

Lecture Slides for

INTRODUCTION TO
Machine Learning
2nd Edition

ETHEM ALPAYDIN
© The MIT Press, 2010

alpaydin@boun.edu.tr
<http://www.cmpe.boun.edu.tr/~ethem/i2ml2e>

CHAPTER 6:

Dimensionality Reduction

Why Reduce Dimensionality?

- Reduces time complexity: Less computation
- Reduces space complexity: Less 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 $x_i, i = 1, \dots, d$ dimensions to new $k < d$ dimensions, $z_j, j = 1, \dots, k$

Principal components analysis (PCA), linear discriminant analysis (LDA), factor analysis (FA)

Subset Selection

- There are 2^d subsets of d features
- Forward search: Add the best feature at each step
 - Set of features F initially \emptyset .
 - At each iteration, find the best new feature
$$j = \operatorname{argmin}_j E (F \cup x_j)$$
 - Add x_j to F if $E (F \cup x_j) < E (F)$
- 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)

Principal Components Analysis (PCA)

- Find a low-dimensional space such that when \mathbf{x} is projected there, information loss is minimized.
- The projection of \mathbf{x} on the direction of \mathbf{w} is: $z = \mathbf{w}^T \mathbf{x}$
- Find \mathbf{w} such that $\text{Var}(z)$ is maximized

$$\begin{aligned}\text{Var}(z) &= \text{Var}(\mathbf{w}^T \mathbf{x}) = E[(\mathbf{w}^T \mathbf{x} - \mathbf{w}^T \boldsymbol{\mu})^2] \\ &= E[(\mathbf{w}^T \mathbf{x} - \mathbf{w}^T \boldsymbol{\mu})(\mathbf{w}^T \mathbf{x} - \mathbf{w}^T \boldsymbol{\mu})] \\ &= E[\mathbf{w}^T (\mathbf{x} - \boldsymbol{\mu})(\mathbf{x} - \boldsymbol{\mu})^T \mathbf{w}] \\ &= \mathbf{w}^T E[(\mathbf{x} - \boldsymbol{\mu})(\mathbf{x} - \boldsymbol{\mu})^T] \mathbf{w} = \mathbf{w}^T \boldsymbol{\Sigma} \mathbf{w}\end{aligned}$$

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

- Maximize $\text{Var}(z)$ subject to $||\mathbf{w}||=1$

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

$\Sigma \mathbf{w}_1 = \alpha \mathbf{w}_1$ that is, \mathbf{w}_1 is an eigenvector of Σ

Choose the one with the largest eigenvalue for $\text{Var}(z)$ to be max

- Second principal component: Max $\text{Var}(z_2)$, s.t., $||\mathbf{w}_2||=1$ and orthogonal to \mathbf{w}_1

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

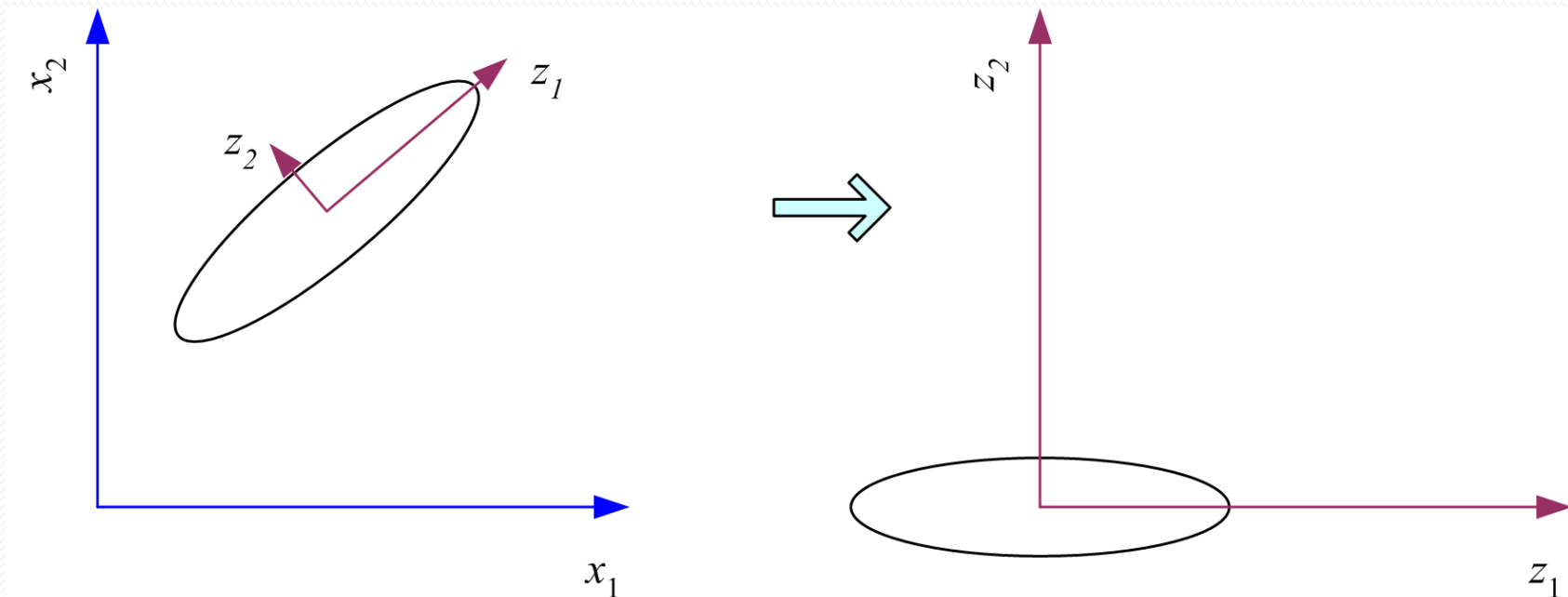
$\Sigma \mathbf{w}_2 = \alpha \mathbf{w}_2$ that is, \mathbf{w}_2 is another eigenvector of Σ
and so on.

What PCA does

$$\mathbf{z} = \mathbf{W}^T(\mathbf{x} - \mathbf{m})$$

where the columns of \mathbf{W} are the eigenvectors of Σ , and \mathbf{m} is sample mean

Centers the data at the origin and rotates the axes



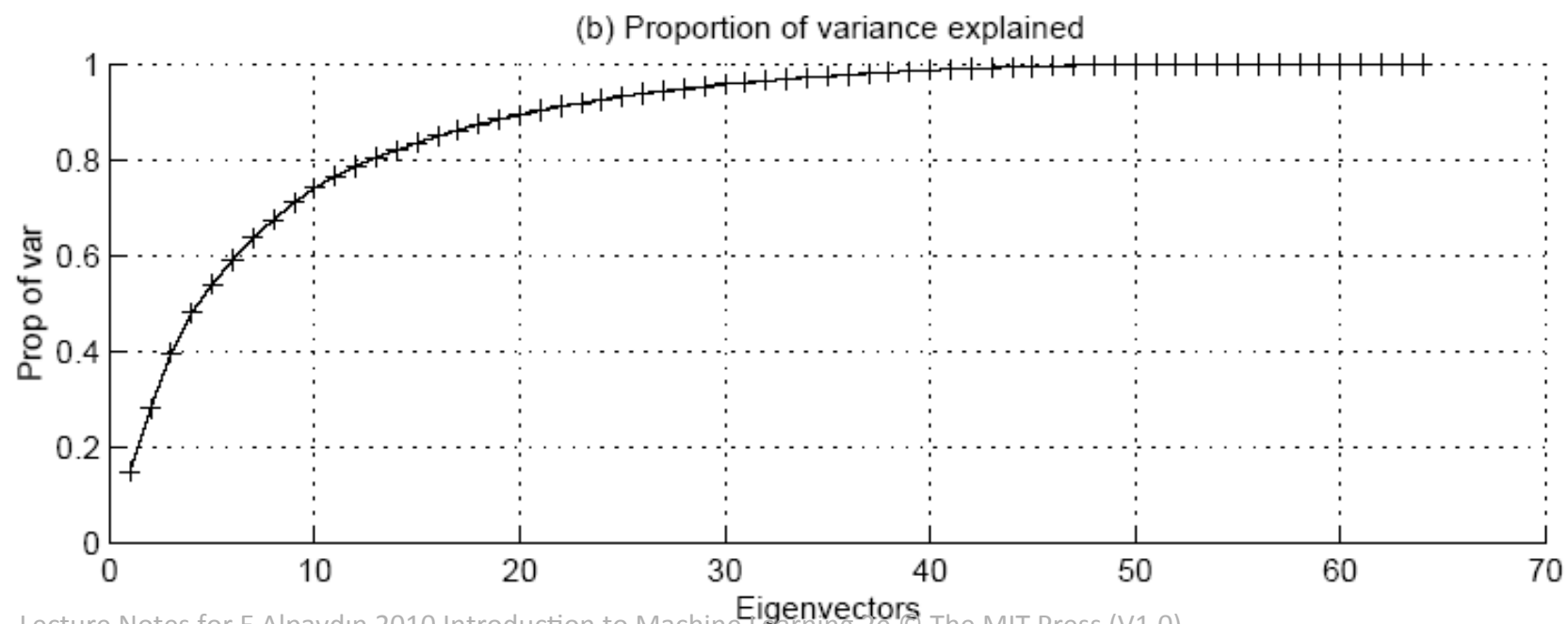
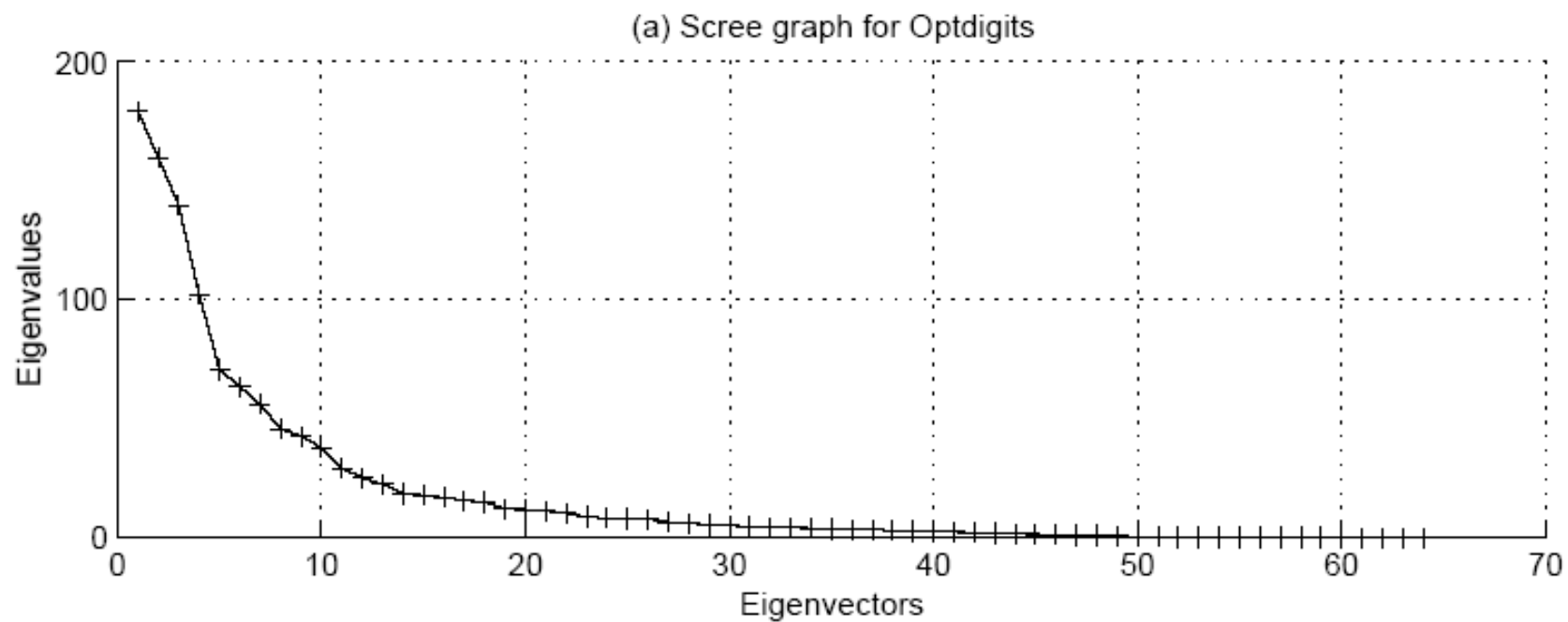
How to choose k ?

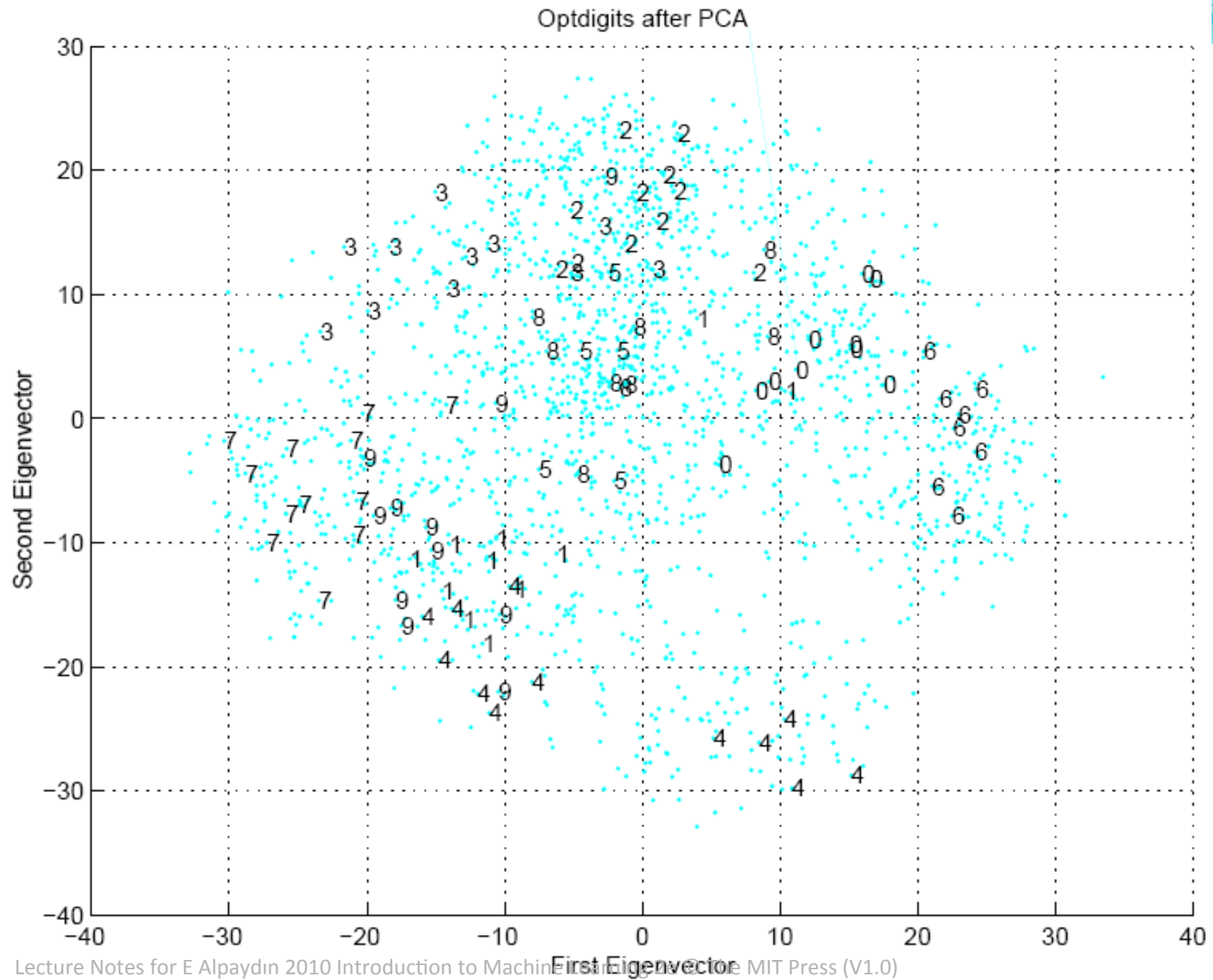
- Proportion of Variance (PoV) explained

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

when λ_i are sorted in descending order

- Typically, stop at PoV>0.9
- Scree graph plots of PoV vs k , stop at “elbow”





Factor Analysis

- Find a small number of factors \mathbf{z} , which when combined generate \mathbf{x} :

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

where $z_j, j = 1, \dots, k$ are the latent factors with

$$E[z_j] = 0, \text{Var}(z_j) = 1, \text{Cov}(z_i, z_j) = 0, i \neq j,$$

ε_i are the noise sources

$$E[\varepsilon_i] = \psi_i, \text{Cov}(\varepsilon_i, \varepsilon_j) = 0, i \neq j, \text{Cov}(\varepsilon_i, z_j) = 0,$$

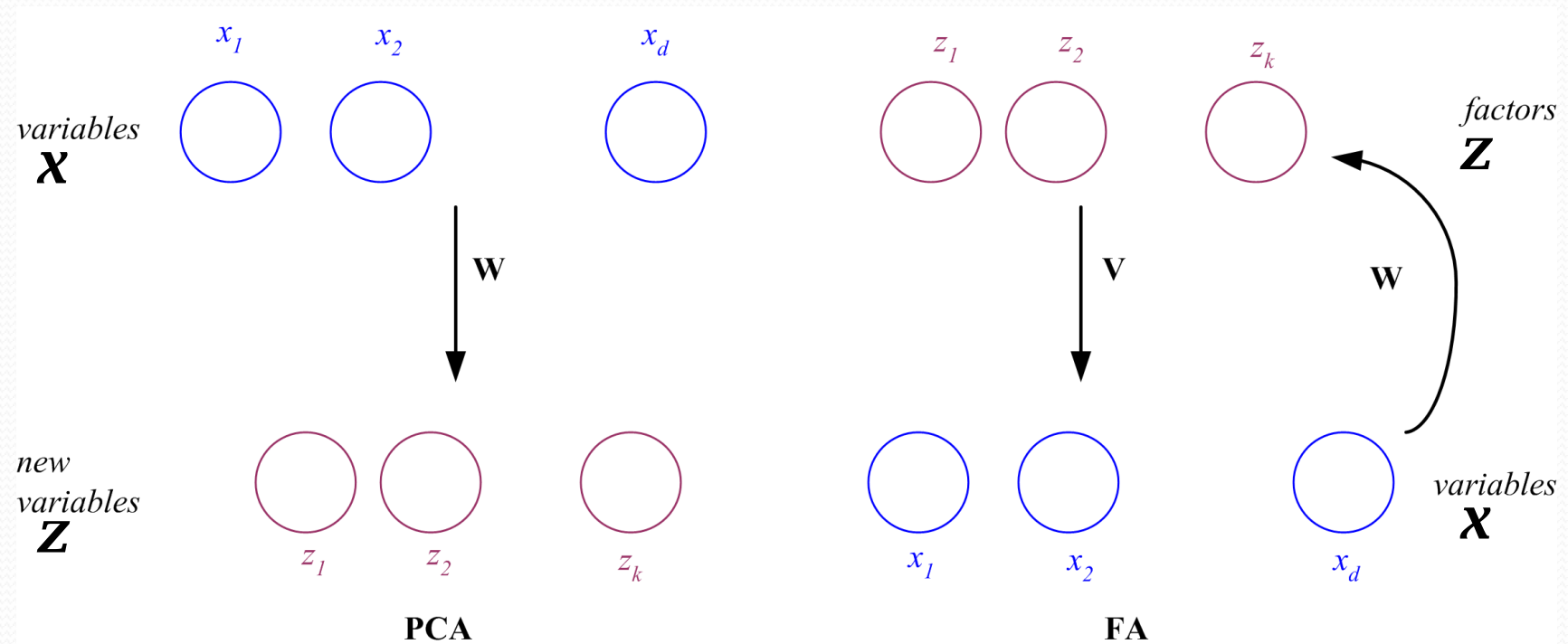
and v_{ij} are the factor loadings

PCA vs FA

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

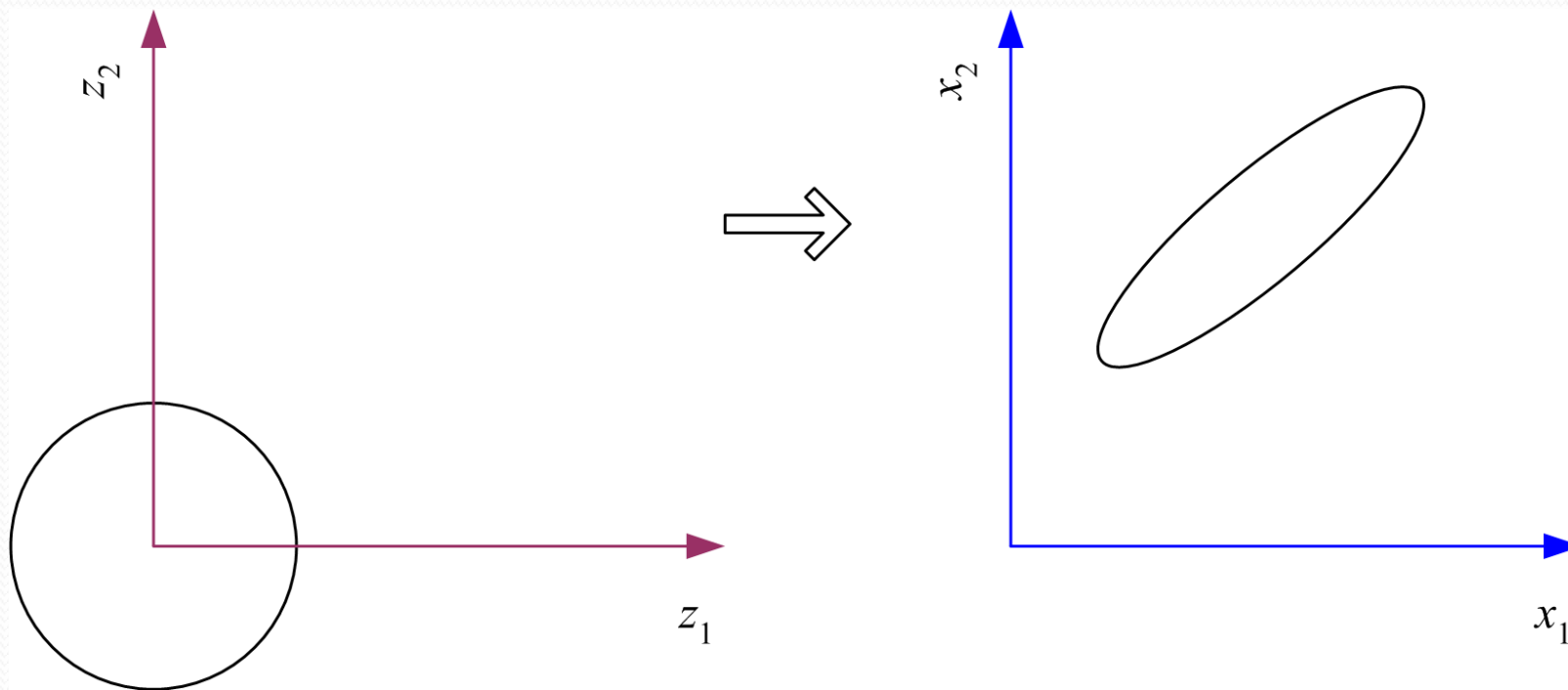
- FA From \mathbf{z} to \mathbf{x}

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



Factor Analysis

- In FA, factors z_j are stretched, rotated and translated to generate \mathbf{x}



Multidimensional Scaling

- Given pairwise distances between N points,

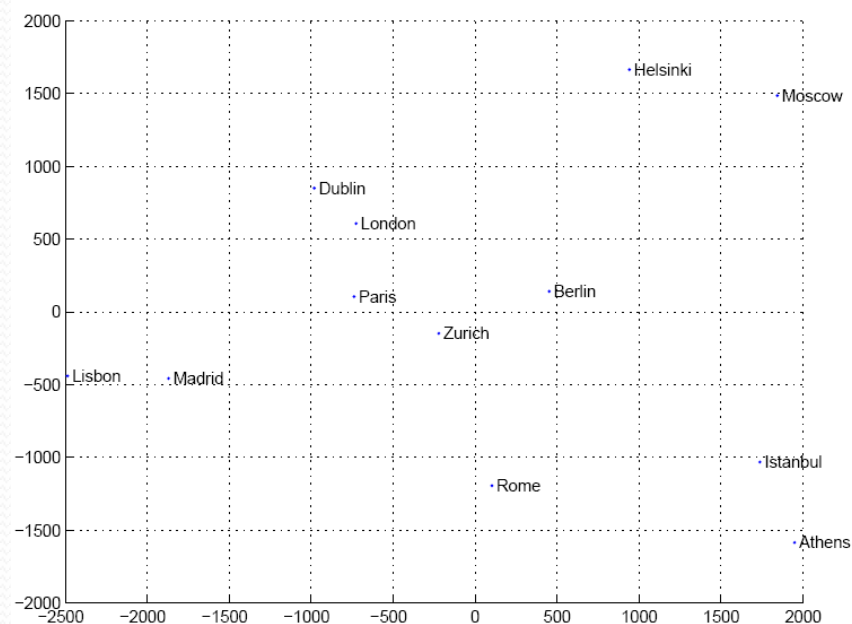
$$d_{ij}, i, j = 1, \dots, N$$

place on a low-dim map s.t. distances are preserved.

- $\mathbf{z} = \mathbf{g}(\mathbf{x} \mid \vartheta)$ Find ϑ that min Sammon stress

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

Map of Europe by MDS



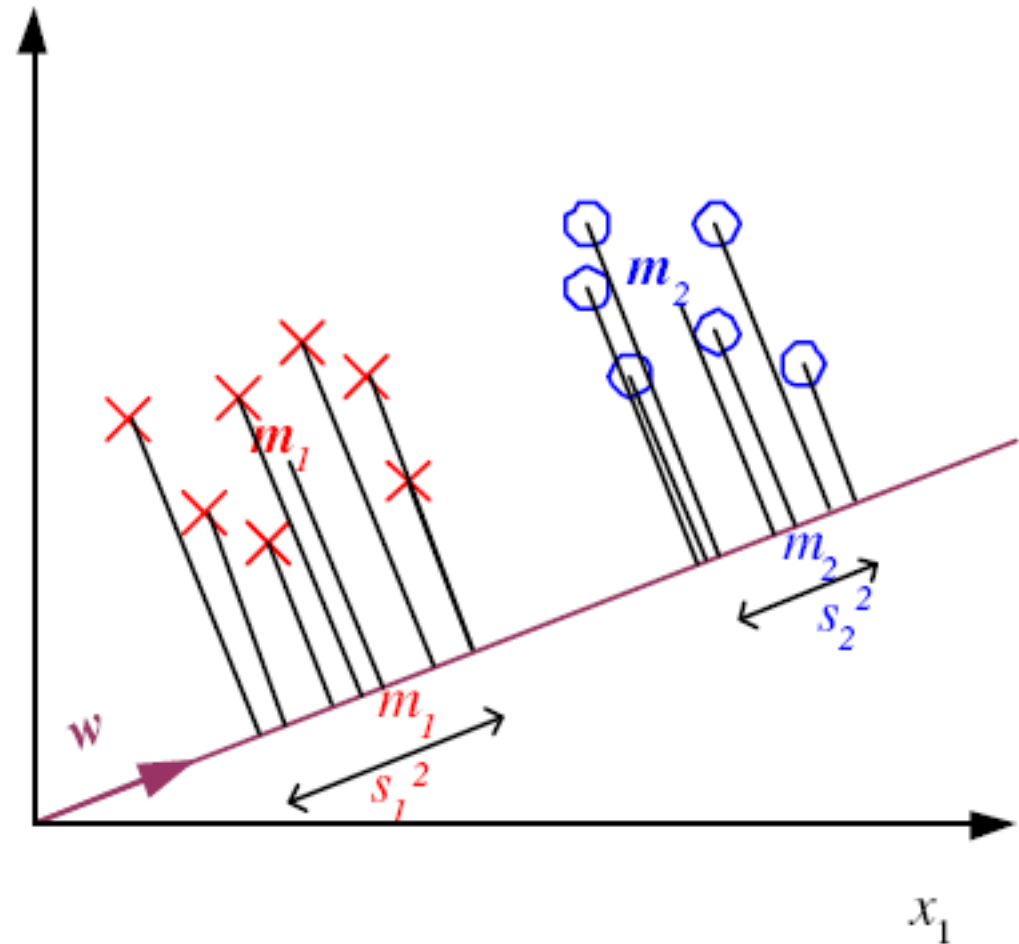
Map from CIA – The World Factbook: <http://www.cia.gov/>

Linear Discriminant Analysis

- Find a low-dimensional space such that when \mathbf{x} is projected, classes are well-separated.
- Find \mathbf{w} that maximizes

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

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



- Between-class scatter:

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

- Within-class scatter:

$$\begin{aligned}s_1^2 &= \sum_t (\mathbf{w}^T \mathbf{x}^t - m_1)^2 r^t \\ &= \sum_t \mathbf{w}^T (\mathbf{x}^t - \mathbf{m}_1) (\mathbf{x}^t - \mathbf{m}_1)^T \mathbf{w} r^t = \mathbf{w}^T \mathbf{S}_1 \mathbf{w}\end{aligned}$$

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

$$s_1^2 + s_2^2 = \mathbf{w}^T \mathbf{S}_W \mathbf{w} \text{ where } \mathbf{S}_W = \mathbf{S}_1 + \mathbf{S}_2$$

Fisher's Linear Discriminant

- Find \mathbf{w} that max

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

- 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:

$$\mathbf{S}_W = \sum_{i=1}^K \mathbf{S}_i \quad \mathbf{S}_i = \sum_t r_i^t (\mathbf{x}^t - \mathbf{m}_i)(\mathbf{x}^t - \mathbf{m}_i)^T$$

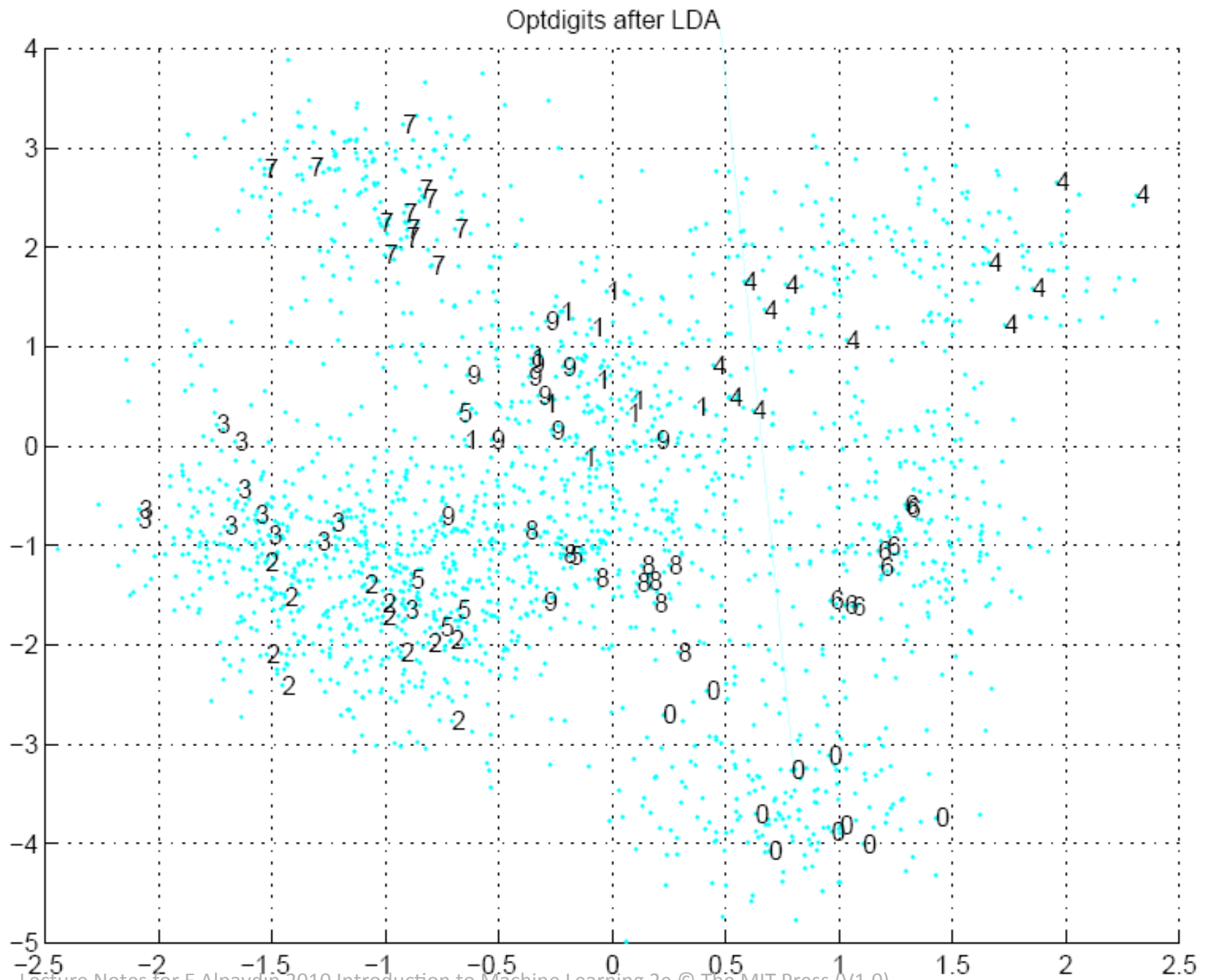
- Between-class scatter:

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

- Find \mathbf{W} that max

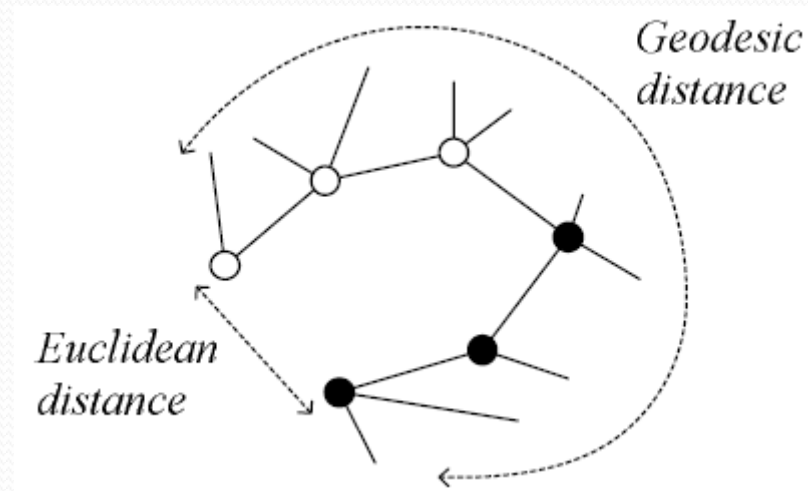
$$J(\mathbf{W}) = \frac{|\mathbf{W}^T \mathbf{S}_B \mathbf{W}|}{|\mathbf{W}^T \mathbf{S}_W \mathbf{W}|}$$

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



Isomap

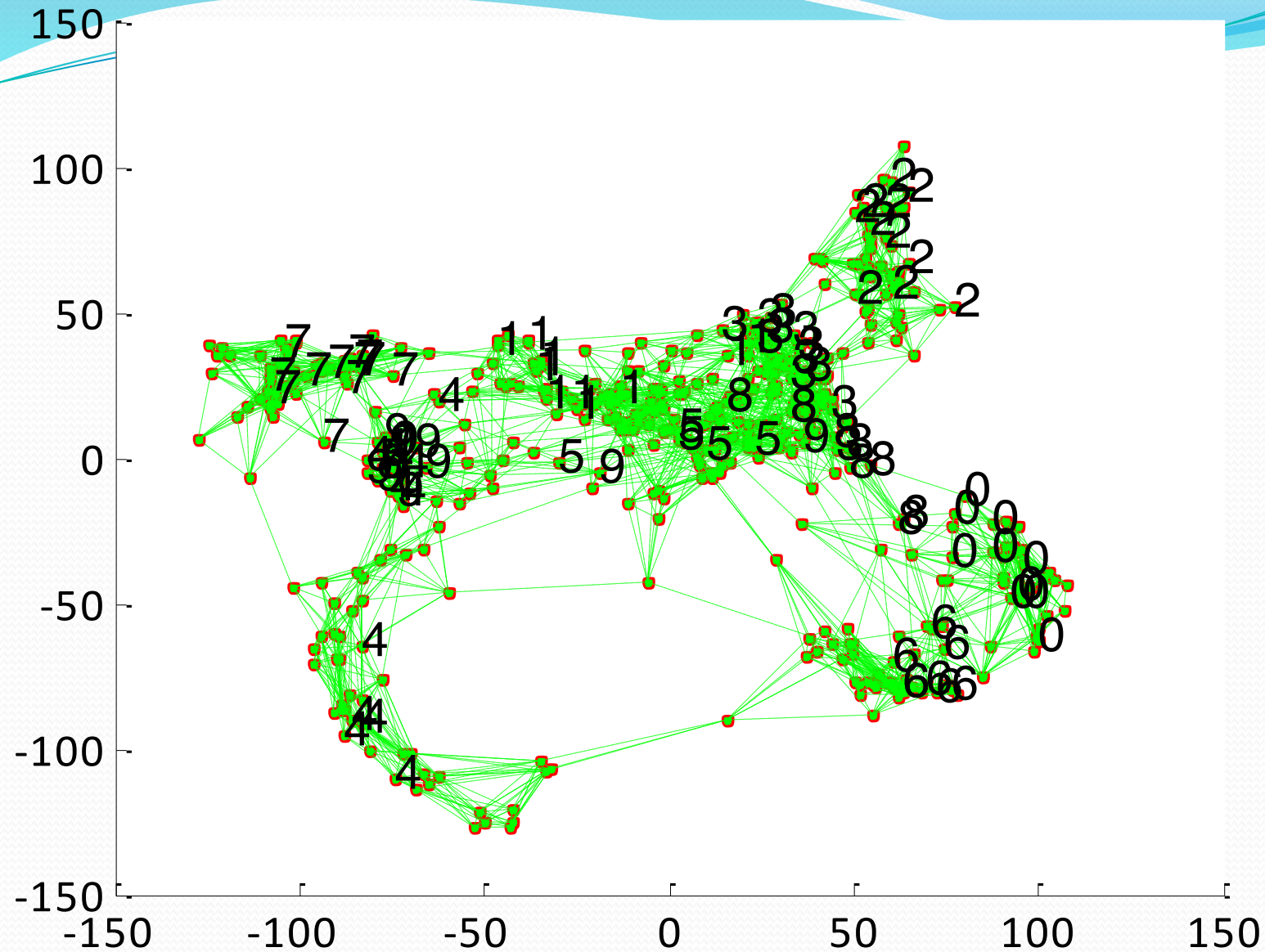
- Geodesic distance is the distance along the manifold that the data lies in, as opposed to the Euclidean distance in the input space



Isomap

- Instances r and s are connected in the graph if $||\mathbf{x}^r - \mathbf{x}^s|| < \epsilon$ or if \mathbf{x}^s is one of the k neighbors of \mathbf{x}^r
The edge length is $||\mathbf{x}^r - \mathbf{x}^s||$
- For two nodes r and s not connected, the distance is equal to the shortest path between them
- Once the $N \times N$ 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>

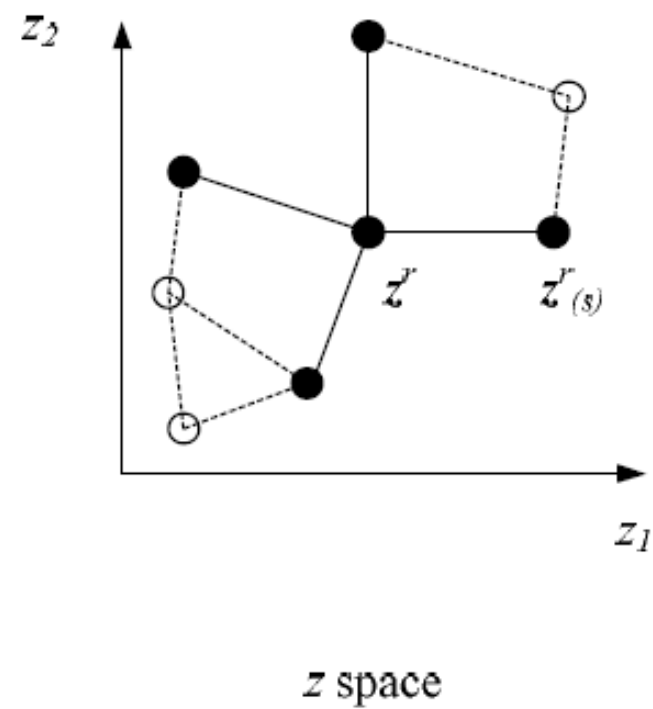
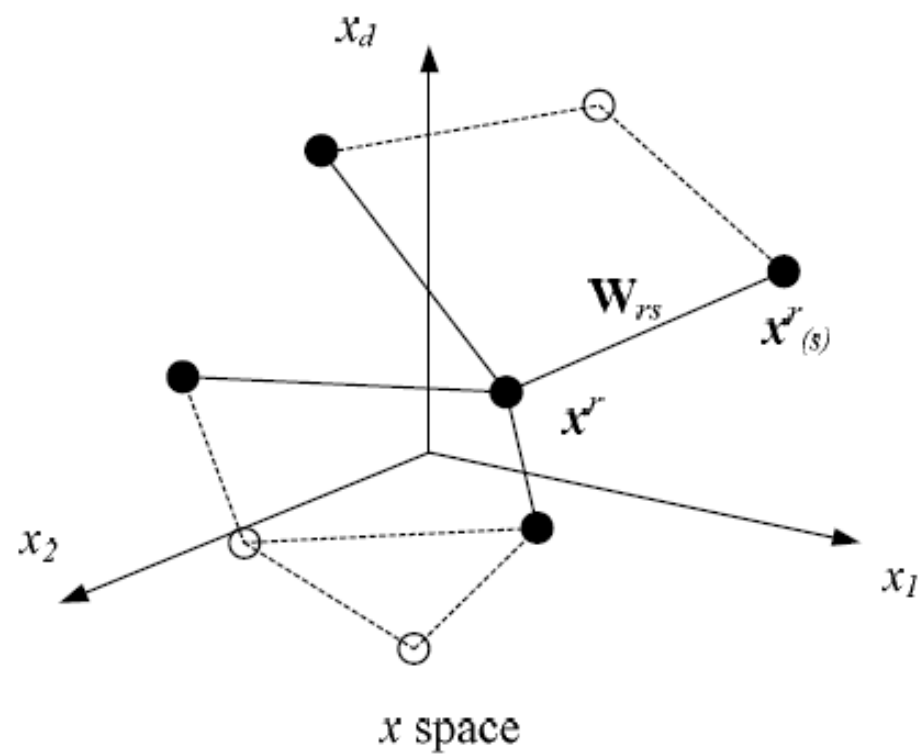
Locally Linear Embedding

1. Given \mathbf{x}^r find its neighbors $\mathbf{x}_{(r)}^s$
2. Find \mathbf{W}_{rs} that minimize

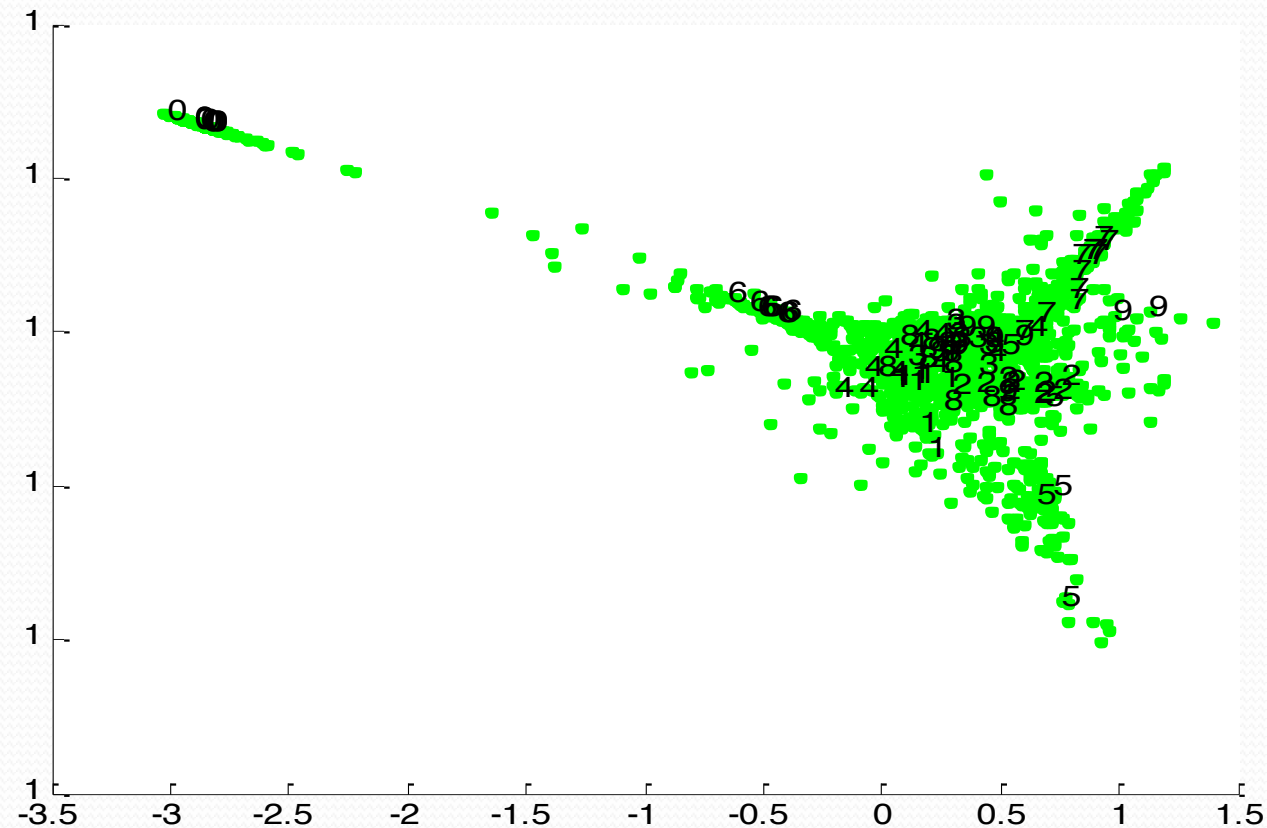
$$E(\mathbf{W} | X) = \sum_r \left\| \mathbf{x}^r - \sum_s \mathbf{W}_{rs} \mathbf{x}_{(r)}^s \right\|^2$$

3. Find the new coordinates \mathbf{z}^r that minimize

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

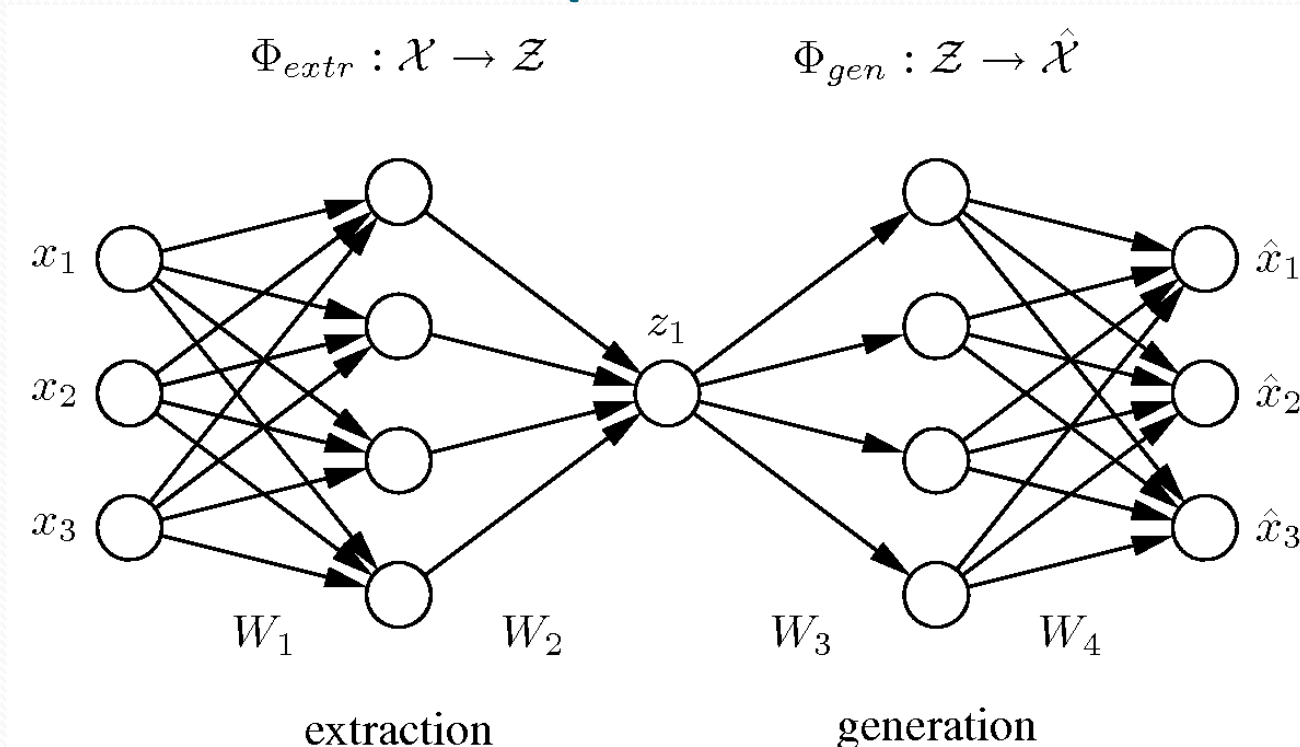


LLE on Optdigits



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

Nonlinear PCA (Nonlinear autoencoder)



You need at least 4 layers of weights for nonlinear dimensionality reduction.
Note that it could be very costly to train the neural network for a large number of features D

Source: Matthias Scholz www.nlpca.de/

Kernel PCA [Scholkopf et.al. 1998]

- Kernel substitution: allows expression of an algorithm only in terms of kernels $k(x, x_n) = \phi(x)^T \phi(x_n)$ (dot product in the $\phi(x)$ space which could be very high dimensional)
- Kernel PCA: extend PCA so that instance vectors only appear in terms of dot products.
- If $k(x, x_n) = x^T x_n$ kernel PCA reduces to PCA

Kernel PCA

Assume both x and $\phi(x)$ [for the time being] have zero mean.

$$Su_i = \lambda_i u_i \quad S_{D \times D} = \frac{1}{N} \sum_{n=1}^N x_n x_n^T \quad u_i^T u_i = 1, \quad Cv_i = \lambda_i v_i \quad C_{M \times M} = \frac{1}{N} \sum_{n=1}^N \phi(x_n) \phi(x_n)^T$$

$$\frac{1}{N} \sum_{n=1}^N \phi(x_n) \phi(x_n)^T v_i = \lambda_i v_i \rightarrow v_i = \sum_{n=1}^N a_{in} \phi(x_n)$$

$$\frac{1}{N} \sum_{n=1}^N \phi(x_n) \phi(x_n)^T \sum_{m=1}^N a_{im} \phi(x_m) = \lambda_i \sum_{n=1}^N a_{in} \phi(x_n)$$

$$\frac{1}{N} \sum_{n=1}^N k(x_l, x_n) \sum_{m=1}^N a_{im} k(x_n, x_m) = \lambda_i \sum_{n=1}^N a_{in} k(x_l, x_n) \quad // \text{multiply both sides by } \phi(x_l)$$

$$K^2 a_i = \lambda_i N K a_i \rightarrow K a_i = \lambda_i N a_i$$

Taking into account unit length of v_i and nonzero mean $\phi(x)$, compute :

$$\tilde{K} = K - 1_N K - K 1_N + 1_N K 1_N \quad \text{compute eigen vectors of } \tilde{K}$$

projection of a point x onto eigenvector i is given by :

$$y_i(x) = \phi(x)^T v_i = \sum_{n=1}^N a_{in} \phi(x)^T \phi(x_n) = \sum_{n=1}^N a_{in} k(x, x_n)$$

mRMR (minimum Redundancy Maximum Relevance) [Peng 2003]

- Measure feature-feature (redundancy) and feature-label (relevance) correlations using mutual information:

$$I(X, Y) = \sum_{y \in Y} \sum_{x \in X} p(x, y) \log \left(\frac{p(x, y)}{p_1(x) p_2(y)} \right)$$

- S: set of selected features, F_i, F_j features, H: labels

$$\text{min Red, } \text{Red} = \frac{1}{|S|^2} \sum_{F_i, F_j \in S} I(F_i, F_j) \quad \text{max Rel, } \text{Rel} = \frac{1}{|S|} \sum_{F_i \in S} I(F_i, H)$$

- MID: max Rel-Red
- MIQ: max Rel/Red

mRMR Algorithm

for all F_i in $i=1..d$ do

 Compute Rel_i between F_i and H using MI

end for

Sort (decreasing) features according to their Rel_i values

Initialize feature subset $S = \{\text{the most relevant feature}\}$

$i = 2$

while $i \leq d$

 Compute MIQ_i (or MID_i) for each unselected feature

 let j = the feature with max MIQ (or MID)

$S \leftarrow S \cup F_j$

$i=j+1$

endwhile

mRMR Notes

Need to discretize features in order to compute MI, or need to use a nonparametric method to compute MI.

Advantages of mRMR:

mRMR is much faster than wrapper methods (i.e. forward-backward selection)

Since it takes into account the label information it is more beneficial for classification than PCA

Since MI is a nonlinear measure of similarity, even if there are nonlinear correlations between features/labels, they are taken into account.

Some Feature Selection Tools

- Weka
- PrTools
- Many methods are easily implemented using Matlab
- mRMR source code is available from Peng's web site.