimensionality reduction and data representation," Neural Computation, vol. 15, no. 6, pp. 1373 - 1396, 2003

流形学习-性能对比



(a) Data



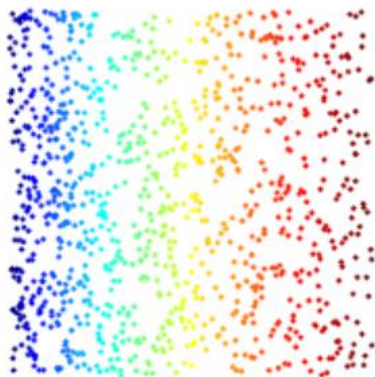
(b) Isomap



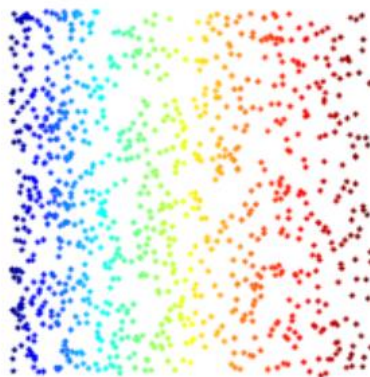
(c) LLE



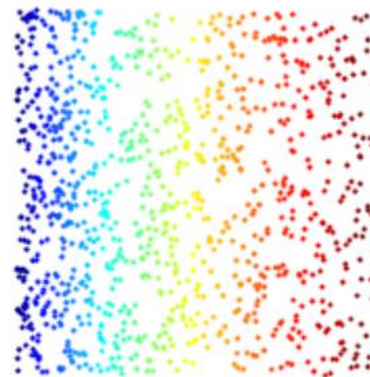
(d) LE



(e) HLLE



(f) LTSA



(g) LSE

流形学习-总结

— 统一的学习模型

- 目标：给定高维数据 $\{x_i\}_{i=1}^n \subset R^m$ 寻找其低维表示。
- 学习模型： $\{y_i\}_{i=1}^n \subset R^d, d < m$ 。

对L 进行特征值分解

$$\begin{aligned} & \min_Y \text{trace}(YMY^T) \\ \text{s.t. } & YY^T = I, Y = [y_1, y_2, \dots, y_n] \in R^{d \times n} \end{aligned}$$

任务：构造 M 矩阵——与数据图构造和局部描述紧密相关！

内容提要

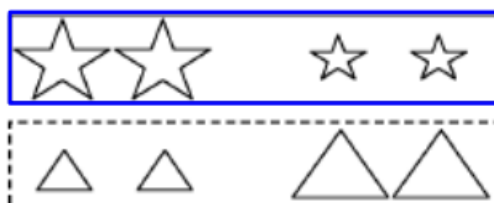
- 引言
- 主成分分析
- 多维缩放
- 流形学习方法
- 距离度量学习

度量学习

● 学习任务——距离度量是模式分类的基础！

—能否：让距离度量反映用户偏好或某些先验知识？

形状
检索



依形状



依大小


图像
检索



度量学习

● 学习任务——距离度量是模式分类的基础！

— 给定一些相似点对和不相似点对，学习一个距离度量来反映用户的某种偏好。

$$d = \sqrt{\sum_{i=1}^m (x_i - y_i)^2} = \left((x - y)^T (x - y) \right)^{1/2}$$


欧氏距离同等地对待每一个特征分量，不能反映出偏好信息！

度量学习

- 学习任务——距离度量是模式分类的基础！
- 考虑马氏距离

单位矩阵

$$d_E(x, y) = \sqrt{(x - y)^T (x - y)} = \sqrt{(x - y)^T \mathbf{I} (x - y)}$$

$$d_A(x, y) = \sqrt{(x - y)^T \mathbf{A} (x - y)} = \sqrt{(x - y)^T \mathbf{W} \mathbf{W}^T (x - y)}$$

度量学习

- 学习任务——距离度量是模式分类的基础！

- 考虑马氏距离

单位矩阵

$$d_E(x, y) = \sqrt{(x - y)^T (x - y)} = \sqrt{(x - y)^T \mathbf{I} (x - y)}$$

$$d_A(x, y) = \sqrt{(x - y)^T \mathbf{A} (x - y)} = \sqrt{(x - y)^T \mathbf{W} \mathbf{W}^T (x - y)}$$

其中 \mathbf{A} 为一个半正定矩阵。令 $\mathbf{A} = \mathbf{W} \mathbf{W}^T$ ，则马氏矩阵度量等价于将数据通过 $\mathbf{y} = \mathbf{W}^T \mathbf{x}$ 变换后在新空间做欧氏距离计算。

度量学习

- 经典学习模型

$$\begin{aligned} & \max_A \sum_{x_i, x_j \in D} d_A(x_i, x_j) \\ & \text{s.t.} \sum_{x_i, x_j \in S} [d_A(x_i, x_j)]^2 \leq 1 \\ & \quad A \succeq 0 \end{aligned}$$

E. P. Xing, A. Y. Ng, M. I. Jordan, S. Russell, Distance metric learning, with application to clustering with side-information, NIPS, 2003

度量学习

● 经典学习模型

$$\begin{aligned} & \max_A \sum_{x_i, x_j \in D} d_A(x_i, x_j) \\ & s. t. \sum_{x_i, x_j \in S} [d_A(x_i, x_j)]^2 \leq 1 \\ & A \succeq 0 \end{aligned}$$

不相似点对马氏距离尽可能小

相似点对马氏距离尽可能小

半正定

E. P. Xing, A. Y. Ng, M. I. Jordan, S. Russell, Distance metric learning, with application to clustering with side-information, NIPS, 2003

度量学习

- 距离准则—迹比值最大化

$$\max_{W^T W = I} g(W) = \frac{\text{tr}(W^T S_b W)}{\text{tr}(W^T S_w W)}$$

所有**不相似**点对
距离平方和

所有**相似**点对
距离平方和

矩阵的迹

度量学习

- 距离准则—迹比值最大化

$$\max_{W^T W = I} g(W) = \frac{\text{tr}(W^T S_b W)}{\text{tr}(W^T S_w W)}$$

所有**不相似**点对
距离平方和

所有**相似**点对
距离平方和

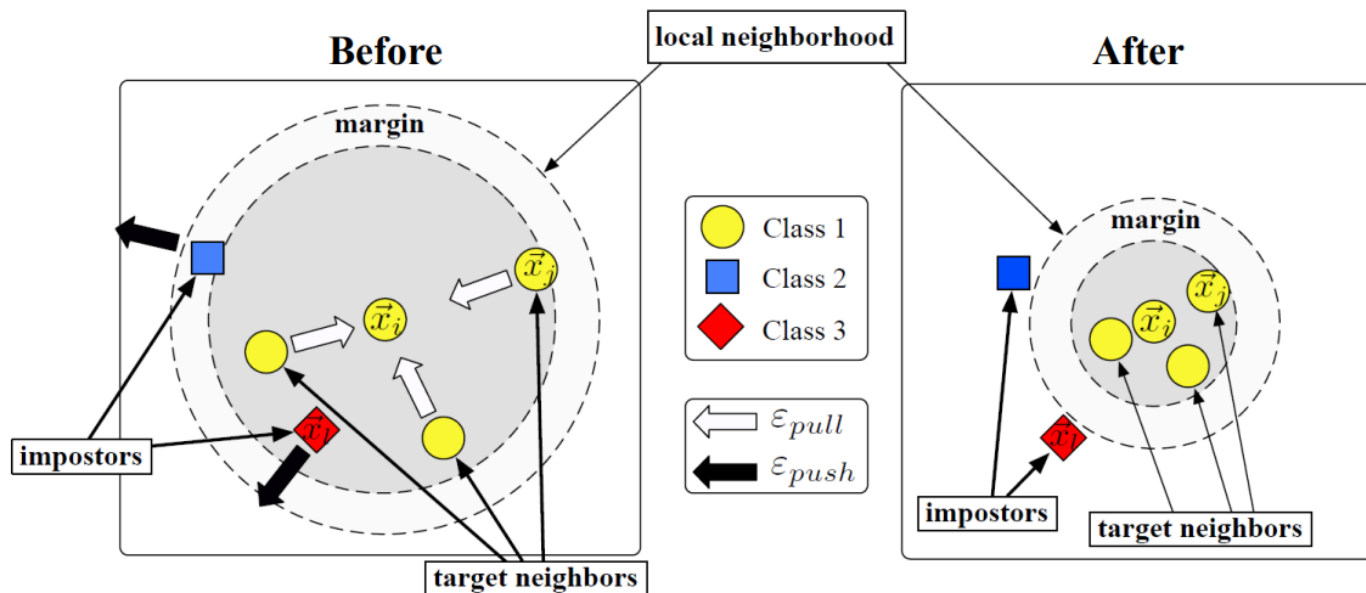
矩阵的迹

$$S_w = \sum_{x_i, x_j \in S} (x_i - x_j)(x_i - x_j)^T, S_b = \sum_{x_i, x_j \in D} (x_i - x_j)(x_i - x_j)^T$$

度量学习

● 最大间隔近邻分类 (LMNN)

— 基本思想：大间隔最近邻居算法的主要想法就是通过学习一种距离度量使得在一个新的转换空间中，对于一个输入 x_i 的 k 个近邻都属于同一类别，而不同类别的样本与 x_i 保持一定大的距离。



度量学习

● 最大间隔近邻分类 (LMNN)

— 学习模型

$$\begin{aligned} \xi_{ijl} &\geq 0 \\ M &\succcurlyeq 0 \end{aligned} \quad \min \sum_{ij} \eta_{ij} (x_i - x_j)^T M (x_i - x_j) + \sum_{ijl} \eta_{ij} (1 - y_{il}) \xi_{ijl}$$
$$s. t. (x_i - x_l)^T M (x_i - x_l) - (x_i - x_j)^T M (x_i - x_j) \geq 1 - \xi_{ijl}$$

$$\eta_{ij} = \begin{cases} 1, & \text{if } (\mathbf{x}_i, \mathbf{x}_j) \in \text{the same class, and } \mathbf{x}_j \in N_k(\mathbf{x}_i) \\ 0, & \text{otherwise} \end{cases}$$
$$y_{il} = \begin{cases} 1, & \text{if } (\mathbf{x}_i, \mathbf{x}_l) \in \text{the same class} \\ 0, & \text{otherwise} \end{cases}$$

度量学习

● 最大间隔近邻分类 (LMNN)

— 学习模型

属于同一类的邻近点之间的距离要小

同一邻域内，不属于同一类的点对

$$\xi_{ijl} \geq 0$$
$$M \succcurlyeq 0$$

$$\min \sum_{ij} \eta_{ij} (x_i - x_j)^T M (x_i - x_j) + \sum_{ijl} \eta_{ij} (1 - y_{il}) \xi_{ijl}$$

Relax variables

$$s. t. \underline{(x_i - x_l)^T M (x_i - x_l) - (x_i - x_j)^T M (x_i - x_j) \geq 1 - \xi_{ijl}}$$

同一邻域内，不属于同一类的点之间的距离要大于“1”

度量学习

- 当前的一些主要的工作

- 基于弱信息（侧信息）的学习模型构建
- 基于信息论的距离度量学习
- 相关/相似关系学习
- 核矩阵学习

参考文献

- 周志华. 《机器学习》
- 李航. 《统计学习方法》

致 谢

- 感谢**向世明**老师的20版PPT作为原始材料
- 感谢**王锐**与**段俊贤**对本PPT的制作与修改

致谢

- 感谢向世明老师的20版PPT作为原始材料
- 感谢王锐与段俊贤对本PPT的制作与修改

Thank All of You!
(Questions?)

赫然

rhe@nlpr.ia.ac.cn

<https://rhe-web.github.io/>

智能感知与计算研究中心 (CRIPAC)

中科院自动化研究所 模式识别国家重点实验室