

# 机器学习引论

彭玺

[pengxi@scu.edu.cn](mailto:pengxi@scu.edu.cn)

[www.pengxi.me](http://www.pengxi.me)

四川大学-机器学习引论

# 提纲

- 一 . Review
- 二 . Locally Linear Embedding
- 三 . Laplacian Eigenmap

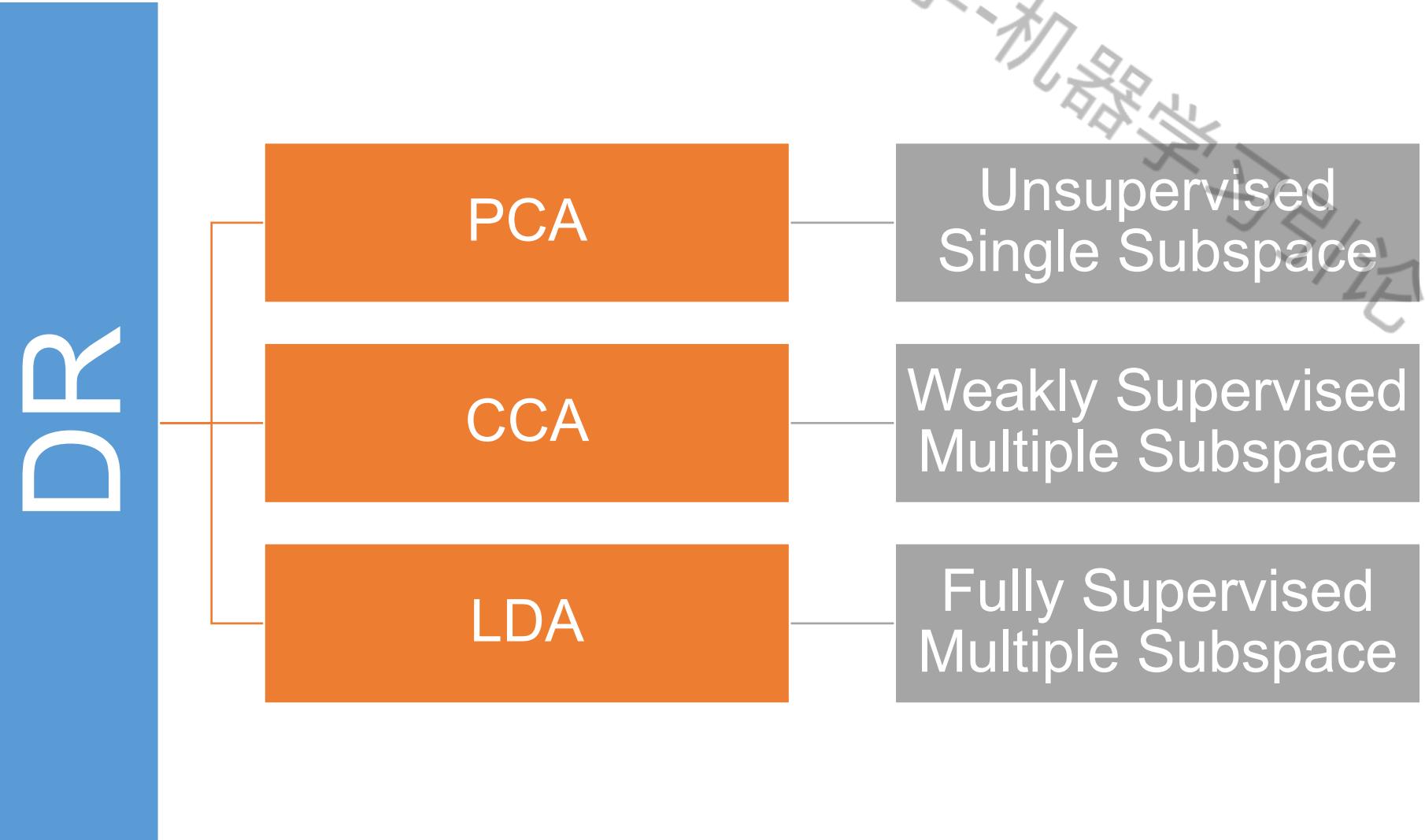
四川大学-机器学习引论

# 提纲

- 一 . Review
- 二 . Locally Linear Embedding
- 三 . Laplacian Eigenmap

四川大学-机器学习引论

# 一、Review



# 一、Review

How to compute redundancy in mathematics?

- Covariance:

If the entries in the column vector

$$\mathbf{X} = \begin{bmatrix} X_1 \\ \vdots \\ X_n \end{bmatrix}$$

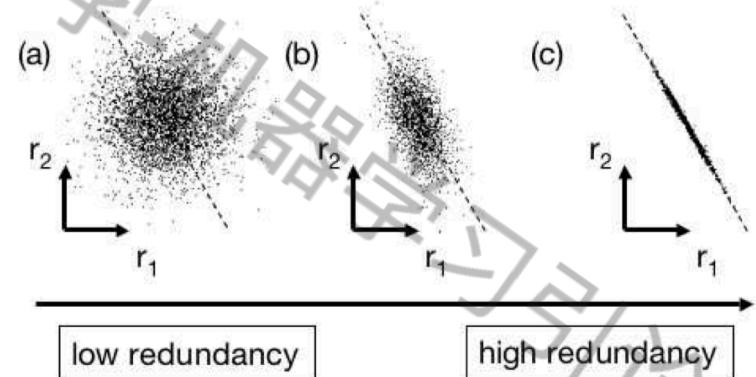
are **random variables**, each with finite **variance**, then the covariance matrix  $\Sigma$  is the matrix whose  $(i, j)$  entry is the **covariance**

$$\Sigma_{ij} = \text{cov}(X_i, X_j) = E[(X_i - \mu_i)(X_j - \mu_j)] = E[X_i X_j] - \mu_i \mu_j$$

where the operator  $E$  denotes the expected (mean) value of its argument, and

$$\mu_i = E(X_i)$$

- $\Sigma_{ij} = 0$  if and only if  $i$  and  $j$  are entirely **uncorrelated**.
- Otherwise,  $i$  and  $j$  are **correlated**.



Identical with **cosine distance**  
with normalized input

Correlated=redundant!

Futher reading: In fact, the variances  $\Sigma_{ii}$  also defines the signal-to-noise ratio.

# 一、Review

As the covariance defines the redundancy, then one could **remove the redundancy** in low dimensional space by **diagonalizing the covariance matrix**.

投影 :  $\mathbf{w}^T \mathbf{x}$

$$\text{方差} : \frac{1}{n} \sum_{i=1}^n (\mathbf{w}^T \mathbf{x}_i)^2 = \mathbf{w}^T S \mathbf{w}$$

$$S = \frac{1}{n} \sum_i \mathbf{x}_i \mathbf{x}_i^T$$

最大方差 :

$$\begin{aligned} \max_{\mathbf{w}} \quad & \mathbf{w}^T S \mathbf{w} \\ s.t. \quad & \|\mathbf{w}\| = 1 \end{aligned}$$

$$\text{Note that: } \frac{1}{n} \sum_{i=1 \dots n} (\mathbf{w}^T \mathbf{x}_i)^2 = \mathbf{w}^T X X^T \mathbf{w}$$

拉格朗日乘数法 :

$$\begin{aligned} L &= \mathbf{w}^T S \mathbf{w} + \lambda(1 - \mathbf{w}^T \mathbf{w}) \\ \frac{\partial L}{\partial \mathbf{w}} &= 2S\mathbf{w} - 2\lambda\mathbf{w} \\ S\mathbf{w} &= \lambda\mathbf{w} \end{aligned}$$

方差 :

$$\mathbf{w}^T S \mathbf{w} = \mathbf{w}^T \lambda \mathbf{w} = \lambda$$

# 一、Review

View 2: minimizing reconstruction error/description length.

正交基 :

$$\mathbf{u}_1, \dots, \mathbf{u}_D$$

原始数据 :

$$\mathbf{x}_i = \sum_{j=1}^D \alpha_{ij} \mathbf{u}_j$$

基坐标 :

$$\alpha_{ij} = \mathbf{u}_j^T \mathbf{x}_i$$

降维重建 :

$$\hat{\mathbf{x}}_i = \sum_{j=1}^d \alpha_{ij} \mathbf{u}_j$$

$$\begin{aligned} \frac{1}{n} \sum_{i=1}^n \|\mathbf{x}_i - \hat{\mathbf{x}}_i\|^2 &= \frac{1}{n} \sum_{i=1}^n \left\| \sum_{j=1}^D \alpha_{ij} \mathbf{u}_j - \sum_{j=1}^d \alpha_{ij} \mathbf{u}_j \right\|^2 \\ &= \frac{1}{n} \sum_{i=1}^n \left\| \sum_{j=d+1}^D \alpha_{ij} \mathbf{u}_j \right\|^2 \\ &= \frac{1}{n} \sum_{i=1}^n \sum_{j=d+1}^D \alpha_{ij}^2 \\ &= \frac{1}{n} \sum_{i=1}^n \sum_{j=d+1}^D \mathbf{u}_j^T \mathbf{x}_i \mathbf{x}_i^T \mathbf{u}_j \\ &= \sum_{j=d+1}^D \mathbf{u}_j^T S \mathbf{u}_j \quad \text{等价方差最小} \end{aligned}$$

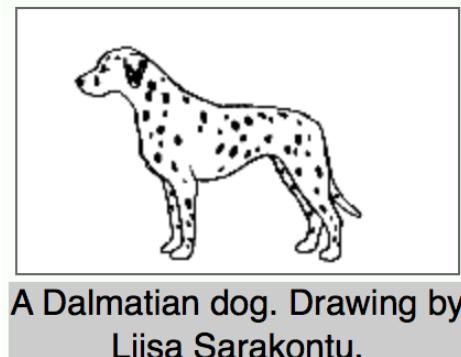
# 一、Review

Canonical correlation analysis (CCA) is a way of measuring the linear relationship between two multidimensional variables. It finds two bases, one for each variable, that are optimal with respect to correlations and, at the same time, it finds the corresponding correlations.

In other words, it finds the two bases in which the correlation matrix between the variables is diagonal and the correlations on the diagonal are maximized.

The latent assumption of CCA is that these two variables describe the same object in different views, e.g., image+text,

- multi-view analysis



# 一、Review

Let  $x = \mathbf{w}_x^T \mathbf{x}$  and  $y = \mathbf{w}_y^T \mathbf{y}$  (canonical variates), the formulation of CCA is as follows:

$$\Sigma_{ij} = \text{cov}(X_i, X_j) = \mathbb{E}[(X_i - \mu_i)(X_j - \mu_j)] = \mathbb{E}[X_i X_j] - \mu_i \mu_j$$

$$\operatorname{argmax}_{\mathbf{w}_x, \mathbf{w}_y} \frac{\text{cov}(\mathbf{w}_x^T \mathbf{x}, \mathbf{w}_y^T \mathbf{y})}{\sqrt{\text{var}(\mathbf{w}_x^T \mathbf{x}) \text{var}(\mathbf{w}_y^T \mathbf{y})}}$$

$$= \operatorname{argmax}_{\mathbf{w}_x, \mathbf{w}_y} \frac{\mathbf{w}_x^T \mathbf{x} (\mathbf{w}_y^T \mathbf{y})^T}{\sqrt{(\mathbf{w}_x^T \mathbf{x})(\mathbf{w}_x^T \mathbf{x})^T (\mathbf{w}_y^T \mathbf{y})(\mathbf{w}_y^T \mathbf{y})^T}}$$

CCA is formulated as

$$\boxed{\operatorname{argmax}_{\mathbf{w}_x, \mathbf{w}_y} \mathbf{w}_x^\top \mathbf{C}_{xy} \mathbf{w}_y,}$$

$$\begin{aligned} \text{s.t. } \mathbf{w}_x^\top \mathbf{C}_{xx} \mathbf{w}_x &= 1, \\ \mathbf{w}_y^\top \mathbf{C}_{yy} \mathbf{w}_y &= 1. \end{aligned}$$

$$x \in \mathcal{R}^d$$

$$\mathbf{x} \in \mathcal{R}^{m_x}$$

$$\mathbf{w}_x \in \mathcal{R}^{m_x \times d}$$

$$y \in \mathcal{R}^d$$

$$\mathbf{y} \in \mathcal{R}^{m_y}$$

$$\mathbf{w}_y \in \mathcal{R}^{m_y \times d}$$

# 一、Review

Incorporating these two constraints, the Lagrangian  $\mathcal{J}$  is given by

$$\mathcal{J} = \mathbf{w}_x^\top \mathbf{C}_{xy} \mathbf{w}_y + \lambda_x (1 - \mathbf{w}_x^\top \mathbf{C}_{xx} \mathbf{w}_x) + \lambda_y (1 - \mathbf{w}_y^\top \mathbf{C}_{yy} \mathbf{w}_y).$$

Let  $\frac{\partial \mathcal{J}}{\partial \mathbf{w}_x} = 0$  and  $\frac{\partial \mathcal{J}}{\partial \mathbf{w}_y} = 0$  lead to

$$\mathbf{C}_{xy} \mathbf{w}_y - 2\lambda_x \mathbf{C}_{xx} \mathbf{w}_x = 0, \quad (1)$$

$$\mathbf{C}_{yx} \mathbf{w}_x - 2\lambda_y \mathbf{C}_{yy} \mathbf{w}_y = 0. \quad (2)$$

Pre-multiply (1) by  $\mathbf{w}_x^\top$  and pre-multiply (2)  $\mathbf{w}_y^\top$  to obtain

$$\mathbf{w}_x^\top \mathbf{C}_{xy} \mathbf{w}_y - 2\lambda_x \underbrace{\mathbf{w}_x^\top \mathbf{C}_{xx} \mathbf{w}_x}_1 = 0,$$

$$\mathbf{w}_y^\top \mathbf{C}_{yx} \mathbf{w}_x - 2\lambda_y \underbrace{\mathbf{w}_y^\top \mathbf{C}_{yy} \mathbf{w}_y}_1 = 0,$$

# 一、Review

$$\mathbf{w}_x^\top \mathbf{C}_{xy} \mathbf{w}_y - 2\lambda_x \underbrace{\mathbf{w}_x^\top \mathbf{C}_{xx} \mathbf{w}_x}_1 = 0,$$

$$\mathbf{w}_y^\top \mathbf{C}_{yx} \mathbf{w}_x - 2\lambda_y \underbrace{\mathbf{w}_y^\top \mathbf{C}_{yy} \mathbf{w}_y}_1 = 0,$$

leading to

$$\mathbf{w}_x^\top \mathbf{C}_{xy} \mathbf{w}_y = 2\lambda_x,$$

$$\mathbf{w}_y^\top \mathbf{C}_{yx} \mathbf{w}_x = 2\lambda_y.$$

Since  $\mathbf{w}_x^\top \mathbf{C}_{xy} \mathbf{w}_y = \mathbf{w}_y^\top \mathbf{C}_{yx} \mathbf{w}_x$ , we have  $\lambda = 2\lambda_x = 2\lambda_y$ .

Then, CCA is solved by following generalized eigenvalue problem

$$\begin{bmatrix} 0 & \mathbf{C}_{xy} \\ \mathbf{C}_{yx} & 0 \end{bmatrix} \begin{bmatrix} \mathbf{w}_x \\ \mathbf{w}_y \end{bmatrix} = \lambda \begin{bmatrix} \mathbf{C}_{xx} & 0 \\ 0 & \mathbf{C}_{yy} \end{bmatrix} \begin{bmatrix} \mathbf{w}_x \\ \mathbf{w}_y \end{bmatrix},$$

where  $\lambda = 2\lambda_x = 2\lambda_y$ .

# 一、Review

Recall the generalized eigenvalue problem:

$$\begin{bmatrix} 0 & C_{xy} \\ C_{yx} & 0 \end{bmatrix} \begin{bmatrix} w_x \\ w_y \end{bmatrix} = \lambda \begin{bmatrix} C_{xx} & 0 \\ 0 & C_{yy} \end{bmatrix} \begin{bmatrix} w_x \\ w_y \end{bmatrix}.$$

This problem has  $m_x + m_y$  eigenvalues  $\{\lambda_1, -\lambda_1, \dots, \lambda_m, -\lambda_m, 0, \dots, 0\}$ , where  $m = \min(m_x, m_y)$ .

The generalized eigenvalue problem can be re-written as

$$\begin{bmatrix} C_{xx} & C_{xy} \\ C_{yx} & C_{yy} \end{bmatrix} \begin{bmatrix} w_x \\ w_y \end{bmatrix} = (1 + \lambda) \begin{bmatrix} C_{xx} & 0 \\ 0 & C_{yy} \end{bmatrix} \begin{bmatrix} w_x \\ w_y \end{bmatrix}.$$

This problem has  $m_x + m_y$  eigenvalues  $\{1 + \lambda_1, 1 - \lambda_1, \dots, 1 + \lambda_m, 1 - \lambda_m, 1, \dots, 1\}$ .

# 一、Review

## CCA: Extension to Multiple Sets of Variables

Consider  $n$  multiple sets of variables,  $\mathbf{x}_1 \in \mathbb{R}^{m_1}, \mathbf{x}_2 \in \mathbb{R}^{m_2}, \dots, \mathbf{x}_n \in \mathbb{R}^{m_n}$ .

Then, CCA is formulated as

$$\arg \max_{\mathbf{w}_1, \dots, \mathbf{w}_n} \sum_{i=1}^n \sum_{j=1}^n \mathbf{w}_i^\top \mathbf{C}_{ij} \mathbf{w}_j,$$

$$\text{s.t. } \sum_{i=1}^n \mathbf{w}_i^\top \mathbf{C}_{ii} \mathbf{w}_j = 1.$$

Why does not separate?

# 一、Review

Incorporating these two constraints, the Lagrangian  $\mathcal{J}$  is given by

$$\mathcal{J} = \sum_{i=1}^n \sum_{j=1}^n \mathbf{w}_i^\top C_{ij} \mathbf{w}_j + \lambda \left( 1 - \sum_{i=1}^n \mathbf{w}_i^\top C_{ii} \mathbf{w}_j \right).$$

It follows from  $\frac{\partial \mathcal{L}}{\partial \mathbf{w}_i} = 0$  for  $i = 1, \dots, n$  that we have

$$\begin{bmatrix} C_{11} & C_{12} & \cdots & C_{1n} \\ C_{21} & C_{22} & \cdots & C_{2n} \\ \vdots & \vdots & & \vdots \\ C_{n1} & C_{n2} & \cdots & C_{nn} \end{bmatrix} \begin{bmatrix} \mathbf{w}_1 \\ \mathbf{w}_2 \\ \vdots \\ \mathbf{w}_n \end{bmatrix} = \lambda \begin{bmatrix} C_{11} & 0 & \cdots & 0 \\ 0 & C_{22} & \cdots & 0 \\ \vdots & & \ddots & 0 \\ 0 & \cdots & 0 & C_{nn} \end{bmatrix} \begin{bmatrix} \mathbf{w}_1 \\ \mathbf{w}_2 \\ \vdots \\ \mathbf{w}_n \end{bmatrix}.$$

The minimal generalized eigenvalue has the fixed range  $[0, 1]$ , whereas the maximal generalized eigenvalue has a range dependent on the dimensions of the variables. The minimal generalized eigenvalue is more convenient.

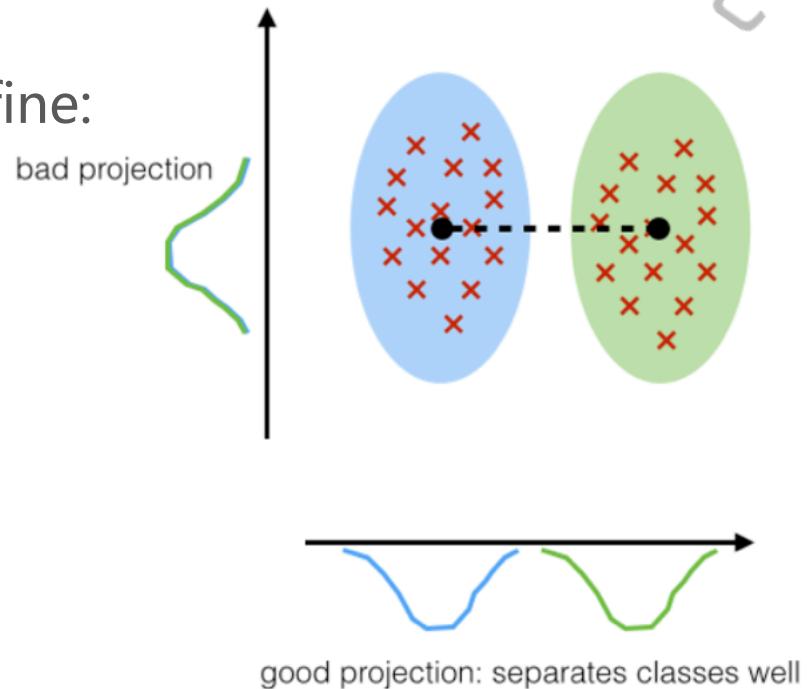
# 一、Review

The basic idea behind LDA:

- Learning a projection matrix  $W$  so that the **within-class data points** are as **close** as possible and **between-class data points** as **far** as possible.

With the help of label, we could define:

- within-class scatter
- between-class scatter matrix



$y, \tilde{\mu}_i \in \mathcal{R}^d$

$\mathbf{x} \in \mathcal{R}^D$

$\mathbf{w} \in \mathcal{R}^{D \times d}$

$y \in \mathcal{R}^d$

$\mathbf{S}_i, \mathbf{S}_B \in \mathcal{R}^{D \times D}$

# 一、Review

Binary Class:

Define the **within-class scatter** for projected samples by  $\tilde{s}_1^2 + \tilde{s}_2^2$ , where

$$\tilde{s}_i^2 = \sum_{y \in \mathcal{Y}_i} (y - \tilde{\mu}_i)^2 = \mathbf{w}^\top \underbrace{\left[ \sum_{\mathbf{x} \in C_i} (\mathbf{x} - \boldsymbol{\mu}_i)(\mathbf{x} - \boldsymbol{\mu}_i)^\top \right]}_{\mathbf{S}_i} \mathbf{w}.$$

and **between-class scatter matrix**:

$$\mathbf{S}_B = (\boldsymbol{\mu}_1 - \boldsymbol{\mu}_2)(\boldsymbol{\mu}_1 - \boldsymbol{\mu}_2)^\top$$

where

$$y = \mathbf{w}^\top \mathbf{x} \quad \boldsymbol{\mu}_i = \frac{1}{N_i} \sum_{\mathbf{x} \in C_i} \mathbf{x}. \quad \tilde{\mu}_i = \frac{1}{N_i} \sum_{y \in \mathcal{Y}_i} y = \frac{1}{N_i} \sum_{\mathbf{x} \in C_i} \mathbf{w}^\top \mathbf{x} = \mathbf{w}^\top \boldsymbol{\mu}_i.$$

# 一、Review

LDA aims to find

$$\mathbf{w} = \arg \max_{\mathbf{w}} \frac{|\tilde{\mu}_1 - \tilde{\mu}_2|^2}{\tilde{s}_1^2 + \tilde{s}_2^2} = \arg \max_{\mathbf{w}} \frac{\mathbf{w}^\top \mathbf{S}_B \mathbf{w}}{\mathbf{w}^\top \mathbf{S}_W \mathbf{w}},$$

$$\boxed{\arg \max_{\mathbf{w}} \frac{\mathbf{w}^\top \mathbf{S}_B \mathbf{w}}{\mathbf{w}^\top \mathbf{S}_W \mathbf{w}}} \Rightarrow \boxed{\mathbf{S}_B \mathbf{w} = \lambda \mathbf{S}_W \mathbf{w}}$$

(generalized eigenvalue problem).

Clearly, the optimal  $\mathbf{w}$  consists of  $K-1$  largest eigenvector of  $\mathbf{S}_W^{-1} \mathbf{S}_B$

where  $K$  denotes the class number.

$y, \tilde{\mu}_i \in \mathcal{R}^d$

$\mathbf{x} \in \mathcal{R}^D$

$\mathbf{w} \in \mathcal{R}^{D \times d}$

$y \in \mathcal{R}^d$

$\mathbf{S}_W, \mathbf{S}_B \in \mathcal{R}^{D \times D}$

$\tilde{\mathbf{S}}_W, \tilde{\mathbf{S}}_B \in \mathcal{R}^{d \times d}$

# 一、Review

## Extension to Multiple Classes:

Within-class scatter matrix

$$\mathbf{S}_W = \sum_{i=1}^K \sum_{\mathbf{x} \in C_i} (\mathbf{x} - \boldsymbol{\mu}_i)(\mathbf{x} - \boldsymbol{\mu}_i)^\top.$$

Between-class scatter matrix

$$\mathbf{S}_B = \sum_{i=1}^K \sum_{C_i} (\boldsymbol{\mu}_i - \boldsymbol{\mu})(\boldsymbol{\mu}_i - \boldsymbol{\mu})^\top = \sum_{i=1}^K N_i (\boldsymbol{\mu}_i - \boldsymbol{\mu})(\boldsymbol{\mu}_i - \boldsymbol{\mu})^\top.$$

Total scatter matrix:  $\mathbf{S}_T = \mathbf{S}_W + \mathbf{S}_B$

$$\mathbf{S}_T = \sum_{\mathbf{x}} (\mathbf{x} - \boldsymbol{\mu})(\mathbf{x} - \boldsymbol{\mu})^\top.$$

$$\text{Rank}(\mathbf{S}_B) \leq K - 1, \quad \text{Rank}(\mathbf{S}_W) \leq N - K, \quad \text{Rank}(\mathbf{S}_T) \leq N - 1.$$

# 一、Review

FLD seeks  $K - 1$  discriminant functions  $\mathbf{W}$  such that  $\mathbf{y} = \mathbf{W}^\top \mathbf{x}$ :

$$\begin{aligned}\mathbf{W} &= \arg \max_{\mathbf{W}} \mathcal{J}_{FLD} \\ &= \arg \max_{\mathbf{W}} \text{tr} \left\{ \tilde{\mathbf{S}}_W^{-1} \tilde{\mathbf{S}}_B \right\} \\ &= \arg \max_{\mathbf{W}} \text{tr} \left\{ \left( \mathbf{W}^\top \mathbf{S}_W \mathbf{W} \right)^{-1} \left( \mathbf{W}^\top \mathbf{S}_B \mathbf{W} \right) \right\},\end{aligned}$$

leading to

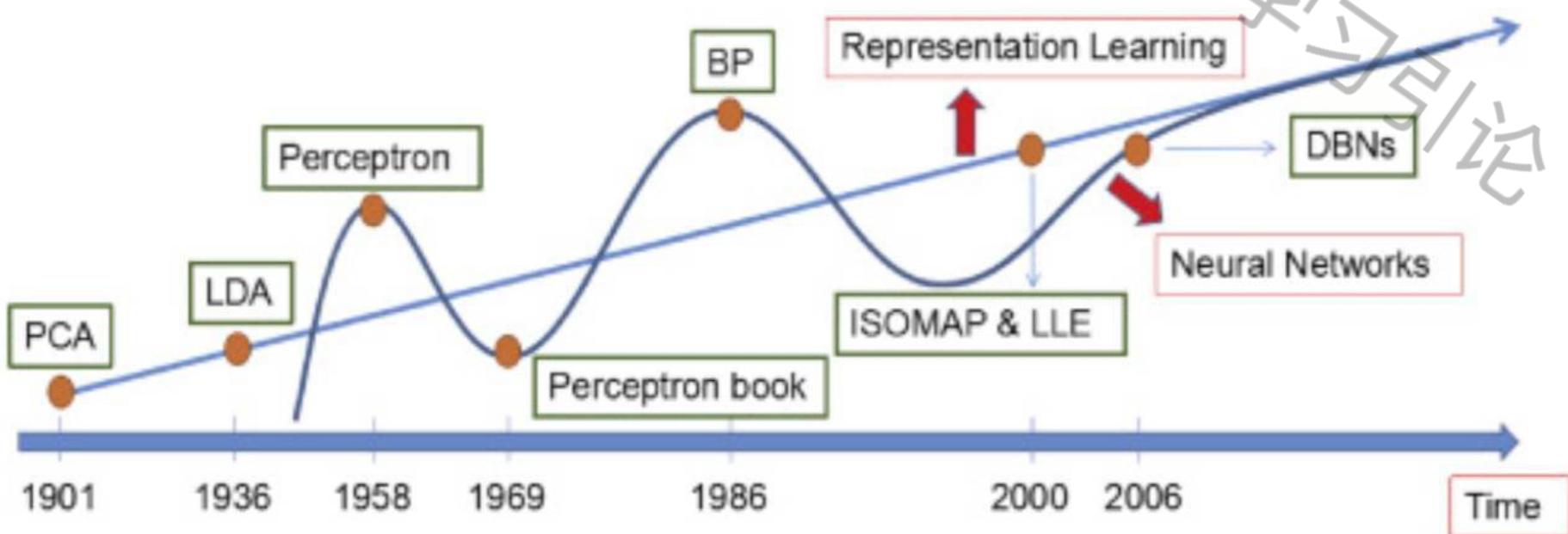
$$\boxed{\arg \max_{\mathbf{W}} \mathcal{J}_{FLD}} \Rightarrow \boxed{\mathbf{S}_B \mathbf{w}_i = \lambda_i \mathbf{S}_W \mathbf{w}_i}.$$

generalized eigenvalue problem

# 提纲

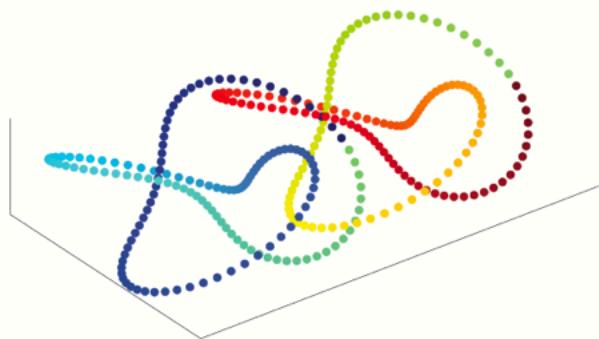
- 一 . Review
- 二 . Locally Linear Embedding
- 三 . Laplacian Eigenmap

## 二、 LLE

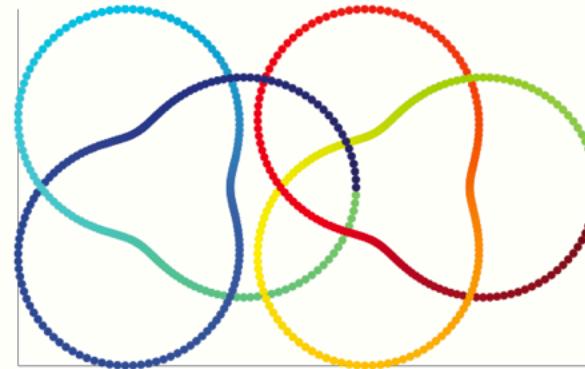


## 二、 LLE

### Nonlinear Dimensionality Reduction



(a)



(b)

- Sam T. Roweis<sup>1</sup> and Lawrence K. Saul, Nonlinear Dimensionality Reduction by Locally Linear Embedding, Science, 2000;
- Lawrence K. Saul and Sam T. Roweis<sup>1</sup>, Think Globally, Fit Locally- Unsupervised Learning of Low Dimensional Manifolds, JMLR2003.

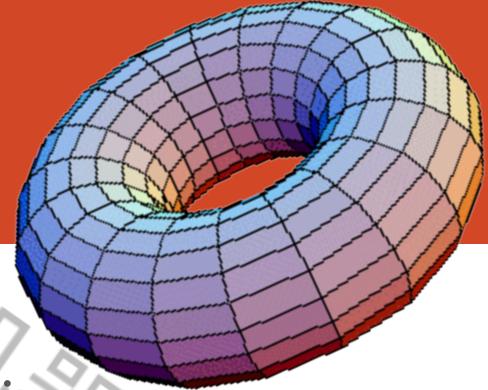
# Tips:

Sam Roweis is an Associate Professor in the Department of Computer Science at New York University. His research interests are in machine learning, data mining, and statistical signal processing. Roweis did his undergraduate degree at the University of Toronto in the Engineering Science program and earned his doctoral degree in 1999 from the California Institute of Technology working with John Hopfield. He did a postdoc with Geoff Hinton and Zoubin Ghahramani at the Gatsby Unit in London. He then was at the University of Toronto from 2001-2009 as an assistant and later associate professor. He was a visiting faculty member at MIT in 2005. He has also worked at several industrial research labs including Google, Bell Labs, Whizbang! Labs and Microsoft. He has won several awards and was the holder of a Canada Research Chair in Statistical Machine Learning, a Sloan Research Fellowship, a Premier's Research Excellence Award, and is still a Scholar of the Canadian Institute for Advanced Research.



Born: April 27, 1972  
Died: January 12, 2010, Washington Square Village

## 二、LLE



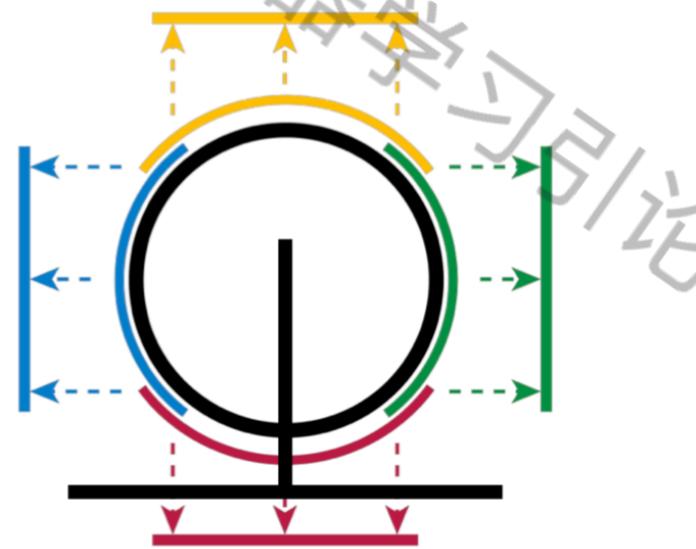
Manifold Learning : 流形，多样体 ( many fold 日文 ) :

- 文天祥《正气歌》：天地有正气，杂然赋流形，下则为河岳，上则为日星<sup>1</sup>
- 易经《彖》：大哉乾元，万物资始，乃统天。云行雨施，品物流形。
- 流形，各种形态，指宇宙间所有的一切<sup>2</sup>。
- 黎曼，一个高维空间靠多重延伸“流”出从而量出一个高维空间的度量进而知其形状（内蕴几何）。

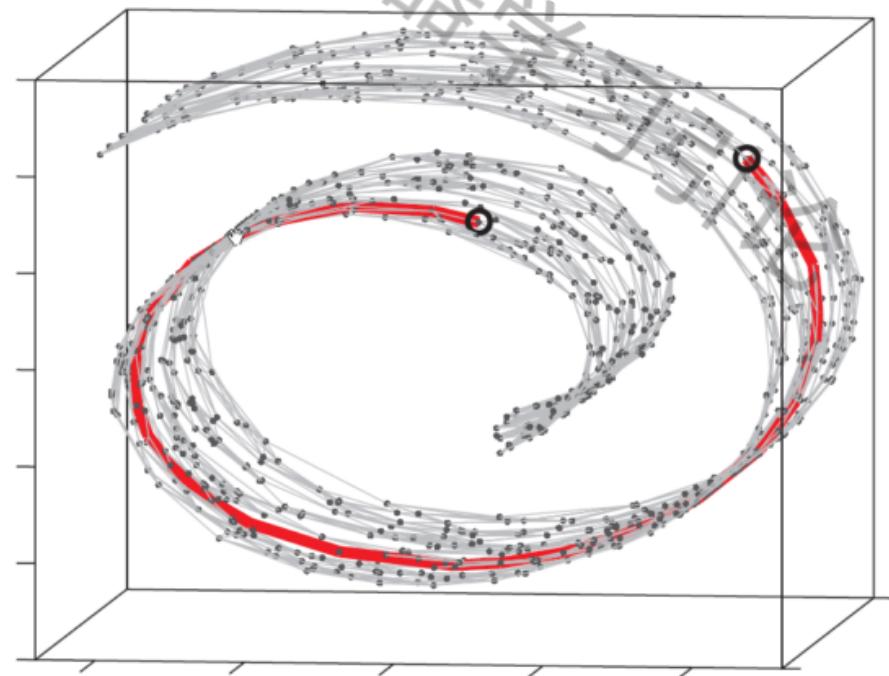
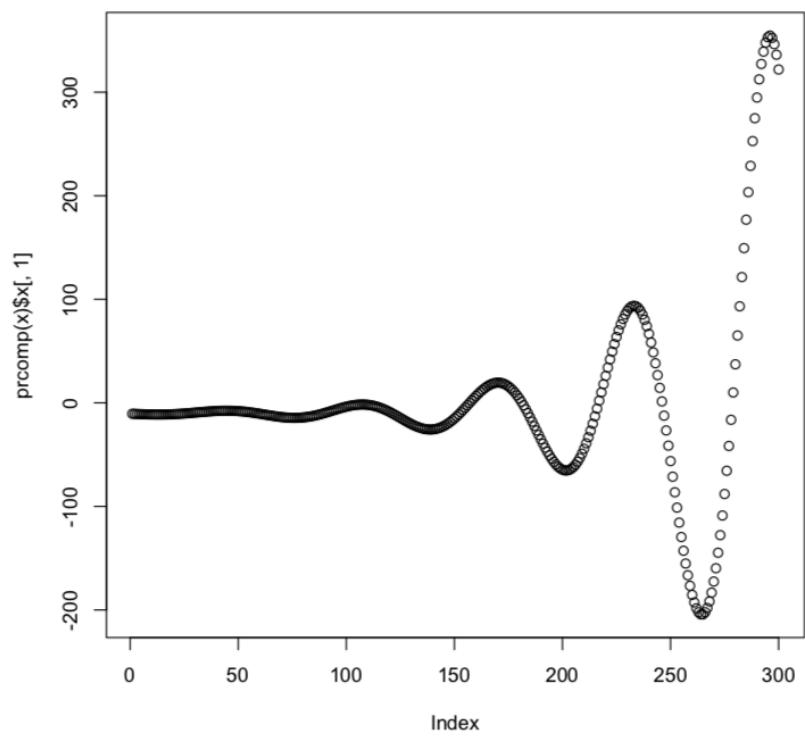
1. 江泽涵：《我国数学名词的早期工作》（《数学通报》1980年12期）
2. 朱东润主编《中国历代文学作品选》

## 二、LLE

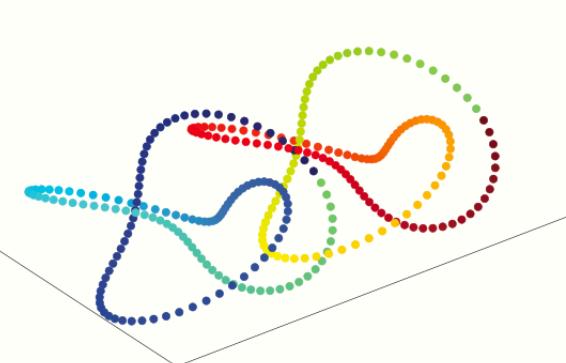
- ▶ Manifold = Many + Fold, 很多曲面片的叠加
- ▶ 叠加但不是拼接，不自交
- ▶ 欧氏空间属于流形
- ▶ 任何一个流形都可以嵌入到足够高维度的欧氏空间中  
(Whitney 嵌入定理)



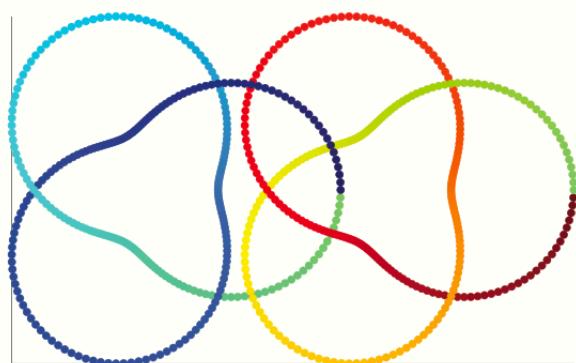
## 二、LLE



## 二、 LLE

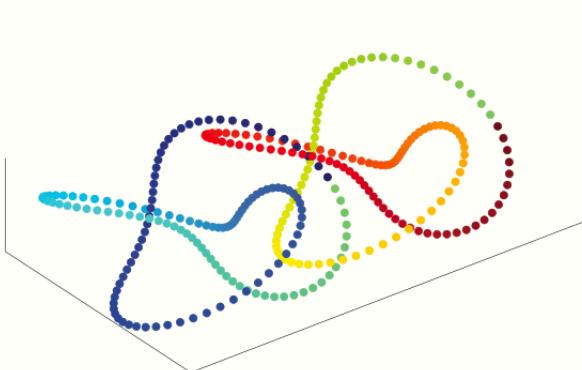


(a)

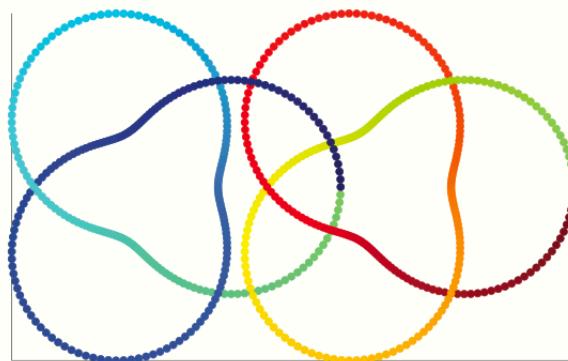


(b)

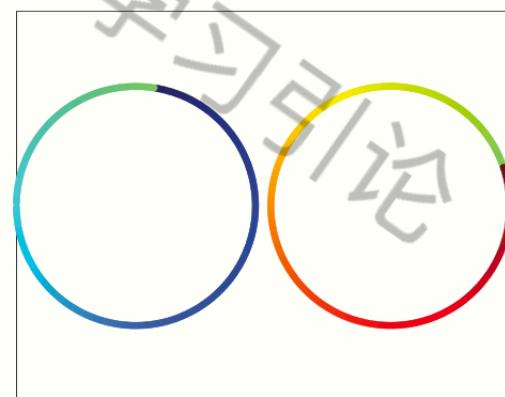
## 二、LLE



(a)



(b)



(c)

## 二、 LLE

- A manifold is a topological space which is **locally Euclidean**.
- Euclidean space is a simplest example of a manifold.
- The dimension of a manifold is the minimum integer number of coordinates necessary to identify each point in that manifold.

$$\mathbf{x}_i \in \mathcal{R}^D,$$

$$\mathbf{y}_i \in \mathcal{R}^d,$$

$$\mathbf{D}_i \in \mathcal{R}^{D \times k},$$

$$\hat{\mathbf{D}}_i \in \mathcal{R}^{d \times k},$$

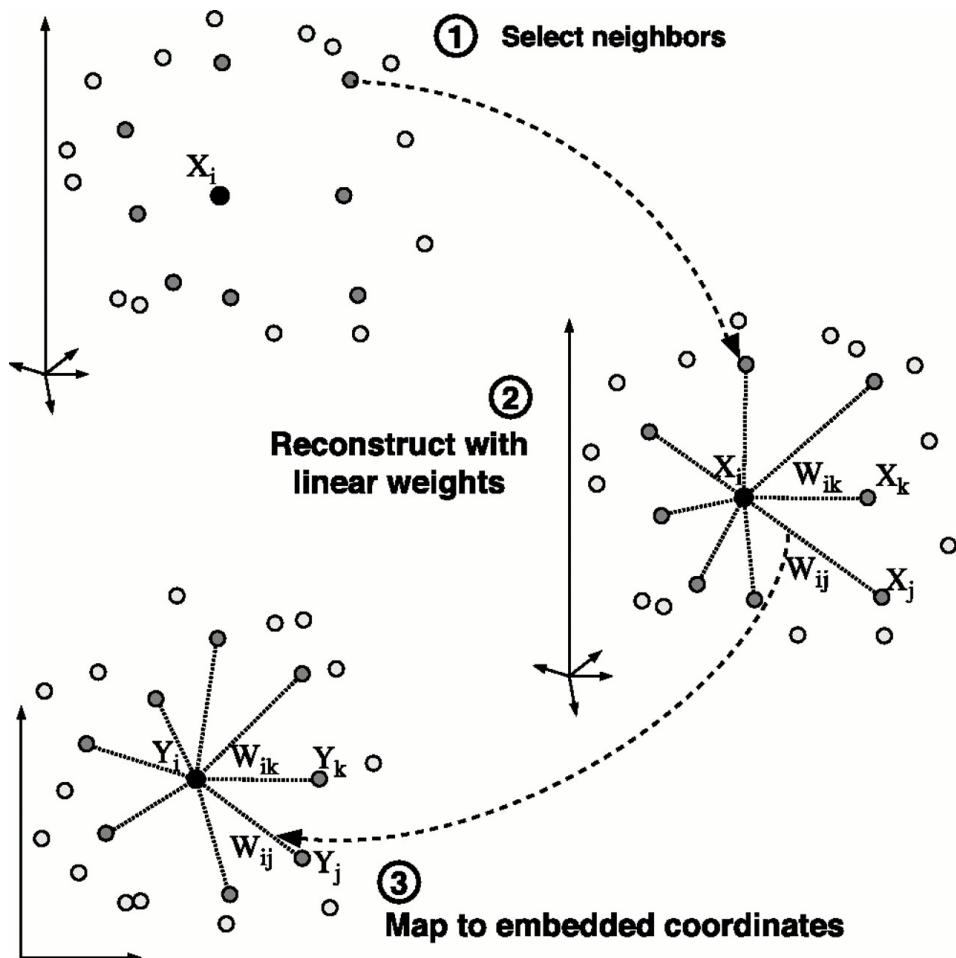
$$\mathbf{W}_{ij} \in \mathcal{R}^1,$$

$$\mathbf{D}_{ij} \in \mathcal{R}^D,$$

$$\mathbf{W}_i \in \mathcal{R}^k$$

$$\hat{\mathbf{D}}_{ij} \in \mathcal{R}^d$$

## 二、 LLE



- **Locally**: for each data point  $\mathbf{x}_i$ , finding its  $k$  nearest neighbors  $\mathbf{D}_i$ 
  - | To find a set of Euclidean space
  - | because a manifold is a topological space which is **locally Euclidean**.

$$\mathbf{x}_i \in \mathcal{R}^D,$$

$$\mathbf{y}_i \in \mathcal{R}^d,$$

$$\mathbf{D}_i \in \mathcal{R}^{D \times k},$$

$$\hat{\mathbf{D}}_i \in \mathcal{R}^{d \times k},$$

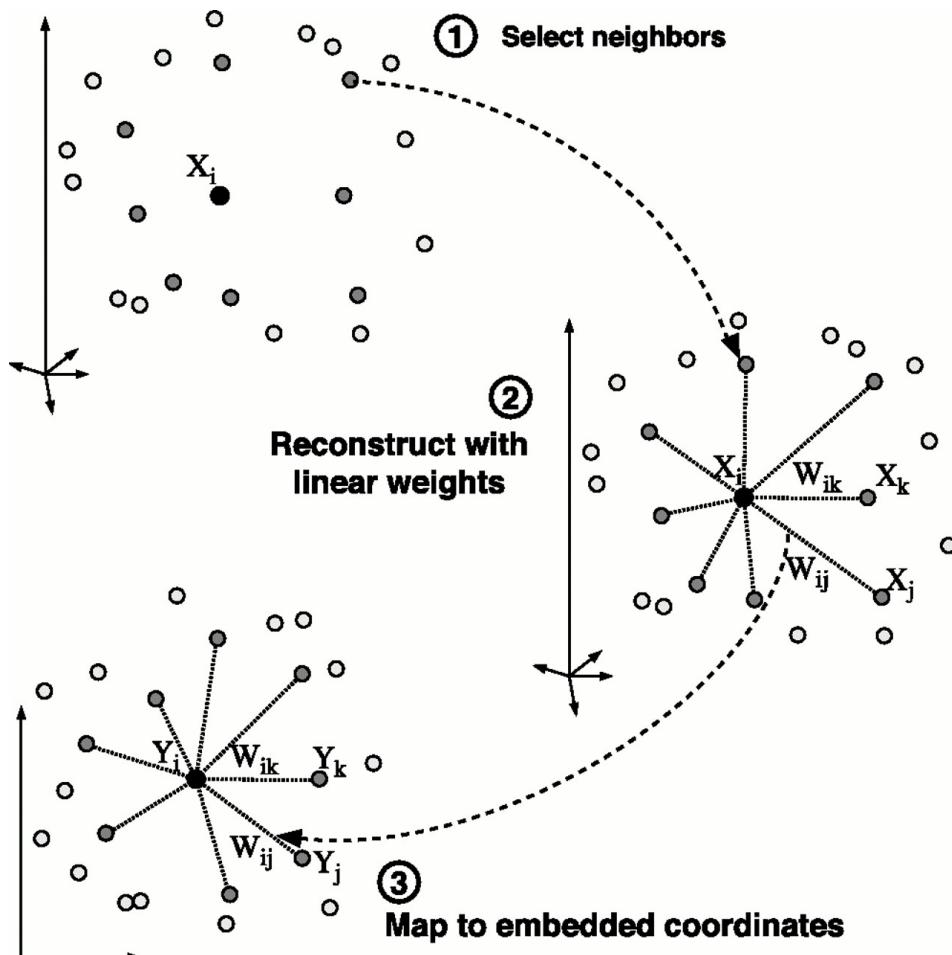
$$\mathbf{W}_{ij} \in \mathcal{R}^1,$$

$$\mathbf{D}_{ij} \in \mathcal{R}^D,$$

$$\mathbf{W}_i \in \mathcal{R}^k$$

$$\hat{\mathbf{D}}_{ij} \in \mathcal{R}^d$$

## 二、 LLE



- Locally: for each data point  $\mathbf{x}_i$ , finding its  $k$  nearest neighbors  $\mathbf{D}_i$
- Linear**: compute the linear reconstruction coefficient  $\mathbf{W}_{ij}$  w.r.t.  $\mathbf{D}_i$  via

$$\min_{\mathbf{W}_{ij}} \left\| \mathbf{x}_i - \sum_{j=1}^k \mathbf{W}_{ij} \mathbf{D}_{ij} \right\|_2^2$$

$$\text{s.t. } \sum_{j=1}^k \mathbf{W}_{ij} = 1$$

As  $\mathbf{x}_i$  and  $\mathbf{D}_i$  lies on the Euclidean space, there could be linearly represented (by  $\mathbf{W}_{ij}$ ) each other thanks to the property of linear space.

Here, the neighborhood size  $k$  should be larger than the intrinsic dimension of manifold.

$$\mathbf{x}_i \in \mathcal{R}^D,$$

$$\mathbf{y}_i \in \mathcal{R}^d,$$

$$\mathbf{D}_i \in \mathcal{R}^{D \times k},$$

$$\hat{\mathbf{D}}_i \in \mathcal{R}^{d \times k},$$

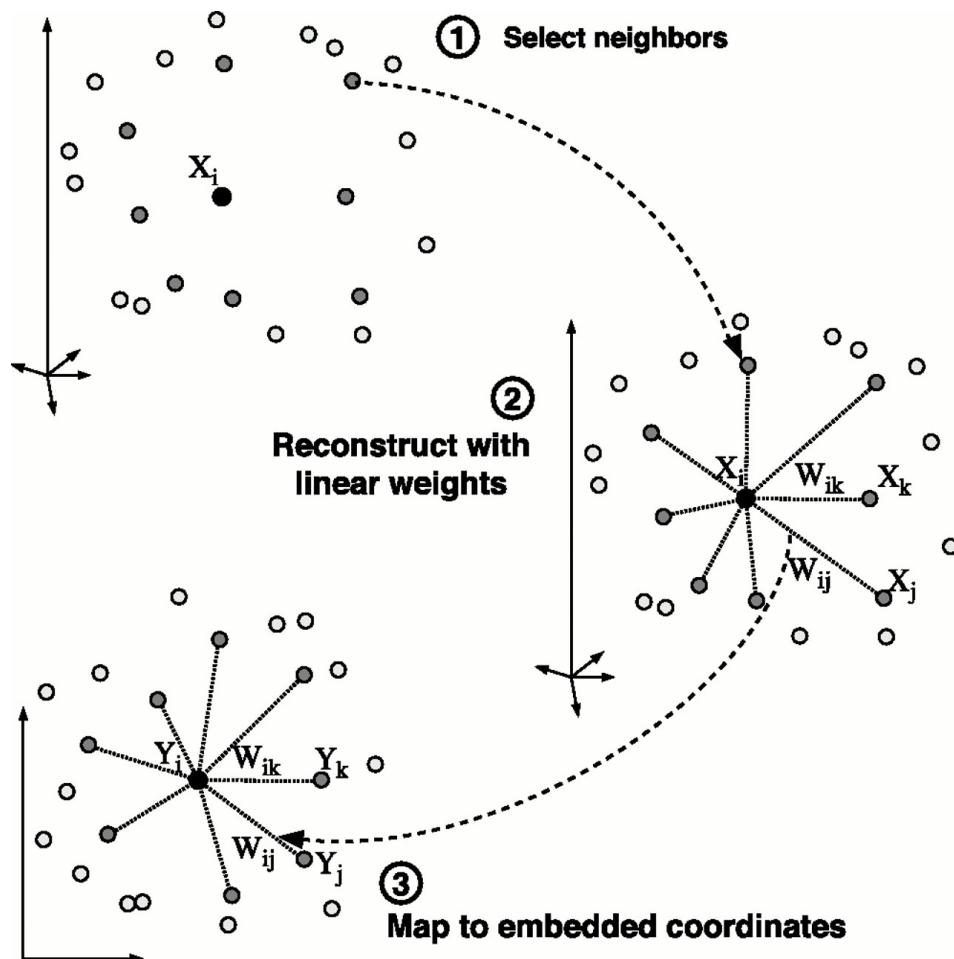
$$\mathbf{W}_{ij} \in \mathcal{R}^1,$$

$$\mathbf{D}_{ij} \in \mathcal{R}^D,$$

$$\mathbf{W}_i \in \mathcal{R}^k$$

$$\hat{\mathbf{D}}_{ij} \in \mathcal{R}^d$$

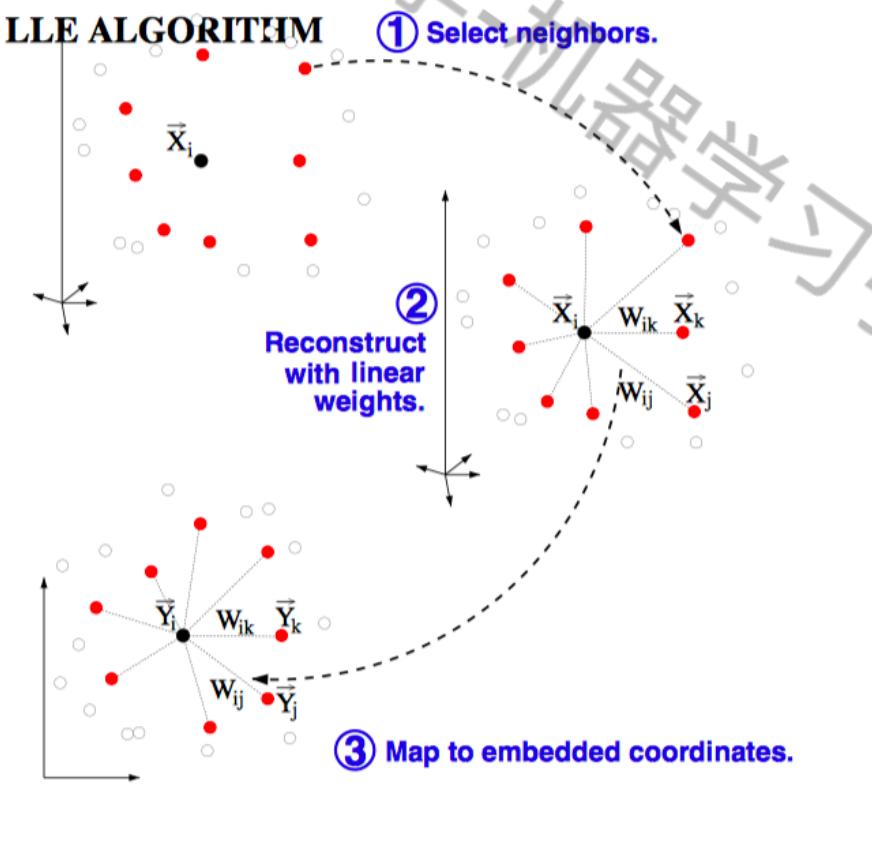
## 二、 LLE



- Locally: for each data point  $\mathbf{x}_i$ , finding its  $k$  nearest neighbors  $\mathbf{D}_i$
  - Linear: compute the linear reconstruction coefficient  $\mathbf{W}_{ij}$  w.r.t.  $\mathbf{D}_i$  via
- $$\min_{\mathbf{W}_{ij}} \left\| \mathbf{x}_i - \sum_{j=1}^k \mathbf{W}_{ij} \mathbf{D}_{ij} \right\|_2^2$$
- $$\text{s.t. } \sum_{j=1}^k \mathbf{W}_{ij} = 1$$
- Embedding**: using  $\mathbf{W}$  as an invariance for DR by embedding it into the manifold via:
- $$\min_{\mathbf{Y}} \left\| \mathbf{y}_i - \sum_{j=1}^k \mathbf{W}_{ij} \hat{\mathbf{D}}_{ij} \right\|_2^2$$
- $$\text{s.t. } \mathbf{y}_i^T \mathbf{y}_i = 1$$
- $\hat{\mathbf{D}}_i$  denotes the neighbors of  $\mathbf{x}_i$  in the projection space.

## 二、LLE

1. Compute the neighbors of each data point,  $\vec{X}_i$ .
2. Compute the weights  $W_{ij}$  that best reconstruct each data point  $\vec{X}_i$  from its neighbors, minimizing the cost in Equation (1) by constrained linear fits.
3. Compute the vectors  $\vec{Y}_i$  best reconstructed by the weights  $W_{ij}$ , minimizing the quadratic form in Equation (2) by its bottom nonzero eigenvectors.



A manifold is a topological space which is **locally Euclidean**, based on which LLE

- Computes the **locally linear** reconstruction coefficient in the input space;
- Using the obtained coefficient as a prior/invariance to find the **low dimensional space** which is also locally Euclidean.

## 二、 LLE

- Locally: for each data point  $\mathbf{x}_i$ , finding its  $k$  nearest neighbors  $\mathbf{D}_i$
- Linear: compute the linear reconstruction coefficient  $\mathbf{W}_{ij}$  w.r.t.  $\mathbf{D}_i$  via

$$\min_W \|\mathbf{x}_i - \sum_{j=1}^k \mathbf{W}_{ij} \mathbf{D}_{ij}\|_2^2 \quad \text{s.t. } \sum_{j=1}^k \mathbf{W}_{ij} = 1$$

The constraint is to enforce the data points lies into the affine space – a special variant of Euclidean space, which enjoys following advantages:

- Rotation/scale invariance: Let the matrix  $\mathbf{R}$  be the rotation/scale operation on  $\mathbf{X}$ ,

$$\mathbf{R}\mathbf{x}_i = \sum_{j=1}^k \mathbf{V}_{ij} \mathbf{R}\mathbf{D}_{ij}$$

$$\rightarrow \mathbf{R}^\dagger \mathbf{R}\mathbf{x}_i = \sum_{j=1}^k \mathbf{R}^\dagger \mathbf{V}_{ij} \mathbf{R}\mathbf{D}_j \rightarrow \mathbf{x}_i = \sum_{j=1}^k \mathbf{W}_{ij} \mathbf{D}_{ij} = \sum_{j=1}^k \mathbf{V}_{ij} \mathbf{D}_{ij} \rightarrow \mathbf{W}_{ij} = \mathbf{V}_{ij}$$

## 二、 LLE

- Locally: for each data point  $\mathbf{x}_i$ , finding its  $k$  nearest neighbors  $\mathbf{D}_i$
- Linear: compute the linear reconstruction coefficient  $\mathbf{W}_{ij}$  w.r.t.  $\mathbf{D}_i$  via

$$\min_{\mathbf{W}} \left\| \mathbf{x}_i - \sum_{j=1}^k \mathbf{W}_{ij} \mathbf{D}_{ij} \right\|_2^2 \quad \text{s.t. } \sum_{j=1}^k \mathbf{W}_{ij} = 1$$

The constraint is to enforce the data points lies into the affine space – a special variant of Euclidean space, which enjoys following advantages:

$$\begin{aligned}\hat{\mathbf{W}} &= \arg \min_{\sum_j w_{ij}=1} \sum_i |(\mathbf{x}_i - \mathbf{a}) - \sum_j w_{ij} (\mathbf{x}_{ij} - \mathbf{a})|^2 \\ &= \arg \min_{\sum_j w_{ij}=1} \sum_i |(\mathbf{x}_i - \mathbf{a}) - (\sum_j w_{ij} \mathbf{x}_{ij}) + (\sum_j w_{ij} \mathbf{a})|^2 \\ &= \arg \min_{\sum_j w_{ij}=1} \sum_i |(\mathbf{x}_i - \mathbf{a}) - (\sum_j w_{ij} \mathbf{x}_{ij}) + \mathbf{a}|^2 \\ &= \arg \min_{\sum_j w_{ij}=1} \sum_i |\mathbf{x}_i - \sum_j w_{ij} \mathbf{x}_{ij}|^2 \\ &= \mathbb{W}\end{aligned}$$

Translation  
invariance

## 二、 LLE

- Locally: for each data point  $\mathbf{x}_i$ , finding its  $k$  nearest neighbors  $\mathbf{D}_i$
- Linear: compute the linear reconstruction coefficient  $\mathbf{W}_{ij}$  w.r.t.  $\mathbf{D}_i$  via

$$\min_{\mathbf{W}} \left\| \mathbf{x}_i - \sum_{j=1}^k \mathbf{W}_{ij} \mathbf{D}_{ij} \right\|_2^2 \quad \text{s.t. } \sum_{j=1}^k \mathbf{W}_{ij} = 1$$

with following optimization process:

$$\begin{aligned} \mathcal{L} &= \left\| \mathbf{x}_i - \sum_{j=1}^k \mathbf{W}_{ij} \mathbf{D}_{ij} \right\|_2^2 + \lambda \left( 1 - \sum_{j=1}^k \mathbf{W}_{ij} \right) = \boxed{\left\| \sum_{j=1}^k \mathbf{W}_{ij} \mathbf{x}_i - \sum_{j=1}^k \mathbf{W}_{ij} \mathbf{D}_{ij} \right\|_2^2 + \lambda \left( 1 - \sum_{j=1}^k \mathbf{W}_{ij} \right)} \\ &= \left\| \sum_{j=1}^k (\mathbf{x}_i - \mathbf{D}_{ij}) \mathbf{W}_{ij} \right\|_2^2 + \lambda \left( 1 - \sum_{j=1}^k \mathbf{W}_{ij} \right) = \mathbf{w}^T (\mathbf{X}_i - \mathbf{D}_i)^T (\mathbf{X}_i - \mathbf{D}_i) \mathbf{w} + \lambda \left( 1 - \mathbf{1}^T \mathbf{w} \right) \end{aligned}$$

## 二、 LLE

- Locally: for each data point  $\mathbf{x}_i$ , finding its  $k$  nearest neighbors  $\mathbf{D}_i$
- Linear: compute the linear reconstruction coefficient  $\mathbf{W}_{ij}$  w.r.t.  $\mathbf{D}_i$  via

$$\min_{\mathbf{W}} \left\| \mathbf{x}_i - \sum_{j=1}^k \mathbf{W}_{ij} \mathbf{D}_{ij} \right\|_2^2 \quad \text{s.t. } \sum_{j=1}^k \mathbf{W}_{ij} = 1$$

with following optimization process:

$$\begin{aligned} \mathcal{L} &= \left\| \mathbf{x}_i - \sum_{j=1}^k \mathbf{W}_{ij} \mathbf{D}_{ij} \right\|_2^2 + \lambda \left( 1 - \sum_{j=1}^k \mathbf{W}_{ij} \right) = \left\| \sum_{j=1}^k \mathbf{W}_{ij} \mathbf{x}_i - \sum_{j=1}^k \mathbf{W}_{ij} \mathbf{D}_{ij} \right\|_2^2 + \lambda \left( 1 - \sum_{j=1}^k \mathbf{W}_{ij} \right) \\ &= \boxed{\left\| \sum_{j=1}^k (\mathbf{x}_i - \mathbf{D}_{ij}) \mathbf{W}_{ij} \right\|_2^2 + \lambda \left( 1 - \sum_{j=1}^k \mathbf{W}_{ij} \right)} = \mathbf{w}^T (\mathbf{X}_i - \mathbf{D}_i)^T (\mathbf{X}_i - \mathbf{D}_i) \mathbf{w} + \lambda \left( 1 - \mathbf{1}^T \mathbf{w} \right) \end{aligned}$$

## 二、 LLE

- Locally: for each data point  $\mathbf{x}_i$ , finding its  $k$  nearest neighbors  $\mathbf{D}_i$
- Linear: compute the linear reconstruction coefficient  $\mathbf{W}_{ij}$  w.r.t.  $\mathbf{D}_i$  via

$$\min_{\mathbf{W}} \left\| \mathbf{x}_i - \sum_{j=1}^k \mathbf{W}_{ij} \mathbf{D}_{ij} \right\|_2^2 \quad \text{s.t. } \sum_{j=1}^k \mathbf{W}_{ij} = 1$$

with following optimization process:

$$\begin{aligned} \mathcal{L} &= \left\| \mathbf{x}_i - \sum_{j=1}^k \mathbf{W}_{ij} \mathbf{D}_{ij} \right\|_2^2 + \lambda \left( 1 - \sum_{j=1}^k \mathbf{W}_{ij} \right) = \left\| \sum_{j=1}^k \mathbf{W}_{ij} \mathbf{x}_i - \sum_{j=1}^k \mathbf{W}_{ij} \mathbf{D}_{ij} \right\|_2^2 + \lambda \left( 1 - \sum_{j=1}^k \mathbf{W}_{ij} \right) \\ &= \left\| \sum_{j=1}^k (\mathbf{x}_i - \mathbf{D}_{ij}) \mathbf{W}_{ij} \right\|_2^2 + \lambda \left( 1 - \sum_{j=1}^k \mathbf{W}_{ij} \right) = \boxed{\mathbf{w}^T (\mathbf{X}_i - \mathbf{D}_i)^T (\mathbf{X}_i - \mathbf{D}_i) \mathbf{w} + \lambda (1 - \mathbf{1}^T \mathbf{w})} \end{aligned}$$

w is a vector whose elements  
are  $\mathbf{W}_{ij}$   $\mathbf{X}_i$  is a matrix whose  
column is  $\mathbf{x}_i$

## 二、 LLE

$$\mathbf{G}_i = (\mathbf{X}_i - \mathbf{D}_i)^T (\mathbf{X}_i - \mathbf{D}_i)$$

- Locally: for each data point  $\mathbf{x}_i$ , finding its  $k$  nearest neighbors  $\mathbf{D}_i$
- Linear**: compute the linear reconstruction coefficient  $\mathbf{W}_{ij}$  w.r.t.  $\mathbf{D}_i$  via

$$\min_{\mathbf{W}} \left\| \mathbf{x}_i - \sum_{j=1}^k \mathbf{W}_{ij} \mathbf{D}_{ij} \right\|_2^2 \quad \text{s.t. } \sum_{j=1}^k \mathbf{W}_{ij} = 1$$

with following optimization process:

$$\begin{aligned} \mathcal{L} &= \left\| \mathbf{x}_i - \sum_{j=1}^k \mathbf{W}_{ij} \mathbf{D}_{ij} \right\|_2^2 + \lambda \left( 1 - \sum_{j=1}^k \mathbf{W}_{ij} \right) = \left\| \sum_{j=1}^k \mathbf{W}_{ij} \mathbf{x}_i - \sum_{j=1}^k \mathbf{W}_{ij} \mathbf{D}_{ij} \right\|_2^2 + \lambda \left( 1 - \sum_{j=1}^k \mathbf{W}_{ij} \right) \\ &= \left\| \sum_{j=1}^k (\mathbf{x}_i - \mathbf{D}_{ij}) \mathbf{W}_{ij} \right\|_2^2 + \lambda \left( 1 - \sum_{j=1}^k \mathbf{W}_{ij} \right) = \boxed{\mathbf{w}^T (\mathbf{X}_i - \mathbf{D}_i)^T (\mathbf{X}_i - \mathbf{D}_i) \mathbf{w} + \lambda (1 - \mathbf{1}^T \mathbf{w})} \end{aligned}$$

Let  $\frac{\partial \mathcal{L}}{\partial \mathbf{W}} = 0$ , we have  $(\mathbf{X}_i - \mathbf{D}_i)^T (\mathbf{X}_i - \mathbf{D}_i) \mathbf{w} = \lambda \mathbf{1}$

$$\rightarrow \mathbf{w} = \frac{\mathbf{G}_i^\dagger \mathbf{1}}{\mathbf{1}^T \mathbf{G}_i^\dagger \mathbf{1}}$$

$\lambda$  is a constant which is used to achieve the constraint.

Note that, a small number will be added onto the main diagonal entries of  $\mathbf{G}_i^\dagger$  for nonsingularity.

## 二、 LLE

- Locally: for each data point  $\mathbf{x}_i$ , finding its  $k$  nearest neighbors  $\mathbf{D}_i$
- Linear: compute the linear reconstruction coefficient  $\mathbf{W}_{ij}$  w.r.t.  $\mathbf{D}_i$  via

$$\min_{\mathbf{W}} \left\| \mathbf{x}_i - \sum_{j=1}^k \mathbf{W}_{ij} \mathbf{D}_{ij} \right\|_2^2 \quad \text{s.t. } \sum_{j=1}^k \mathbf{W}_{ij} = 1$$

- **Embedding**: using  $\mathbf{W}$  as an invariance for DR by embedding it into a low dimensional space via:

$$\begin{aligned} & \min_{\mathbf{y}_i} \sum_{i=1}^i \left\| \mathbf{y}_i - \sum_{j=1}^k \mathbf{W}_{ij} \hat{\mathbf{D}}_{ij} \right\|_2^2 \quad \text{s.t. } \mathbf{y}_i^T \mathbf{y}_i = 1 \\ & \rightarrow \min_{\mathbf{Y}} \left\| \mathbf{Y} - \mathbf{Y} \mathbf{W} \right\|_F^2 \quad \text{s.t. } \text{tr}(\mathbf{Y}^T \mathbf{Y}) = 1 \end{aligned}$$

## 二、 LLE

- Locally: for each data point  $\mathbf{x}_i$ , finding its  $k$  nearest neighbors  $\mathbf{D}_i$
- Linear: compute the linear reconstruction coefficient  $\mathbf{W}_{ij}$  w.r.t.  $\mathbf{D}_i$  via

$$\min_{\mathbf{W}} \left\| \mathbf{x}_i - \sum_{j=1}^k \mathbf{W}_{ij} \mathbf{D}_{ij} \right\|_2^2 \quad \text{s.t. } \sum_{j=1}^k \mathbf{W}_{ij} = 1$$

- Embedding: using  $\mathbf{W}$  as an invariance for DR by embedding it into a low dimensional space via:

$$\min_{\mathbf{y}_i} \sum_{i=1}^i \left\| \mathbf{y}_i - \sum_{j=1}^k \mathbf{W}_{ij} \hat{\mathbf{D}}_{ij} \right\|_2^2 \quad \text{s.t. } \mathbf{y}_i^T \mathbf{y}_i = 1$$

$$\rightarrow \min_{\mathbf{Y}} \left\| \mathbf{Y} - \mathbf{Y}\mathbf{W} \right\|_F^2 \quad \text{s.t. } \text{tr}(\mathbf{Y}^T \mathbf{Y}) = 1$$

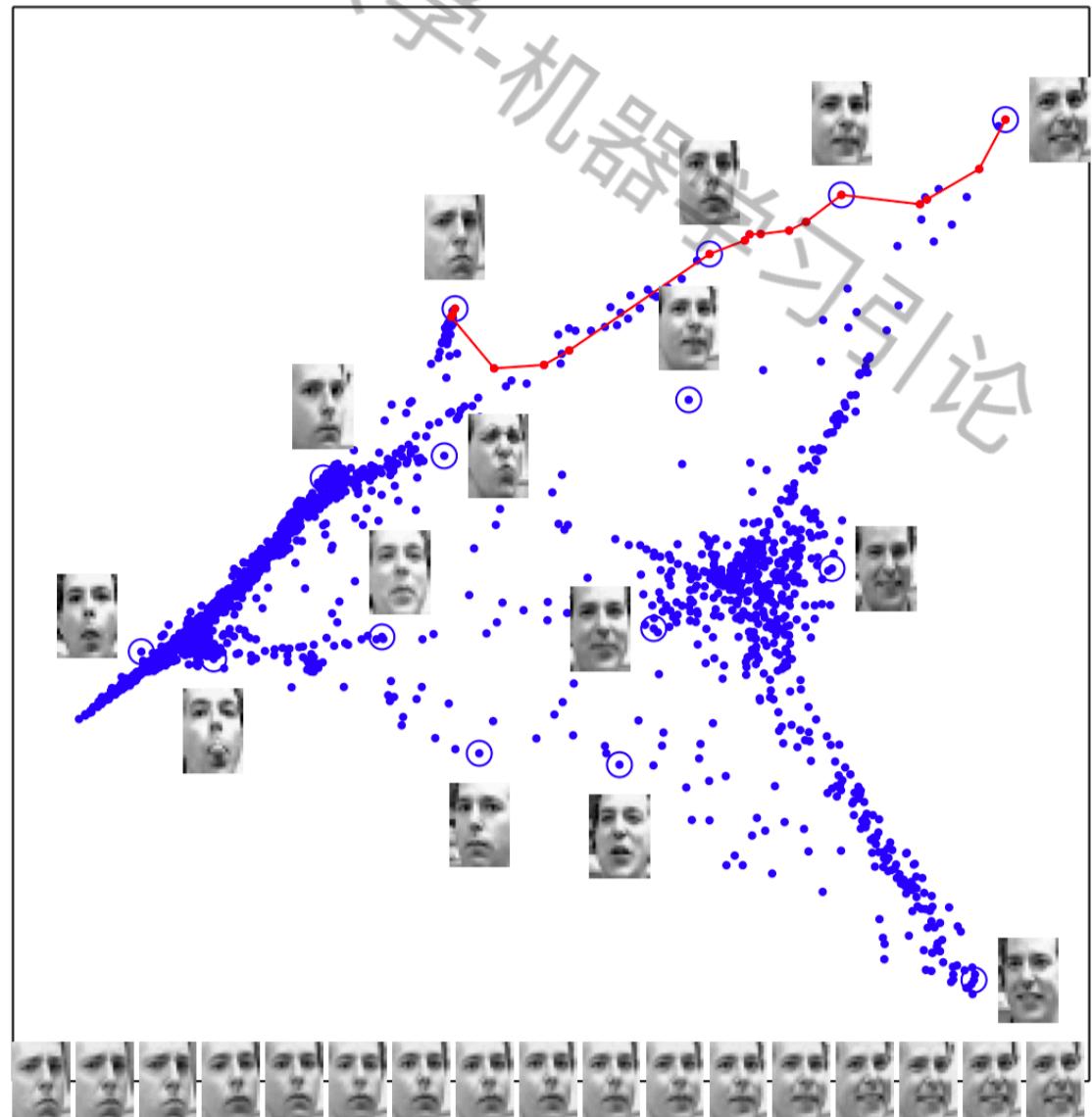
Let  $\mathcal{L} = \left\| \mathbf{Y} - \mathbf{Y}\mathbf{W} \right\|_F^2 + \lambda \text{trace}(\mathbf{I} - \mathbf{Y}^T \mathbf{Y})$  and its derivative w.r.t.  $\mathbf{Y}$  be zero, then we have

$$2(\mathbf{I} - \mathbf{W})(\mathbf{I} - \mathbf{W})^T \mathbf{Y}^T = 2\lambda \mathbf{Y}^T$$

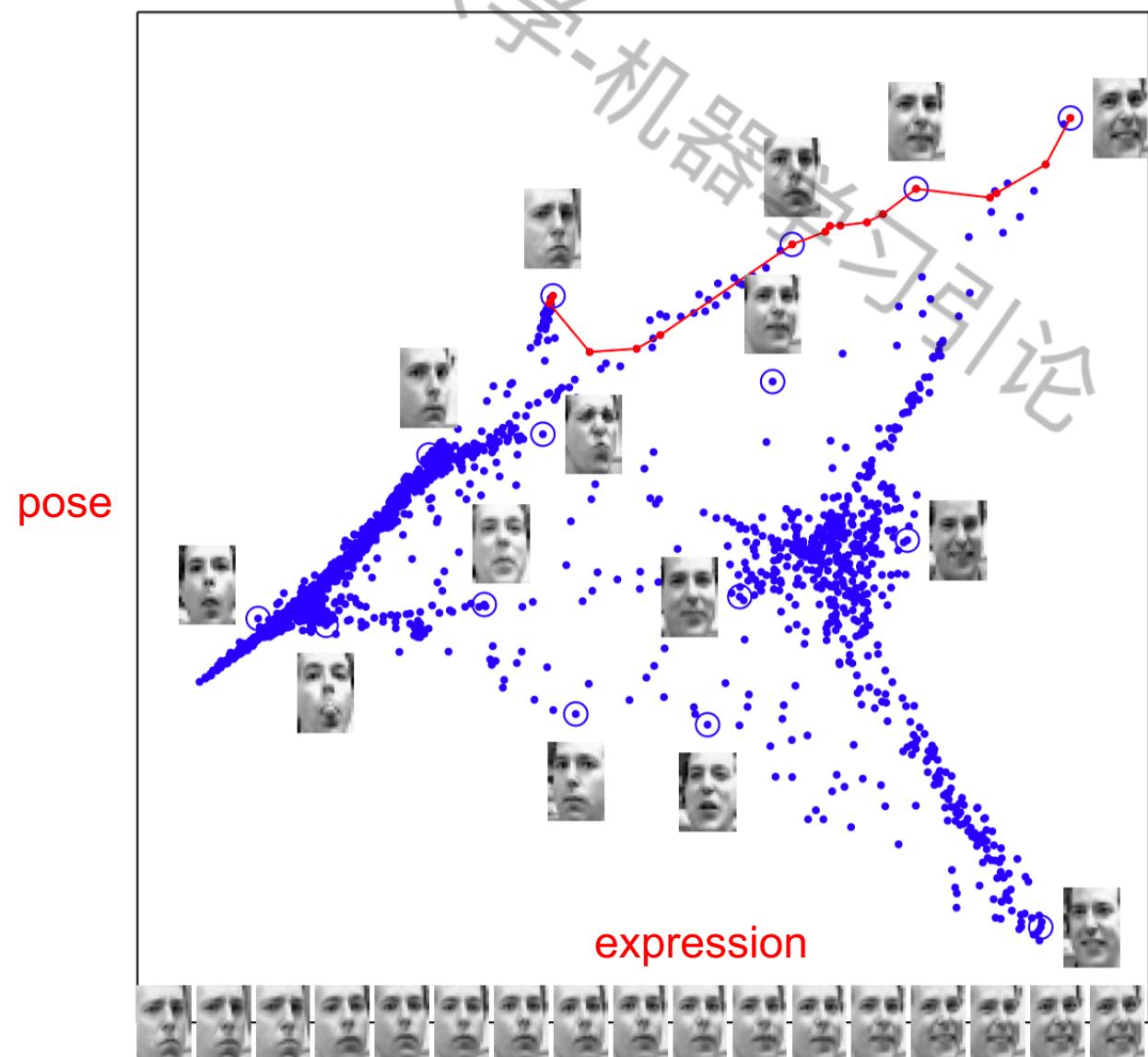
Clearly, the optimal  $\mathbf{Y}$  consists of  $d$  eigenvectors corresponding to  $d$  smallest nonzero eigenvalue of  $(\mathbf{I} - \mathbf{W})(\mathbf{I} - \mathbf{W})^T$

## 二、LLE

How many manifold learned in this example? And what are them?



## 二、LLE



# 提纲

- . Review
- . Locally Linear Embedding
- . Laplacian Eigenmap

M. Belkin and P. Niyogi. “Laplacian eigenmaps for dimensionality reduction and data representation,” *Neural Comput.*, 15(6):1373–1396, 2003.

$$\mathbf{W}_{ij} \in \mathcal{R}^1$$

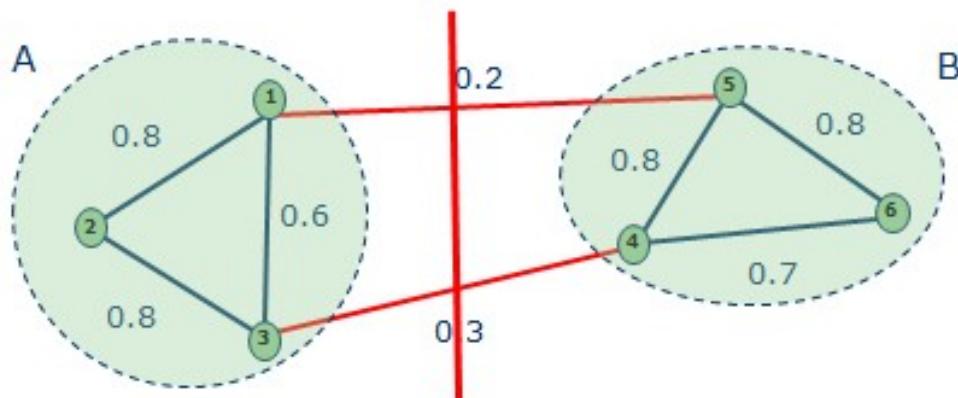
$$\mathbf{x}_i, \mathbf{x}_j \in \mathcal{R}^D$$

$$\epsilon, t > 0$$

### 三、 Laplacian Eigenmap

- Step 1: find k nearest neighbors for each data point
- Step 2: obtain a local invariance by constructing a similarity graph via

$$\mathbf{W}_{ij} = \begin{cases} \exp^{-\frac{\|\mathbf{x}_i - \mathbf{x}_j\|_2^2}{t}} & \text{if } \|\mathbf{x}_i - \mathbf{x}_j\| < \epsilon \text{ or they are knn} \\ 0 & \text{otherwise} \end{cases}$$



$$\mathbf{W}_{ij} \in \mathcal{R}^1$$

$$\mathbf{x}_i, \mathbf{x}_j \in \mathcal{R}^D$$

$$\epsilon, t > 0$$

### 三、 Laplacian Eigenmap

- Step 1: find k nearest neighbors for each data point
- Step 2: obtain a local invariance by constructing a similarity graph via

$$\mathbf{W}_{ij} = \begin{cases} \exp^{-\frac{\|\mathbf{x}_i - \mathbf{x}_j\|_2^2}{t}} & \text{if } \|\mathbf{x}_i - \mathbf{x}_j\| < \epsilon \text{ or they are knn} \\ 0 & \text{otherwise} \end{cases}$$

- Step 3: embed  $\mathbf{W}$  into a low-dimensional space by

$$\min_{\mathbf{Y}} \sum_i \sum_j \mathbf{W}_{ij} \|\mathbf{y}_i - \mathbf{y}_j\|_2^2$$

$$\text{s.t. } \text{tr}(\mathbf{Y} \mathbf{D} \mathbf{Y}^T) = 1$$

$$\mathbf{W}_{ij} \in \mathcal{R}^1$$

$$\mathbf{x}_i, \mathbf{x}_j \in \mathcal{R}^D$$

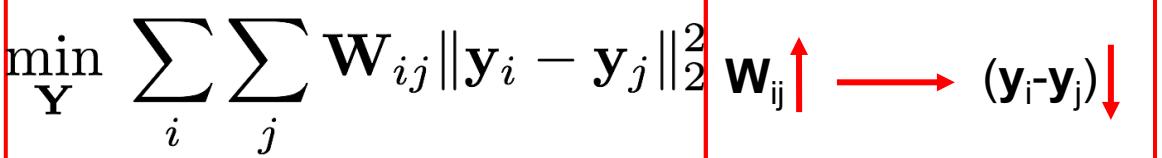
$$\epsilon, t > 0$$

### 三、 Laplacian Eigenmap

- Step 1: find k nearest neighbors for each data point
- Step 2: obtain a local invariance by constructing a similarity graph via

$$\mathbf{W}_{ij} = \begin{cases} \exp^{-\frac{\|\mathbf{x}_i - \mathbf{x}_j\|_2^2}{t}} & \text{if } \|\mathbf{x}_i - \mathbf{x}_j\| < \epsilon \text{ or they are knn} \\ 0 & \text{otherwise} \end{cases}$$

- Step 3: embed  $\mathbf{W}$  into a low-dimensional space by

$$\min_{\mathbf{Y}} \sum_i \sum_j \mathbf{W}_{ij} \|\mathbf{y}_i - \mathbf{y}_j\|_2^2$$


$$\text{s.t. } \text{tr}(\mathbf{Y} \mathbf{D} \mathbf{Y}^T) = 1$$

$$\mathbf{W}_{ij} \in \mathcal{R}^1$$

$$\mathbf{x}_i, \mathbf{x}_j \in \mathcal{R}^D$$

$$\epsilon, t > 0$$

### 三、 Laplacian Eigenmap

- Step 1: find k nearest neighbors for each data point
- Step 2: obtain a local invariance by constructing a similarity graph via

$$\mathbf{W}_{ij} = \begin{cases} \exp^{-\frac{\|\mathbf{x}_i - \mathbf{x}_j\|_2^2}{t}} & \text{if } \|\mathbf{x}_i - \mathbf{x}_j\| < \epsilon \text{ or they are knn} \\ 0 & \text{otherwise} \end{cases}$$

- Step 3: embed  $\mathbf{W}$  into a low-dimensional space by

$$\min_{\mathbf{Y}} \sum_i \sum_j \mathbf{W}_{ij} \|\mathbf{y}_i - \mathbf{y}_j\|_2^2$$

$$\text{s.t. } \text{tr}(\mathbf{YD}\mathbf{Y}^T) = 1$$

The constraint removes an arbitrary scaling factor in the embedding. Matrix D provides a natural measure on the vertices of the graph. The bigger the value  $D_{ii}$  is, the more important is that vertex.

$$\mathbf{D}_{ii} = \sum_j \mathbf{W}_{ij}, i = j$$

$$\mathbf{D}_{ij} = 0, i \neq j$$

### 三、 Laplacian Eigenmap

$$\begin{aligned} & \sum_{i=1}^k \sum_{j=1}^k \| \mathbf{y}_i - \mathbf{y}_j \|^2 W_{ij} \\ &= \sum_{i=1}^k \sum_{j=1}^k (\mathbf{y}_i^\top \mathbf{y}_i - 2\mathbf{y}_i^\top \mathbf{y}_j + \mathbf{y}_j^\top \mathbf{y}_j) W_{ij} \\ &= \sum_{i=1}^k \left( \sum_{j=1}^k W_{ij} \right) \mathbf{y}_i^\top \mathbf{y}_i + \sum_{j=1}^k \left( \sum_{i=1}^k W_{ij} \right) \mathbf{y}_j^\top \mathbf{y}_j \\ &\quad - 2 \sum_{i=1}^k \sum_{j=1}^k \mathbf{y}_i^\top \mathbf{y}_j W_{ij} \\ &= 2 \sum_{i=1}^k D_{ii} \mathbf{y}_i^\top \mathbf{y}_i - 2 \sum_{i=1}^k \sum_{j=1}^k \mathbf{y}_i^\top \mathbf{y}_j W_{ij} \\ &= 2 \sum_{i=1}^k (\sqrt{D_{ii}} \mathbf{y}_i)^\top (\sqrt{D_{ii}} \mathbf{y}_i) - 2 \sum_{i=1}^k \mathbf{y}_i^\top \left( \sum_{j=1}^k \mathbf{y}_j W_{ij} \right) \\ &= 2 \text{Tr}[(\mathbf{Y} \sqrt{\mathbf{D}})(\mathbf{Y} \sqrt{\mathbf{D}})^\top] - 2 \sum_{i=1}^k \mathbf{y}_i^\top (\mathbf{Y} \mathbf{W})_i \\ &= 2 \text{Tr}[\mathbf{Y}(\mathbf{D} - \mathbf{W})\mathbf{Y}^\top] = 2 \text{Tr}[\mathbf{YLY}^\top] \end{aligned}$$

### 三、 Laplacian Eigenmap

$$\begin{aligned} & \sum_{i=1}^k \sum_{j=1}^k \| \mathbf{y}_i - \mathbf{y}_j \|^2 W_{ij} \\ &= \sum_{i=1}^k \sum_{j=1}^k (\mathbf{y}_i^\top \mathbf{y}_i - 2\mathbf{y}_i^\top \mathbf{y}_j + \mathbf{y}_j^\top \mathbf{y}_j) W_{ij} \\ &= \sum_{i=1}^k \left( \sum_{j=1}^k W_{ij} \right) \mathbf{y}_i^\top \mathbf{y}_i + \sum_{j=1}^k \left( \sum_{i=1}^k W_{ij} \right) \mathbf{y}_j^\top \mathbf{y}_j \\ &\quad - 2 \sum_{i=1}^k \sum_{j=1}^k \mathbf{y}_i^\top \mathbf{y}_j W_{ij} \\ &= 2 \sum_{i=1}^k D_{ii} \mathbf{y}_i^\top \mathbf{y}_i - 2 \sum_{i=1}^k \sum_{j=1}^k \mathbf{y}_i^\top \mathbf{y}_j W_{ij} \\ &= 2 \sum_{i=1}^k (\sqrt{D_{ii}} \mathbf{y}_i)^\top (\sqrt{D_{ii}} \mathbf{y}_i) - 2 \sum_{i=1}^k \mathbf{y}_i^\top \left( \sum_{j=1}^k \mathbf{y}_j W_{ij} \right) \\ &= 2 \text{Tr}[(\mathbf{Y} \sqrt{\mathbf{D}})(\mathbf{Y} \sqrt{\mathbf{D}})^\top] - 2 \sum_{i=1}^k \mathbf{y}_i^\top (\mathbf{Y} \mathbf{W})_i \\ &= 2 \text{Tr}[\mathbf{Y}(\mathbf{D} - \mathbf{W})\mathbf{Y}^\top] = 2 \text{Tr}[\mathbf{YLY}^\top] \end{aligned}$$

### 三、 Laplacian Eigenmap

$$\begin{aligned} & \sum_{i=1}^k \sum_{j=1}^k \| \mathbf{y}_i - \mathbf{y}_j \|^2 W_{ij} \\ &= \sum_{i=1}^k \sum_{j=1}^k (\mathbf{y}_i^\top \mathbf{y}_i - 2\mathbf{y}_i^\top \mathbf{y}_j + \mathbf{y}_j^\top \mathbf{y}_j) W_{ij} \\ &= \sum_{i=1}^k \left( \sum_{j=1}^k W_{ij} \right) \mathbf{y}_i^\top \mathbf{y}_i + \sum_{j=1}^k \left( \sum_{i=1}^k W_{ij} \right) \mathbf{y}_j^\top \mathbf{y}_j \\ &\quad - 2 \sum_{i=1}^k \sum_{j=1}^k \mathbf{y}_i^\top \mathbf{y}_j W_{ij} \\ &= 2 \sum_{i=1}^k D_{ii} \mathbf{y}_i^\top \mathbf{y}_i - 2 \sum_{i=1}^k \sum_{j=1}^k \mathbf{y}_i^\top \mathbf{y}_j W_{ij} \\ &= 2 \sum_{i=1}^k (\sqrt{D_{ii}} \mathbf{y}_i)^\top (\sqrt{D_{ii}} \mathbf{y}_i) - 2 \sum_{i=1}^k \mathbf{y}_i^\top \left( \sum_{j=1}^k \mathbf{y}_j W_{ij} \right) \\ &= 2 \text{Tr}[(\mathbf{Y} \sqrt{\mathbf{D}})(\mathbf{Y} \sqrt{\mathbf{D}})^\top] - 2 \sum_{i=1}^k \mathbf{y}_i^\top (\mathbf{Y} \mathbf{W})_i \\ &= 2 \text{Tr}[\mathbf{Y}(\mathbf{D} - \mathbf{W})\mathbf{Y}^\top] = 2 \text{Tr}[\mathbf{YLY}^\top] \end{aligned}$$

### 三、 Laplacian Eigenmap

$$\begin{aligned} & \sum_{i=1}^k \sum_{j=1}^k \| \mathbf{y}_i - \mathbf{y}_j \|^2 W_{ij} \\ &= \sum_{i=1}^k \sum_{j=1}^k (\mathbf{y}_i^\top \mathbf{y}_i - 2\mathbf{y}_i^\top \mathbf{y}_j + \mathbf{y}_j^\top \mathbf{y}_j) W_{ij} \\ &= \sum_{i=1}^k \left( \sum_{j=1}^k W_{ij} \right) \mathbf{y}_i^\top \mathbf{y}_i + \sum_{j=1}^k \left( \sum_{i=1}^k W_{ij} \right) \mathbf{y}_j^\top \mathbf{y}_j \\ &\quad - 2 \sum_{i=1}^k \sum_{j=1}^k \mathbf{y}_i^\top \mathbf{y}_j W_{ij} \\ &= 2 \sum_{i=1}^k D_{ii} \mathbf{y}_i^\top \mathbf{y}_i - 2 \sum_{i=1}^k \sum_{j=1}^k \mathbf{y}_i^\top \mathbf{y}_j W_{ij} \\ &= 2 \sum_{i=1}^k (\sqrt{D_{ii}} \mathbf{y}_i)^\top (\sqrt{D_{ii}} \mathbf{y}_i) - 2 \sum_{i=1}^k \mathbf{y}_i^\top \left( \sum_{j=1}^k \mathbf{y}_j W_{ij} \right) \\ &= 2 \text{Tr}[(\mathbf{Y} \sqrt{\mathbf{D}})(\mathbf{Y} \sqrt{\mathbf{D}})^\top] - 2 \sum_{i=1}^k \mathbf{y}_i^\top (\mathbf{Y} \mathbf{W}^\top)_i \\ &= 2 \text{Tr}[\mathbf{Y}(\mathbf{D} - \mathbf{W})\mathbf{Y}^\top] = 2 \text{Tr}[\mathbf{YLY}^\top] \end{aligned}$$

### 三、 Laplacian Eigenmap

$$\begin{aligned} & \sum_{i=1}^k \sum_{j=1}^k \| \mathbf{y}_i - \mathbf{y}_j \|^2 W_{ij} \\ &= \sum_{i=1}^k \sum_{j=1}^k (\mathbf{y}_i^\top \mathbf{y}_i - 2\mathbf{y}_i^\top \mathbf{y}_j + \mathbf{y}_j^\top \mathbf{y}_j) W_{ij} \\ &= \sum_{i=1}^k \left( \sum_{j=1}^k W_{ij} \right) \mathbf{y}_i^\top \mathbf{y}_i + \sum_{j=1}^k \left( \sum_{i=1}^k W_{ij} \right) \mathbf{y}_j^\top \mathbf{y}_j \\ &\quad - 2 \sum_{i=1}^k \sum_{j=1}^k \mathbf{y}_i^\top \mathbf{y}_j W_{ij} \\ &= 2 \sum_{i=1}^k D_{ii} \mathbf{y}_i^\top \mathbf{y}_i - 2 \sum_{i=1}^k \sum_{j=1}^k \mathbf{y}_i^\top \mathbf{y}_j W_{ij} \\ &= 2 \sum_{i=1}^k (\sqrt{D_{ii}} \mathbf{y}_i)^\top (\sqrt{D_{ii}} \mathbf{y}_i) - 2 \sum_{i=1}^k \mathbf{y}_i^\top \left( \sum_{j=1}^k \mathbf{y}_j W_{ij} \right) \\ &= 2 \text{Tr}[(\mathbf{Y} \sqrt{\mathbf{D}})(\mathbf{Y} \sqrt{\mathbf{D}})^\top] - 2 \sum_{i=1}^k \mathbf{y}_i^\top (\mathbf{Y} \mathbf{W}^\top)_i \\ &= 2 \text{Tr}[\mathbf{Y}(\mathbf{D} - \mathbf{W})\mathbf{Y}^\top] = 2 \text{Tr}[\mathbf{YLY}^\top] \end{aligned}$$

### 三、 Laplacian Eigenmap

$$\begin{aligned} & \sum_{i=1}^k \sum_{j=1}^k \| \mathbf{y}_i - \mathbf{y}_j \|^2 W_{ij} \\ &= \sum_{i=1}^k \sum_{j=1}^k (\mathbf{y}_i^\top \mathbf{y}_i - 2\mathbf{y}_i^\top \mathbf{y}_j + \mathbf{y}_j^\top \mathbf{y}_j) W_{ij} \\ &= \sum_{i=1}^k \left( \sum_{j=1}^k W_{ij} \right) \mathbf{y}_i^\top \mathbf{y}_i + \sum_{j=1}^k \left( \sum_{i=1}^k W_{ij} \right) \mathbf{y}_j^\top \mathbf{y}_j \\ &\quad - 2 \sum_{i=1}^k \sum_{j=1}^k \mathbf{y}_i^\top \mathbf{y}_j W_{ij} \\ &= 2 \sum_{i=1}^k D_{ii} \mathbf{y}_i^\top \mathbf{y}_i - 2 \sum_{i=1}^k \sum_{j=1}^k \mathbf{y}_i^\top \mathbf{y}_j W_{ij} \\ &= 2 \sum_{i=1}^k (\sqrt{D_{ii}} \mathbf{y}_i)^\top (\sqrt{D_{ii}} \mathbf{y}_i) - 2 \sum_{i=1}^k \mathbf{y}_i^\top \left( \sum_{j=1}^k \mathbf{y}_j W_{ij} \right) \\ &= 2 \text{Tr}[(\mathbf{Y} \sqrt{\mathbf{D}})(\mathbf{Y} \sqrt{\mathbf{D}})^\top] - 2 \sum_{i=1}^k \mathbf{y}_i^\top (\mathbf{Y} \mathbf{W}^\top)_i \\ &= 2 \text{Tr}[\mathbf{Y}(\mathbf{D} - \mathbf{W})\mathbf{Y}^\top] = 2 \text{Tr}[\mathbf{YLY}^\top] \end{aligned}$$

Then, the loss is as blow

$$\begin{aligned} & \min_{\mathbf{Y}} \text{Tr}(\mathbf{YLY}^\top) \\ \text{s.t. } & \mathbf{YDY}^\top = \mathbf{I} \end{aligned}$$

### 三、 Laplacian Eigenmap

Let  $\mathcal{L} = \text{Tr}(\mathbf{Y}\mathbf{L}\mathbf{Y}^T + \boldsymbol{\Lambda}(\mathbf{I} - \mathbf{Y}\mathbf{D}\mathbf{Y}^T))$ ,

where  $\boldsymbol{\Lambda}$  is a diagonal matrix whose entries are Lagrange multipliers.

We compute the derivative of  $\zeta$  with respect to  $\mathbf{Y}$  as

$$\frac{\partial \zeta}{\partial \mathbf{Y}} = \mathbf{LY} - \mathbf{DY}\boldsymbol{\Lambda}$$

The optimal  $\mathbf{Y}$  satisfies

$$\mathbf{LY} - \mathbf{DY}\boldsymbol{\Lambda} = \mathbf{0} \tag{5}$$

which is a generalized eigenvalue problem, we turn Equa. (6) into a simple eigenvalue problem by post-multiplying  $\mathbf{D}^{-1}$ , The optimal  $\mathbf{Y}$  satisfies

$$\mathbf{D}^{-1}\mathbf{LY} = \mathbf{Y}\boldsymbol{\Lambda} \tag{6}$$

Note that:  $\mathbf{L}$  is a symmetric matrix

# Taking Home & Test Questions

- The definition of manifold
- Why manifold learning could achieve nonlinear dimension reduction?
- What means Locally, Linear and embedding?
- What is the assumption adopted in these three steps?
- What is the local consistency proposed and used by LLE/LE? And why could be?
- What is the key (key idea) of LLE/LE? In one sentence.
- Implement LLE/LE
- Could LLE/LE address multiple subspace dimension reduction?
- What advantages and disadvantages of LLE/LE?
- Could PCA learn a manifold? What is the difference among PCA, CCA, LDA, LLE, and LE?
- Compare the objective function of PCA, CCA, LDA, LLE, and LE, and summary their difference and similarity.

Others...

四川大学-机器学习引论

Q&A

THANKS!