

机器学习之无监督学习

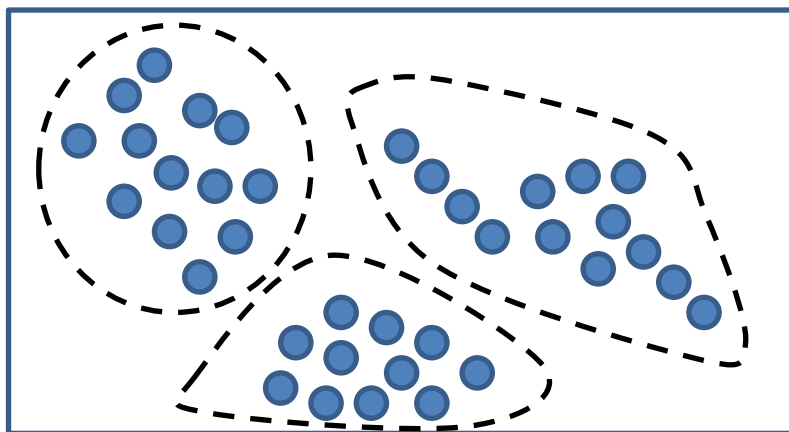
数据降维算法

倪冰冰

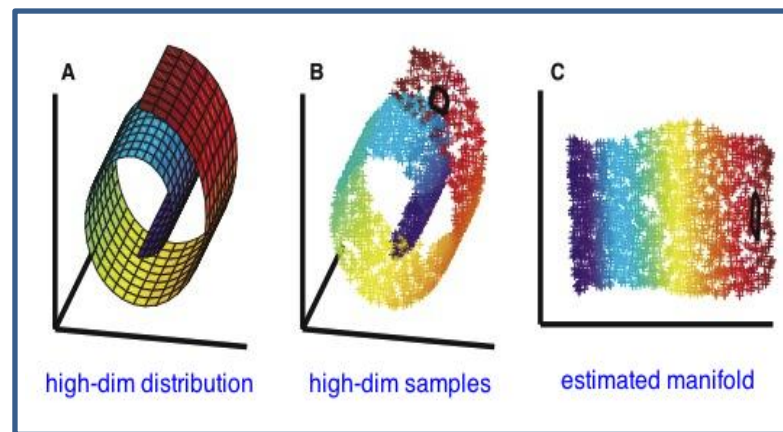
上海交通大学

引言

- 什么是无监督学习？
 1. 数据没有明确的标签信息。
 2. 我们希望仅依赖数据本身来探索其具有的内在结构信息。
- 无监督学习的种类有哪些？



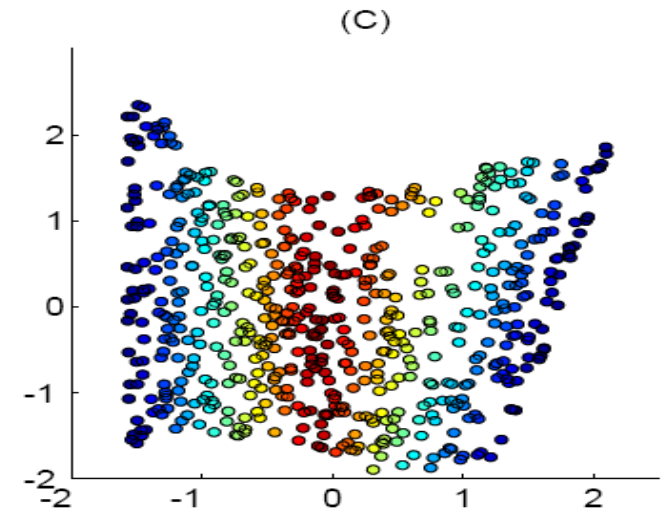
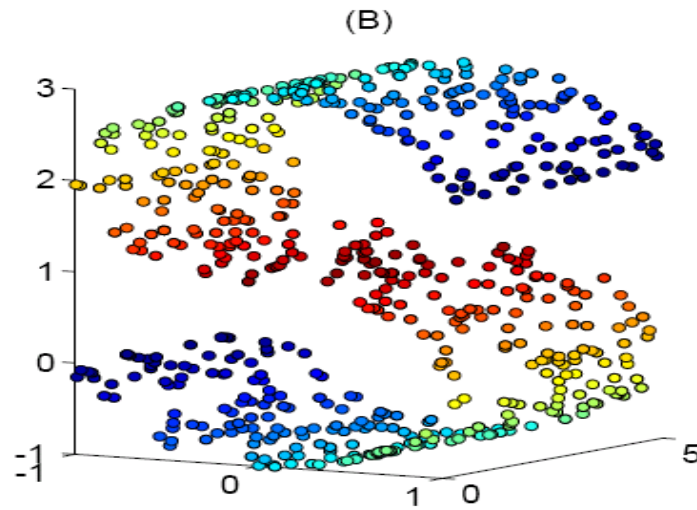
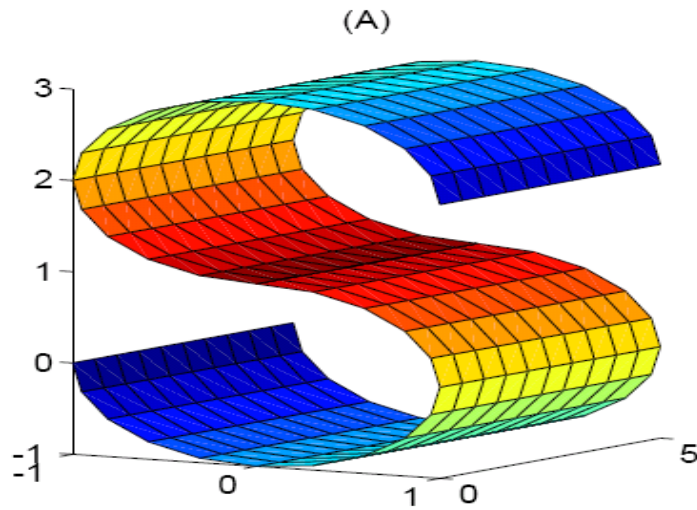
聚类学习



表征学习(降维)

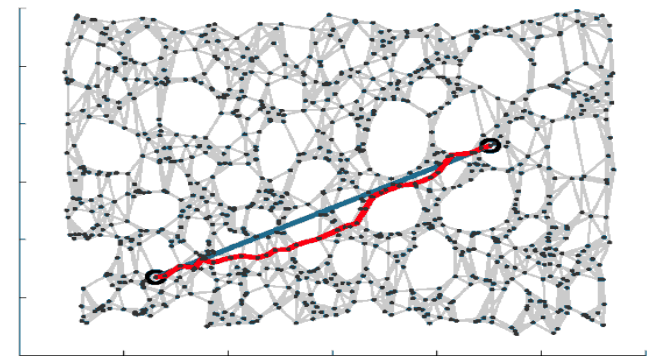
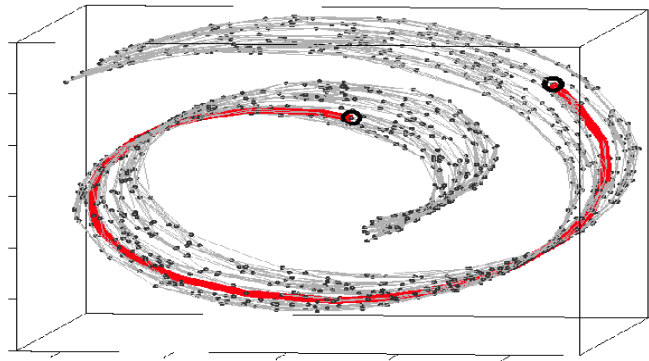
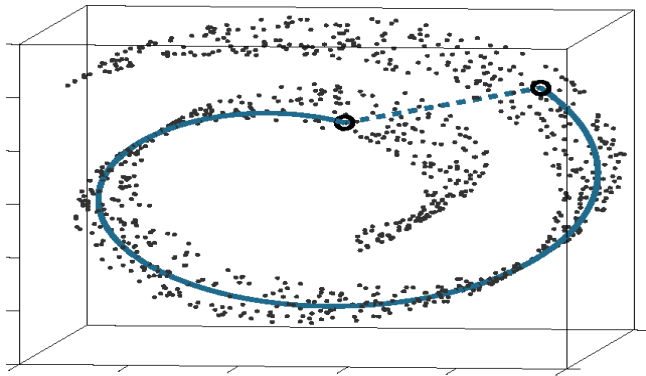
非线性数据降维

- 如果数据集拥有非线性的结构，那么线性方法如PCA和MDS等均不能有效工作
- 非线性的数据降维方法因此应运而生



ISOMAP

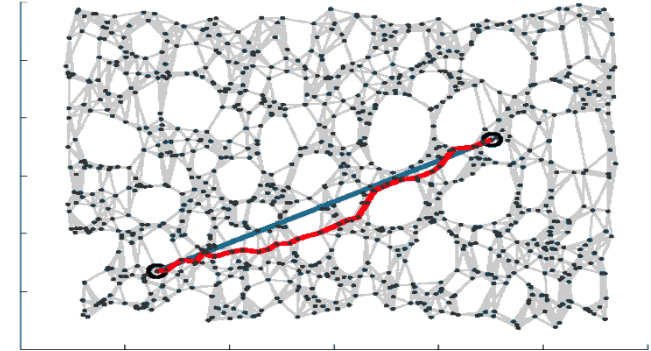
- 非线性结构举例（瑞士卷）
只有geodesic距离才能反映该数据的真正低维流形结构
- ISOMAP (Isometric feature Mapping)
保留数据集的本征几何结构
每两个数据点之间计算geodesic距离（区别与Euclidean距离）



ISOMAP (算法描述)

● 步骤 1

- 在原空间（高维度空间）找到每个点的“邻居” $d_x(i, j)$
- 点与点之间邻接关系（邻接度）可以由图来表示

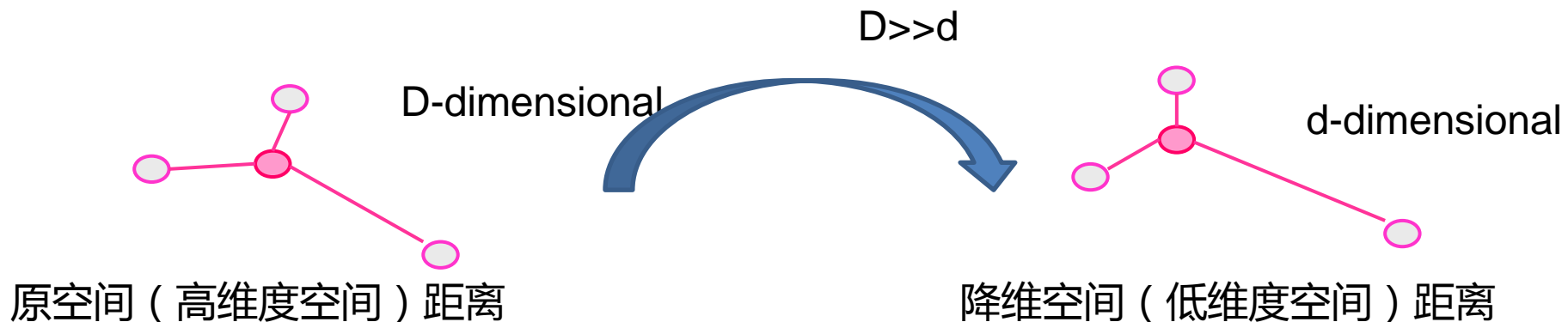


● 步骤 2

- 估计任意两数据点之间的geodesic距离，即在数据流形上的最短路径 $d_G(i, j)$

● 步骤 3

- 构建一个低维空间中的数据流形，使得点与点之间的距离保持之前计算的原始空间中的geodesic距离

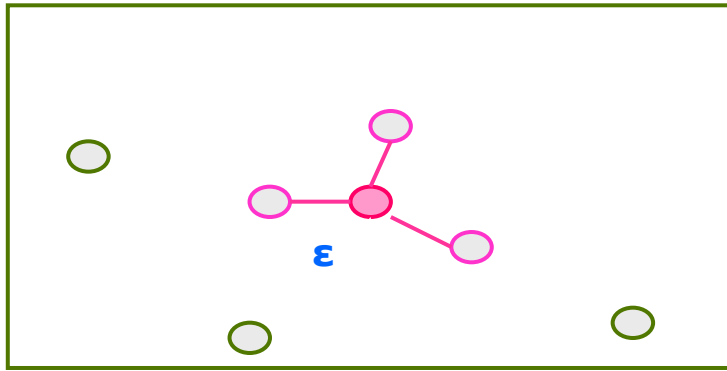


ISOMAP (算法描述)

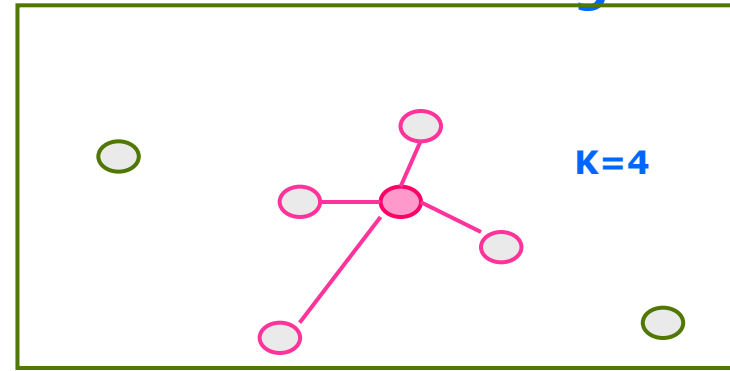
- 步骤1

- 两种邻接点表示方式 $d_x(i, j)$

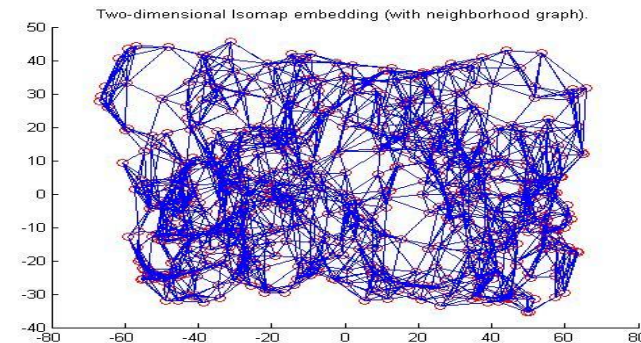
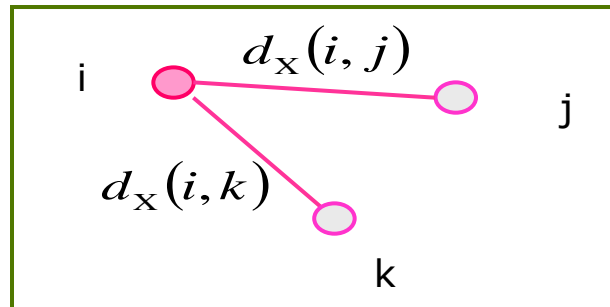
ϵ -radius



K-nearest neighbors



- 构建邻接图



ISOMAP (算法描述)

- 步骤2

- 构建每两个数据点之间geodesic距离：最短路径算法.
- 可以使用Floyd' s算法或者Dijkstra' s算法

$$d_G(i, j)$$

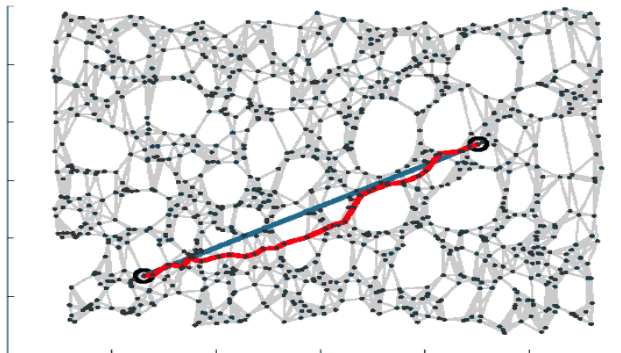
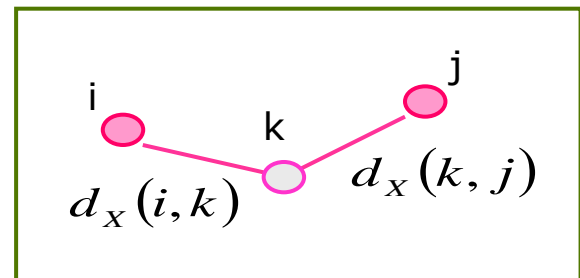
Dijkstra' s算法

$$d_G(i, j) = d_X(i, j) \text{ neighboring } i, j$$

$$d_G(i, j) = \infty \text{ otherwise}$$

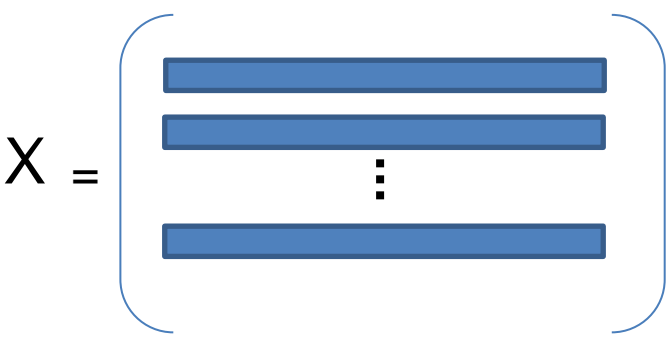
for $k = 1, 2, \dots, N$

$$d_G(i, j) = \min\{d_X(i, j), d_X(i, k) + d_X(k, j)\}$$

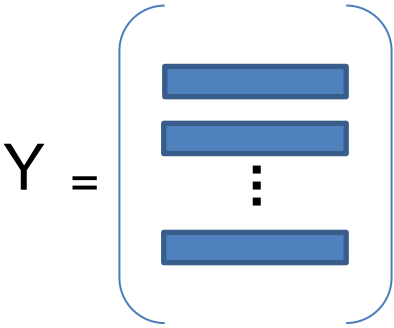
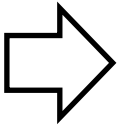


ISOMAP (算法描述)

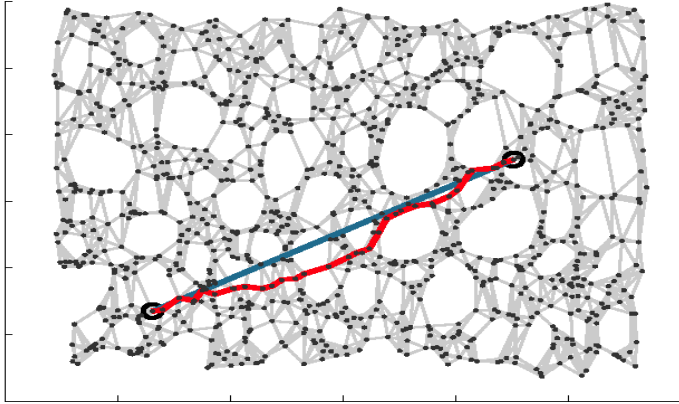
- 步骤3
 - 构建一个低维度特征集Y来最佳保留原始空间中的数据流形结构
 - 使用类似MDS的算法解出特征集Y



原特征空间



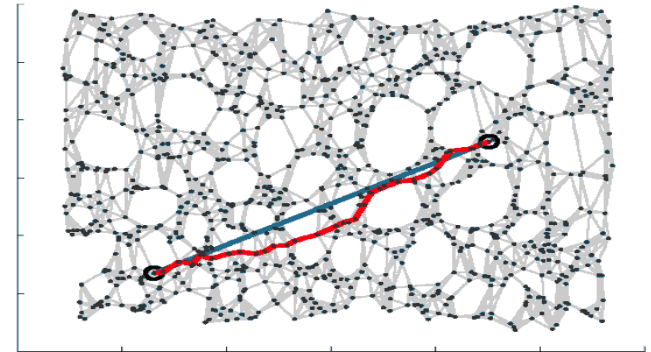
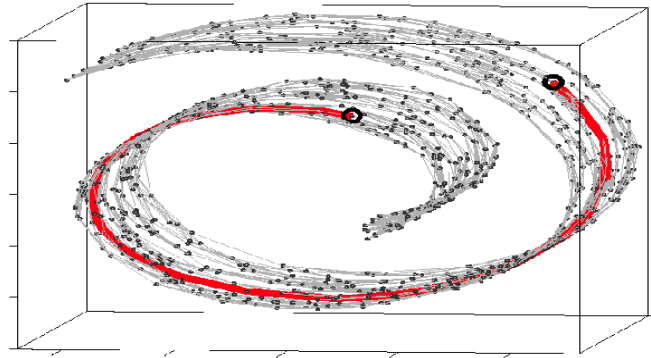
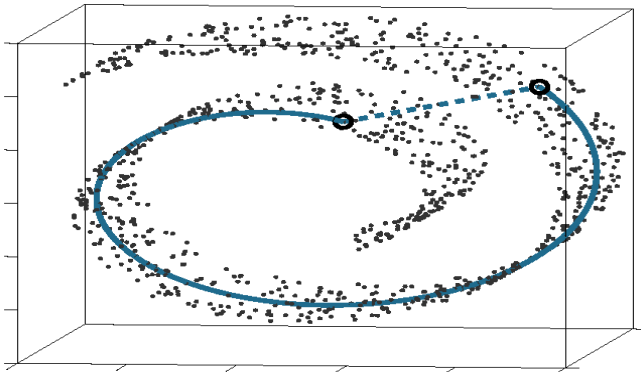
降维后特征空间



算法步骤	MDS	ISOMAP
1	计算两两数据样本间欧式距离 d_{ij}	计算两两数据样本间geodesic距离 d_{ij}
2	通过 d_{ij} 构建T矩阵 ($T=YY^T$)	通过 d_{ij} 构建T矩阵 ($T=YY^T$)
3	通过矩阵特征值分解获得低维表示Y	通过矩阵特征值分解获得低维表示Y

ISOMAP低维度流形结构恢复

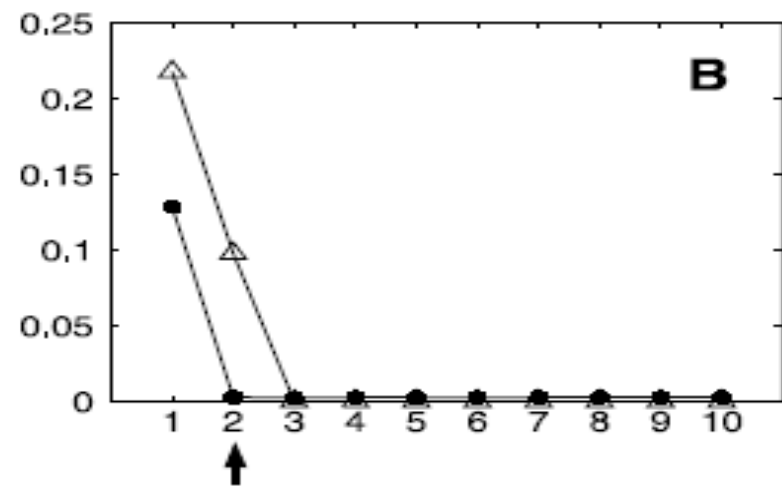
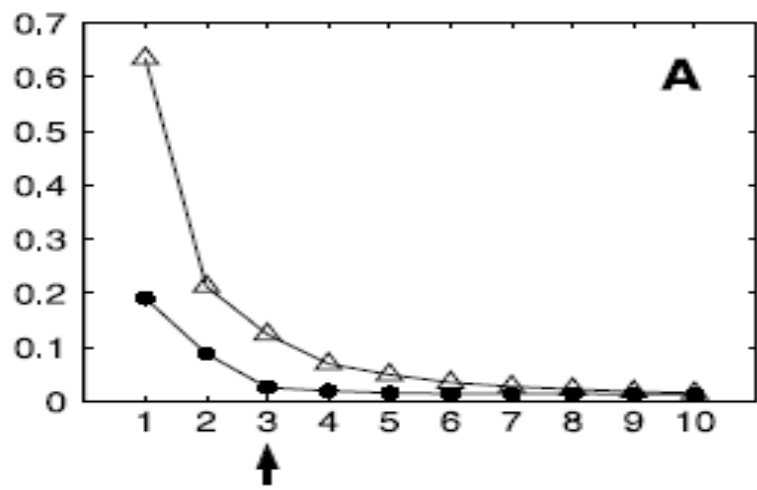
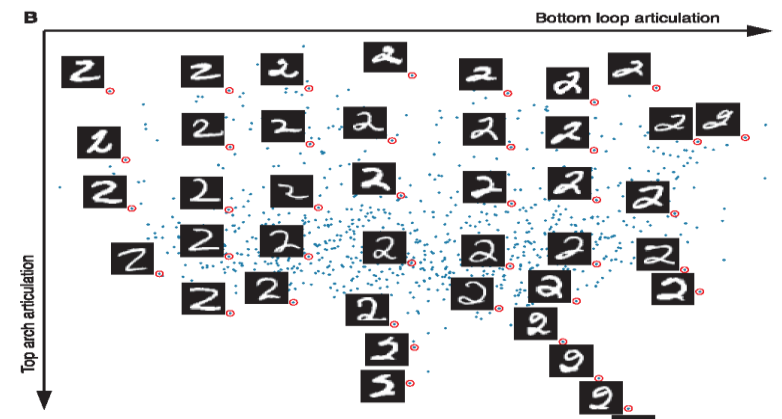
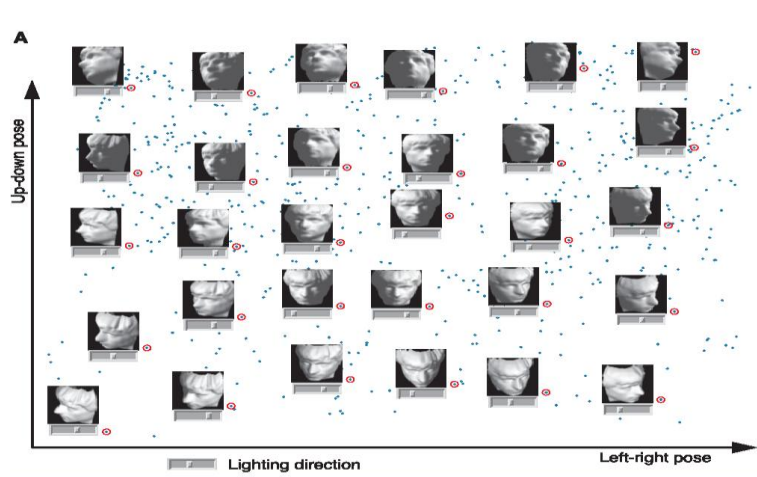
- ISOMAP可以保证所恢复的低维度特征集很好地反映原空间中的非线性流形结构
- 当采样密度更大时（采样点更多），所恢复的低维度流形结构更加贴近原空间中结构



低维非线性流形例子(ISOMAP)

Face

Hand writing



MDS : open triangles
ISOMAP : filled circles

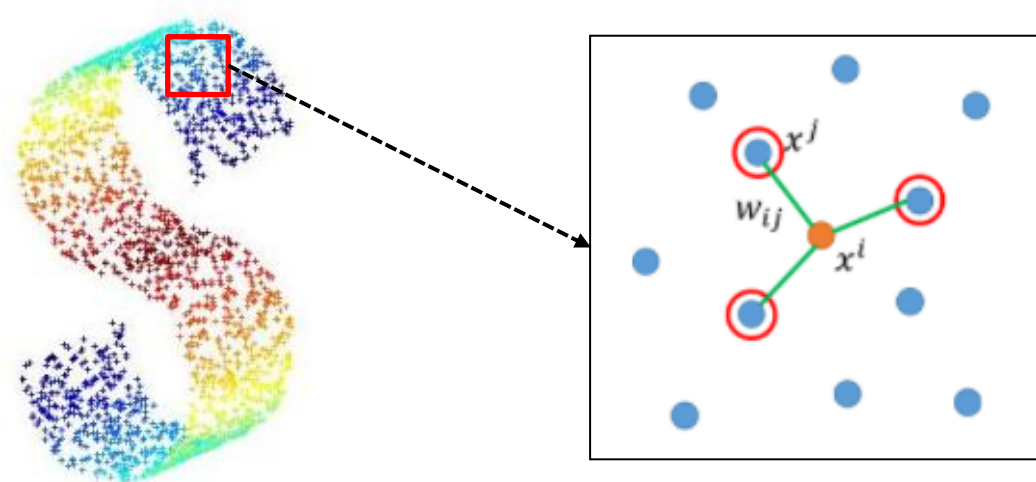
Locally Linear Embedding

- LLE (Locally Linear Embedding)
 - 保留邻接几何结构
 - 可以全局地映射到低维空间，通过局部的线性关系建模，恢复全局的非线性结构
 - **假设**：每个数据点和它的邻居点处在同一个线性的小块上，因此每个数据点可以被它的邻居们所线性重建：

$$\hat{X}_i = \sum_j W_{ij} X_j$$

$$W_{ij} = 0 \text{ if } X_j \text{ is not a neighbor of } X_i$$

- W_{ij} 表示第j-th 数据点对i-th 数据点的重建系数



LLE (算法描述)

- 最小化误差函数：

$$\mathcal{E}(W) = \sum_i \left| X_i - \sum_j W_{ij} X_j \right|^2$$

约束条件：

$$W_{ij} = 0 \quad \text{if } X_j \text{ is not a neighbor of } X_i$$

$$\sum_j W_{ij} = 1$$

- 显式解(使用拉格朗日乘子):

$$C_{jk} = (x_i - x_j) \cdot (x_i - x_k)$$

$$w_j = \frac{\sum_k C_{jk}^{-1}}{\sum_{lm} C_{lm}^{-1}}$$

LLE (算法描述)

- 低维度空间解Y

最小化:

$$\phi(Y) = \sum_i \left| Y_i - \sum_j W_{ij} Y_j \right|^2$$

约束:

$$\sum_i Y_i = 0, \quad \frac{1}{N} \sum_i Y Y^T = I$$

Quadratic form, $\min Tr(YMY^T)$

where:

$$M = (I - W)^T (I - W)$$

- 求解：找出M的d个最小特征值对应的特征向量.

提示：拉格朗日算子： $\min Tr(YMY^T) - \lambda(YY^T - nI)$

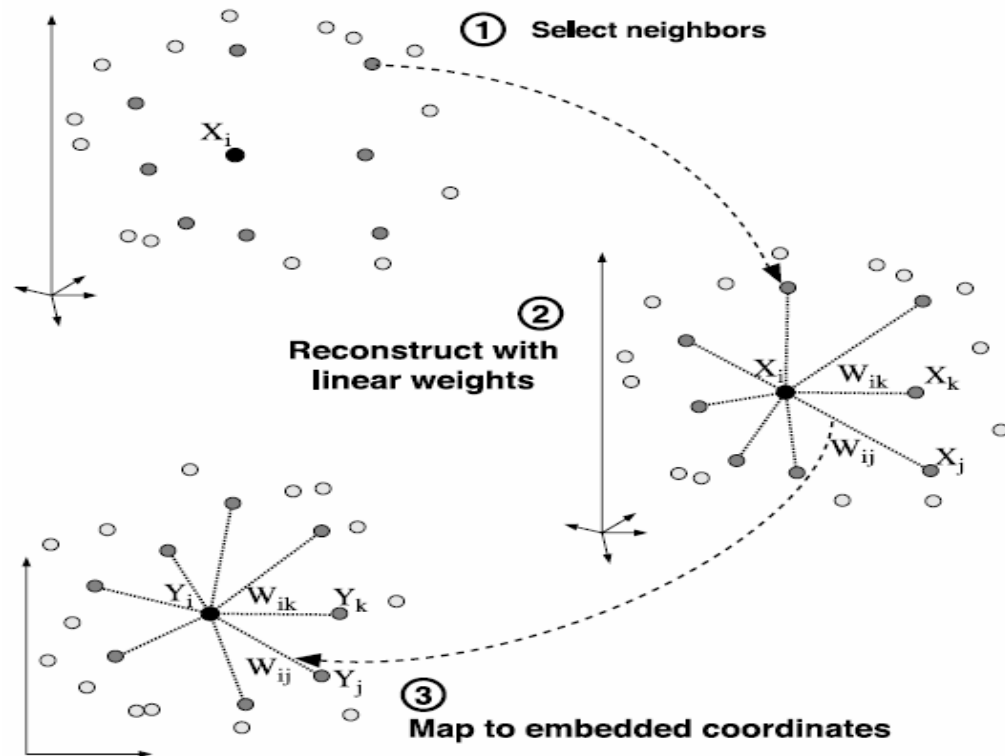
LLE (算法总结)

- Step 1
 - 为每个数据 X_i 找到K个近邻
- Step 2
 - 计算权重系数 W_{ij} 表示数据点 X_i 的邻居点对它的重建作用, 通过最小化原空间重建误差 (约束条件下) :

$$\varepsilon(W) = \sum_i \left| X_i - \sum_j W_{ij} X_j \right|^2$$

- Step 3
 - 计算低维表示 Y , 保留重建权重 W_{ij} , 最小化低维空间重建误差 (约束条件下) :

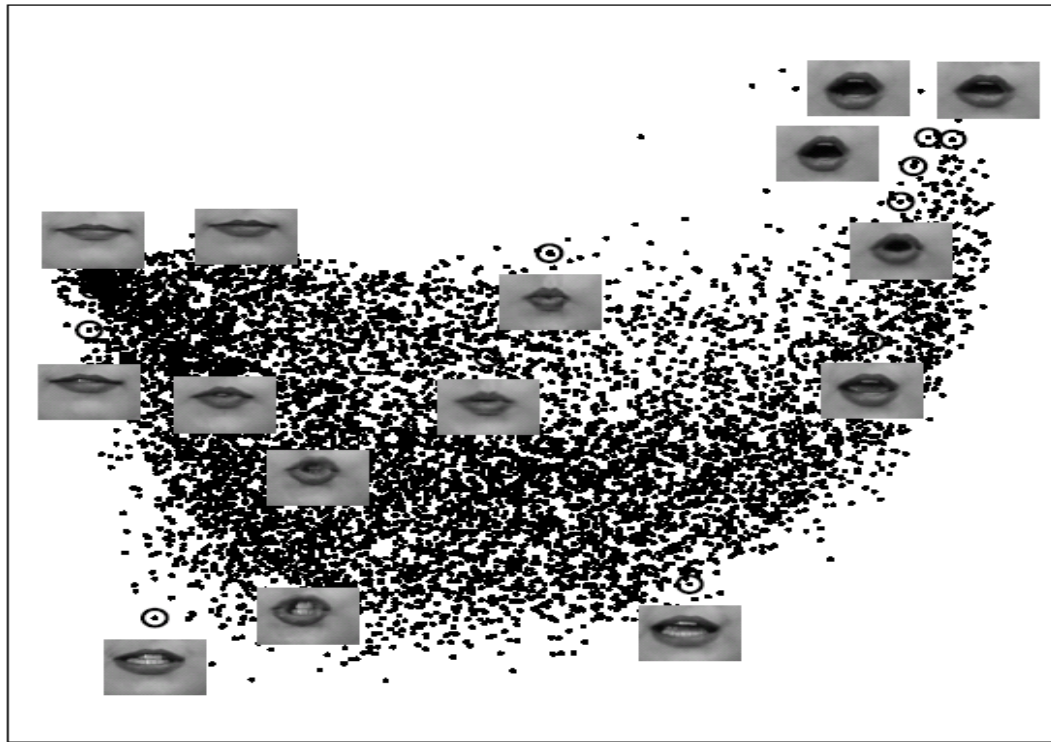
$$\phi(Y) = \sum_i \left| Y_i - \sum_j W_{ij} Y_j \right|^2$$



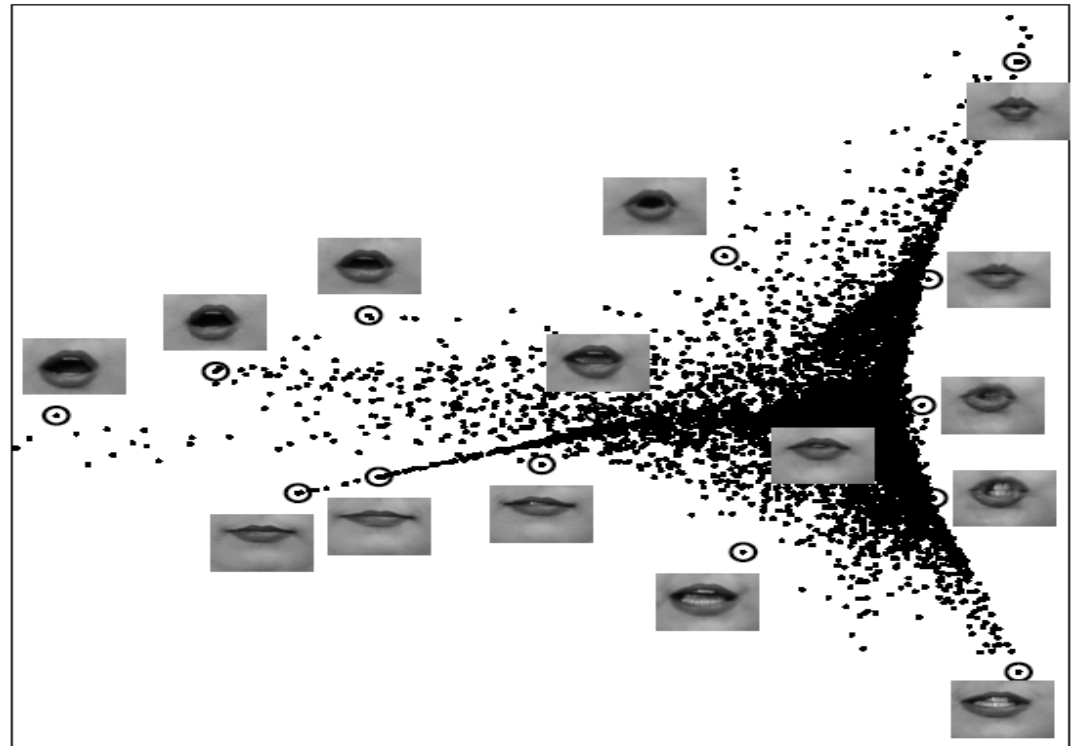
低维度流形例子(LLE)

- Lips

PCA



LLE



ISOMAP LLE小结

- ISOMAP
 - 两两点之间使用geodesic manifold distances代替Euclidean Distance
- LLE
 - 从局部的线性结构关系，恢复全局的非线性流形
- ISOMAP vs LLE
 - 保留了邻接的几何结构
 - LLE需要更多的训练数据
 - ISOMAP计算效率高，实用性更高
 - LLE与ISOMAP都没有显式的映射函数，故使用比较麻烦

THANK YOU

AI300学院

