Lecture Slides for

# INTRODUCTION TO MACHINE LEARNING

## 3RD EDITION

ETHEM ALPAYDIN
© The MIT Press, 2014

*alpaydin@boun.edu.tr*
*http://www.cmpe.boun.edu.tr/~ethem/i2ml3e*

CHAPTER 6:

# DIMENSIONALITY REDUCTION

# Why Reduce Dimensionality?

- Reduces time complexity: Less computation

- Reduces space complexity: Fewer parameters

- Saves the cost of observing the feature

- Simpler models are more robust on small datasets

- More interpretable; simpler explanation

- Data visualization (structure, groups, outliers, etc) if plotted in 2 or 3 dimensions

# Feature Selection vs Extraction

- Feature selection: Choosing $k<d$ important features, ignoring the remaining $d - k$
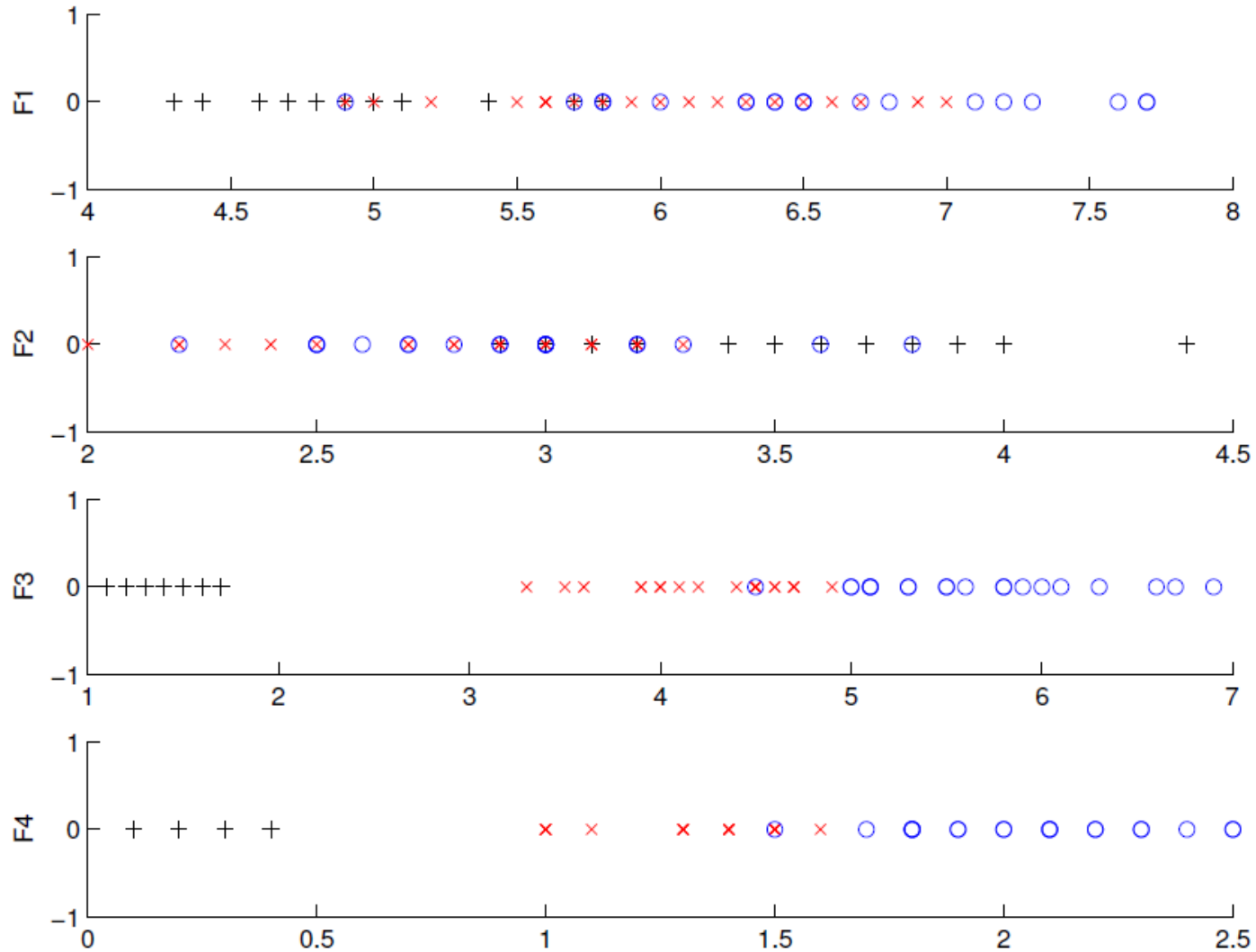
    Subset selection algorithms

- Feature extraction: Project the

    original $x_i$ , $i = 1,...,d$ dimensions to

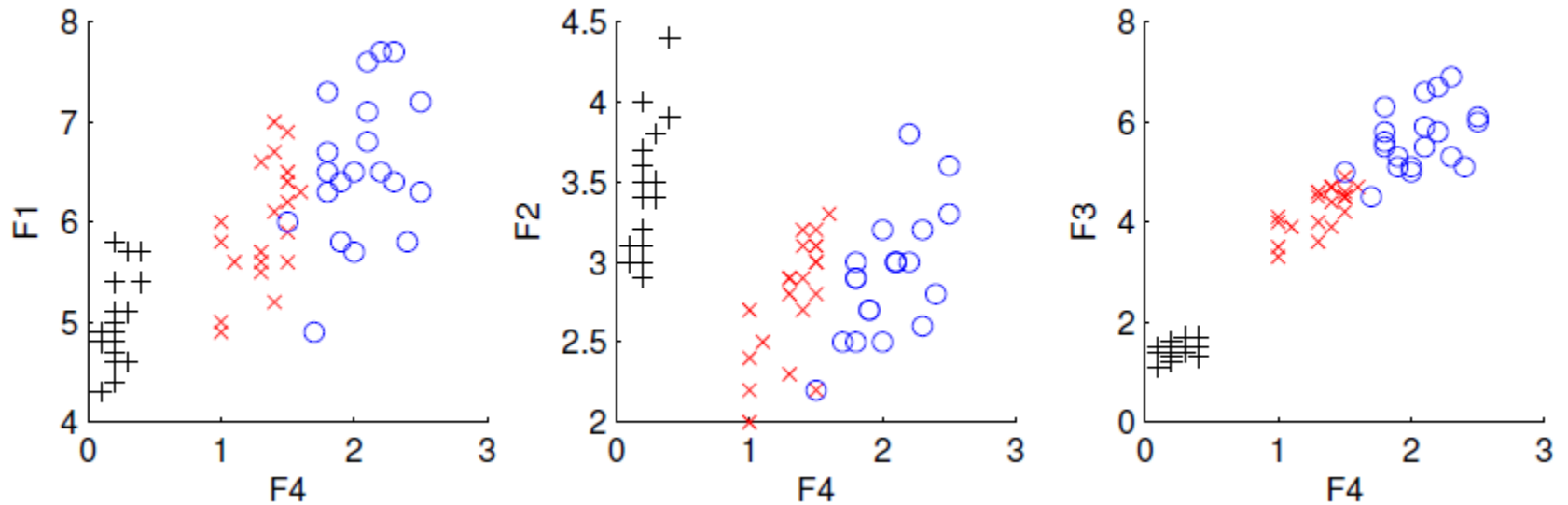    new $k<d$ dimensions, $z_j$ , $j = 1,...,k$

# Subset Selection

- There are $2^d$ subsets of $d$ features
- Forward search: Add the best feature at each step
  - Set of features $F$ initially $\varnothing$.
  - At each iteration, find the best new feature
    $i = \text{argmin}_i \; E \, ( \, F \cup x_i \, )$
  - Add $x_i$ to $F$ if $E \, ( \, F \cup x_i \, ) < E \, ( \, F \, )$

- Hill-climbing $O(d^2)$ algorithm
- Backward search: Start with all features and remove one at a time, if possible.
- Floating search (Add $k$, remove $l$)

# Iris data: Single feature



Chosen

# Iris data: Add one more feature to F4



Chosen

# Principal Components Analysis

- Find a low-dimensional space such that when $\boldsymbol{x}$ is projected there, information loss is minimized.

- The projection of $\boldsymbol{x}$ on the direction of $\boldsymbol{w}$ is: $z = \boldsymbol{w}^T\boldsymbol{x}$

- Find $\boldsymbol{w}$ such that Var(z) is maximized

$$\text{Var}(z) = \text{Var}(\boldsymbol{w}^T\boldsymbol{x}) = \text{E}[(\boldsymbol{w}^T\boldsymbol{x} - \boldsymbol{w}^T\boldsymbol{\mu})^2]$$

$$= \text{E}[(\boldsymbol{w}^T\boldsymbol{x} - \boldsymbol{w}^T\boldsymbol{\mu})(\boldsymbol{w}^T\boldsymbol{x} - \boldsymbol{w}^T\boldsymbol{\mu})]$$

$$= \text{E}[\boldsymbol{w}^T(\boldsymbol{x} - \boldsymbol{\mu})(\boldsymbol{x} - \boldsymbol{\mu})^T\boldsymbol{w}]$$

$$= \boldsymbol{w}^T \text{E}[(\boldsymbol{x} - \boldsymbol{\mu})(\boldsymbol{x} - \boldsymbol{\mu})^T]\boldsymbol{w} = \boldsymbol{w}^T \sum \boldsymbol{w}$$

where $\text{Var}(\boldsymbol{x}) = \text{E}[(\boldsymbol{x} - \boldsymbol{\mu})(\boldsymbol{x} - \boldsymbol{\mu})^T] = \sum$

- Maximize Var($z$) subject to $||w||=1$

$$\max_{\mathbf{w}_1} \mathbf{w}_1^T \Sigma \mathbf{w}_1 - \alpha \left( \mathbf{w}_1^T \mathbf{w}_1 - 1 \right)$$

$\sum w_1 = \alpha w_1$ that is, $w_1$ is an eigenvector of $\sum$

Choose the one with the largest eigenvalue for Var($z$) to be max

- Second principal component: Max Var($z_2$), s.t., $||w_2||=1$ and orthogonal to $w_1$

$$\max_{\mathbf{w}_2} \mathbf{w}_2^T \Sigma \mathbf{w}_2 - \alpha \left( \mathbf{w}_2^T \mathbf{w}_2 - 1 \right) - \beta \left( \mathbf{w}_2^T \mathbf{w}_1 - 0 \right)$$

$\sum w_2 = \alpha \, w_2$ that is, $w_2$ is another eigenvector of $\sum$

and so on.

# What PCA does

$$z = W^T(x - m)$$

where the columns of **W** are the eigenvectors of $\sum$ and **m** is sample mean

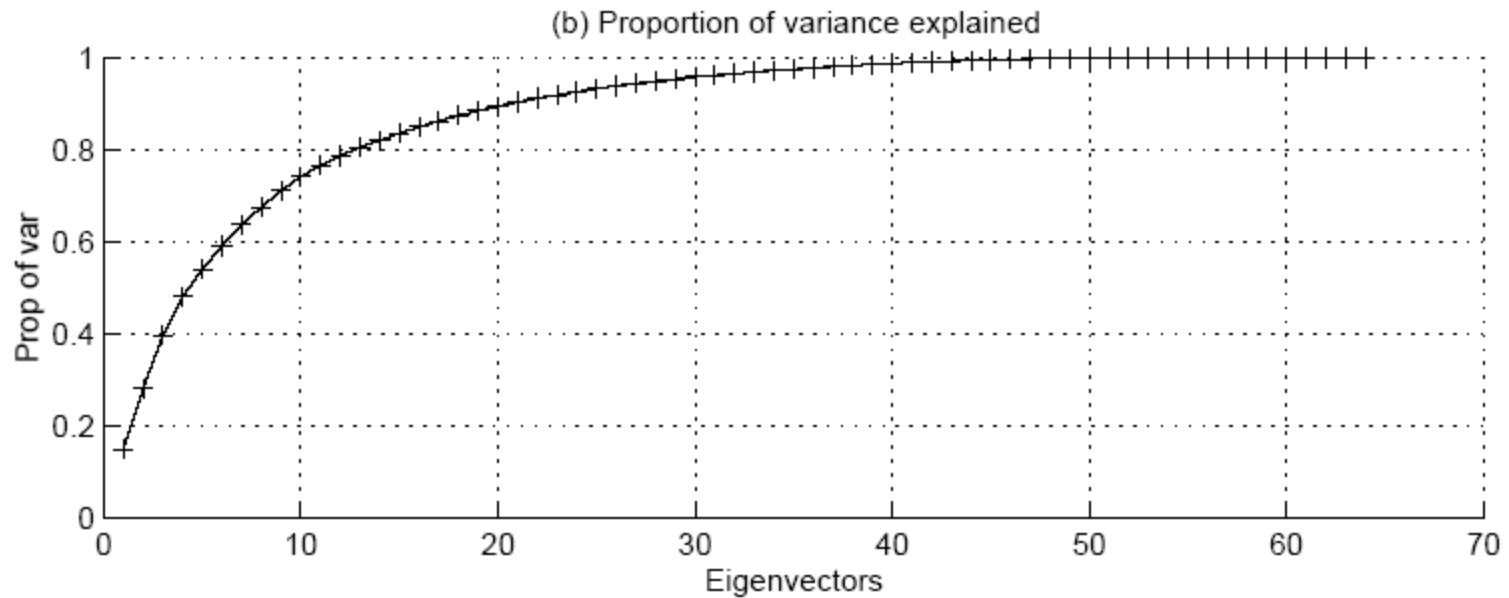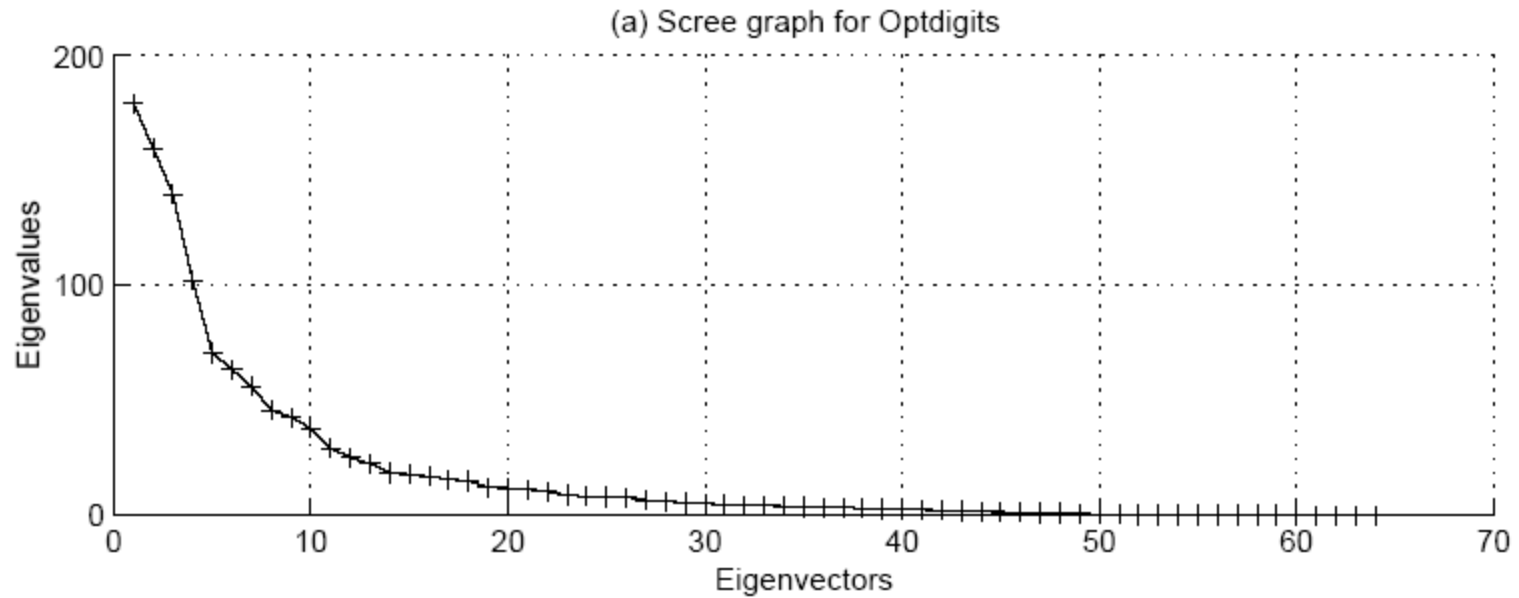Centers the data at the origin and rotates the axes
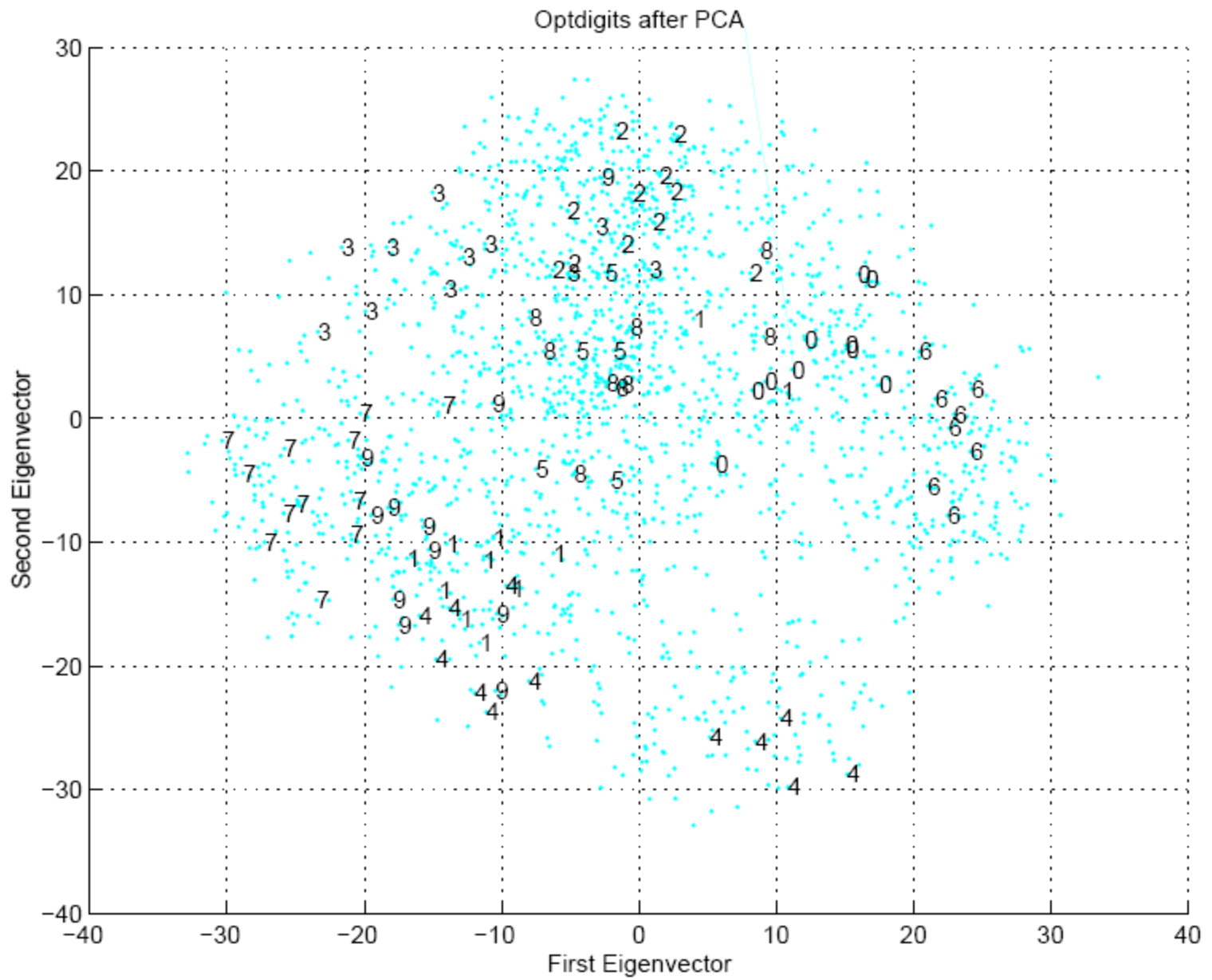
# How to choose k ?

□ Proportion of Variance (PoV) explained

$$\frac{\lambda_1 + \lambda_2 + \cdots + \lambda_k}{\lambda_1 + \lambda_2 + \cdots + \lambda_k + \cdots + \lambda_d}$$

when $\lambda_i$ are sorted in descending order

□ Typically, stop at PoV>0.9

□ Scree graph plots of PoV vs *k*, stop at "elbow"

(a) Scree graph for Optdigits

(b) Proportion of variance explained

Optdigits after PCA

13

# Feature Embedding

- When $X$ is the $N$x$d$ data matrix,

$X^TX$ is the dxd matrix (covariance of features, if mean-centered)

$XX^T$ is the $N$x$N$ matrix (pairwise similarities of instances)

- PCA uses the eigenvectors of $X^TX$ which are $d$-dim and can be used for projection

- Feature embedding uses the eigenvectors of $XX^T$ which are $N$-dim and which give directly the coordinates after projection

- Sometimes, we can define pairwise similarities (or distances) between instances, then we can use feature embedding without needing to represent instances as vectors.

# Factor Analysis

☐ Find a small number of factors **z**, which when combined generate **x** :

$$x_i - \mu_i = v_{i1}z_1 + v_{i2}z_2 + \ldots + v_{ik}z_k + \varepsilon_i$$

where $z_j$, $j = 1, \ldots, k$ are the latent factors with
$$E[\, z_i\,] = 0,\ Var(z_i) = 1,\ Cov(z_{i\,,}\, z_j) = 0,\ i \neq j\,,$$
$\varepsilon_i$ are the noise sources
$$E[\, \varepsilon_i\,] = \psi_i,\ Cov(\varepsilon_i\,,\, \varepsilon_j) = 0,\ i \neq j,\ Cov(\varepsilon_i\,,\, z_j) = 0\,,$$
and $v_{ij}$ are the factor loadings

# PCA vs FA

☐ PCA     From $x$ to $z$     $z = W^T(x - \mu)$

☐ FA     From $z$ to $x$     $x - \mu = Vz + \varepsilon$



**PCA**                    **FA**

# Factor Analysis

- In FA, factors $z_i$ are stretched, rotated and translated to generate $x$

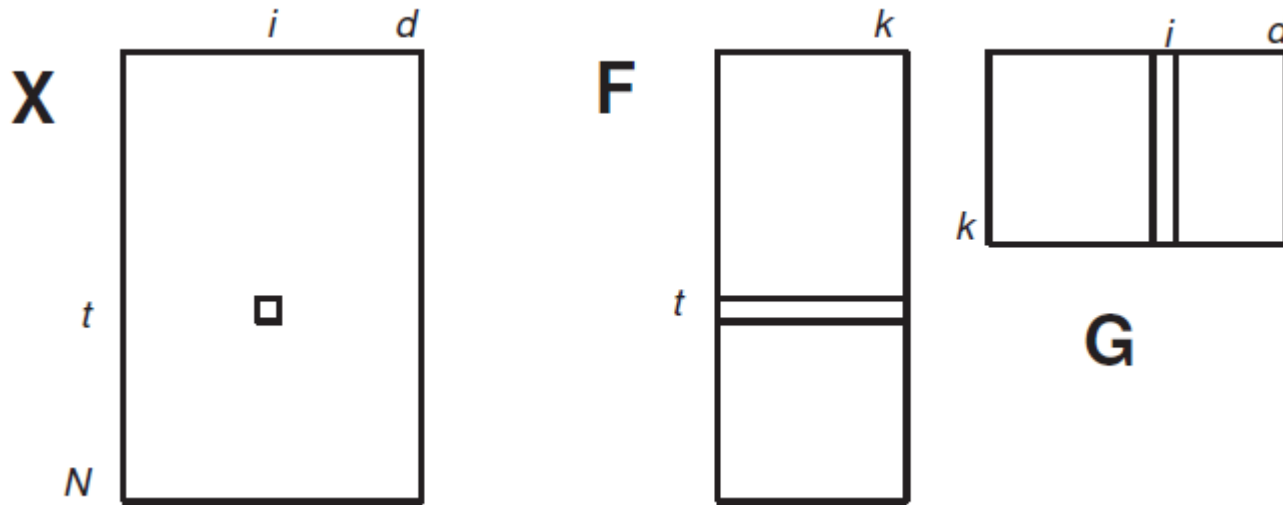# Singular Value Decomposition and Matrix Factorization

- Singular value decomposition: $X = VAW^T$

  $V$ is $N$x$N$ and contains the eigenvectors of $XX^T$

  $W$ is $d$x$d$ and contains the eigenvectors of $X^TX$

  and $A$ is $N$x$d$ and contains singular values on its first $k$ diagonal

- $X = u_1 a_1 v_1^T + \ldots + u_k a_k v_k^T$ where $k$ is the rank of $X$

# Matrix Factorization

- Matrix factorization: **X=FG**

  **F** is *Nxk* and **G** is *kxd*



$$X_{ti} = F_t^T G_i = \sum_{j=1}^{k} F_{tj} G_{ji}$$

*Latent semantic indexing*

# Multidimensional Scaling

☐ Given pairwise distances between *N* points,

$$d_{ij},\ i,j = 1,...,N$$

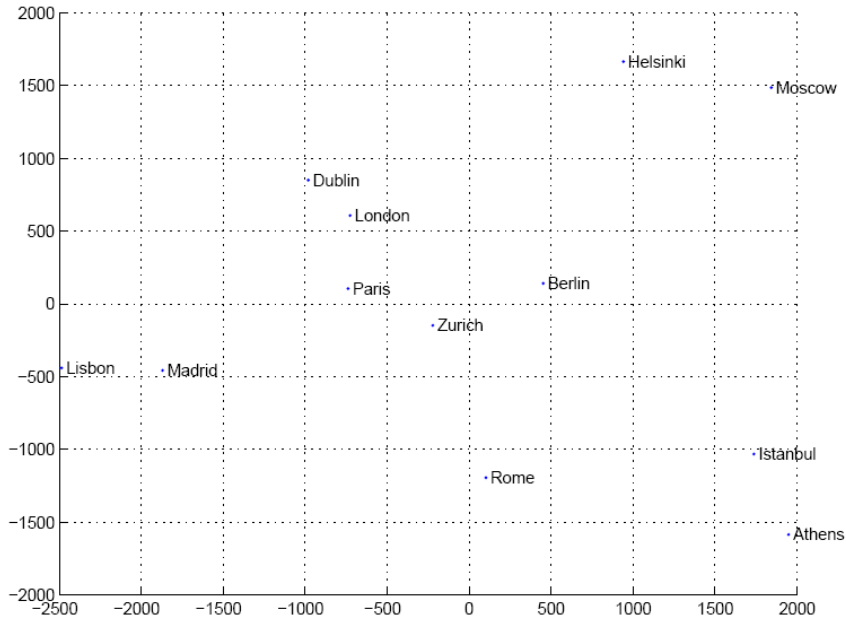place on a low-dim map s.t. distances are preserved (by feature embedding)

☐ $z = g\,(x\mid\theta)$     Find $\theta$ that min Sammon stress

$$E(\theta\mid\mathcal{X}) = \sum_{r,s} \frac{\left(\left\|\mathbf{z}^r - \mathbf{z}^s\right\| - \left\|\mathbf{x}^r - \mathbf{x}^s\right\|\right)^2}{\left\|\mathbf{x}^r - \mathbf{x}^s\right\|^2}$$

$$= \sum_{r,s} \frac{\left(\left\|\mathbf{g}(\mathbf{x}^r\mid\theta) - \mathbf{g}(\mathbf{x}^s\mid\theta)\right\| - \left\|\mathbf{x}^r - \mathbf{x}^s\right\|\right)^2}{\left\|\mathbf{x}^r - \mathbf{x}^s\right\|^2}$$

# Map of Europe by MDS

Map from CIA – The World Factbook: http://www.cia.gov/

# Linear Discriminant Analysis

□ Find a low-dimensional space such that when **x** is projected, classes are well-separated.

□ Find **w** that maximizes

$$J(\mathbf{w}) = \frac{(m_1 - m_2)^2}{s_1^2 + s_2^2}$$

$$m_1 = \frac{\sum_t \mathbf{w}^T \mathbf{x}^t r^t}{\sum_t r^t} \qquad s_1^2 = \sum_t \left(\mathbf{w}^T \mathbf{x}^t - m_1\right)^2 r^t$$

□ Between-class scatter:

$$(m_1 - m_2)^2 = \left(\mathbf{w}^T \mathbf{m}_1 - \mathbf{w}^T \mathbf{m}_2\right)^2$$

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

$$= \mathbf{w}^T \mathbf{S}_B \mathbf{w} \ \text{where} \ \mathbf{S}_B = (\mathbf{m}_1 - \mathbf{m}_2)(\mathbf{m}_1 - \mathbf{m}_2)^T$$

□ Within-class scatter:

$$s_1^2 = \sum_t \left(\mathbf{w}^T \mathbf{x}^t - m_1\right)^2 r^t$$

$$= \sum_t \mathbf{w}^T (\mathbf{x}^t - \mathbf{m}_1)(\mathbf{x}^t - \mathbf{m}_1)^T \mathbf{w} r^t = \mathbf{w}^T \mathbf{S}_1 \mathbf{w}$$

$$\text{where} \ \mathbf{S}_1 = \sum_t (\mathbf{x}^t - \mathbf{m}_1)(\mathbf{x}^t - \mathbf{m}_1)^T r^t$$

$$s_1^2 + s_1^2 = \mathbf{w}^T \mathbf{S}_W \mathbf{w} \ \text{where} \ \mathbf{S}_W = \mathbf{S}_1 + \mathbf{S}_2$$

# Fisher's Linear Discriminant

☐ Find **w** that max

$$J(\mathbf{w}) = \frac{\mathbf{w}^T \mathbf{S}_B \mathbf{w}}{\mathbf{w}^T \mathbf{S}_W \mathbf{w}} = \frac{\left| \mathbf{w}^T (\mathbf{m}_1 - \mathbf{m}_2) \right|^2}{\mathbf{w}^T \mathbf{S}_W \mathbf{w}}$$

☐ LDA soln:

$$\mathbf{w} = c \cdot \mathbf{S}_W^{-1} (\mathbf{m}_1 - \mathbf{m}_2)$$

☐ Parametric soln:

$$\mathbf{w} = \Sigma^{-1} (\mu_1 - \mu_2)$$
$$\text{when } p(\mathbf{x} | C_i) \sim \mathcal{N}(\mu_i, \Sigma)$$

# K>2 Classes

- Within-class scatter:

$$\mathbf{S}_W = \sum_{i=1}^{K} \mathbf{S}_i \qquad \mathbf{S}_i = \sum_{t} r_i^t \left( \mathbf{x}^t - \mathbf{m}_i \right) \left( \mathbf{x}^t - \mathbf{m}_i \right)^T$$
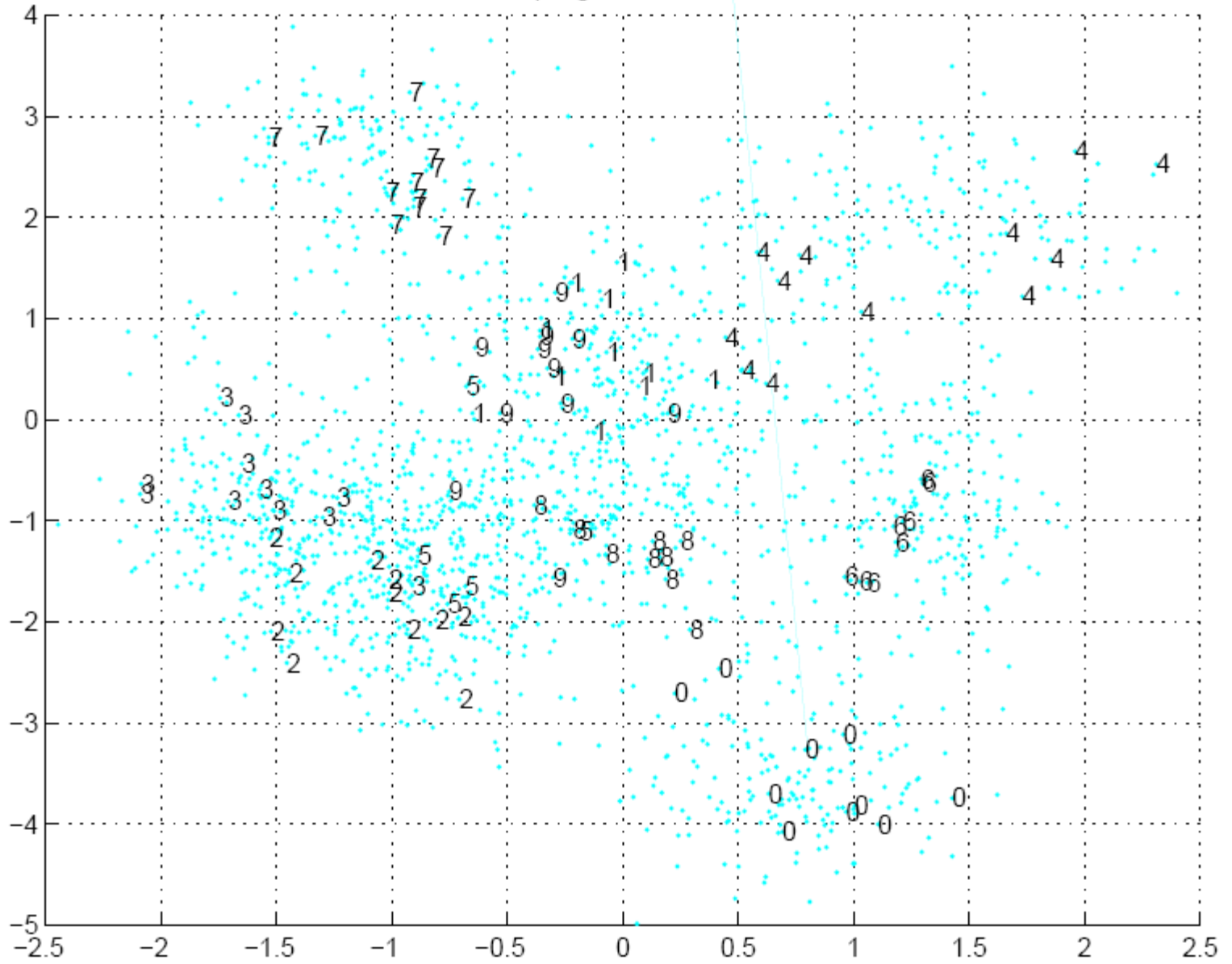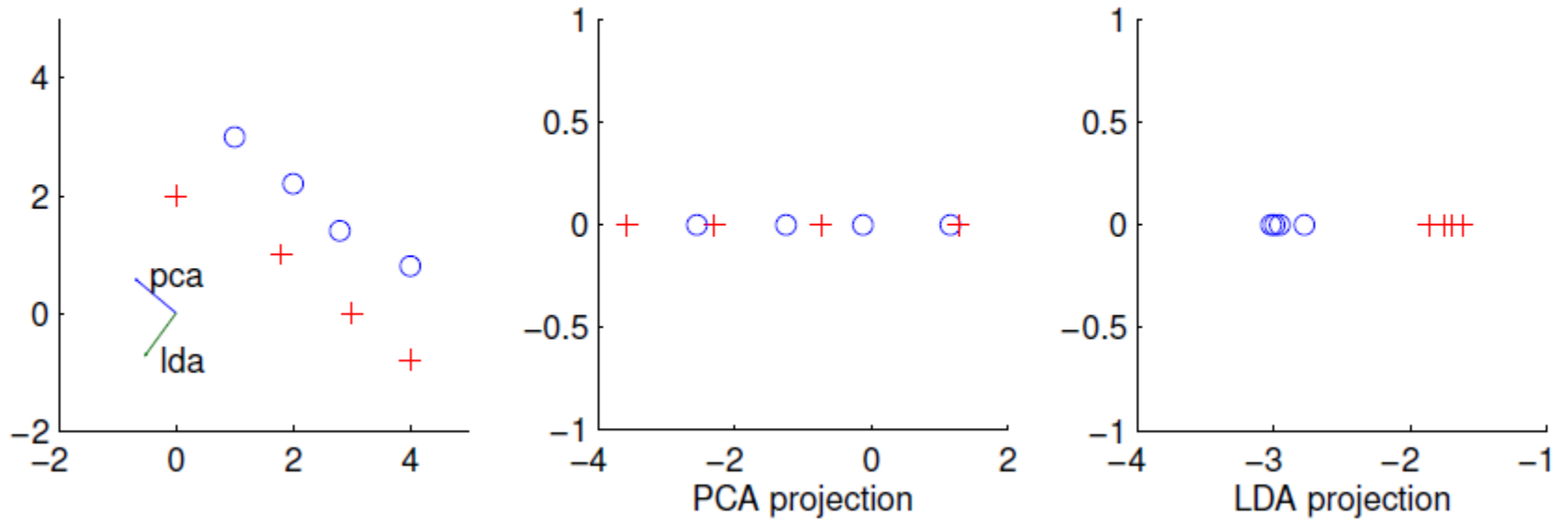
- Between-class scatter:

$$\mathbf{S}_B = \sum_{i=1}^{K} N_i \left( \mathbf{m}_i - \mathbf{m} \right) \left( \mathbf{m}_i - \mathbf{m} \right)^T \qquad \mathbf{m} = \frac{1}{K} \sum_{i=1}^{K} \mathbf{m}_i$$

- Find **W** that max

$$J(\mathbf{W}) = \frac{\left| \mathbf{W}^T \mathbf{S}_B \mathbf{W} \right|}{\left| \mathbf{W}^T \mathbf{S}_W \mathbf{W} \right|}$$

The largest eigenvectors of $\mathbf{S}_W^{-1}\mathbf{S}_B$; maximum rank of $K$-1

Optdigits after LDA
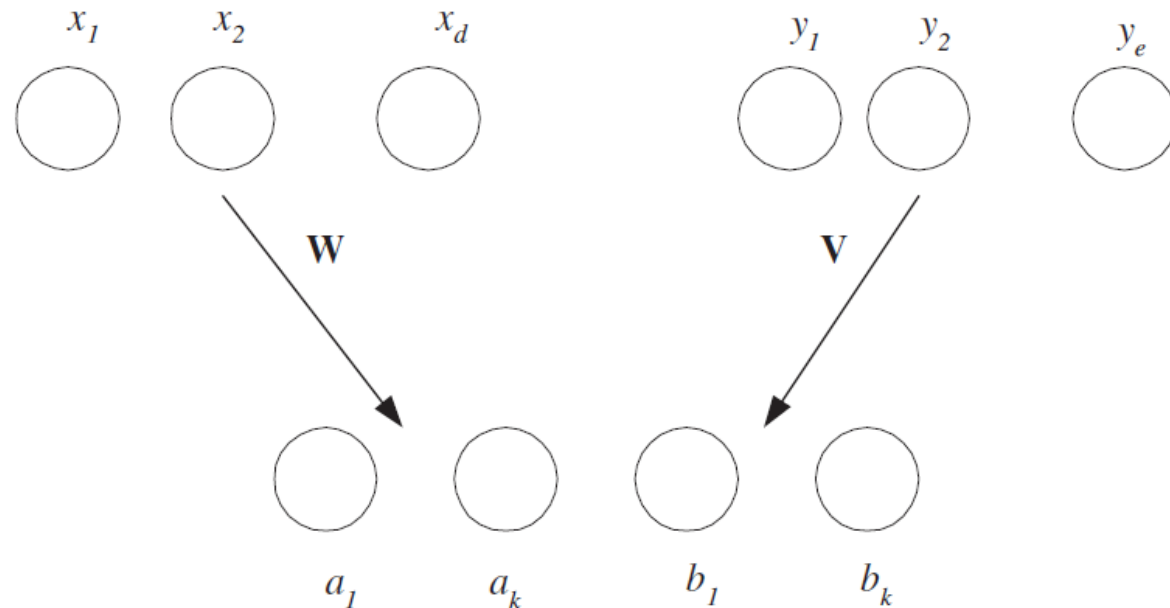
# PCA vs LDA

# Canonical Correlation Analysis

- $X = \{x^t, y^t\}_t$ ; two sets of variables $x$ and $y$ x

- We want to find two projections $w$ and $v$ st when $x$ is projected along $w$ and $y$ is projected along $v$, the correlation is maximized:

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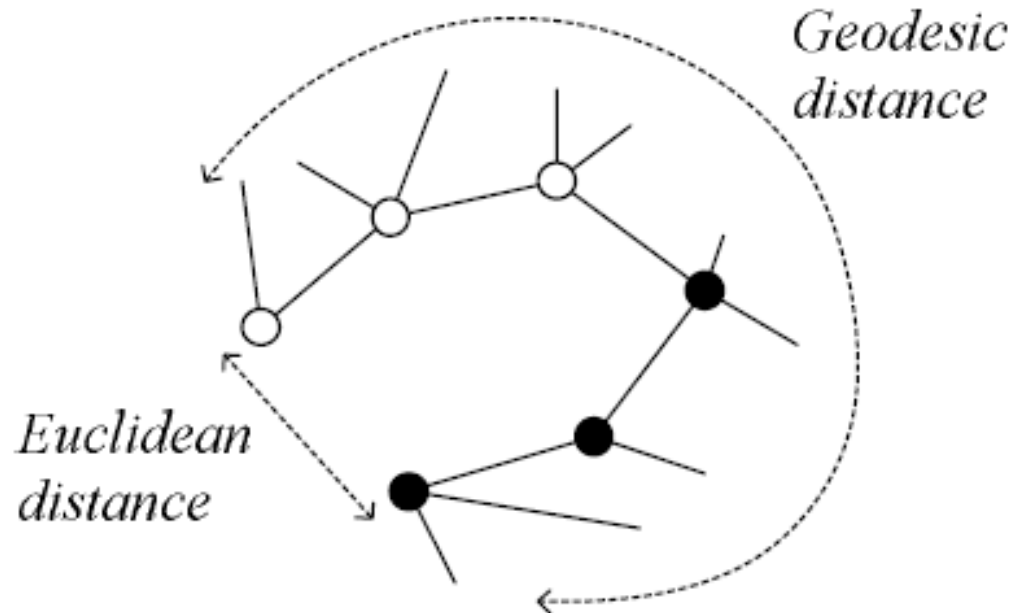\end{aligned}
$$

# CCA

- *x* and *y* may be two different views or modalities; e.g., image and word tags, and CCA does a joint mapping

# Isomap

- Geodesic distance is the distance along the manifold that the data lies in, as opposed to the Euclidean distance in the input space
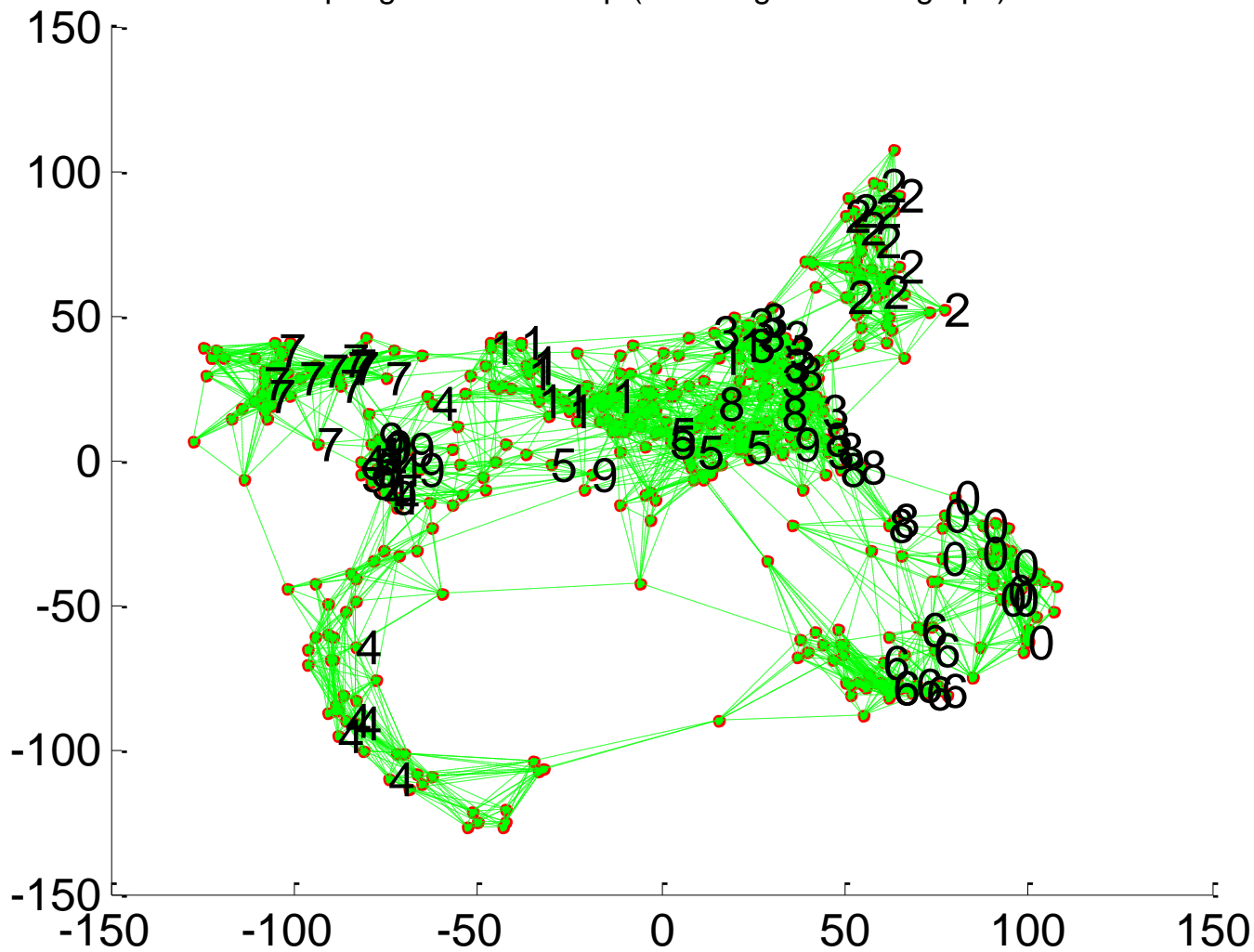
# Isomap

- Instances r and s are connected in the graph if $||x^r\text{-}x^s|| < \varepsilon$ or if $x^s$ is one of the *k* neighbors of $x^r$

  The edge length is $||x^r\text{-}x^s||$

- For two nodes r and s not connected, the distance is equal to the shortest path between them

- Once the *NxN* distance matrix is thus formed, use MDS to find a lower-dimensional mapping

Optdigits after Isomap (with neighborhood graph).

Matlab source from http://web.mit.edu/cocosci/isomap/isomap.html
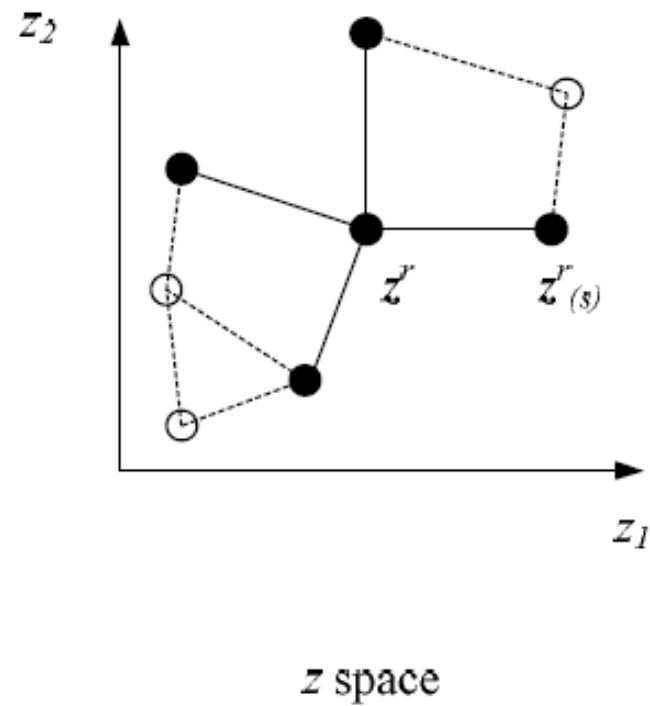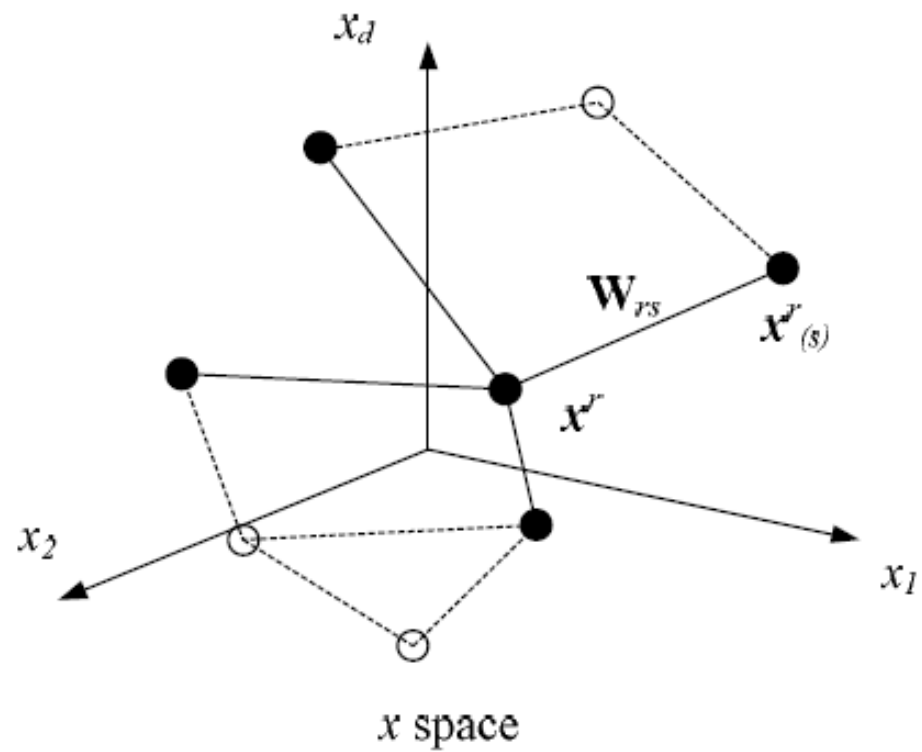
# Locally Linear Embedding

1. Given $\mathbf{x}^r$ find its neighbors $\mathbf{x}^s_{(r)}$

2. Find $\mathbf{W}_{rs}$ that minimize

$$E(\mathbf{W} \mid X) = \sum_r \left\| \mathbf{x}^r - \sum_s \mathbf{W}_{rs} \mathbf{x}^s_{(r)} \right\|^2$$
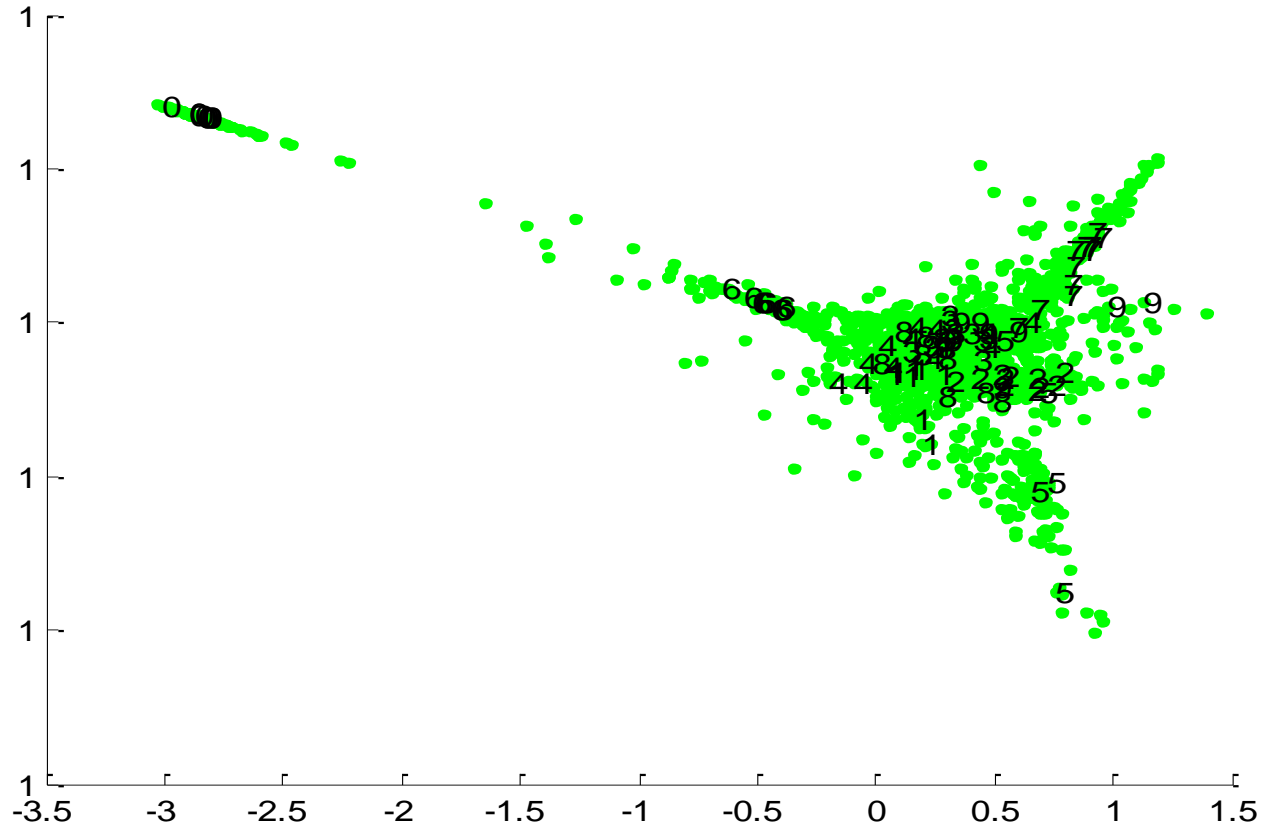
3. Find the new coordinates $\mathbf{z}^r$ that minimize

$$E(\mathbf{z} \mid \mathbf{W}) = \sum_r \left\| \mathbf{z}^r - \sum_s \mathbf{W}_{rs} \mathbf{z}^s_{(r)} \right\|^2$$

$x$ space

$z$ space

# LLE on Optdigits

Matlab source from http://www.cs.toronto.edu/~roweis/lle/code.html

# Laplacian Eigenmaps

□ Let *r* and *s* be two instances and $B_{rs}$ is their similarity, we want to find $\mathbf{z}^r$ and $\mathbf{z}^s$ that

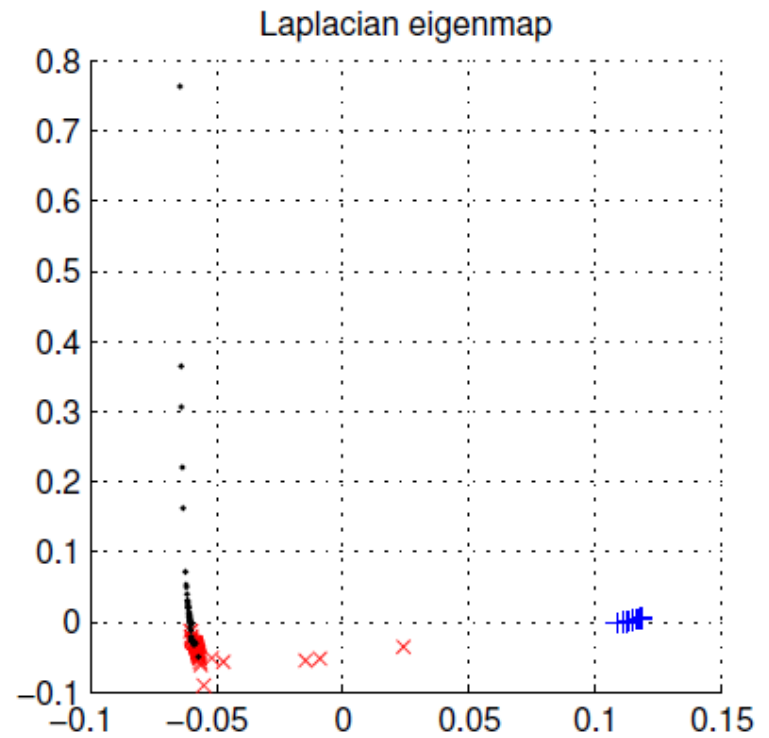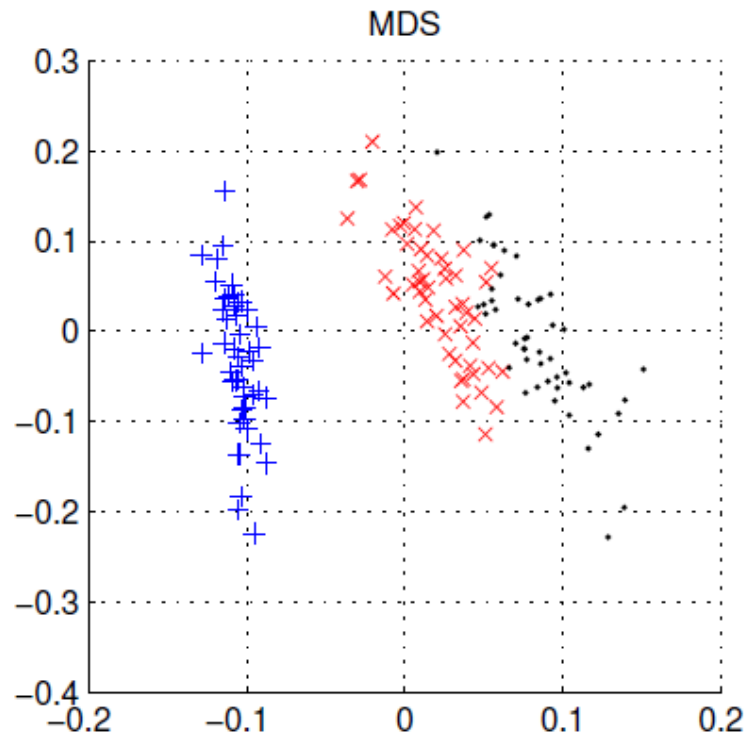$$\min \sum_{r,s} \|\mathbf{z}^r - \mathbf{z}^s\|^2 B_{rs}$$

□ $B_{rs}$ can be defined in terms of similarity in an original space: 0 if $\mathbf{x}^r$ and $\mathbf{x}^s$ are too far, otherwise

$$B_{rs} = \exp\left[-\frac{\|\mathbf{x}^r - \mathbf{x}^s\|^2}{2\sigma^2}\right]$$

□ Defines a graph Laplacian, and feature embedding returns $\mathbf{z}^r$

# Laplacian Eigenmaps on Iris

*Spectral clustering (chapter 7)*