



*Lecture Slides for*

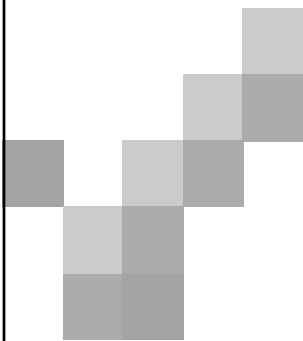
INTRODUCTION TO

*Machine Learning*

ETHEM ALPAYDIN  
© The MIT Press, 2004

**Edited for CS 536 Fall 2005 – Rutgers University  
Ahmed Elgammal**

*alpaydin@boun.edu.tr*  
*<http://www.cmpe.boun.edu.tr/~ethem/i2ml>*



CHAPTER 6:

*Dimensionality  
Reduction*

## *Why Reduce Dimensionality?*

1. Reduces time complexity: Less computation
2. Reduces space complexity: Less parameters
3. Saves the cost of observing the feature
4. Simpler models are more robust on small datasets
5. More interpretable; simpler explanation
6. 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  
 $j = \operatorname{argmin}_i E(F \cup x_i)$
  - 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  $\operatorname{Var}(z)$  is maximized

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

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

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

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

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

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  $\mathbf{w}_1$

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

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

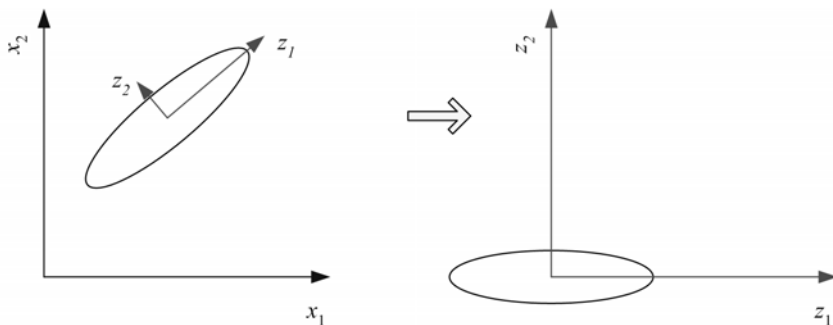
7

## What PCA does

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

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

Centers the data at the origin and rotates the axes



8

## How to choose $k$ ?

- Proportion of Variance (PoV) explained

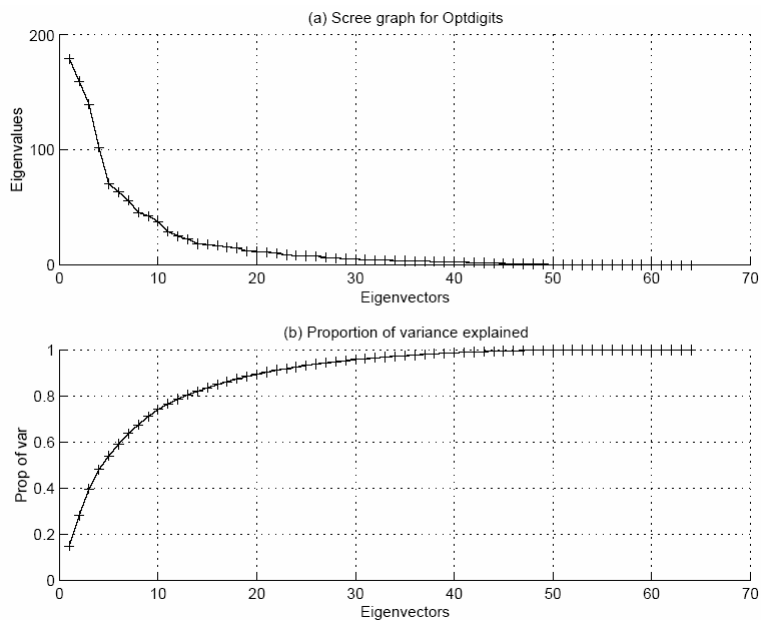
$$\frac{\lambda_1 + \lambda_2 + \dots + \lambda_k}{\lambda_1 + \lambda_2 + \dots + \lambda_k + \dots + \lambda_d}$$

when  $\lambda_i$  are sorted in descending order

- Typically, stop at  $\text{PoV} > 0.9$
- Scree graph plots of PoV vs  $k$ , stop at “elbow”

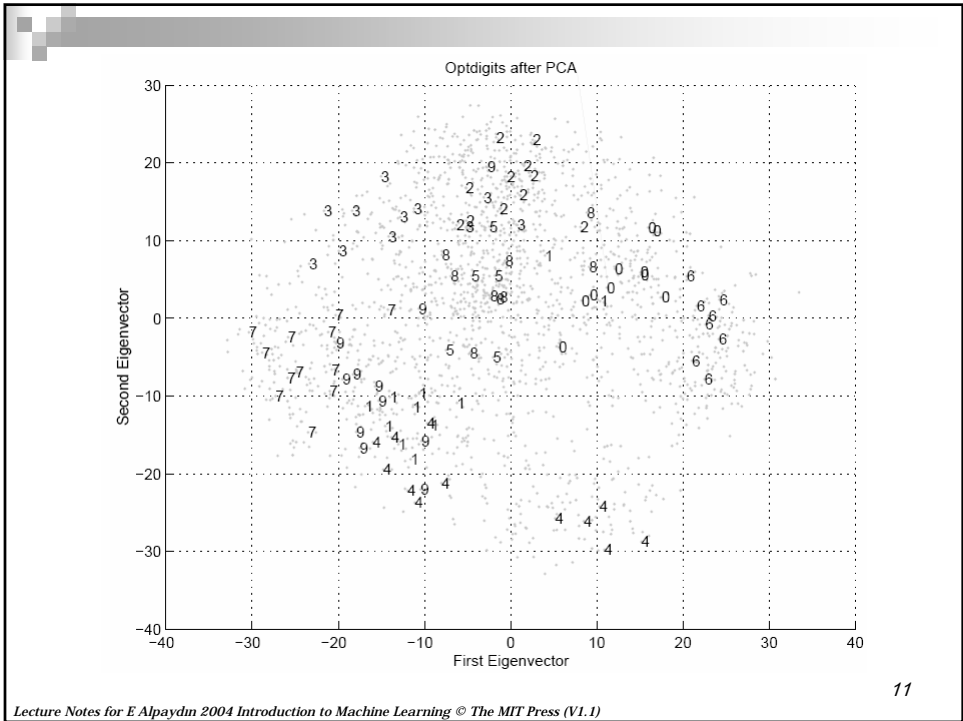
9

Lecture Notes for E Alpaydm 2004 Introduction to Machine Learning © The MIT Press (V1.1)



10

Lecture Notes for E Alpaydm 2004 Introduction to Machine Learning © The MIT Press (V1.1)



## Principle Component Analysis

### PCA

- Given a set of points  $\{x_1, x_2, \dots, x_N\}, x_i \in R^d$
- We are looking for a linear projection: a linear combination of orthogonal basis vectors

$x \approx A \cdot c$

$R^d$

$R^m, m \ll d$

$x$

$\approx$

$c^1$

+

$c^2$

+

$c^3$

+

...

+

$c^m$

$A$

$x$

$\approx$

...

$c$

What is the projection that minimizes the reconstruction error ?

$$E = \sum_i \|x_i - Ac_i\|^2$$

12

Lecture Notes for E Alpaydm 2004 Introduction to Machine Learning © The MIT Press (V1.1)

# Principle Component Analysis

## PCA

- Given a set of points

$$\{x_1, x_2, \dots, x_N\}, x_i \in R^d$$

- Center the points: compute

$$\mu = \frac{1}{N} \sum_i x_i$$

$$P = [x_1 - \mu, x_2 - \mu, \dots, x_N - \mu], x_i \in R^d$$

- Compute covariance matrix  $Q = PP^T$
- Compute the eigen vectors for  $Q \longrightarrow Qe_k = \lambda_k e_k$
- Eigenvectors are the orthogonal basis we are looking for*

13

## Singular Value Decomposition

- SVD: If  $A$  is a real  $m$  by  $n$  matrix then there exist orthogonal matrices  $U$  ( $m \cdot m$ ) and  $V$  ( $n \cdot n$ ) such that

$$U^t A V = \Sigma = \text{diag}(\sigma_1, \sigma_2, \dots, \sigma_p) \quad p = \min\{m, n\}$$

$$U^t A V = \Sigma \quad A = U \Sigma V^t$$

- Singular values:** Non negative square roots of the eigenvalues of  $A^t A$ . Denoted  $\sigma_i, i=1, \dots, n$
- $A^t A$  is symmetric  $\Rightarrow$  eigenvalues and singular values are real.
- Singular values arranged in decreasing order.

$$A^t A = (U \Sigma V^t)^t (U \Sigma V^t) = V \Sigma^t U^t U \Sigma V^t = V \Sigma^t \Sigma V^t = V \Sigma^2 V^{-1}$$

$$(A^t A)V = V \Sigma^2$$

$$(A^t A)v = v \lambda$$

$$\begin{array}{|c|} \hline A \\ \hline m \times n \\ \hline \end{array} = \begin{array}{|c|} \hline U \\ \hline m \times m \\ \hline \end{array} \begin{array}{|c|} \hline \Sigma \\ \hline m \times n \\ \hline \end{array} \begin{array}{|c|} \hline V^t \\ \hline n \times n \\ \hline \end{array}$$

14

## SVD for PCA

- SVD can be used to efficiently compute the image basis

$$PP^t = (U \Sigma V^t)(U \Sigma V^t)^t = U \Sigma V^t V \Sigma^t U^t = U \Sigma^t \Sigma U^t = U \Sigma^2 U^{-1}$$

$$(PP^t)U = U \Sigma^2$$

$$(PP^t)v = v\lambda$$

- $U$  are the eigen vectors (basis)
- Most important thing to notice: Distance in the eigen-space is an approximation to the correlation in the original space

$$\|x_i - x_j\| \approx \|c_i - c_j\|$$

## PCA

$$\begin{array}{ccc} & x \approx Ac & \\ R^d \swarrow & & \nwarrow R^m, m \ll d \\ & c \approx A^T x & \end{array}$$

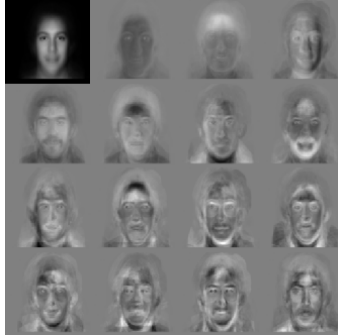
- Most important thing to notice: Distance in the eigen-space is an approximation to the correlation in the original space

$$\|x_i - x_j\| \approx \|c_i - c_j\|$$



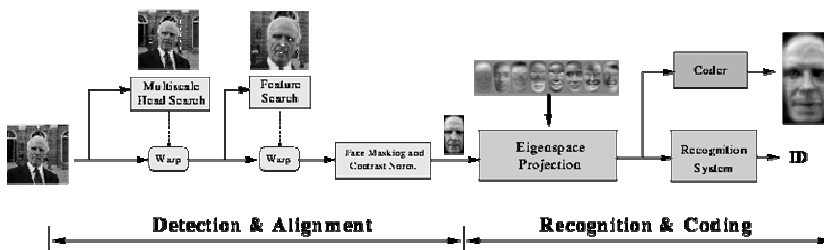
# Eigenface

- Use PCA and subspace projection to perform face recognition
- How to describe a face as a linear combination of face basis
- Matthew Turk and Alex Pentland “Eigenfaces for Recognition” 1991



# Face Recognition - Eigenface

- MIT Media Lab - Face Recognition demo page  
<http://vismod.media.mit.edu/vismod/demos/facerec/>



# 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

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

$\varepsilon_i$  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

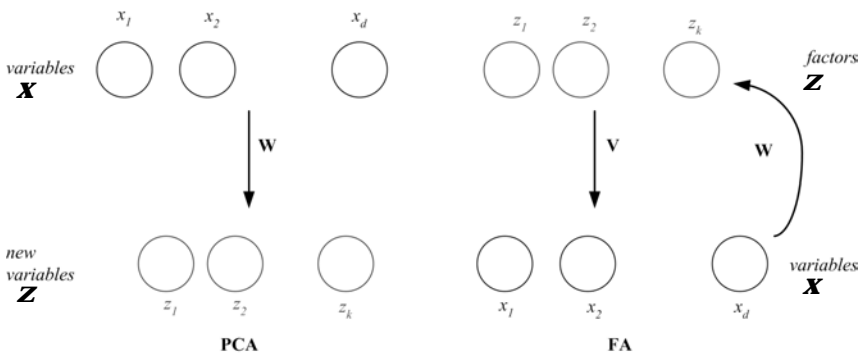
# PCA vs FA

- PCA From  $\mathbf{x}$  to  $\mathbf{z}$

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

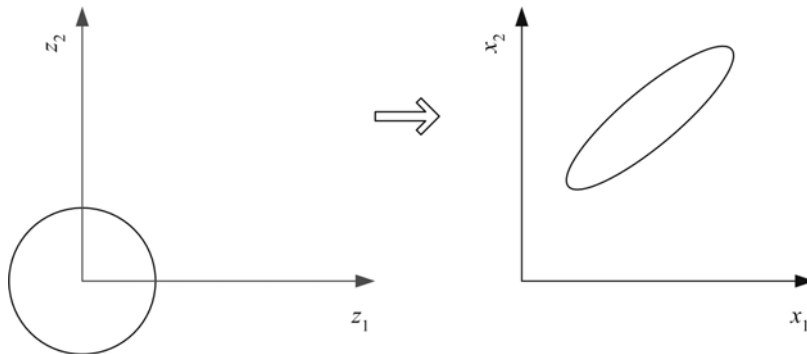
$$\mathbf{z} = \mathbf{W}^T(\mathbf{x} - \boldsymbol{\mu})$$

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



## Factor Analysis

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



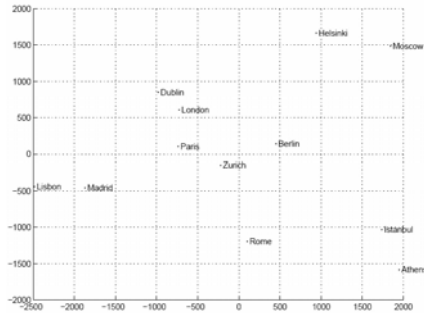
## Multidimensional Scaling

- Given pairwise distances between  $N$  points,  $d_{ij}$ ,  $i, j = 1, \dots, N$   
place on a low- dim map s.t. distances are preserved.
- $\mathbf{z} = \mathbf{g}(\mathbf{x} | \theta)$  Find  $\theta$  that min Sammon stress

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

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

## Map of Europe by MDS



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

23

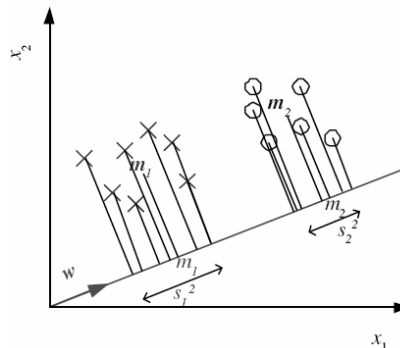
Lecture Notes for E Alpaydm 2004 Introduction to Machine Learning © The MIT Press (V1.1)

## Linear Discriminant Analysis

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

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

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



24

Lecture Notes for E Alpaydm 2004 Introduction to Machine Learning © The MIT Press (V1.1)

- Between-class scatter:

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

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

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

25

## *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:

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

- Parametric soln:

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

26

## $K > 2$ Classes

- Within-class scatter:

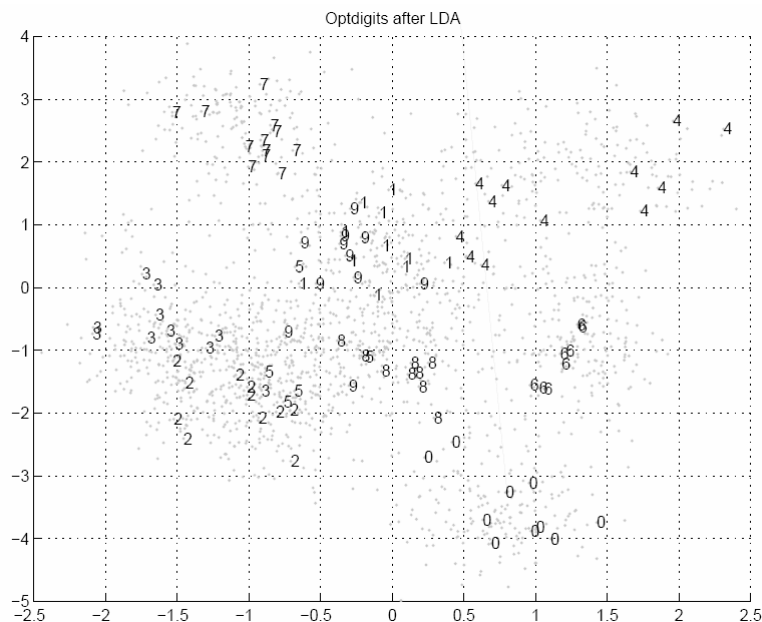
$$\mathbf{S}_W = \sum_{i=1}^K \mathbf{S}_i \quad \mathbf{S}_i = \sum_t r_i^t (\mathbf{x}^t - \mathbf{m}_i)(\mathbf{x}^t - \mathbf{m}_i)^T$$

- 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

$$J(\mathbf{W}) = \frac{|\mathbf{W}^T \mathbf{S}_B \mathbf{W}|}{|\mathbf{W}^T \mathbf{S}_W \mathbf{W}|} \quad \begin{array}{l} \text{The largest eigenvectors of } \mathbf{S}_W^{-1} \mathbf{S}_B \\ \text{Maximum rank of } K-1 \end{array}$$



# Separating Style and Content

- Objective: Decomposing two factors using linear methods
  - Content: which character
  - Style : which font
- “Bilinear models”
- J. Tenenbaum and W. Freeman “Separating Style and Content with Bilinear Models” Neural computation 2000

A Classification

A	B	C	D	E
A	B	C	D	E
A	B	C	D	E
A	B	C	D	E
A	B	C	D	E
B	C	A	E	D

B Extrapolation

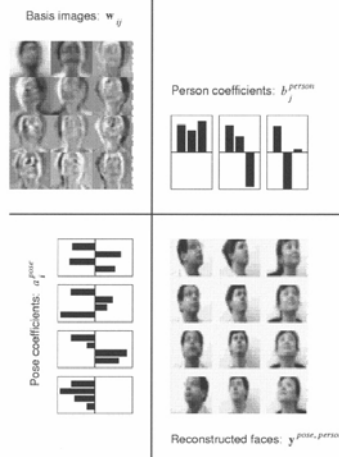
A	B	C	D	E
A	B	C	D	E
A	B	C	D	E
A	B	C	D	E
A	B	C	D	E
?	?	C	D	E

Figures from J. Tenenbaum and W. Freeman “Separating Style and Content with Bilinear Models” Neural computation 2000

# Bilinear Models

- Symmetric bilinear model

$$y^{sc} = \sum_{i,j} w_{ij} a_i^s b_j^c$$

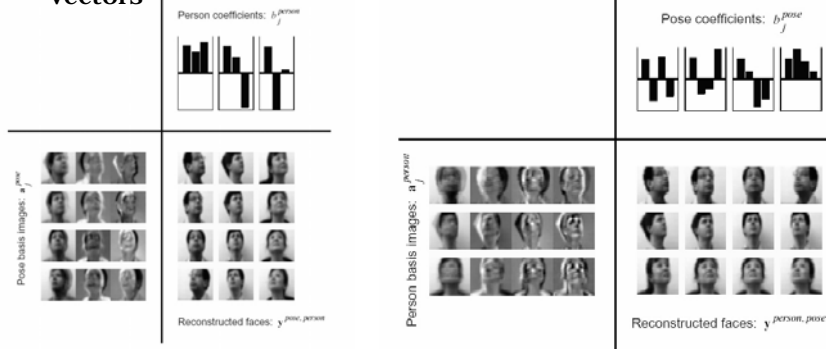


Figures from J. Tenenbaum and W. Freeman “Separating Style and Content with Bilinear Models” Neural computation 2000

# Bilinear models

$$y^{sc} = A^s b^c$$

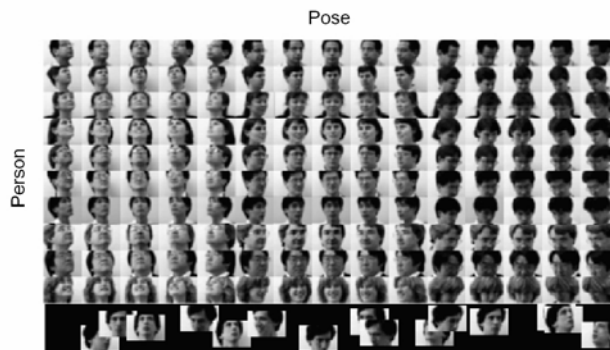
- Asymmetric bilinear model: use style dependent basis vectors



Head pose as style factor  
person as content

Person as style factor  
pose as content

Figures from J. Tenenbaum and W. Freeman "Separating Style and Content with Bilinear Models" Neural computation 2000



Figures from J. Tenenbaum and W. Freeman "Separating Style and Content with Bilinear Models" Neural computation 2000