

## CH. 6: Dimensionality Reduction

- **Objectives:** Reduces space complexity,  
Reduces time complexity,  
Data visualization
- **Two main methods of dimensionality reduction:**
  - 1) **Feature selection:** Choosing  $k < d$  important features,
  - 2) **Feature extraction:** Mapping data points from  $d$ -D to  $k$ -D space, where  $k < d$ , while preserving as many properties of the data as possible.

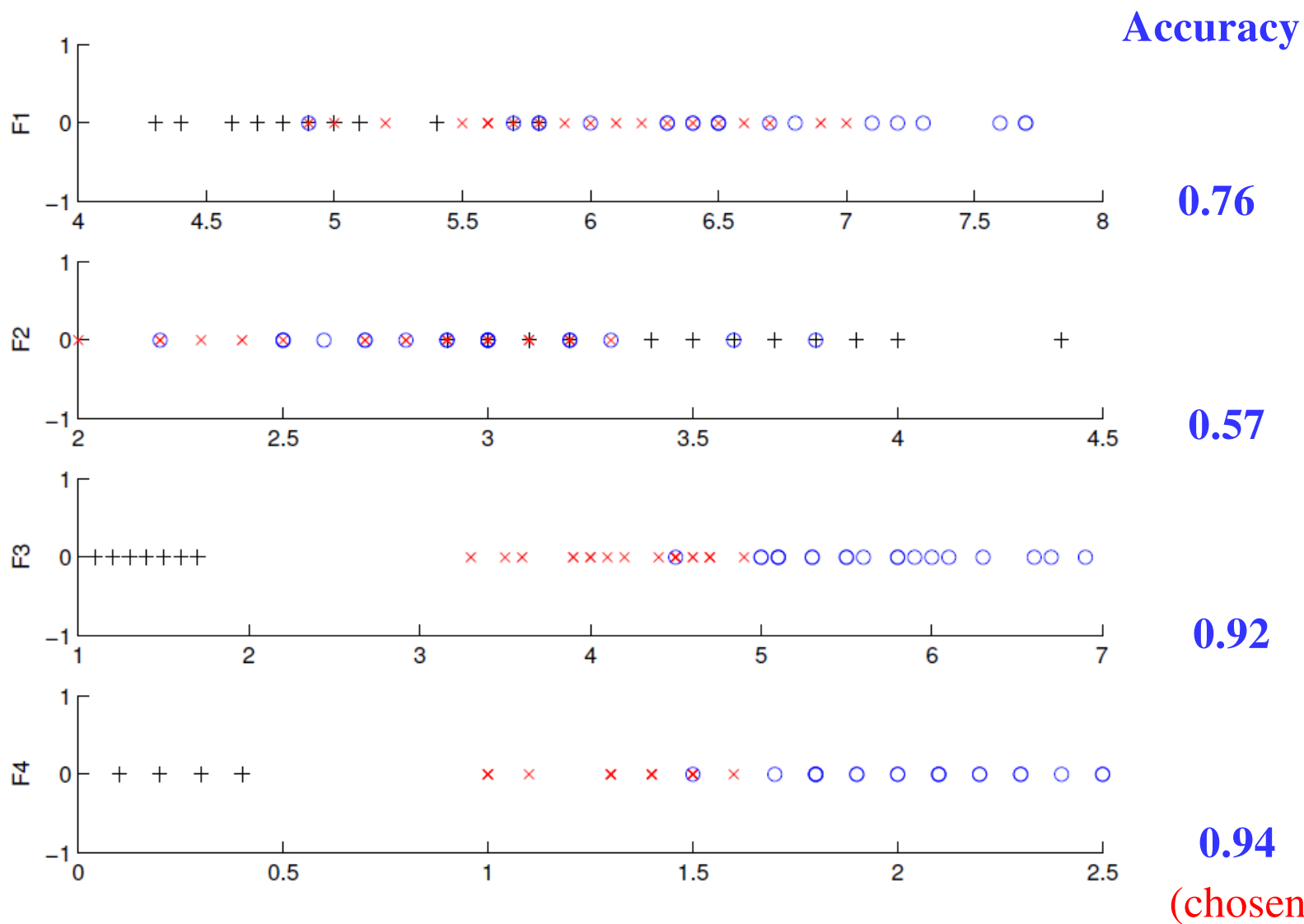
## 6.1 Feature Selection

i) **Forward search:** Add the best feature at each step

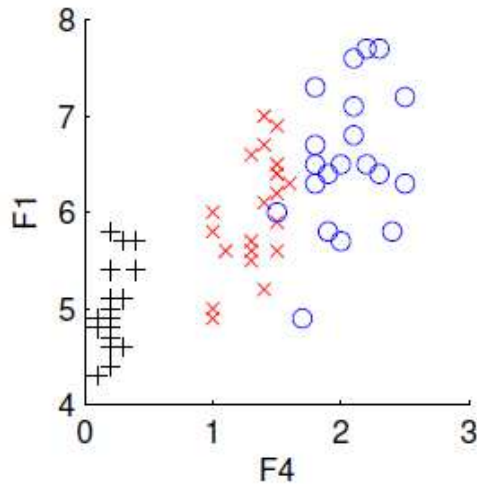
- ▣ Initially,  $F = \phi$  ( $F$ : feature set)
- ▣ At each iteration, find the best new feature using a sample  $j = \arg \max_i P(F \cup x_i)$  where  
 $P()$ : performance function
- ▣ Add  $x_j$  to  $F$  if  $P(F \cup x_i) > P(F)$

Example: Iris data (3 classes: (+, ×, ○),  
4 features: (F1,F2,F3,F4) )

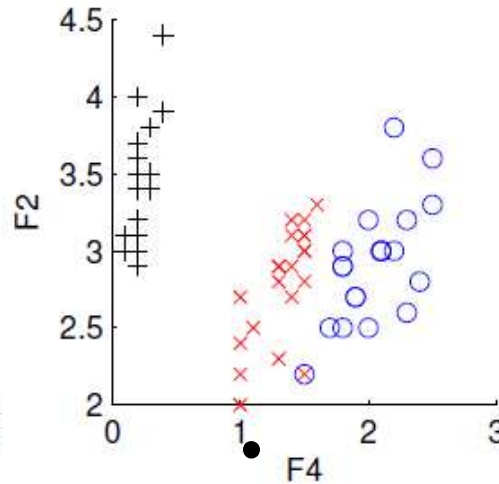
# Single feature



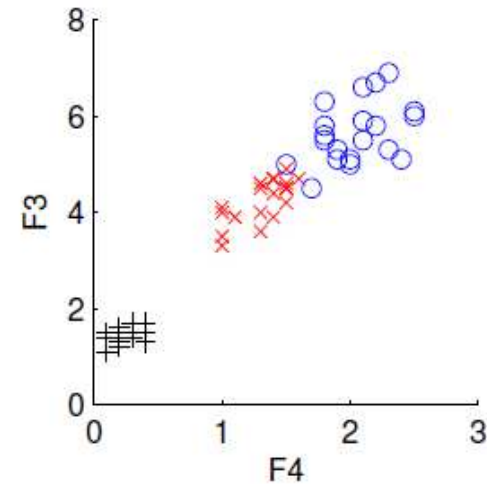
Add one more feature to F4



Accuracy      **0.87**



**0.92**



**0.96**  
(chosen)

Since the accuracies of (F1,F3,F4) and (F2,F3,F4) are both 0.94 smaller than 0.96.

Stop the feature selection process at (F3,F4) .

**ii) Backward search:** Start with all features and  
remove one at a time.

Remove  $x_j$  from  $F$  if  $j = \arg \min_i E(F - \{x_i\})$

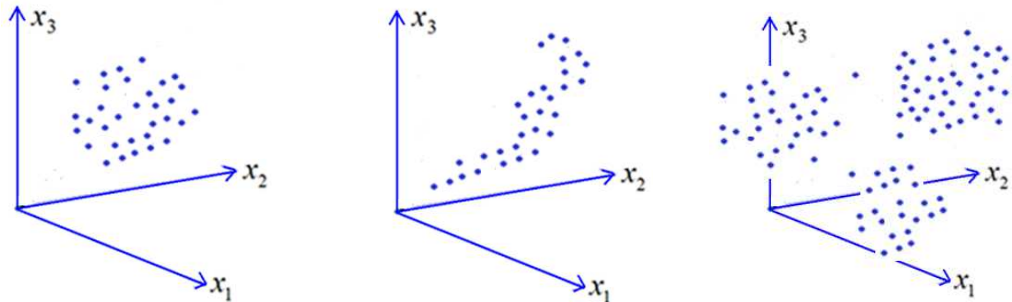
$E()$ : error function

**iii) Floating search:** The numbers of added and  
removed features can change at each step.

## 6.2 Feature Extraction

- Graphical representation ([visualization](#))

A data set  $X = \{\mathbf{x}_i\}_{i=1}^N$  can be represented as a set of points in a space.



The data set may possess certain properties. Feature extraction (FE) attempts to preserve or even improve the properties during an FE process.

## i) Principal Components Analysis (PCA)

Data points:  $\mathbf{x}_i = (x_{i1}, x_{i2}, \dots, x_{id})^T$ ,  $i = 1, 2, \dots, N$

Matrix representation:

$$\mathbf{X} = (\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_N)^T = \begin{pmatrix} x_{11} & x_{12} & \cdots & x_{1d} \\ x_{21} & x_{22} & \cdots & x_{2d} \\ \cdots & \cdots & \cdots & \cdots \\ x_{N1} & x_{N2} & \cdots & x_{Nd} \end{pmatrix}$$

Mean vector:  $\boldsymbol{\mu}_x = \frac{1}{N} \sum_{i=1}^N \mathbf{x}_i$ ,

Covariance matrix:  $\mathbf{C}_x = \frac{1}{N} \sum_{i=1}^N (\mathbf{x}_i - \boldsymbol{\mu}_x)(\mathbf{x}_i - \boldsymbol{\mu}_x)^T$

Let  $\lambda_i$  and  $\mathbf{e}_i$ ,  $i = 1, \dots, d$ , be the eigenvalues and eigenvectors of  $C_x$ , i.e.,  $C_x \mathbf{e}_i = \lambda_i \mathbf{e}_i$ .

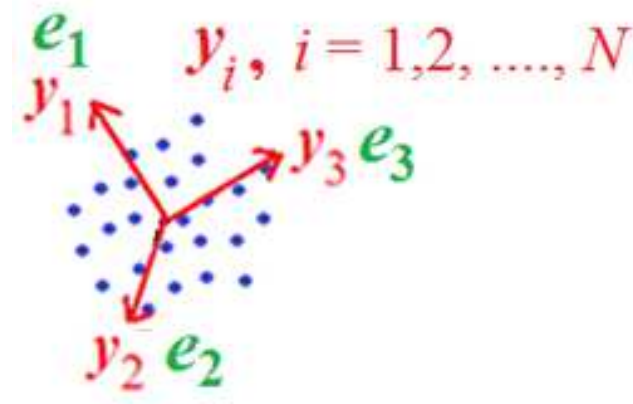
Suppose  $\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_d$ .

Let  $A_{d \times d} = [\mathbf{e}_1 \ \mathbf{e}_2 \ \dots \ \mathbf{e}_d]^T$ .

Compute  $\mathbf{y}_i = A(\mathbf{x}_i - \boldsymbol{\mu}_x)$ ,  $i = 1, 2, \dots, N$ .

y-axes corresponding to eigenvectors  $\mathbf{e}$ 's are orthogonal, i.e., uncorrelated.

The variances over y-axes  
 $\approx$  eigenvalues





For dimensionality reduction,

$$\text{Let } A_k = [\mathbf{e}_1 \cdots \mathbf{e}_k]^T, \quad k < d. \quad \hat{\mathbf{y}} = A_k (\mathbf{x} - \boldsymbol{\mu}_x)$$
$$\hat{\mathbf{y}}_{k \times 1} = (A_k)_{k \times d} (\mathbf{x} - \boldsymbol{\mu}_x)_{d \times 1}$$

Let  $\hat{\mathbf{x}}$  be the reconstruction of  $\mathbf{x}$  from  $\hat{\mathbf{y}}$ .

The reconstruction error  $\|\mathbf{x}_i - \hat{\mathbf{x}}_i\|^2$  depends on

$$\sum_{i=k+1}^d \lambda_i, \text{ which are relatively smaller than } \sum_{i=1}^k \lambda_i.$$

# Eigen faces for face recognition



Figure 3.23: 32 original images of a boy's face, each  $321 \times 261$  pixels.

$$\begin{array}{ccccccccc}
 \mathbf{i} & & \mathbf{e}_1 & & \mathbf{e}_2 & & \mathbf{e}_3 & & \mathbf{e}_4 \\
 \text{[Image]} & = q_1 & \text{[Image]} & + q_2 & \text{[Image]} & + q_3 & \text{[Image]} & + q_4 & \text{[Image]} \\
 \end{array}
 \quad
 \begin{array}{l}
 q_i = \mathbf{i}^T \mathbf{e}_i, \\
 i = 1, 2, 3, 4
 \end{array}$$

Figure 3.24: Reconstruction of the image from four basis vectors  $\mathbf{b}_i$ ,  $i = 1, \dots, 4$  which can be displayed as images. The linear combination was computed as  $q_1 \mathbf{b}_1 + q_2 \mathbf{b}_2 + q_3 \mathbf{b}_3 + q_4 \mathbf{b}_4 = 0.078 \mathbf{b}_1 + 0.062 \mathbf{b}_2 - 0.182 \mathbf{b}_3 + 0.179 \mathbf{b}_4$ .

## ii) Feature Embedding (FE)

FE places  $d$ -D data points in a  $k$ -D space ( $k < d$ ) such that pairwise similarities in the new space respect the original pairwise similarities.

Let  $X_{N \times d}$  be data matrix,  $\lambda_i$  and  $\mathbf{w}_i$  be the eigenvalues and eigenvectors of correlation matrix  $(X^T X)_{d \times d}$  of features, i.e.,  $(X^T X)\mathbf{w}_i = \lambda_i \mathbf{w}_i$ . Multiply both sides  $X$ ,  $X(X^T X)\mathbf{w}_i = (XX^T)X\mathbf{w}_i = \lambda_i X\mathbf{w}_i$ , i.e.,  $\lambda_i$ ,  $X\mathbf{w}_i$  are eigenvalues and eigenvectors of similarity matrix  $(XX^T)_{N \times N}$  of instances.

Let  $\mathbf{v}_i = X\mathbf{w}_i$ ,  $i = 1, \dots, k$  ( $< d$ ) corresponding to  $k$  leading eigenvalues, which form the coordinates of the new space.

Since  $\mathbf{v}_i$  are the eigenvectors of similarity matrix  $(XX^T)_{N \times N}$ . Pairwise similarities between instances will be preserved in the new space.

### iii) Factor Analysis (FA)

- In PCA, from the original features  $x_i, i = 1, \dots, d$  to form a new set of features  $z_j = \sum_{i=1}^d w_{ji} x_i, j = 1, \dots, k$

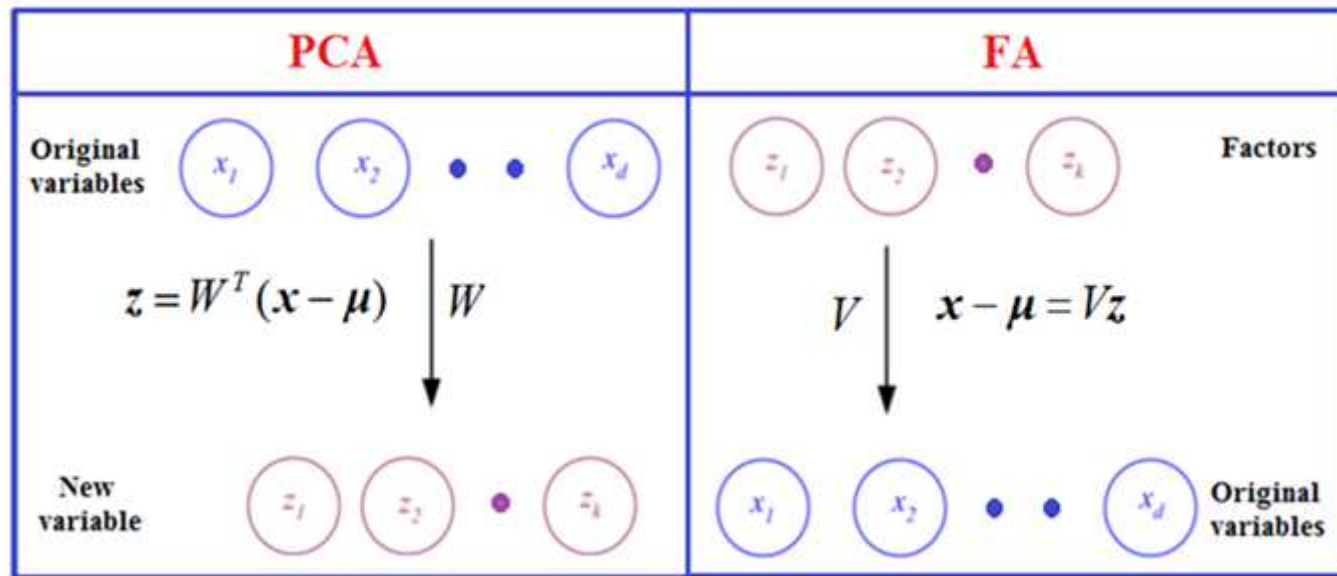
For dimension reduction,  $k < d$ . Mathematically,

$$\mathbf{z} = \mathbf{W}^T (\mathbf{x} - \boldsymbol{\mu}).$$

- In FA, a set of unobservable (latent) factors  $z_j, j = 1, \dots, k$  that combine to generate  $x_i, i = 1, \dots, d$ .

$$x_i = \sum_{j=1}^k v_{ij} z_j, \quad i = 1, \dots, d. \quad \text{Mathematically,}$$

$$\mathbf{x} - \boldsymbol{\mu} = \mathbf{V}\mathbf{z}.$$



- Given a sample  $X = \{\mathbf{x}^t\}_{t=1}^N$ ,  $\mathbf{x} = (x_1, x_2, \dots, x_d)$ , find a small number of factors  $z_i$ ,  $i = 1, \dots, k$  ( $k < d$ ), s.t. each  $x_i$  can be written as a weighted sum of  $z_i$ ,
- $$x_i - \mu_i = v_{i1}z_1 + v_{i2}z_2 + \dots + v_{ik}z_k + \varepsilon_i, \quad i = 1, \dots, d$$
- In vector-matrix form,  $\mathbf{x} - \boldsymbol{\mu} = \mathbf{V}\mathbf{z} + \boldsymbol{\varepsilon}$ .

where  $z_i$  : **latent factors** ( $\sim N(0,1)$ ,  $\text{Cov}(z_i, z_j) = 0$ ,  $i \neq j$ )

$v_{ij}$  : **factor loadings**

$\varepsilon_i$  : errors ( $E[\varepsilon_i] = 0$ ,  $\text{Var}(\varepsilon_i) = \psi_i$ ,

$\text{Cov}(\varepsilon_i, z_j) = 0$ ,  $\forall i, j$      $\text{Cov}(\varepsilon_i, \varepsilon_j) = 0$ ,  $i \neq j$ ).

- **Example:** Let  $s_c, s_e, s_m, s_p$  and  $s_{ch}$  be the **score variables** of Chinese( $c$ ), English( $e$ ), Mathematics( $m$ ), Physics( $p$ ), and Chemistry( $ch$ ), respectively, which are observable. Let  $z_m, z_i$  and  $z_o$  be the **talent variables** of memory( $m$ ), inference( $i$ ), organization( $o$ ), which are latent.

Specifically, given the scores of a student

$$s_c = 78, s_e = 82, s_m = 94, s_p = 89, s_{ch} = 92,$$

what are **loadings**  $v_{ij}$  ( $i = c, e, m, p, ch; j = m, i, o$ )

of **factors**  $z_m, z_i$  and  $z_o$  of the student?

$$s = (s_c \ s_e \ s_m \ s_p \ s_{ch}), \quad s_i = \sum_{j \in \{m, i, o\}} v_{ij} z_j \quad (i = c, e, m, p, ch)$$

**Two uses of factor analysis:**

- i) Knowledge extraction,
- ii) Dimensionality reduction



Knowledge Extraction – Given  $X$  and  $Z$ , find  $V$

From  $\mathbf{x} - \boldsymbol{\mu} = V\mathbf{z} + \boldsymbol{\varepsilon}$ , for simplicity, let  $\boldsymbol{\mu} = \mathbf{0}$ ,

$$\Rightarrow \mathbf{x} = V\mathbf{z} + \boldsymbol{\varepsilon}.$$

$$\begin{aligned}\Sigma &= \text{Cov}(\mathbf{x}) = \text{Cov}(V\mathbf{z} + \boldsymbol{\varepsilon}) = \text{Cov}(V\mathbf{z}) + \text{Cov}(\boldsymbol{\varepsilon}) \\ &= V\text{Cov}(\mathbf{z})V^T + \boldsymbol{\psi} = VIV^T + \boldsymbol{\psi} = VV^T + \boldsymbol{\psi}\end{aligned}$$

( $\because z_i \sim N(0,1)$ ,  $\text{Cov}(z_i, z_j) = 0$ ,  $i \neq j$ ,  $\therefore \text{Cov}(\mathbf{z}) = I$ )

$\boldsymbol{\psi}$ : diagonal matrix with  $\psi_i$  on the diagonals.

Ignoring  $\boldsymbol{\psi}$ ,  $\Sigma = VV^T$ .

Let  $S$  be the estimator  $\Sigma$  of sample  $X$ ,  $S = VV^T$ .

Spectral decomposition of  $S$ :

$$S = E\Lambda E^T = (E\Lambda^{1/2})(E\Lambda^{1/2})^T = VV^T,$$

$$\therefore V = E\Lambda^{1/2}$$

Dimensionality Reduction – Given  $X$ , find  $Z$

$$\text{Let } z_j = \sum_{i=1}^d w_{ji} x_i, \quad j = 1, \dots, k$$

$$z_1 = \sum_{i=1}^d w_{1i} x_i = w_{11}x_1 + w_{12}x_2 + \dots + w_{1d}x_d$$

.....

$$z_k = \sum_{i=1}^d w_{ki} x_i = w_{k1}x_1 + w_{k2}x_2 + \dots + w_{kd}x_d$$

$$\begin{pmatrix} z_1 \\ z_2 \\ \dots \\ z_k \end{pmatrix}_{k \times 1} = \begin{bmatrix} w_{11} & w_{12} & \dots & w_{1d} \\ w_{21} & w_{22} & \dots & w_{2d} \\ \dots & \dots & \dots & \dots \\ w_{k1} & w_{k2} & \dots & w_{kd} \end{bmatrix}_{k \times d} \begin{pmatrix} x_1 \\ x_2 \\ \dots \\ x_d \end{pmatrix}_{d \times 1}$$

In vector-matrix form,  $\mathbf{z} = \mathbf{W}^T \mathbf{x}$ .

Given a sample  $\mathbf{X} = \{\mathbf{x}^t\}_{t=1}^N$ ,  $\mathbf{z}^i = \mathbf{W}^T \mathbf{x}^i$ ,  $i = 1, \dots, N$

In matrix form,  $\mathbf{Z}_{N \times k} = \mathbf{X}_{N \times d} \mathbf{W}_{d \times k}$ .

Solve for  $\mathbf{W} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{Z}$

$S = X^T X$  is the estimated covariance matrix  $\Sigma$   
of sample  $X$ .

$$S = VV^T, \quad V = E\Lambda^{1/2}.$$

$\mathbf{x} - \boldsymbol{\mu} = V\mathbf{z} + \boldsymbol{\varepsilon}$ . Ignore  $\boldsymbol{\mu}$  and  $\boldsymbol{\varepsilon}$ ,  $\Rightarrow \mathbf{x} = V\mathbf{z}$ .

$\mathbf{x}\mathbf{z}^T = V\mathbf{z}\mathbf{z}^T$ . Given a sample  $X = \{\mathbf{x}^t\}_{t=1}^N$ ,

$$\mathbf{x}^i \mathbf{z}^T = V\mathbf{z}^i (\mathbf{z}^i)^T = V, \quad i = 1, \dots, N.$$

In matrix form,  $X^T Z = V$ .

$$W = (X^T X)^{-1} X^T Z = S^{-1}V = (VV^T)^{-1}V, \quad V = E\Lambda^{1/2}$$

$$Z = XW.$$

#### iv) Matrix Factorization (MF)

$$X_{N \times d} = F_{N \times k} G_{k \times d}$$

$k \ll d$

$G$  defines  $k$  new **factors** in terms of the attributes of data  $X$ .

$F$  defines **instances** in terms of the new factors in  $G$ .

Objective: solve  $X = FG$  for  $F = G^+ X$

## v) Linear Discriminant Analysis (LDA)

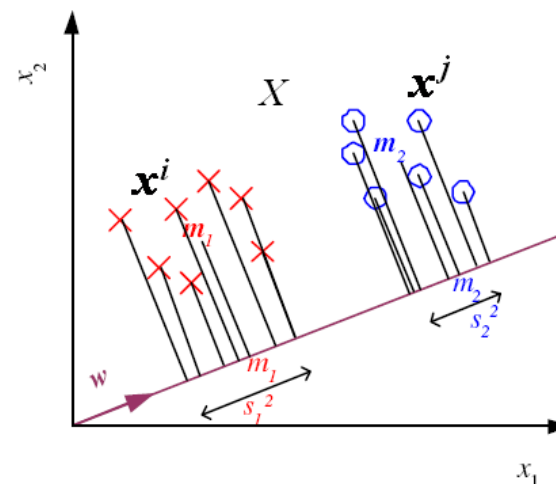
-- Find a low dimension space such that when data are projected onto it, the examples of different **classes** are as **well separated** as possible.

- In **2-Class ( $d$ -D to 1-D)** case, find a direction  $w$ , such that when data are projected onto  $w$ , the examples of different classes are well-separated.

Given a sample

$$X = \{\mathbf{x}^t, r^t\}_{t=1}^N \text{ s.t.}$$

$$r^t = \begin{cases} 1 & \text{if } \mathbf{x}^t \in C_1 \\ 0 & \text{if } \mathbf{x}^t \in C_2 \end{cases}$$



Means:

$$m_1 = \frac{\sum_t \mathbf{w}^T \mathbf{x}^t r^t}{\sum_t r^t} = \mathbf{w}^T \mathbf{m}_1, \quad m_2 = \frac{\sum_t \mathbf{w}^T \mathbf{x}^t (1 - r^t)}{\sum_t (1 - r^t)} = \mathbf{w}^T \mathbf{m}_2$$

Scatters:

$$s_1^2 = \sum_t (\mathbf{w}^T \mathbf{x}^t - m_1)^2 r^t, \quad s_2^2 = \sum_t (\mathbf{w}^T \mathbf{x}^t - m_2)^2 (1 - r^t)$$

□ Find  $\mathbf{w}$  that maximizes  $J(\mathbf{w}) = \frac{(m_1 - m_2)^2}{s_1^2 + s_2^2}$  ---- (A)

$$\begin{aligned} (m_1 - m_2)^2 &= (\mathbf{w}^T \mathbf{m}_1 - \mathbf{w}^T \mathbf{m}_2)^2 \\ &= \mathbf{w}^T (\mathbf{m}_1 - \mathbf{m}_2) (\mathbf{m}_1 - \mathbf{m}_2)^T \mathbf{w} = \mathbf{w}^T S_B \mathbf{w} \end{aligned}$$

where  $S_B = (\mathbf{m}_1 - \mathbf{m}_2) (\mathbf{m}_1 - \mathbf{m}_2)^T$

(Between-class scatter matrix)

$$\begin{aligned}
s_1^2 &= \sum_{t=1}^N (\mathbf{w}^T \mathbf{x}^t - m_1)^2 r^t = \sum_{t=1}^N \mathbf{w}^T r^t (\mathbf{x}^t - \mathbf{m}_1) (\mathbf{x}^t - \mathbf{m}_1)^T \mathbf{w} \\
&= \mathbf{w}^T S_1 \mathbf{w}, \quad \text{where } S_1 = \sum_t r^t (\mathbf{x}^t - \mathbf{m}_1) (\mathbf{x}^t - \mathbf{m}_1)^T
\end{aligned}$$

Similarly,  $s_2^2 = \mathbf{w}^T S_2 \mathbf{w}$

where  $S_2 = \sum_t (1 - r^t) (\mathbf{x}^t - \mathbf{m}_2) (\mathbf{x}^t - \mathbf{m}_2)^T$

$$s_1^2 + s_2^2 = \mathbf{w}^T S_1 \mathbf{w} + \mathbf{w}^T S_2 \mathbf{w} = \mathbf{w}^T (S_1 + S_2) \mathbf{w} = \mathbf{w}^T S_W \mathbf{w}$$

where  $S_W = S_1 + S_2$  (Within-class scatter matrix)

$$\begin{aligned}
(\text{A}) \Rightarrow J(\mathbf{w}) &= \frac{(m_1 - m_2)^2}{s_1^2 + s_2^2} = \frac{\mathbf{w}^T S_B \mathbf{w}}{\mathbf{w}^T S_W \mathbf{w}} \\
&= \frac{\mathbf{w}^T (\mathbf{m}_1 - \mathbf{m}_2) (\mathbf{m}_1 - \mathbf{m}_2)^T \mathbf{w}}{\mathbf{w}^T S_W \mathbf{w}}
\end{aligned}$$

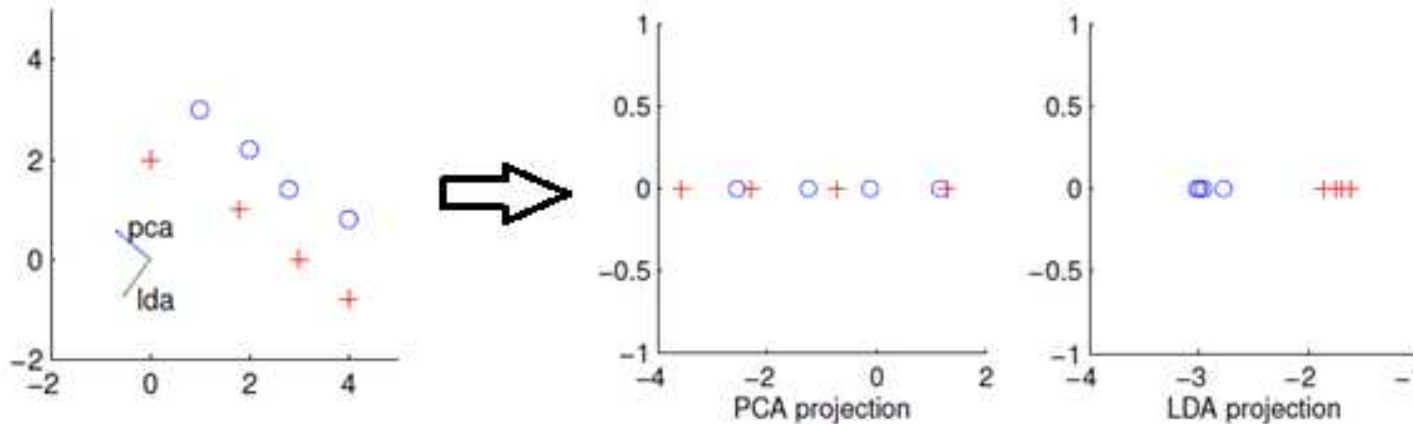


$$\begin{aligned}
\frac{dJ(\mathbf{w})}{d\mathbf{w}} &= \frac{d}{d\mathbf{w}} \left( \frac{\mathbf{w}^T (\mathbf{m}_1 - \mathbf{m}_2) (\mathbf{m}_1 - \mathbf{m}_2)^T \mathbf{w}}{\mathbf{w}^T S_W \mathbf{w}} \right) \quad \left( d\left(\frac{f}{g}\right) = \frac{gdf - fdg}{g^2} \right) \\
&= \mathbf{w}^T S_W \mathbf{w} \frac{2\mathbf{w}^T (\mathbf{m}_1 - \mathbf{m}_2) (\mathbf{m}_1 - \mathbf{m}_2)}{(\mathbf{w}^T S_W \mathbf{w})^2} - \frac{\mathbf{w}^T (\mathbf{m}_1 - \mathbf{m}_2) (\mathbf{m}_1 - \mathbf{m}_2)^T \mathbf{w}}{(\mathbf{w}^T S_W \mathbf{w})^2} 2\mathbf{w}^T S_W \\
&= \frac{2\mathbf{w}^T (\mathbf{m}_1 - \mathbf{m}_2) (\mathbf{m}_1 - \mathbf{m}_2)}{\mathbf{w}^T S_W \mathbf{w}} - \frac{\mathbf{w}^T (\mathbf{m}_1 - \mathbf{m}_2) (\mathbf{m}_1 - \mathbf{m}_2)^T \mathbf{w}}{(\mathbf{w}^T S_W \mathbf{w})^2} 2\mathbf{w}^T S_W \\
&= 2 \frac{\mathbf{w}^T (\mathbf{m}_1 - \mathbf{m}_2)}{\mathbf{w}^T S_W \mathbf{w}} \left( (\mathbf{m}_1 - \mathbf{m}_2) - \frac{\mathbf{w}^T (\mathbf{m}_1 - \mathbf{m}_2)}{\mathbf{w}^T S_W \mathbf{w}} S_W \mathbf{w} \right) = 0 \text{ ----- (B)}
\end{aligned}$$

$$\text{Let } \frac{\mathbf{w}^T (\mathbf{m}_1 - \mathbf{m}_2)}{\mathbf{w}^T S_W \mathbf{w}} = c. \quad (\text{B}) \Rightarrow c \left( (\mathbf{m}_1 - \mathbf{m}_2) - c S_W \mathbf{w} \right) = 0$$

$$\Rightarrow \mathbf{w} = c S_W^{-1} (\mathbf{m}_1 - \mathbf{m}_2).$$

## Example:



- In  $n > 2$  Class ( $d$ -D to  $k$ -D) case,

Within-class scatter matrix:

$$S_W = \sum_{i=1}^n S_i, \text{ where } S_i = \sum_t r_i^t (\mathbf{x}^t - \mathbf{m}_i)(\mathbf{x}^t - \mathbf{m}_i)^T$$
$$r_i^t = 1 \text{ if } \mathbf{x}^t \in C_i \text{ and } 0 \text{ otherwise}$$

Between-class scatter matrix:

$$S_B = \sum_{i=1}^n N_i (\mathbf{m}_i - \mathbf{m})(\mathbf{m}_i - \mathbf{m})^T, \quad \mathbf{m} = \frac{1}{n} \sum_{i=1}^n \mathbf{m}_i, \quad N_i = \sum_t r_i^t$$

Let  $W_{d \times k}$  be the projection matrix from the  $d$ -D space to the  $k$ -D space ( $k < d$ ), then

$(W^T S_W W)_{k \times k}$ ,  $(W^T S_B W)_{k \times k}$  : projections of  $(S_W)_{d \times d}$ ,  $(S_B)_{d \times d}$

A spread measure of a scatter matrix is its determinant.

Find  $W$  such that  $J(W) = \frac{|W^T S_B W|}{|W^T S_W W|}$  is maximized.

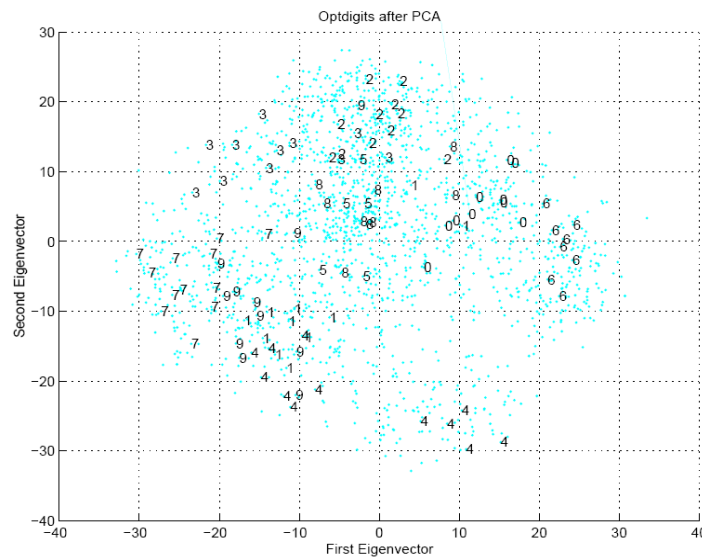
The determinant of matrix  $A_{n \times n}$  is the product of its eigenvalues, i.e.,  $|A| = \lambda_1 \cdot \lambda_2 \cdots \lambda_n$ .

$$\begin{aligned}
 J(W) &= \frac{|W^T S_B W|}{|W^T S_W W|} = \frac{|\lambda_1^B \cdot \lambda_2^B \cdots \lambda_k^B|}{|\lambda_1^W \cdot \lambda_2^W \cdots \lambda_k^W|} = \left| \frac{\lambda_1^B \cdot \lambda_2^B \cdots \lambda_k^B}{\lambda_1^W \cdot \lambda_2^W \cdots \lambda_k^W} \right| \\
 &= |(\lambda_1^W \cdot \lambda_2^W \cdots \lambda_k^W)^{-1} (\lambda_1^B \cdot \lambda_2^B \cdots \lambda_k^B)| \\
 &= \left| |W^T S_W W|^{-1} |W^T S_B W| \right| \underset{\substack{\uparrow \\ |AB|=|A||B|}}{=} |(W^T S_W W)^{-1} (W^T S_B W)|
 \end{aligned}$$

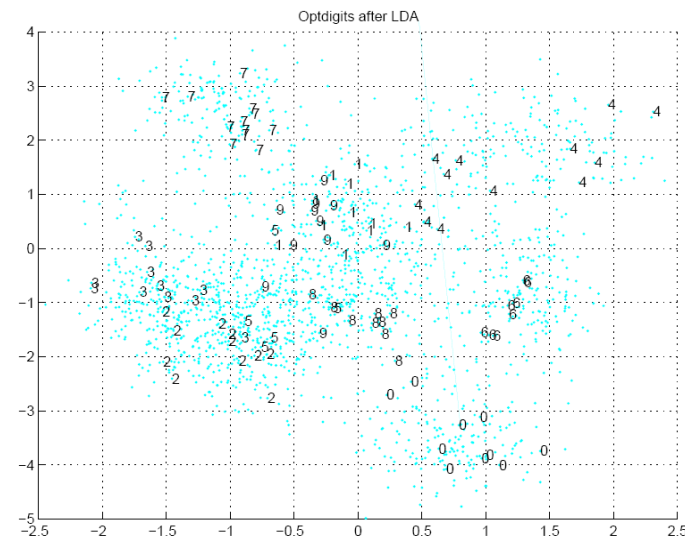
$$\begin{aligned}
 \text{Let } \frac{dJ(W)}{dW} &= \frac{d}{dW} |(W^T S_W W)^{-1} (W^T S_B W)| \\
 &= \frac{d}{dW} |(W^{-1} S_W^{-1} (W^T)^{-1} W^T S_B W)| = \frac{d}{dW} |(W^{-1} S_W^{-1} S_B W)| = 0.
 \end{aligned}$$

The solution of  $W$  is formed by the  $k$  largest eigenvectors of  $S_W^{-1} S_B$ .

**PCA**



**LDA**



Fisher Discriminant Analysis with Kernels, S. Mika, G. Ratsch, J. Weaton, B. Scholkopf, and K.R. Muller, IEEE, 1999.

## vi) Laplacian Eigenmaps (LE)

Let  $\mathbf{x}^r$  and  $\mathbf{x}^s$  be any two out of  $N$  data instances and  $b_{rs}$  is their similarity. Find  $\mathbf{y}^r$  and  $\mathbf{y}^s$  that

$\min \sum_{r,s} \|\mathbf{y}^r - \mathbf{y}^s\|^2 b_{rs}$ , i.e., two similar instances

(large  $b_{rs}$ ) should be close in the new space

(small  $\|\mathbf{y}^r - \mathbf{y}^s\|$ ). Define  $b_{rs} = \exp \left[ -\frac{\|\mathbf{x}^r - \mathbf{x}^s\|^2}{2\sigma^2} \right]$

if  $\mathbf{x}^r$  and  $\mathbf{x}^s$  are in the predefined neighborhood, and 0 otherwise, i.e., only local similarities are cared.

Consider the 1-D new space

$$\begin{aligned}
\sum_{r,s} \|\mathbf{y}^r - \mathbf{y}^s\|^2 b_{rs} &= \frac{1}{2} \sum_{r,s} (y_r - y_s)^2 b_{rs} \\
&= \frac{1}{2} \left( \sum_{r,s} b_{rs} y_r^2 - 2 \sum_{r,s} b_{rs} y_r y_s + \sum_{r,s} b_{rs} y_s^2 \right) \\
&= \frac{1}{2} \left( \sum_r d_r y_r^2 - 2 \sum_{r,s} b_{rs} y_r y_s + \sum_s d_s y_s^2 \right) \quad \left( d_r = \sum_s b_{rs}, \right. \\
&\quad \left. d_s = \sum_r b_{rs} \right) \\
&= \sum_r d_r y_r^2 - \sum_r \sum_s b_{rs} y_r y_s = \mathbf{y}^T D \mathbf{y} - \mathbf{y}^T B \mathbf{y} \\
&= \mathbf{y}^T (D - B) \mathbf{y} = \mathbf{y}^T L \mathbf{y},
\end{aligned}$$

where  $B = [b_{rs}]$ ,  $D = \text{diag}[d_r]$ ,  $L$  : **Laplacian matrix**

$\mathbf{y}$ :  $N$ - $D$  vector,  $y_r$  : the new coordinate of  $\mathbf{x}^r$ .

The solution to  $\min\{\mathbf{y}^T L \mathbf{y}\}$  subject to  $\|\mathbf{y}\| = 1$

$$\approx \frac{d(\mathbf{y}^T L \mathbf{y})}{d\mathbf{y}} = 0 \text{ subject to } \|\mathbf{y}\| = 1$$

$$\approx L\mathbf{y} = 0 \text{ subject to } \|\mathbf{y}\| = 1$$

- The method of **Lagrange multipliers**

$$\text{Error: } E = \|L\mathbf{y} - \mathbf{0}\|^2 = \|L\mathbf{y}\|^2 = \mathbf{y}^T (L^T L) \mathbf{y}$$

$$\text{Constraint: } \|\mathbf{y}\| = 1$$



Minimize  $F(\mathbf{y}) = \mathbf{y}^T (L^T L) \mathbf{y} + \lambda (\|\mathbf{y}\|^2 - 1)$

where  $\lambda$  : Lagrange multiplier

Let  $\frac{dF(\mathbf{y})}{d\mathbf{y}} = 2(L^T L) \mathbf{y} + 2\lambda \mathbf{y} = \mathbf{0}$

We obtain  $(L^T L) \mathbf{y} = -\lambda \mathbf{y}$

The solution  $\mathbf{y}$  is an eigenvector of  $L^T L$  with eigenvalue  $-\lambda$

The associated error  $E = \mathbf{y}^T (L^T L) \mathbf{y} = -\lambda \mathbf{y}^T \mathbf{y} = -\lambda$

The  $\mathbf{y}$  with the smallest  $\lambda$  is the least square error solution of  $L\mathbf{y} = \mathbf{0}$ .

**LE Algorithm:** Given  $N$  points  $\mathbf{x}_i \in R^d$ ,  $i = 1, \dots, N$

1. Put an edge  $e_{ij}$  between nodes  $n_i$  and  $n_j$  if

$$\|\mathbf{x}_i - \mathbf{x}_j\| < \varepsilon.$$

2. Weight the edge by  $b_{ij} = \exp\left(-\frac{\|\mathbf{x}_i - \mathbf{x}_j\|^2}{t}\right)$

3. For each connected component of  $G$ , compute eigenvalues  $\lambda$  and eigenvectors  $\mathbf{y}$  of  $L$ , i.e.,

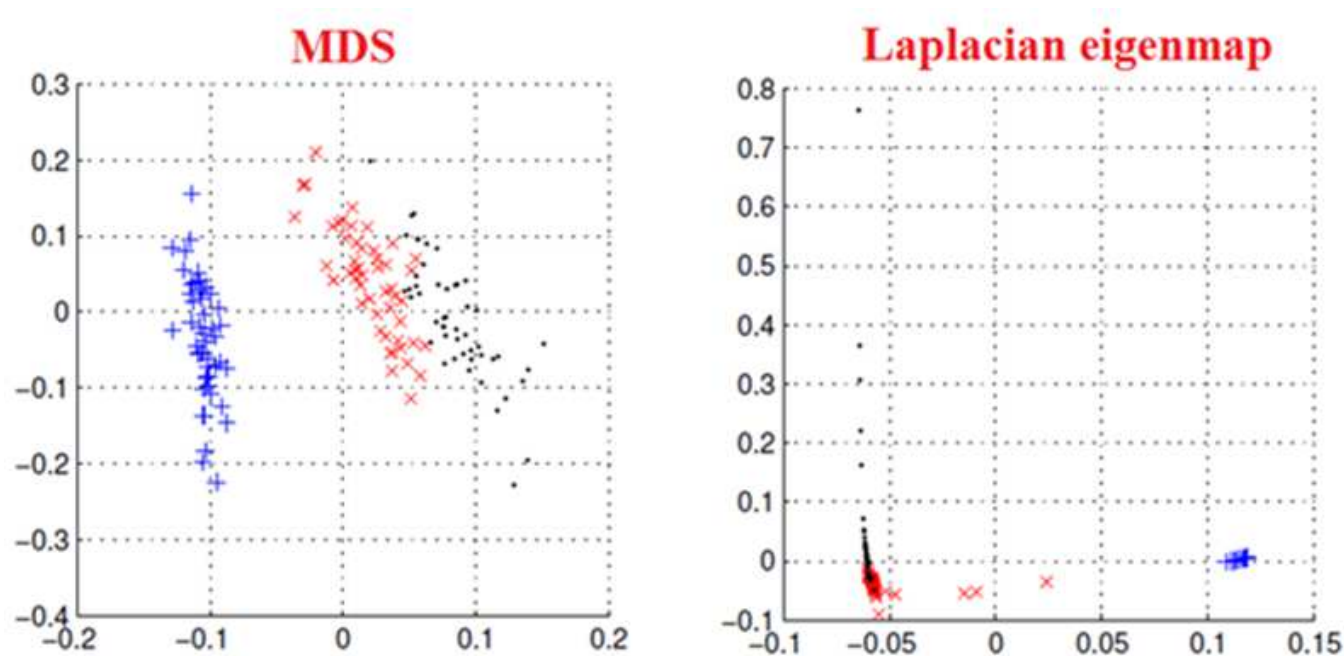
$L\mathbf{y} = \lambda \mathbf{y}$ , where  $L = D - B$ ,  $D$ : diagonal matrix,

$d_{ii} = \sum_j b_{ji}$ . Let  $\mathbf{y}_i$ ,  $i = 0, \dots, k-1$  be the solutions

of (A), ordered from small to large eigenvalues.

## Example: Iris data

LE lead to denser data than MDS.



Laplace Eigenmaps for Dimensionality Reduction and Data Representation, M. Belkin, Neural Computing, 15, pp. 1373-1396, 2003.

## ix) Canonical Correlation Analysis (CCA)

- Given a sample  $X = \{\mathbf{x}^t, \mathbf{y}^t\}_{t=1}^N$ , both  $x$  and  $y$  are inputs, e.g., (1) acoustic information and visual information in speech recognition, (2) image data and text annotations in image retrieval application.
- Take the **correlation** of  $(\mathbf{x}^t, \mathbf{y}^t)$  into account while reducing dimensionality to a joint space, i.e., find two vectors  $\mathbf{w}$  and  $\mathbf{v}$  s.t. when  $\mathbf{x}$  is projected along  $\mathbf{w}$  and  $\mathbf{y}$  is projected along  $\mathbf{v}$ , their correlation  $\rho$  is maximized, where

$$\begin{aligned}\rho &= \frac{\text{Cov}(\mathbf{w}^T \mathbf{x}, \mathbf{v}^T \mathbf{y})}{\sqrt{\text{Var}(\mathbf{w}^T \mathbf{x})} \sqrt{\text{Var}(\mathbf{v}^T \mathbf{y})}} \\ &= \frac{\mathbf{w}^T \text{Cov}(\mathbf{x}, \mathbf{y}) \mathbf{v}}{\sqrt{\mathbf{w}^T \text{Var}(\mathbf{x}) \mathbf{w}} \sqrt{\mathbf{v}^T \text{Var}(\mathbf{y}) \mathbf{v}}} = \frac{\mathbf{w}^T S_{xy} \mathbf{v}}{\sqrt{\mathbf{w}^T S_{xx} \mathbf{w}} \sqrt{\mathbf{v}^T S_{yy} \mathbf{v}}}\end{aligned}$$

where **Covariance matrices:**

$$S_{xx} = \text{Var}(\mathbf{x}) = E[(\mathbf{x} - \boldsymbol{\mu}_x)^2]$$

$$S_{yy} = \text{Var}(\mathbf{y}) = E[(\mathbf{y} - \boldsymbol{\mu}_y)^2]$$

**Cross-covariance matrices:**

$$S_{xy} = \text{Cov}(\mathbf{x}, \mathbf{y}) = E[(\mathbf{x} - \boldsymbol{\mu}_x)(\mathbf{y} - \boldsymbol{\mu}_y)^T]$$

$$S_{yx} = \text{Cov}(\mathbf{y}, \mathbf{x}) = E[(\mathbf{y} - \boldsymbol{\mu}_y)(\mathbf{x} - \boldsymbol{\mu}_x)^T]$$

Let  $\frac{\partial \rho}{\partial \mathbf{w}} = 0$  and  $\frac{\partial \rho}{\partial \mathbf{v}} = 0$

Solutions:  $\mathbf{w}$  is an eigenvector of  $S_{xx}^{-1} S_{xy} S_{yy}^{-1} S_{yx}$ ;  
 $\mathbf{v}$  is an eigenvector of  $S_{yy}^{-1} S_{yx} S_{xx}^{-1} S_{xy}$ .

Choose  $(\mathbf{w}, \mathbf{v})$  with largest eigenvalues as the solution.

$\therefore \rho \propto$  shared eigenvalue of  $\lambda_{\mathbf{w}}$  and  $\lambda_{\mathbf{v}}$

- Look for  $k$  pairs  $(\mathbf{w}_i, \mathbf{v}_i)$ ,  $i = 1, \dots, k$

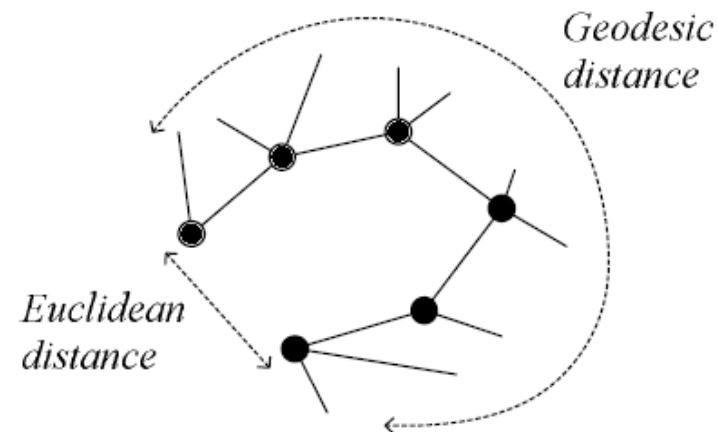
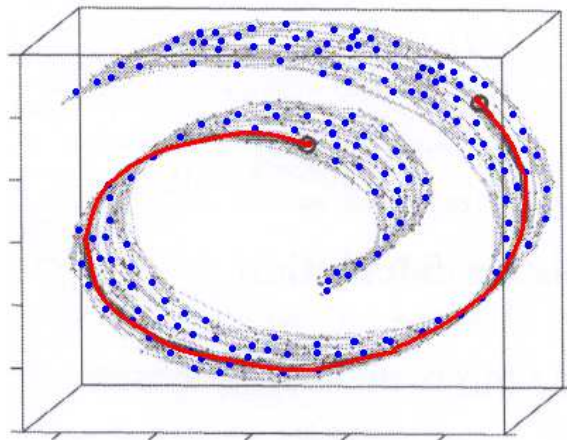
Let  $W = [\mathbf{w}_1, \mathbf{w}_2, \dots, \mathbf{w}_k]$ ,  $V = [\mathbf{v}_1, \mathbf{v}_2, \dots, \mathbf{v}_k]$

$\mathbf{r}^t = W^T \mathbf{x}^t$ ,  $\mathbf{s}^t = V^T \mathbf{y}^t$ .

$(\mathbf{r}^t, \mathbf{s}^t)$ : lower-dimensional representation of  $(\mathbf{x}^t, \mathbf{y}^t)$ .

Canonical Correlation Analysis: An Overview with Application to Learning Methods, D. R. Hardoon, S. Szedmak, and J. Shawe-Taylor, Neural Computation, 16, 2004.

## vi) Isometric Feature Mapping (Isomap)



- Estimates the geodesic distance and applies multidimensional scaling for dimensionality reduction.

$$D = \begin{matrix} & \mathbf{x}_1 & \mathbf{x}_2 & \cdots & \mathbf{x}_N \\ \begin{matrix} \mathbf{x}_1 \\ \mathbf{x}_2 \\ \cdots \\ \mathbf{x}_N \end{matrix} & \begin{pmatrix} d_{11} & d_{12} & \cdots & d_{1N} \\ d_{21} & d_{22} & \cdots & d_{2N} \\ \cdots & \cdots & \cdots & \cdots \\ d_{N1} & d_{N2} & \cdots & d_{NN} \end{pmatrix} \end{matrix}$$



## Multidimensional Scaling (MDS):

Given distance matrix  $D = [d_{rs}]_{N \times N}$ ,

1. Calculate  $B = [b_{rs}]_{N \times N}$ , where

$$b_{rs} = \frac{1}{2}(d_{r\cdot}^2 + d_{\cdot s}^2 - d_{\cdot\cdot}^2 - d_{rs}^2), \quad d_{r\cdot}^2 = \frac{1}{N} \sum_s d_{rs}^2,$$
$$d_{\cdot s}^2 = \frac{1}{N} \sum_r d_{rs}^2, \quad d_{\cdot\cdot}^2 = \frac{1}{N^2} \sum_r \sum_s d_{rs}^2$$

2. Find the spectral decomposition of  $B$ ,

$$B = E \Lambda E^T.$$

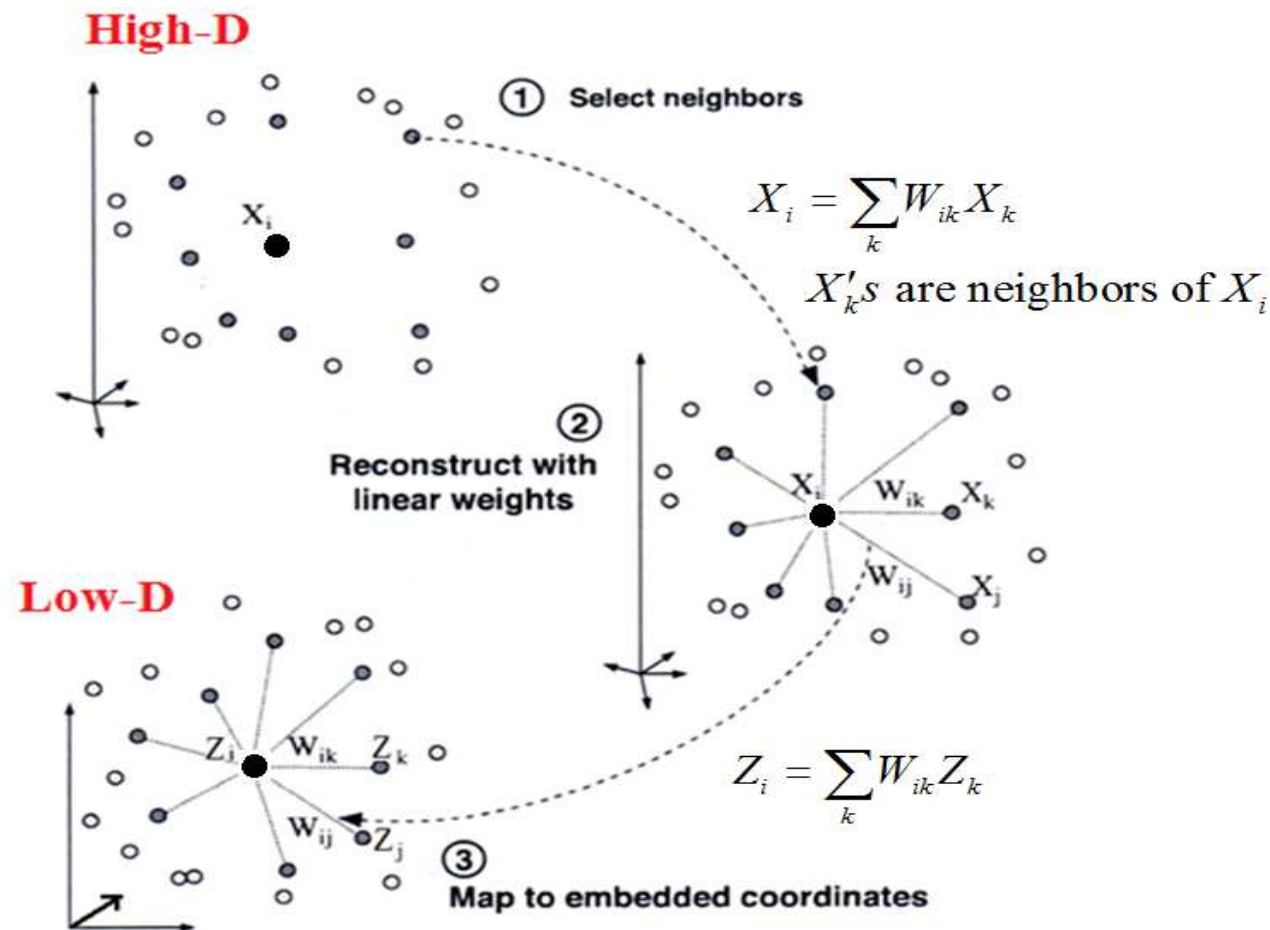
3. For dimensionality reduction, discard from  $\Lambda$  the  $p$  small eigenvalues and from  $E$  the corresponding eigenvectors to form  $\Lambda'$  and  $E'$ , respectively.

4. Find  $Z_{(N-p) \times (N-p)} = E' \Lambda'^{1/2}$ .

\* The coordinates of the points are the rows of  $Z$ .

## vii) Locally Linear Embedding (LLE)

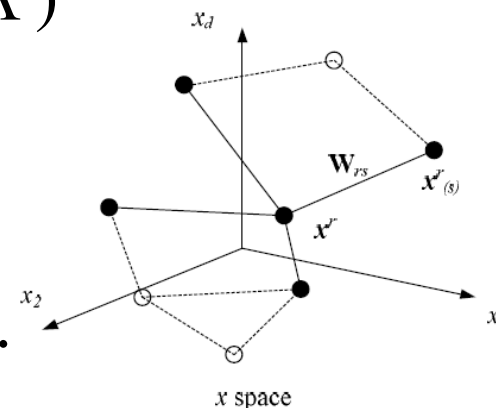
-- Recovers global nonlinear structure from locally linear fit.



1. Given  $\mathbf{x}^r$  and its neighbors  $\mathbf{x}_{(r)}^s$ , find  $W_{rs}$  by minimizing the error function  $E(W | X)$

$$\min_{W_{rs}} E(W | X) = \sum_r \left\| \mathbf{x}^r - \sum_s W_{rs} \mathbf{x}_{(r)}^s \right\|^2$$

subject to  $W_{rr} = 0, \forall r$  and  $\sum_r W_{rs} = 1$ .

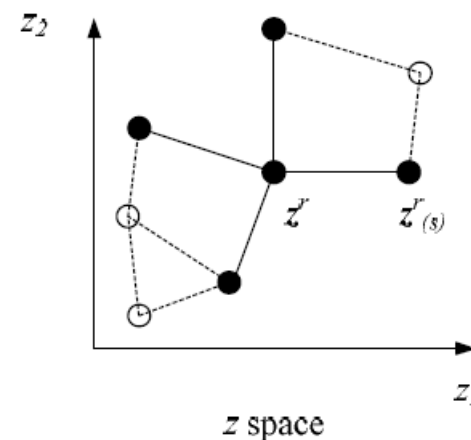


2. Find new coordinates  $\mathbf{z}^r$  that respect the constraints

given by  $W_{rs}$ ,

$$\min_{\mathbf{z}^r} E(Z | W) = \sum_r \left\| \mathbf{z}^r - \sum_s W_{rs} \mathbf{z}_{(r)}^s \right\|^2$$

subject to  $E[\mathbf{z}] = 0, \text{Cov}(\mathbf{z}) = I$



$\text{Dim}(\mathbf{z}\text{-space}) < \text{Dim}(\mathbf{x}\text{-space})$

Nonlinear Dimensionality Reduction by Locally Linear Embedding, S.T. Roweis and L.K. Saul, Science, 290, 2000.

**viii) *t*-Distributed Stochastic Neighbor Embedding**

-- Position the data points in the new space such that local neighborhood statistics are as similar as possible

The probability that  $\mathbf{x}^s$  is a neighbor of  $\mathbf{x}^r$

$$p_{s|r} = \frac{\exp[-\|\mathbf{x}^r - \mathbf{x}^s\|^2 / 2\sigma_r^2]}{\sum_{l \neq r} \exp[-\|\mathbf{x}^r - \mathbf{x}^l\|^2 / 2\sigma_r^2]}$$

$t$ -SNE is the symmetrized version of SNE by

defining  $p_{rs} = \frac{p_{s|r} + p_{r|s}}{2N}$

The probability in the lower-dimensional space is calculated as

$$q_{rs} = \frac{(1 + \|\mathbf{z}^l - \mathbf{z}^m\|^2)^{-1}}{\sum_{l \neq r} \sum_{m \neq l} (1 + \|\mathbf{z}^l - \mathbf{z}^m\|^2)^{-1}} (t - \text{distribution})$$

The aim is to learn  $\mathbf{z}^r$  so that for all pairs  $(r,s)$ ,

$q_{rs}$  can be as close as possible to  $p_{rs}$

Gradient descent for finding optimal  $\mathbf{z}^r$ :

1. Start from small random  $\mathbf{z}^r$
2. Update iteratively in the direction that decreases the KL distance in small steps

KL distance between the probability distributions  $P$  and  $Q$ , from which  $p_{rs}$  and  $q_{rs}$  are drawn

$$\text{KL}(P \parallel Q) = \sum_r \sum_s p_{rs} \log \frac{p_{rs}}{q_{rs}}$$