

Data Analysis and Dimension Reduction

- PCA, LDA and LSA



Berlin Chen

Department of Computer Science & Information Engineering
National Taiwan Normal University

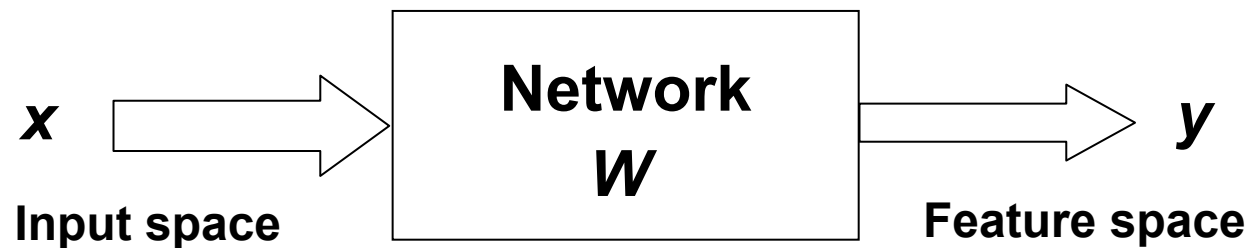


References:

1. *Introduction to Machine Learning* , Chapter 6
2. *Data Mining: Concepts, Models, Methods and Algorithms*, Chapter 3

Introduction (1/3)

- Goal: discover significant patterns or features from the input data
 - Salient feature selection or dimensionality reduction



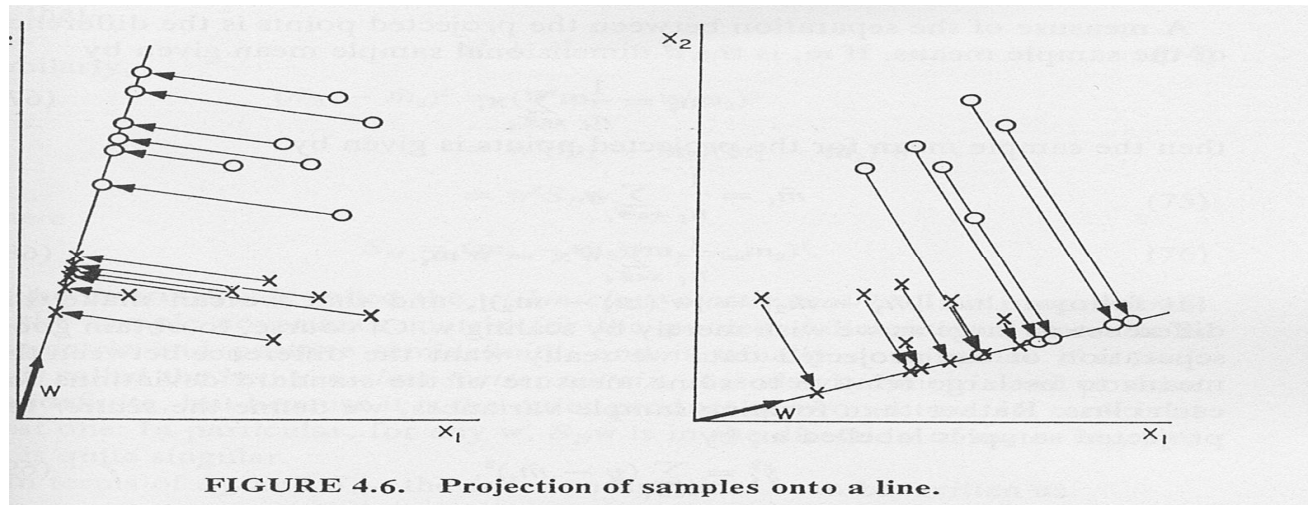
- Compute an input-output mapping based on some desirable properties

Introduction (2/3)

- Principal Component Analysis (PCA)
- Linear Discriminant Analysis (LDA)
- Latent Semantic Analysis (LSA)
- Heteroscedastic Discriminant Analysis (HDA)
- Probabilistic Latent Semantic Analysis (PLSA)

Introduction (3/3)

- Formulation for feature extraction and dimension reduction
 - Model-free (nonparametric)
 - Without prior information: e.g., PCA
 - With prior information: e.g., LDA
 - Model-dependent (parametric), e.g.,
 - HLDA with Gaussian cluster distributions
 - PLSA with multinomial latent cluster distributions

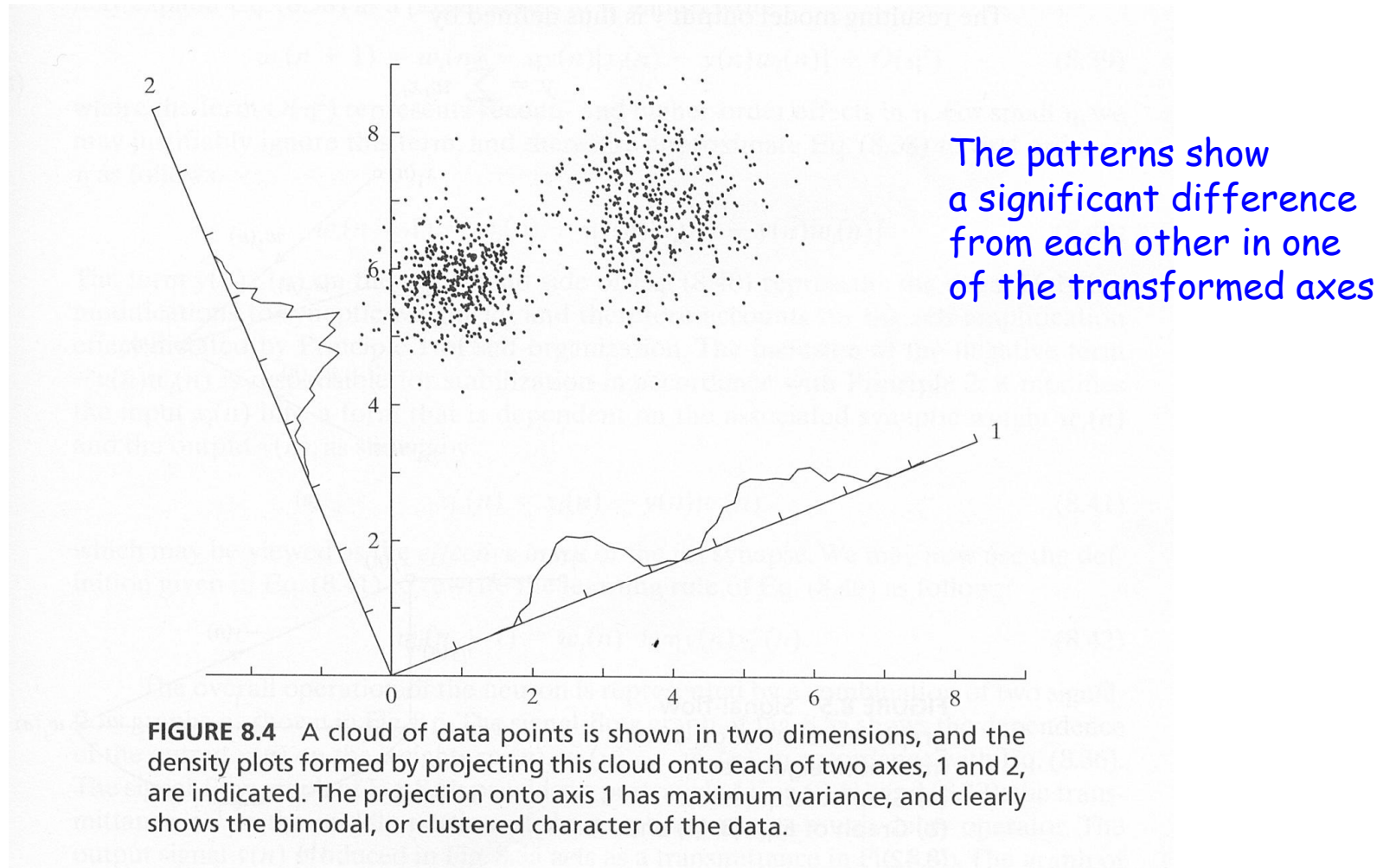


Principal Component Analysis (PCA) (1/2)

Pearson, 1901

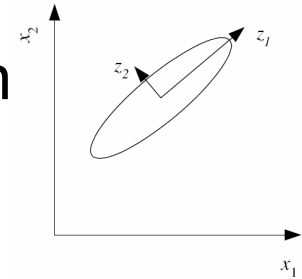
- Known as Karhunen-Loève Transform (1947, 1963)
 - Or Hotelling Transform (1933)
- A standard technique commonly used for data reduction in statistical pattern recognition and signal processing
- A transform by which the data set can be represented by **reduced number of effective features** and still **retain the most intrinsic information content**
 - A small set of features to be found to represent the data samples accurately
- Also called “Subspace Decomposition”, “Factor Analysis” ..

Principal Component Analysis (PCA) (2/2)



PCA Derivations (1/13)

- Suppose \mathbf{x} is an n -dimensional zero mean random vector, $\boldsymbol{\mu} = \mathbf{E} \{ \mathbf{x} \} = \mathbf{0}$
 - If \mathbf{x} is not zero mean, we can subtract the mean before processing the following analysis



- \mathbf{x} can be represented without error by the summation of n linearly independent vectors

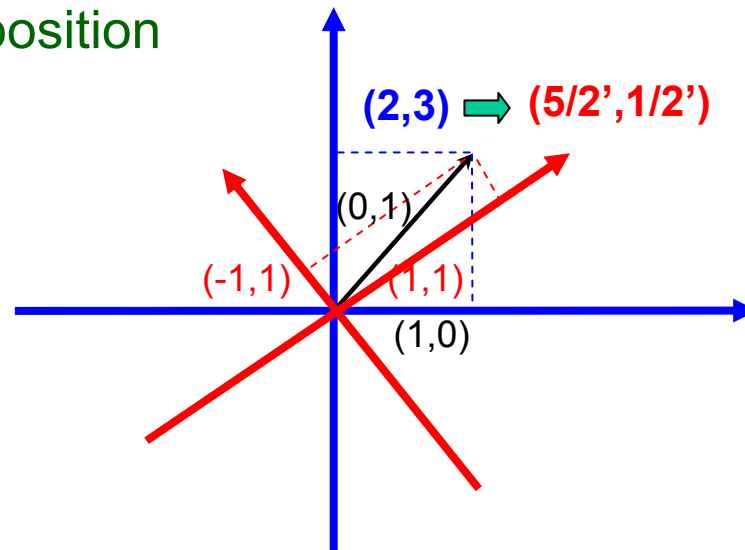
$$\mathbf{x} = \sum_{i=1}^n y_i \boldsymbol{\varphi}_i = \boldsymbol{\Phi} \mathbf{y} \quad \text{where} \quad \mathbf{y} = [y_1 \quad \cdot \quad y_i \quad \cdot \quad y_n]^T$$
$$\boldsymbol{\Phi} = [\boldsymbol{\varphi}_1 \quad \cdot \quad \boldsymbol{\varphi}_i \quad \cdot \quad \boldsymbol{\varphi}_n]$$

The i -th component
in the feature (mapped) space

The basis vectors

PCA Derivations (2/13)

Subspace Decomposition



$$\begin{bmatrix} 2 \\ 3 \end{bmatrix} = 2 \begin{bmatrix} 1 \\ 0 \end{bmatrix} + 3 \begin{bmatrix} 0 \\ 1 \end{bmatrix} = \frac{5}{2} \begin{bmatrix} 1 \\ 1 \end{bmatrix} + \frac{1}{2} \begin{bmatrix} -1 \\ 1 \end{bmatrix} = 1 \begin{bmatrix} 1 \\ 2 \end{bmatrix} + 1 \begin{bmatrix} 1 \\ 1 \end{bmatrix}$$

orthogonal basis sets

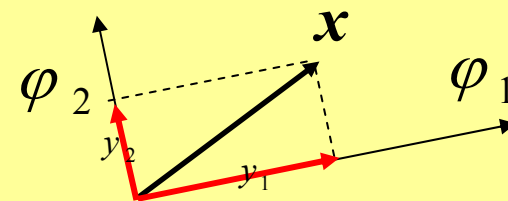
PCA Derivations (3/13)

- Further assume the column (basis) vectors of the matrix Φ form an orthonormal set

$$\varphi_i^T \varphi_j = \begin{cases} 1 & \text{if } i = j \\ 0 & \text{if } i \neq j \end{cases}$$

- Such that y_i is equal to the projection of \mathbf{x} on φ_i

$$\forall_i \quad y_i = \mathbf{x}^T \varphi_i = \varphi_i^T \mathbf{x}$$



$$y_1 = \|\mathbf{x}\| \cos \theta_1 = \|\mathbf{x}\| \frac{\varphi_1^T \mathbf{x}}{\|\mathbf{x}\| \|\varphi_1\|} = \varphi_1^T \mathbf{x}$$

, where $\|\varphi_1\| = 1$

PCA Derivations (4/13)

– Further assume the column (basis) vectors of the matrix Φ form an orthonormal set

- y_i also has the following properties

– Its mean is zero, too

$$\mathbf{E}\{y_i\} = \mathbf{E}\{\varphi_i^T \mathbf{x}\} = \varphi_i^T \mathbf{E}\{\mathbf{x}\} = \varphi_i^T \boldsymbol{\theta} = 0$$

– Its variance is

$$\begin{aligned} \sigma_i^2 &= \mathbf{E}\{y_i^2\} - [\mathbf{E}\{y_i\}]^2 = \mathbf{E}\{y_i^2\} = \mathbf{E}\{\varphi_i^T \mathbf{x} \mathbf{x}^T \varphi_i\} = \varphi_i^T \mathbf{E}\{\mathbf{x} \mathbf{x}^T\} \varphi_i \\ &= \varphi_i^T \mathbf{R} \varphi_i \quad [\mathbf{R} \text{ is the (auto-)correlation matrix of } \mathbf{x}] \end{aligned}$$

- The correlation between two projections y_i and y_j is

$$\begin{aligned} \mathbf{E}\{y_i y_j\} &= \mathbf{E}\{(\varphi_i^T \mathbf{x})(\varphi_j^T \mathbf{x})^T\} = \mathbf{E}\{\varphi_i^T \mathbf{x} \mathbf{x}^T \varphi_j\} \\ &= \varphi_i^T \mathbf{E}\{\mathbf{x} \mathbf{x}^T\} \varphi_j = \varphi_i^T \mathbf{R} \varphi_j \end{aligned}$$

$$\begin{aligned} \Sigma &= \mathbf{E}\{(\mathbf{x} - \boldsymbol{\mu})(\mathbf{x} - \boldsymbol{\mu})^T\} \\ &\approx \left(\frac{1}{N} \sum_{i=1}^N \mathbf{x}_i \mathbf{x}_i^T \right) - \boldsymbol{\mu} \boldsymbol{\mu}^T \\ \mathbf{R} &= \mathbf{E}\{\mathbf{x} \mathbf{x}^T\} \approx \frac{1}{N} \sum_i \mathbf{x}_i \mathbf{x}_i^T \end{aligned}$$

PCA Derivations (5/13)

- **Minimum Mean-Squared Error Criterion**
 - We want to choose only m of $\boldsymbol{\varphi}_i$'s that we still can approximate \mathbf{x} well in **mean-squared error criterion**

original vector $\mathbf{x} = \sum_{i=1}^n y_i \boldsymbol{\varphi}_i = \sum_{i=1}^m y_i \boldsymbol{\varphi}_i + \sum_{j=m+1}^n y_j \boldsymbol{\varphi}_j$

reconstructed vector $\hat{\mathbf{x}}(m) = \sum_{i=1}^m y_i \boldsymbol{\varphi}_i$

$$\bar{\varepsilon}(m) = \mathbf{E} \left\{ \left\| \hat{\mathbf{x}}(m) - \mathbf{x} \right\|^2 \right\} = \mathbf{E} \left\{ \left(\sum_{j=m+1}^n y_j \boldsymbol{\varphi}_j^T \right) \left(\sum_{k=m+1}^n y_k \boldsymbol{\varphi}_k \right) \right\}$$

$$= \mathbf{E} \left\{ \sum_{j=m+1}^n \sum_{k=m+1}^n y_j y_k \boldsymbol{\varphi}_j^T \boldsymbol{\varphi}_k \right\}$$

$$\begin{aligned} E\{y_j\} &= 0 \\ \sigma_j^2 &= E\{y_j^2\} - [E\{y_j\}]^2 \\ &= E\{y_j^2\} \end{aligned}$$

$$= \sum_{j=m+1}^n \mathbf{E} \{ y_j^2 \}$$

$$\because \boldsymbol{\varphi}_j^T \boldsymbol{\varphi}_k = \begin{cases} 1 & \text{if } j = k \\ 0 & \text{if } j \neq k \end{cases}$$

$$= \sum_{j=m+1}^n \sigma_j^2 = \sum_{j=m+1}^n \boldsymbol{\varphi}_j^T \mathbf{R} \boldsymbol{\varphi}_j$$

We should discard the bases where the projections have lower variances

PCA Derivations (6/13)

- Minimum Mean-Squared Error Criterion

- If the orthonormal (basis) set φ_i 's is selected to be the eigenvectors of the correlation matrix \mathbf{R} , associated with eigenvalues λ_i 's

- They will have the property that:

is real and symmetric, therefore its eigenvectors \mathbf{R} form a orthonormal set

$$\mathbf{R} \varphi_j = \lambda_j \varphi_j$$

\mathbf{R} is positive definite ($\mathbf{x}^T \mathbf{R} \mathbf{x} > 0$)
=> all eigenvalues are positive

- Such that the mean-squared error mentioned above will be

$$\begin{aligned} \bar{\varepsilon}(m) &= \sum_{j=m+1}^n \sigma_j^2 \\ &= \sum_{j=m+1}^n \varphi_j^T \mathbf{R} \varphi_j = \sum_{j=m+1}^n \varphi_j^T \lambda_j \varphi_j = \sum_{j=m+1}^n \lambda_j \end{aligned}$$

PCA Derivations (7/13)

- Minimum Mean-Squared Error Criterion

- If the eigenvectors are retained associated with the m largest eigenvalues, the mean-squared error will be

$$\bar{\mathcal{E}}_{eigen}(m) = \sum_{j=m+1}^n \lambda_j \quad (\text{where } \lambda_1 \geq \dots \geq \lambda_m \geq \dots \geq \lambda_n \geq 0)$$

- Any two projections y_i and y_j will be mutually uncorrelated

$$\begin{aligned} E\{y_i y_j\} &= E\left\{(\boldsymbol{\varphi}_i^T \mathbf{x})(\boldsymbol{\varphi}_j^T \mathbf{x})^T\right\} = E\left\{\boldsymbol{\varphi}_i^T \mathbf{x} \mathbf{x}^T \boldsymbol{\varphi}_j\right\} \\ &= \boldsymbol{\varphi}_i^T E\left\{\mathbf{x} \mathbf{x}^T\right\} \boldsymbol{\varphi}_j = \boldsymbol{\varphi}_i^T \mathbf{R} \boldsymbol{\varphi}_j = \lambda_j \boldsymbol{\varphi}_i^T \boldsymbol{\varphi}_j = 0 \end{aligned}$$

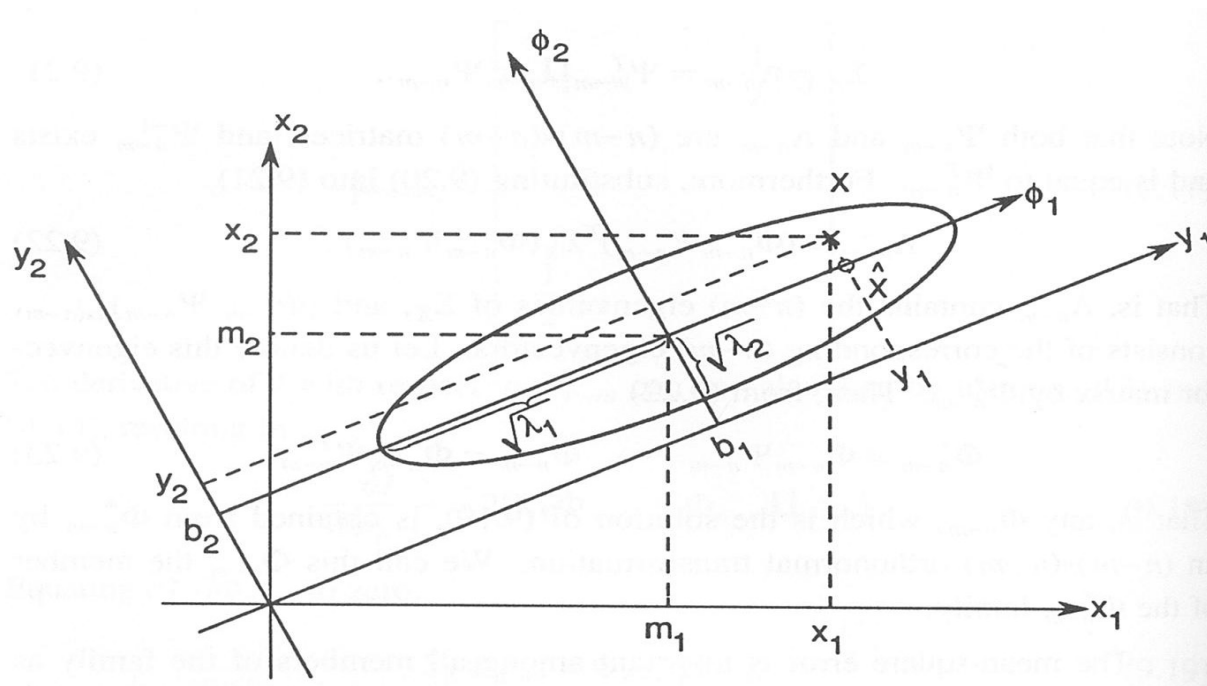
- Good news for most statistical modeling approaches
 - Gaussians and diagonal matrices

$$\begin{aligned} \boldsymbol{\Sigma} &= E[(\mathbf{x} - \boldsymbol{\mu})(\mathbf{x} - \boldsymbol{\mu})^T] \\ &\approx \left(\frac{1}{N} \sum_{i=1}^N \mathbf{x}_i \mathbf{x}_i^T\right) - \boldsymbol{\mu} \boldsymbol{\mu}^T \\ \mathbf{R} &= E\{\mathbf{x} \mathbf{x}^T\} = \frac{1}{N} \sum_i \mathbf{x}_i \mathbf{x}_i^T \end{aligned}$$

$$\begin{bmatrix} \sigma_{11} & \sigma_{12} & \cdot & \cdot & \sigma_{1n} \\ \sigma_{22} & \sigma_{22} & \cdot & \cdot & \cdot \\ \cdot & \cdot & \cdot & \cdot & \cdot \\ \cdot & \cdot & \cdot & \cdot & \cdot \\ \sigma_{n1} & \sigma_{n2} & \cdot & \cdot & \sigma_{nn} \end{bmatrix}$$

PCA Derivations (8/13)

- A Two-dimensional Example of Principle Component Analysis



PCA Derivations (9/13)

- Minimum Mean-Squared Error Criterion

- It can be proved that $\bar{\varepsilon}_{eigen}(m)$ is the optimal solution under the mean-squared error criterion

$$\delta_{jk} = \begin{cases} 1 & \text{if } j = k \\ 0 & \text{if } j \neq k \end{cases}$$

Objective function

To be minimized

Constraints

Define: $J = \sum_{j=m+1}^n \boldsymbol{\varphi}_j^T \mathbf{R} \boldsymbol{\varphi}_j - \sum_{j=m+1}^n \sum_{k=m+1}^n u_{jk} (\boldsymbol{\varphi}_j^T \boldsymbol{\varphi}_k - \delta_{jk})$

$$\frac{\partial \boldsymbol{\varphi}^T \mathbf{R} \boldsymbol{\varphi}}{\partial \boldsymbol{\varphi}} = 2\mathbf{R}\boldsymbol{\varphi}$$

Partial Differentiation

$$\Rightarrow \nabla_{m+1 \leq j \leq n} \frac{\partial J}{\partial \boldsymbol{\varphi}_j} = 2\mathbf{R}\boldsymbol{\varphi}_j - 2 \sum_{k=m+1}^n u_{jk} \boldsymbol{\varphi}_k = \mathbf{0} \quad \left(\text{where } \mathbf{u}_j^T = [u_{j, m+1} \dots u_{j, n}] \right)$$

$$\Rightarrow \nabla_{m+1 \leq j \leq n} \mathbf{R}\boldsymbol{\varphi}_j = \boldsymbol{\Phi}_{n-m} \mathbf{u}_j \quad \left(\text{where } \boldsymbol{\Phi}_{n-m} = [\boldsymbol{\varphi}_{m+1} \dots \boldsymbol{\varphi}_n] \right)$$

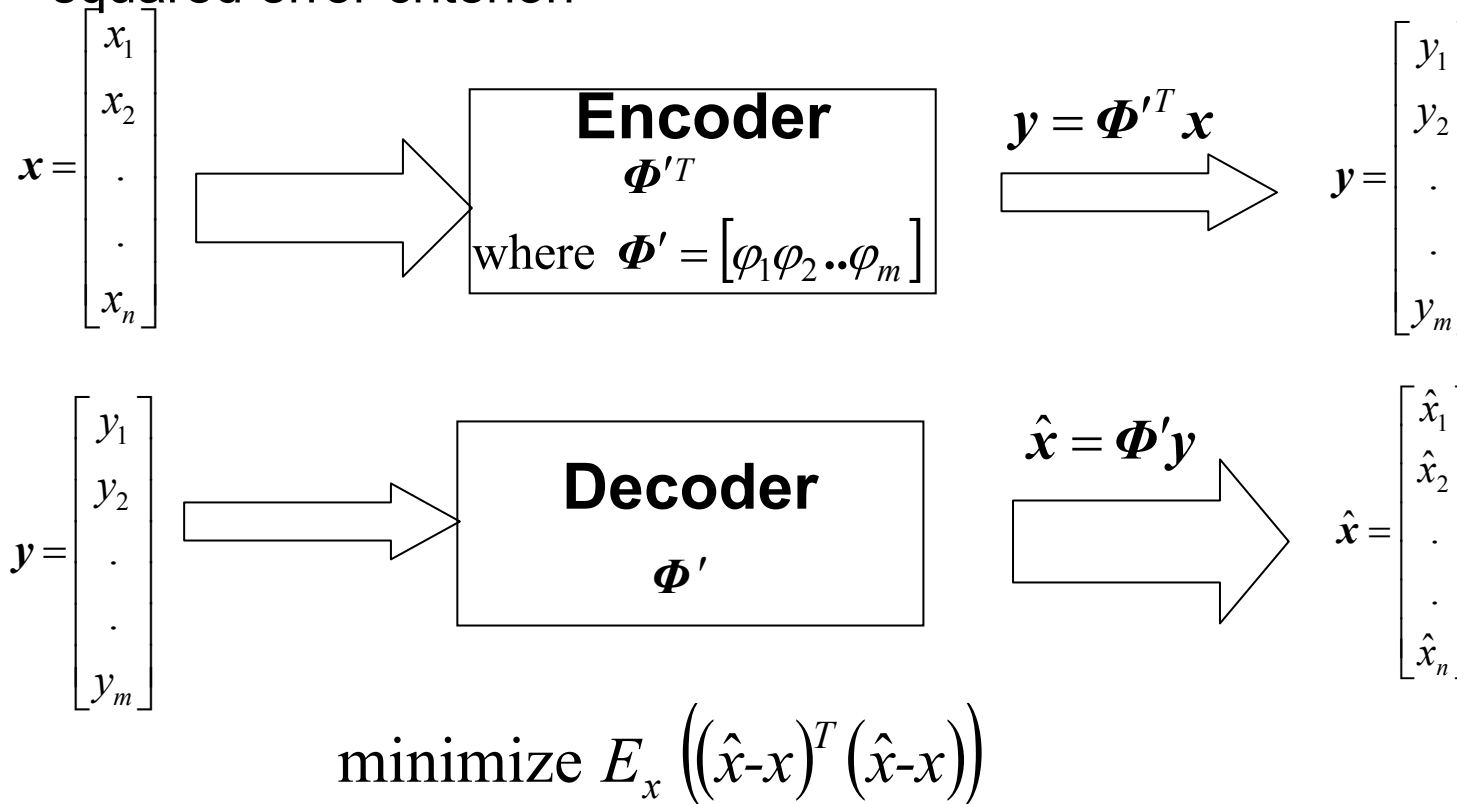
$$\Rightarrow \mathbf{R}[\boldsymbol{\varphi}_{m+1} \dots \boldsymbol{\varphi}_n] = \boldsymbol{\Phi}_{n-m} [\mathbf{u}_{m+1} \dots \mathbf{u}_n]$$

$$\Rightarrow \mathbf{R}\boldsymbol{\Phi}_{n-m} = \boldsymbol{\Phi}_{n-m} \mathbf{U}_{n-m} \quad \left(\text{where } \mathbf{U}_{n-m} = [\mathbf{u}_{m+1} \dots \mathbf{u}_n] \right)$$

Have a particular solution if \mathbf{U}_{n-m} is a diagonal matrix and its diagonal elements is the eigenvalues $\lambda_{m+1} \dots \lambda_n$ of \mathbf{R} and $\boldsymbol{\varphi}_{m+1} \dots \boldsymbol{\varphi}_n$ is their corresponding eigenvectors

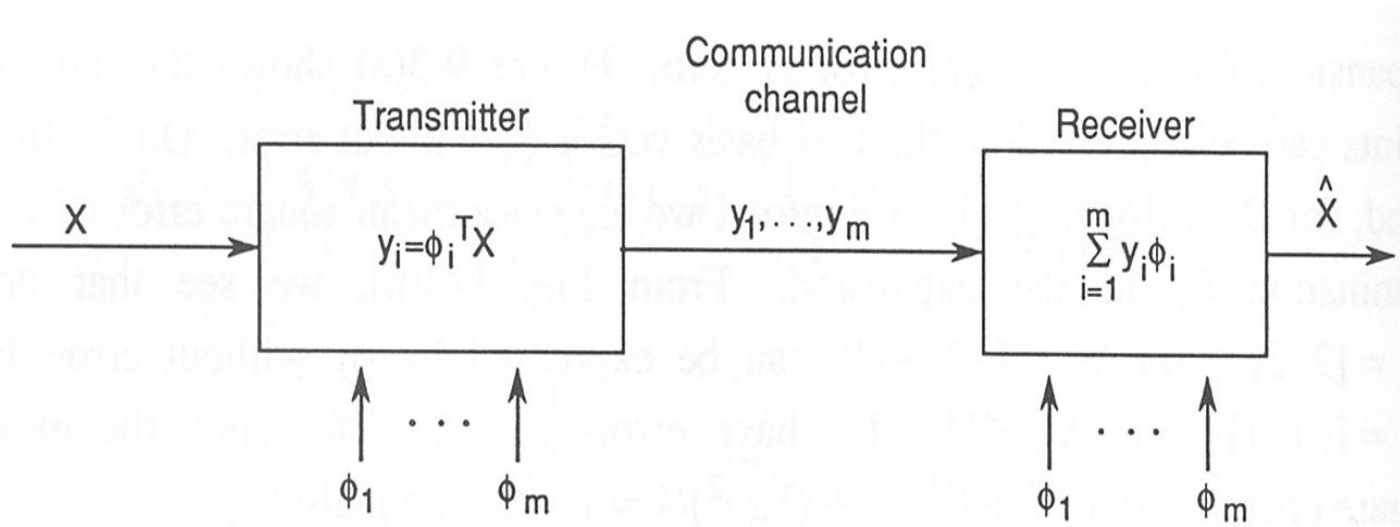
PCA Derivations (10/13)

- Given an input vector \mathbf{x} with dimension m
 - Try to construct a linear transform Φ' (Φ' is an $n \times m$ matrix $m < n$) such that the truncation result, $\Phi'^T \mathbf{x}$, is optimal in mean-squared error criterion



PCA Derivations (11/13)

- Data compression in communication



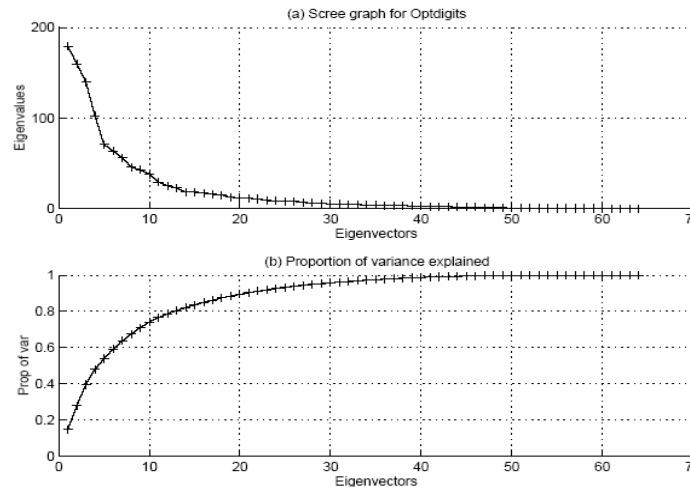
- PCA is an optimal transform for signal representation and dimensional reduction, but not necessary for classification tasks, such as speech recognition **? (To be discussed later on)**
- PCA needs no prior information (e.g. class distributions of output information) of the sample patterns

PCA Derivations (12/13)

- Scree Graph

- The plot of variance as a function of the number of eigenvectors kept

- Select m such that $\frac{\lambda_1 + \lambda_2 + \dots + \lambda_m}{\lambda_1 + \lambda_2 + \dots + \lambda_m + \dots + \lambda_n} \geq \text{Threshold}$



- Or select those eigenvectors with eigenvalues larger than the average input variance (average eigvalue)

$$\lambda_m \geq \frac{1}{n} \sum_{i=1}^n \lambda_i$$

PCA Derivations (13/13)

- PCA finds a linear transform \mathbf{W} such that the **sum of average between-class variation and average within-class variation** is maximal

$$J(\mathbf{W}) = |\tilde{\mathbf{S}}| \stackrel{?}{=} |\tilde{\mathbf{S}}_w + \tilde{\mathbf{S}}_b| = |\mathbf{W}^T \mathbf{S}_w \mathbf{W} + \mathbf{W}^T \mathbf{S}_b \mathbf{W}|$$

$$\mathbf{S} = \frac{1}{N} \sum_i (\mathbf{x}_i - \bar{\mathbf{x}})(\mathbf{x}_i - \bar{\mathbf{x}})^T$$

i ← sample index

$$\mathbf{S}_w = \frac{1}{N} \sum_j N_j \boldsymbol{\Sigma}_j$$

j ← class index

$$\mathbf{S}_b = \frac{1}{N} \sum_j N_j (\bar{\mathbf{x}}_j - \bar{\mathbf{x}})(\bar{\mathbf{x}}_j - \bar{\mathbf{x}})^T$$

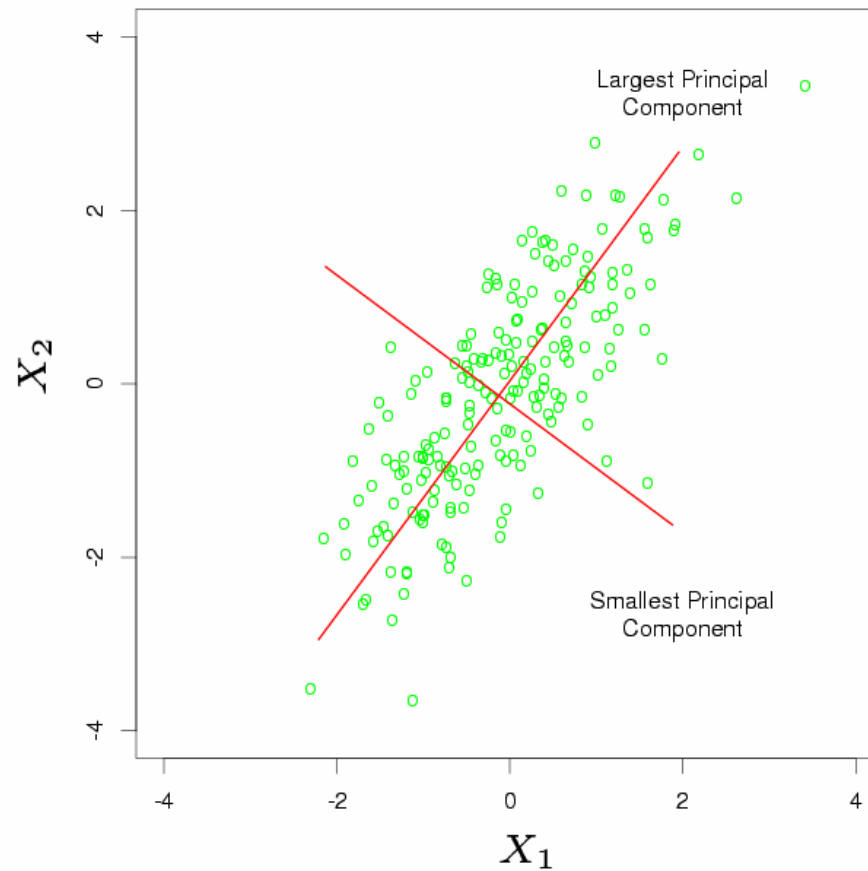
$$\tilde{\mathbf{S}}_w = \mathbf{W}^T \mathbf{S}_w \mathbf{W}$$

$$\tilde{\mathbf{S}}_b = \mathbf{W}^T \mathbf{S}_b \mathbf{W}$$

Try to show that:
 $\mathbf{S} = \mathbf{S}_w + \mathbf{S}_b$

PCA Examples: Data Analysis

- Example 1: principal components of some data points



PCA Examples: Feature Transformation

- Example 2: feature transformation and selection

**Correlation matrix
for old feature
dimensions**

TABLE 3.2 The correlation matrix for Iris data

	Feature 1	Feature 2	Feature 3	Feature 4
Feature 1	1.0000	-0.1094	0.8718	0.8180
Feature 2	-0.1094	1.0000	-0.4205	-0.3565
Feature 3	0.8718	-0.4205	1.0000	0.9628
Feature 4	0.8180	-0.3565	0.9628	1.0000

New feature dimensions

TABLE 3.3 The eigenvalues for Iris data

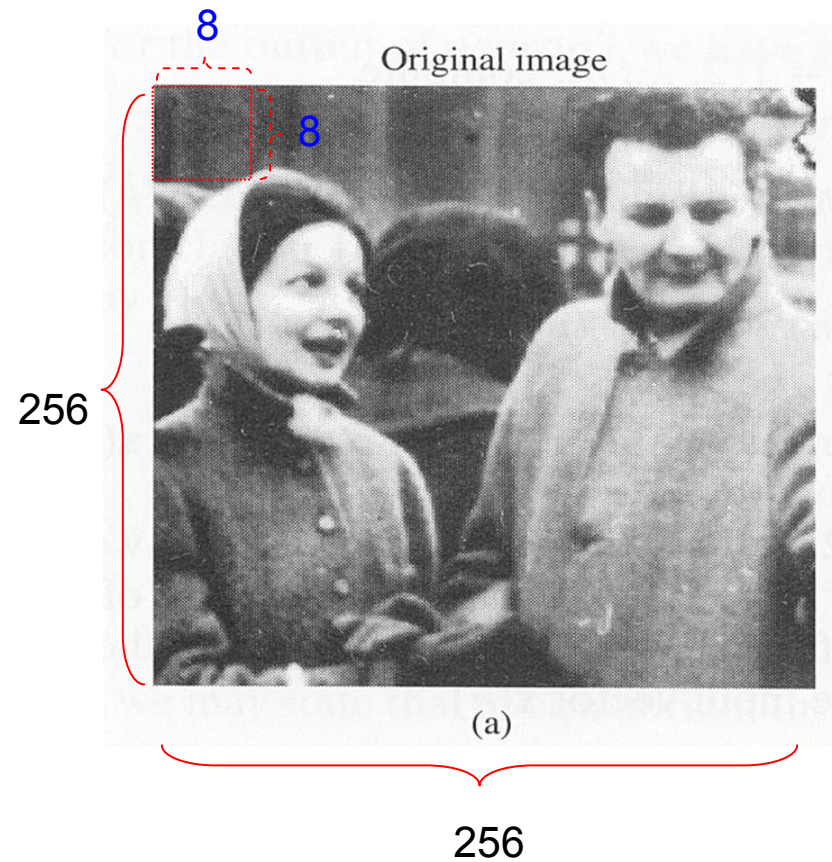
Feature	Eigenvalue
Feature 1	2.91082
Feature 2	0.92122
Feature 3	0.14735
Feature 4	0.02061

$$R = (2.91082 + 0.92122) / (2.91082 + 0.92122 + 0.14735 + 0.02061) \\ = 0.958 > 0.95$$

threshold for information content reserved

PCA Examples: Image Coding (1/2)

- Example 3: Image Coding

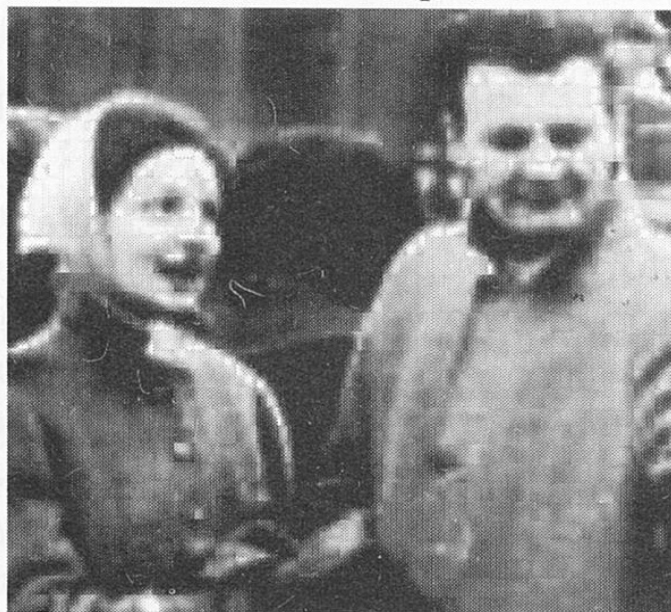


PCA Examples: Image Coding (2/2)

- Example 3: Image Coding (cont.)

Using first 8 components (feature reduction)

15 to 1 compression (value reduction)



(c)



(d)

FIGURE 8.9 (a) An image of parents used in the image coding experiment. (b) 8×8 masks representing the synaptic weights learned by the GHA. (c) Reconstructed image of parents obtained using the dominant 8 principal components without quantization. (d) Reconstructed image of parents with 15 to 1 compression ratio using quantization.

PCA Examples: Eigenface (1/4)

- Example 4: Eigenface in face recognition (Turk and Pentland, 1991)
 - Consider an individual image to be a linear combination of a small number of face components or “eigenfaces” derived from a set of reference images

$$\mathbf{x}_1 = \begin{bmatrix} x_{1,1} \\ x_{1,2} \\ \cdot \\ \cdot \\ x_{1,n} \end{bmatrix}, \mathbf{x}_2 = \begin{bmatrix} x_{2,1} \\ x_{2,2} \\ \cdot \\ \cdot \\ x_{2,n} \end{bmatrix}, \dots, \mathbf{x}_L = \begin{bmatrix} x_{L,1} \\ x_{L,2} \\ \cdot \\ \cdot \\ x_{L,n} \end{bmatrix}$$

- Steps
 - Convert each of the L reference images into a vector of floating point numbers representing light intensity in each pixel
 - Calculate the covariance/correlation matrix between these reference vectors
 - Apply Principal Component Analysis (PCA) find the eigenvectors of the matrix: the eigenfaces
 - Besides, the vector obtained by averaging all images are called “eigenface 0”. The other eigenfaces from “eigenface 1” onwards model the variations from this average face

PCA Examples: Eigenface (2/4)

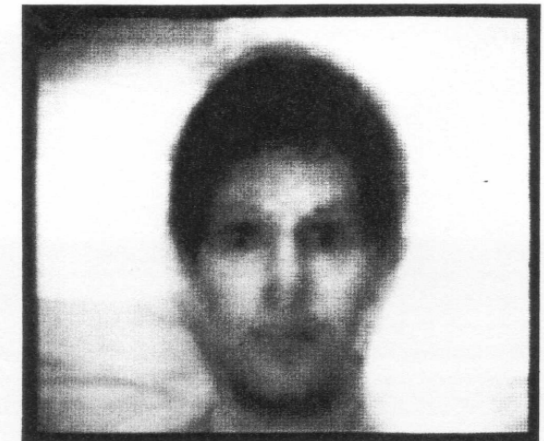
- Example 4: Eigenface in face recognition (cont.)
 - Steps
 - Then the faces are then represented as eigenvoice 0 plus a linear combination of the remain K ($K \leq L$) eigenfaces
 - The Eigenface approach persists the minimum mean-squared error criterion
 - Incidentally, the eigenfaces are not only themselves usually plausible faces, but also directions of variations between faces

$$\hat{\mathbf{x}}_i = \bar{\mathbf{x}} + w_{i,1} \mathbf{e}(1) + w_{i,2} \mathbf{e}(2) + \dots + w_{i,K} \mathbf{e}(K)$$
$$\Rightarrow \mathbf{y}_i = [1, w_{i,1}, w_{i,2}, \dots, w_{i,K}]$$

Feature vector of a person i

PCA Examples: Eigenface (3/4)

Face images as the training set



The averaged face

PCA Examples: Eigenface (4/4)

Seven eigenfaces derived from the training set



(Indicate directions of variations between faces)

A projected face image

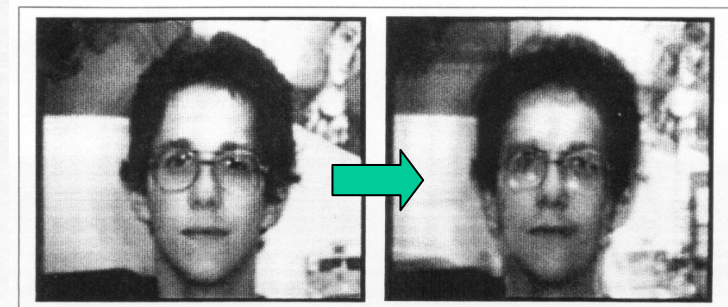
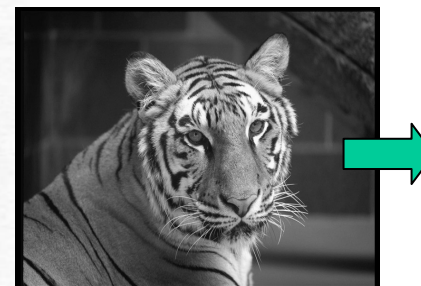
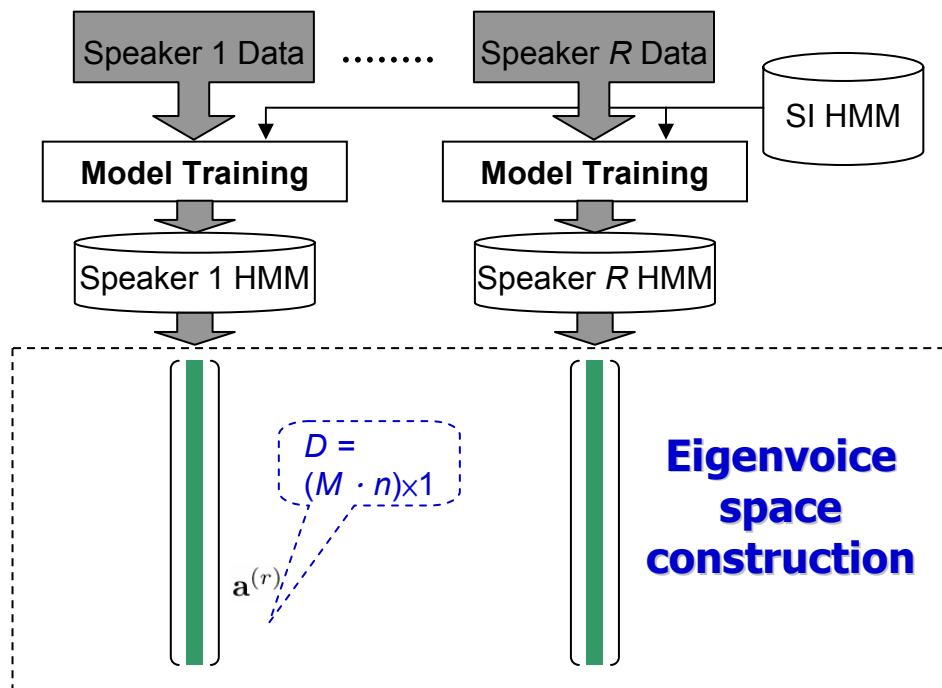


Figure 3. An original face image and its projection onto the face space defined by the eigenfaces of Figure 2.



PCA Examples: Eigenvoice (1/3)

- Example 5: Eigenvoice in speaker adaptation (PSTL, 2000)
 - Steps
 - Concatenating the regarded parameters for each speaker r to form a huge vector $\mathbf{a}^{(r)}$ (a supervectors)
 - SD HMM model mean parameters (μ)



Each new speaker S is represented by a point P in K -space

$$\mathbf{P}_i = \mathbf{e}(0) + w_{i,1} \mathbf{e}(1) + w_{i,2} \mathbf{e}(2) + \dots + w_{i,K} \mathbf{e}(K)$$

SI HMM model

Principal Component Analysis

PCA Examples: Eigenvoice (2/3)

- Example 4: Eigenvoice in speaker adaptation (cont.)

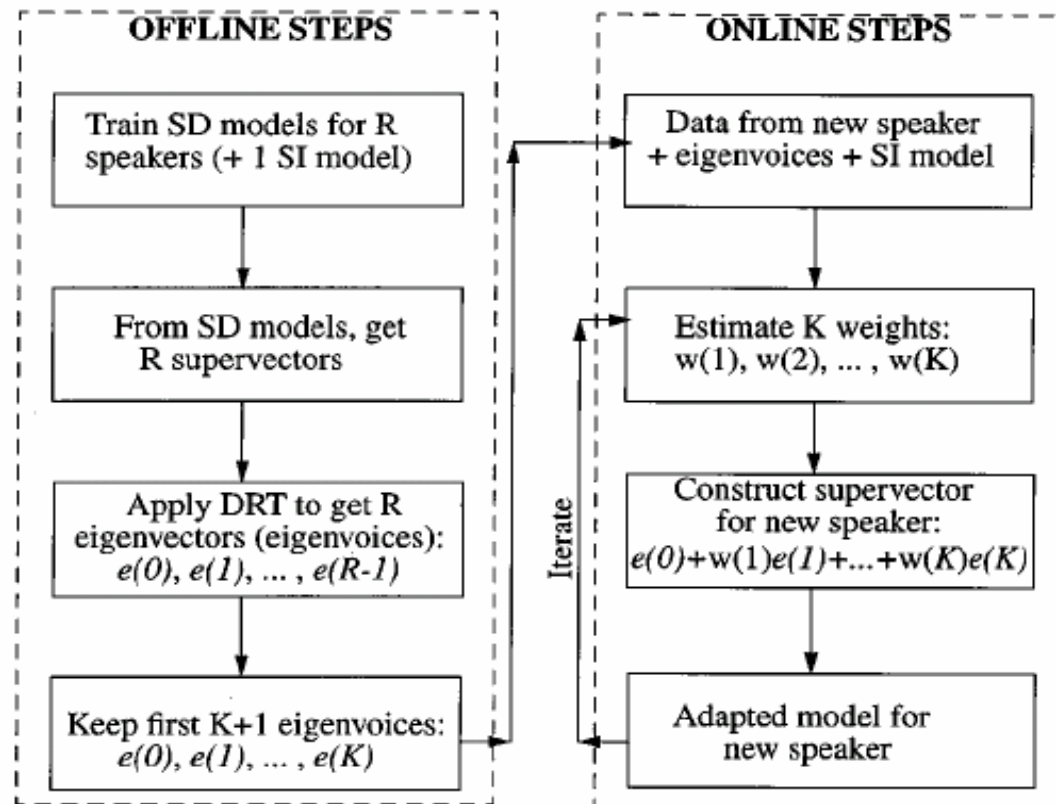


Fig. 1. Block diagram for eigenvoice speaker adaptation

PCA Examples: Eigenvoice (3/3)

- Example 5: Eigenvoice in speaker adaptation (cont.)
 - Dimension 1 (eigenvoice 1):
 - Correlate with pitch or sex
 - Dimension 2 (eigenvoice 2):
 - Correlate with amplitude
 - Dimension 3 (eigenvoice 3):
 - Correlate with second-formant movement

**Note that:
Eigenface performs on feature space
while eigenvoice performs
on model space**

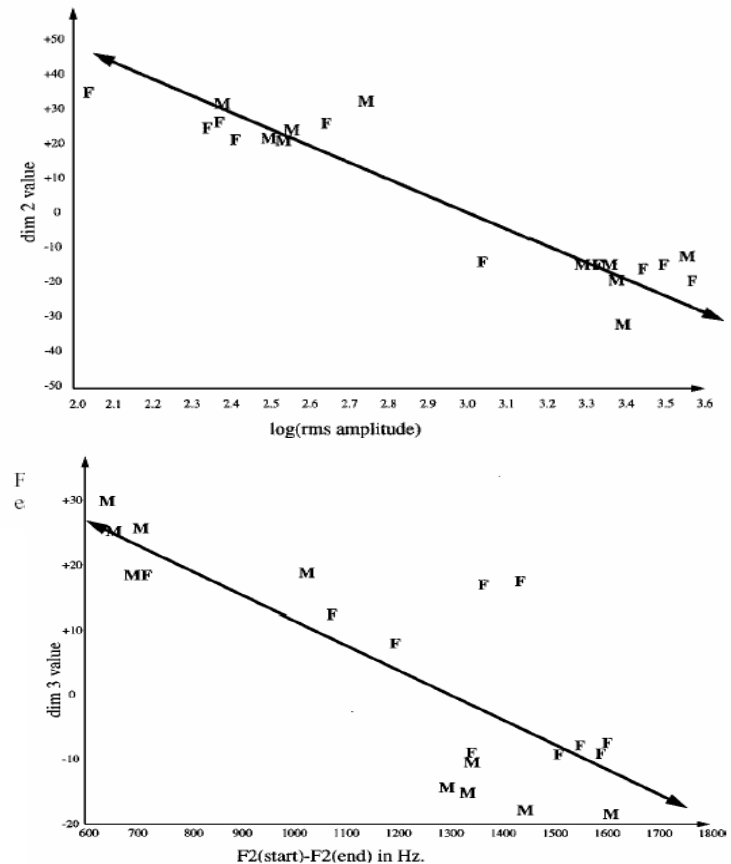


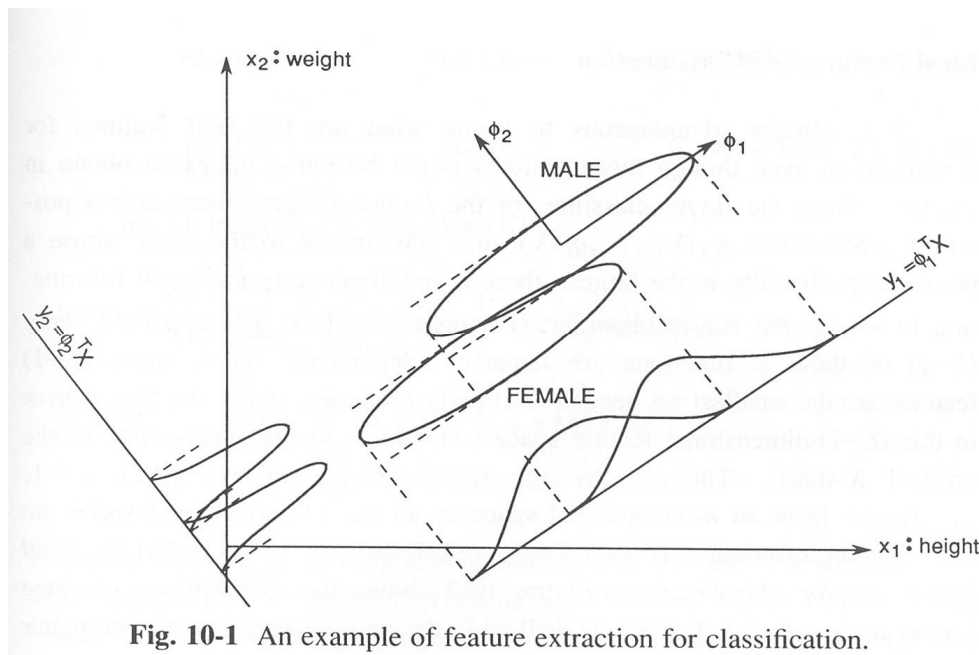
Fig. 4. Dimension 3 versus F2(start)-F2(end) for "U," extreme *M* and *F* in each speaker set

Linear Discriminant Analysis (LDA) (1/2)

- Also called
 - Fisher's Linear Discriminant Analysis, Fisher-Rao Linear Discriminant Analysis
 - Fisher (1936): introduced it for two-class classification
 - Rao (1965): extended it to handle multiple-class classification

Linear Discriminant Analysis (LDA) (2/2)

- Given a set of sample vectors with labeled (class) information, try to find a linear transform \mathbf{W} such that the ratio of **average between-class variation** over **average within-class variation** is maximal



Within-class distributions are assumed here to be Gaussians
With equal variance in the two-dimensional sample space

LDA Derivations (1/4)

- Suppose there are N sample vectors \mathbf{x}_i with dimensionality n , each of them belongs to one of the J classes $g(\mathbf{x}_i) = j$, $j \in \{1, 2, \dots, J\}$, $g(\cdot)$ is class index

- The sample mean is: $\bar{\mathbf{x}} = \frac{1}{N} \sum_{i=1}^N \mathbf{x}_i$

- The class sample means are: $\bar{\mathbf{x}}_j = \frac{1}{N_j} \sum_{g(\mathbf{x}_i)=j} \mathbf{x}_i$

- The class sample covariances are: $\Sigma_j = \frac{1}{N_j} \sum_{g(\mathbf{x}_i)=j} (\mathbf{x}_i - \bar{\mathbf{x}}_j)(\mathbf{x}_i - \bar{\mathbf{x}}_j)^T$

- The **average within-class variation** before transform

$$\mathbf{S}_w = \frac{1}{N} \sum_j N_j \Sigma_j$$

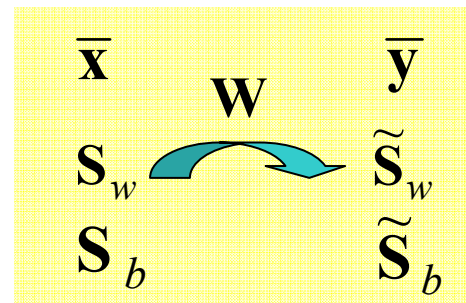
- The **average between-class variation** before transform

$$\mathbf{S}_b = \frac{1}{N} \sum_j N_j (\bar{\mathbf{x}}_j - \bar{\mathbf{x}})(\bar{\mathbf{x}}_j - \bar{\mathbf{x}})^T$$

LDA Derivations (2/4)

- If the transform $\mathbf{W} = [\mathbf{w}_1 \mathbf{w}_2 \dots \mathbf{w}_m]$ is applied
 - The sample vectors will be $\mathbf{y}_i = \mathbf{W}^T \mathbf{x}_i$
 - The sample mean will be $\bar{\mathbf{y}} = \frac{1}{N} \sum_{i=1}^N \mathbf{W}^T \mathbf{x}_i = \mathbf{W}^T \left(\frac{1}{N} \sum_{i=1}^N \mathbf{x}_i \right) = \mathbf{W}^T \bar{\mathbf{x}}$
 - The class sample means will be $\bar{\mathbf{y}}_j = \frac{1}{N_j} \sum_{g(\mathbf{x}_i)=j} \mathbf{W}^T \mathbf{x}_i = \mathbf{W}^T \bar{\mathbf{x}}_j$
 - The **average within-class variation** will be

$$\begin{aligned} \tilde{\mathbf{S}}_w &= \frac{1}{N} \sum_j N_j \left\{ \frac{1}{N_j} \cdot \sum_{g(\mathbf{x}_i)=j} \left(\mathbf{W}^T \mathbf{x}_i - \frac{1}{N_j} \sum_{g(\mathbf{x}_i)=j} (\mathbf{W}^T \mathbf{x}_i) \right) \left(\mathbf{W}^T \mathbf{x}_i - \frac{1}{N_j} \sum_{g(\mathbf{x}_i)=j} (\mathbf{W}^T \mathbf{x}_i) \right)^T \right\} \\ &= \mathbf{W}^T \left\{ \frac{1}{N} \sum_j N_j \mathbf{\Sigma}_j \right\} \mathbf{W} \\ &= \mathbf{W}^T \mathbf{S}_w \mathbf{W} \end{aligned}$$



LDA Derivations (3/4)

- If the transform $\mathbf{W} = [\mathbf{w}_1 \mathbf{w}_2 \dots \mathbf{w}_m]$ is applied
 - Similarly, the **average between-class variation** will be

$$\tilde{\mathbf{S}}_b = \mathbf{W}^T \mathbf{S}_b \mathbf{W}$$

- Try to find optimal \mathbf{W} such that the following objective function is maximized

$$J(\mathbf{W}) = \frac{|\tilde{\mathbf{S}}_b|}{|\tilde{\mathbf{S}}_w|} = \frac{|\mathbf{W}^T \mathbf{S}_b \mathbf{W}|}{|\mathbf{W}^T \mathbf{S}_w \mathbf{W}|}$$

- A closed-form solution: the column vectors of an optimal matrix \mathbf{W} are the generalized eigenvectors corresponding to the largest eigenvalues in

$$\mathbf{S}_b \mathbf{w}_i = \lambda_i \mathbf{S}_w \mathbf{w}_i$$

- That is, \mathbf{w}_i 's are the eigenvectors corresponding to the largest eigenvalues of $\mathbf{S}_w^{-1} \mathbf{S}_b$

$$\mathbf{S}_w^{-1} \mathbf{S}_b \mathbf{w}_i = \lambda_i \mathbf{w}_i$$

LDA Derivations (4/4)

- Proof:

determinant

$$\because \hat{W} = \arg \max_{\hat{W}} J(W) = \arg \max_{\hat{W}} \frac{|\tilde{S}_b|}{|\tilde{S}_w|} = \arg \max_{\hat{W}} \frac{|W^T S_b W|}{|W^T S_w W|}$$

Or equivalently, for each column vector w_i of W , we want to find that :

$$\text{The quadratic form has optimal solution : } \lambda_i = \frac{w_i^T S_b w_i}{w_i^T S_w w_i}$$

$$\left(\frac{F}{G} \right)' = \frac{F'G - GF'}{G^2}$$

$$\Rightarrow \frac{\partial \lambda_i}{\partial w_i} = \frac{2S_b w_i (w_i^T S_w w_i) - 2S_w w_i (w_i^T S_b w_i)}{(w_i^T S_w w_i)^2} = 0$$

$$\frac{d(x^T C x)}{dx} = (C + C^T)x$$

$$\Rightarrow \frac{S_b w_i (w_i^T S_w w_i)}{(w_i^T S_w w_i)^2} - \frac{S_w w_i (w_i^T S_b w_i)}{(w_i^T S_w w_i)^2} = 0$$

$$\frac{S_b w_i}{w_i^T S_w w_i} - \frac{S_w w_i}{w_i^T S_w w_i} \lambda_i = 0 \quad \left(\because \lambda_i = \frac{w_i^T S_b w_i}{w_i^T S_w w_i} \right)$$

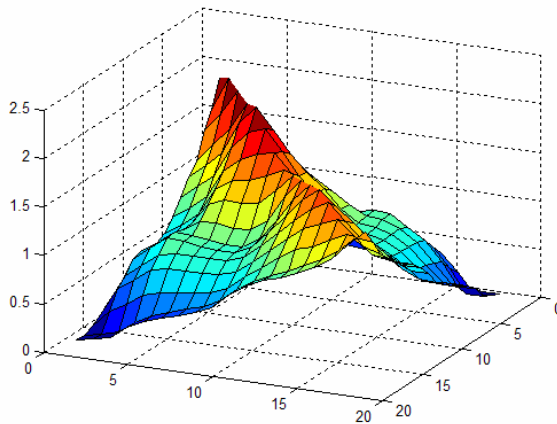
$$\Rightarrow S_b w_i - \lambda_i S_w w_i = 0 \Rightarrow S_b w_i = \lambda_i S_w w_i$$

$$\Rightarrow S_w^{-1} S_b w_i = \lambda_i w_i$$

LDA Examples: Feature Transformation (1/2)

- Example 1: Experiments on Speech Signal Processing

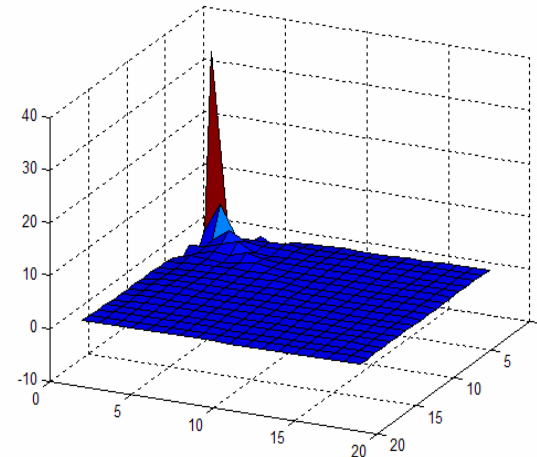
Covariance Matrix of the 18-Mel-filter-bank vectors



Calculated using Year-99's 5471 files

$$\Sigma = \frac{1}{N} \sum_{\mathbf{x}_i} (\mathbf{x}_i - \bar{\mathbf{x}})(\mathbf{x}_i - \bar{\mathbf{x}})^T$$

Covariance Matrix of the 18-cepstral vectors



Calculated using Year-99's 5471 files

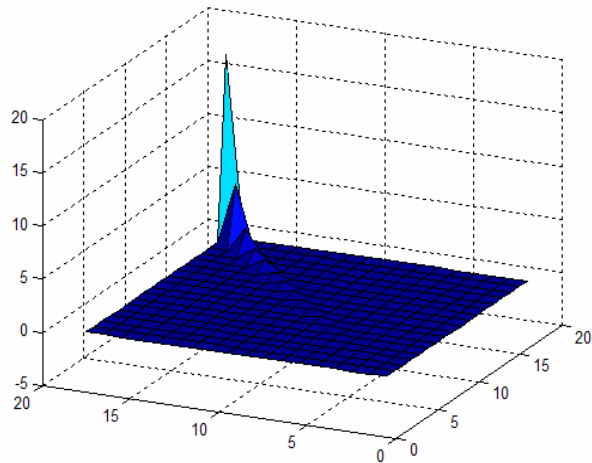
$$\Sigma' = \frac{1}{N} \sum_{\mathbf{y}_i} (\mathbf{y}_i - \bar{\mathbf{y}})(\mathbf{y}_i - \bar{\mathbf{y}})^T$$

After Cosine Transform

LDA Examples: Feature Transformation (2/2)

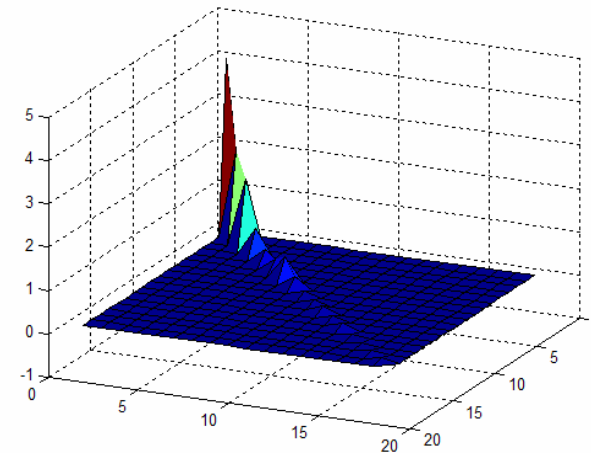
- Example1: Experiments on Speech Signal Processing (cont.)

Covariance Matrix of the 18-PCA-cepstral vectors Covariance Matrix of the 18-LDA-cepstral vectors



Calculated using Year-99's 5471 files

After PCA Transform



Calculated using Year-99's 5471 files

After LDA Transform

	Character Error Rate	
	TC	WG
MFCC	26.32	22.71
LDA-1	23.12	20.17
LDA-2	23.11	20.11

PCA vs. LDA (1/2)

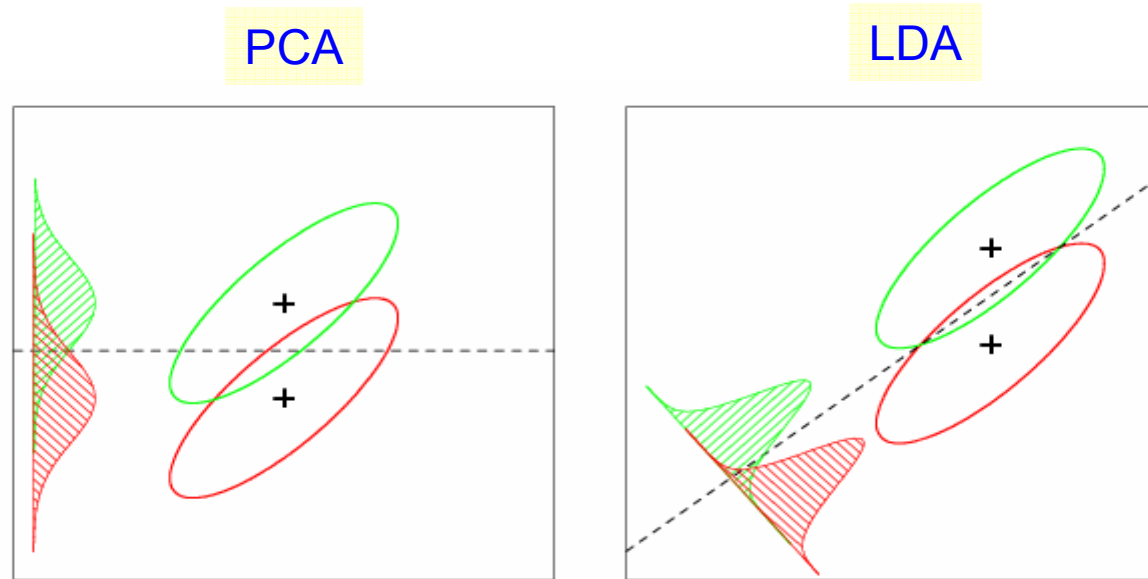


Figure 4.9: *Although the line joining the centroids defines the direction of greatest centroid spread, the projected data overlap because of the covariance (left panel). The discriminant direction minimizes this overlap for Gaussian data (right panel).*

Heteroscedastic Discriminant Analysis (HDA)

- HDA: Heteroscedastic Discriminant Analysis
- The difference in the projections obtained from LDA and HDA for 2-class case

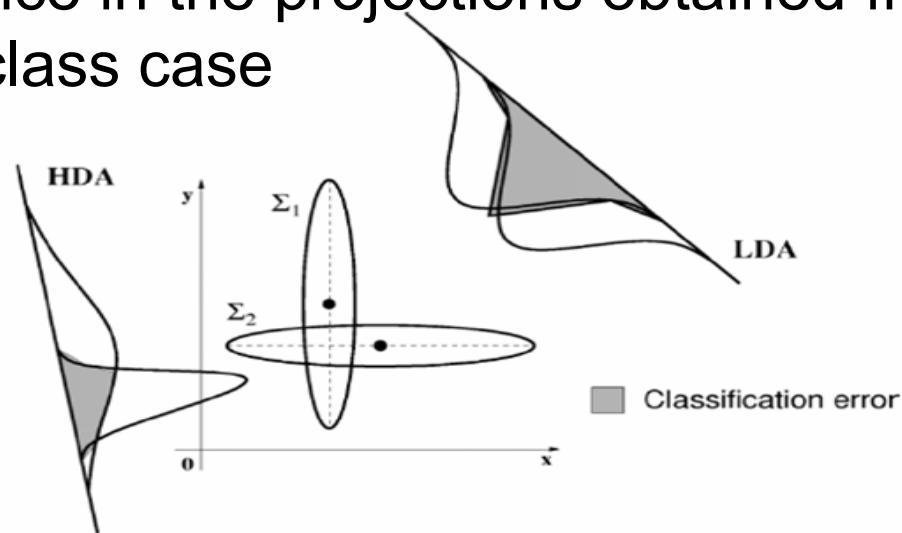


Fig. 1. Difference between LDA and HDA.

- Clearly, the HDA provides a much lower classification error than LDA theoretically
 - However, most statistical modeling assume data samples are Gaussian and have **diagonal** covariance matrices

HW: Feature Transformation (1/4)

- Given two data sets ([MaleData](#), [Female Data](#)) in which each row is a sample with 39 features, please perform the following operations:
 1. Merge these two data sets and find/plot the covariance matrix for the merged data set.
 2. Apply PCA and LDA transformations to the merged data set, respectively. Also, find/plot the covariance matrices for transformations, respectively. Describe the phenomena that you have observed.
 3. Use the first two principal components of PCA as well as the first two eigenvectors of LDA to represent the merged data set. Selectively plot portions of samples from MaleData and FemaleData, respectively. Describe the phenomena that you have observed.

<http://berlin.csie.ntnu.edu.tw/PastCourses/2004S-MachineLearningandDataMining/Homework/HW-1/MaleData.txt>

HW: Feature Transformation (3/4)

- Plot Covariance Matrix

```
CoVar=[
    3.0    0.5    0.4;
    0.9    6.3    0.2;
    0.4    0.4    4.2;
];
colormap('default');
surf(CoVar);
```

- Eigen Decomposition

```
BE=[
    3.0    3.5    1.4;
    1.9    6.3    2.2;
    2.4    0.4    4.2;
];
WI=[
    4.0    4.1    2.1;
    2.9    8.7    3.5;
    4.4    3.2    4.3;
];
```

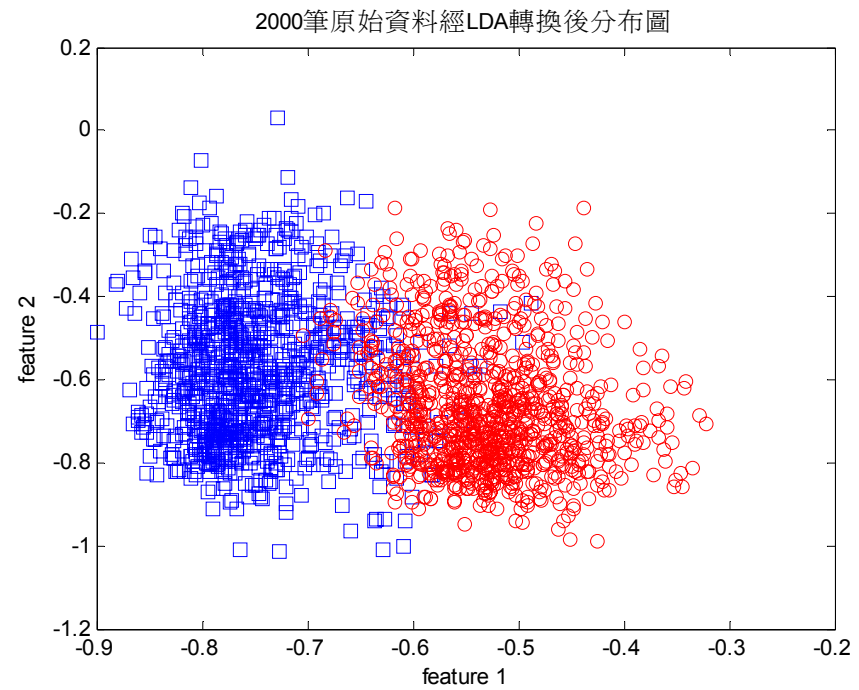
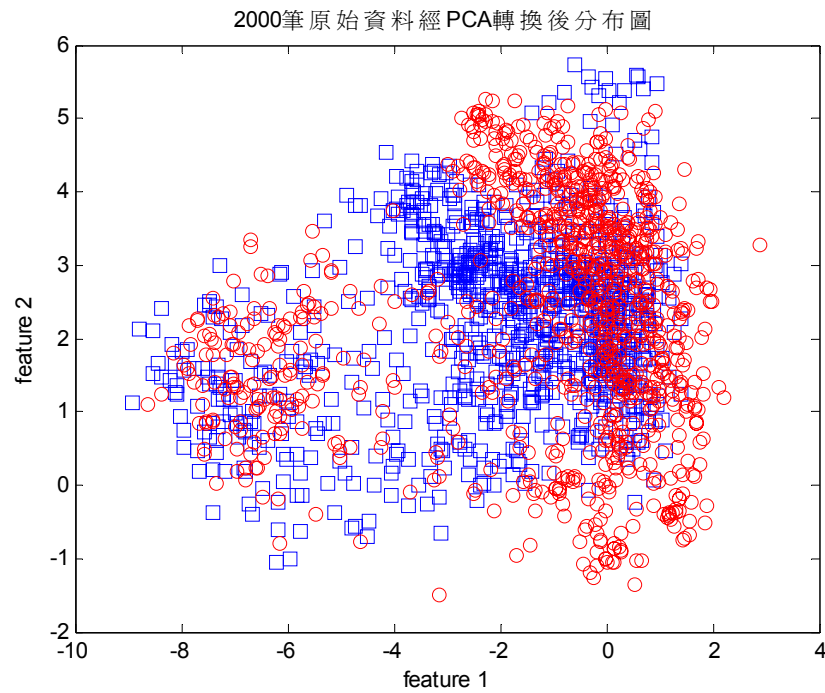
```
%LDA
IWI=inv(WI);
A=IWI*BE;
%PCA
A=BE+WI; % why ?? ( Prove it! )

[V,D]=eig(A);
[V,D]=eigs(A,3);

fid=fopen('Basis','w');
for i=1:3 % feature vector length
    for j=1:3 % basis number
        fprintf(fid,'%10.10f ',V(i,j));
    end
    fprintf(fid,'\n');
end
fclose(fid);
```

HW: Feature Transformation (4/4)

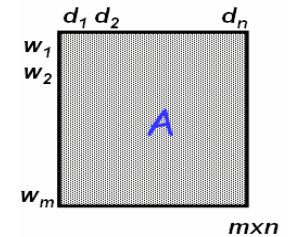
- Examples



Latent Semantic Analysis (LSA) (1/7)

- Also called **Latent Semantic Indexing (LSI)**, **Latent Semantic Mapping (LSM)**
- A technique originally proposed for Information Retrieval (IR), which projects queries and docs into a space with “latent” semantic dimensions

- **Co-occurring terms are projected onto the same dimensions**

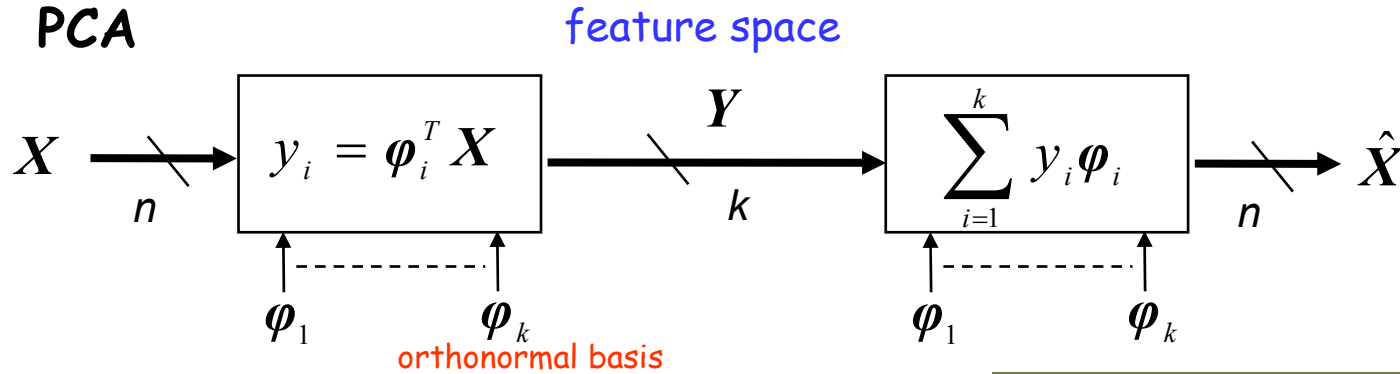


- In the latent semantic space (with fewer dimensions), a query and doc can have high cosine similarity even if they do not share any terms
 - Dimensions of the reduced space correspond to the axes of greatest variation
 - **Closely related to Principal Component Analysis (PCA)**

LSA (2/7)

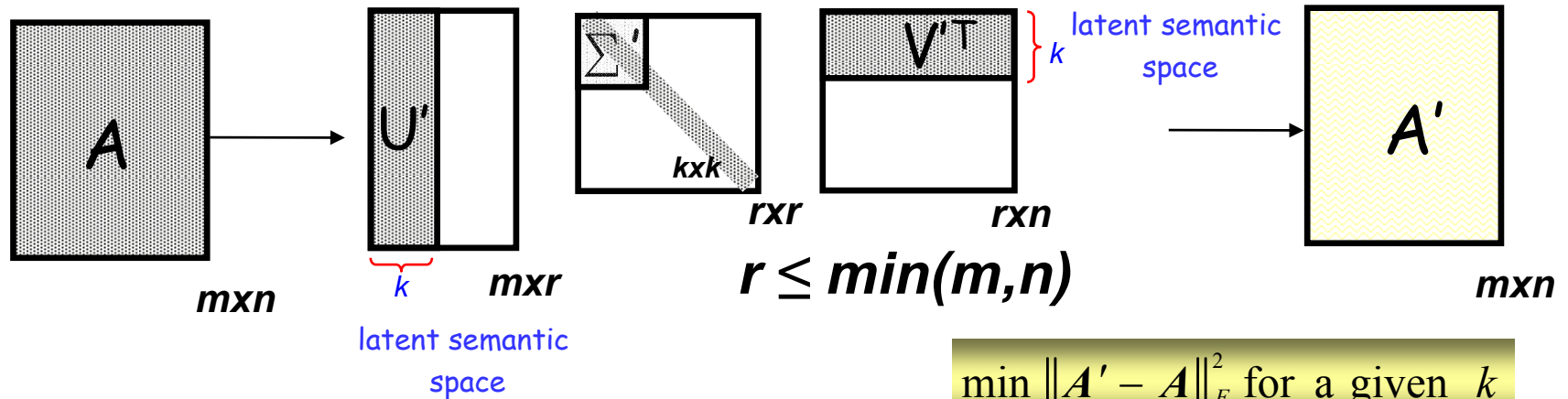
- Dimension Reduction and Feature Extraction

- PCA



- SVD (in LSA)

$$\min \|\hat{X} - X\|^2 \text{ for a given } k$$



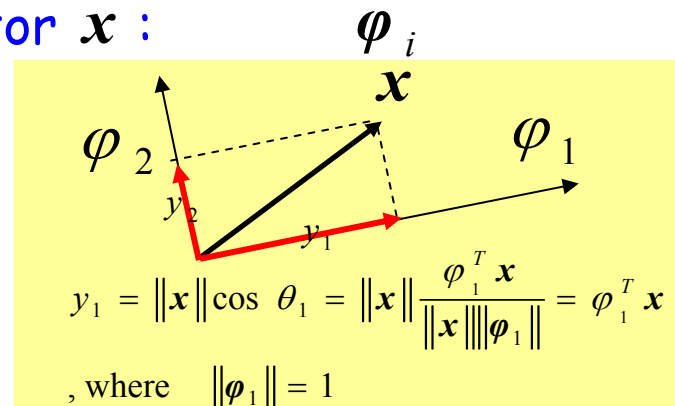
$$\min \|A' - A\|_F^2 \text{ for a given } k$$

LSA (3/7)

- Singular Value Decomposition (SVD) used for the word-document matrix
 - A least-squares method for dimension reduction

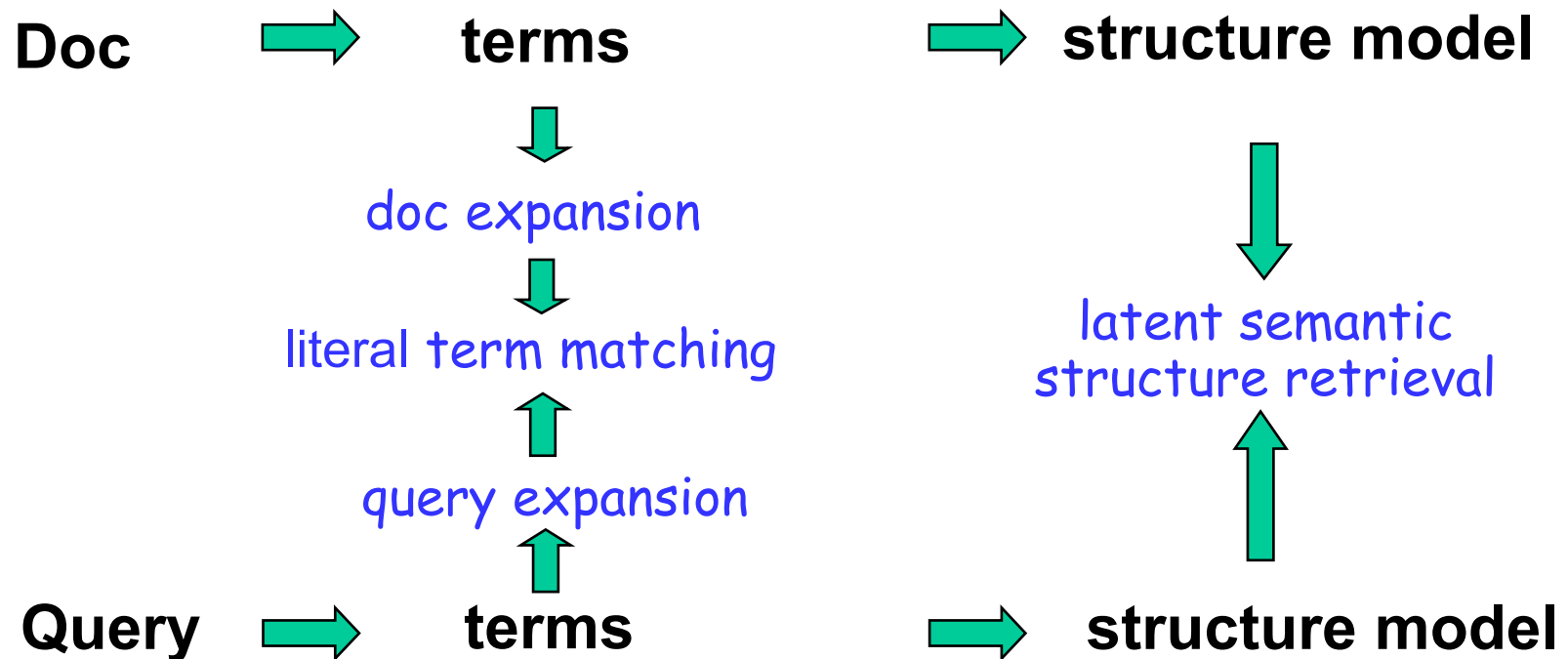
	Term 1	Term 2	Term 3	Term 4
Query	user	interface		
Document 1	user	interface	HCI	interaction
Document 2			HCI	interaction

Projection of a Vector \mathbf{x} :



LSA (4/7)

- Frameworks to circumvent vocabulary mismatch



LSA (5/7)

Titles

- c1: *Human machine interface for Lab ABC computer applications*
- c2: *A survey of user opinion of computer system response time*
- c3: *The EPS user interface management system*
- c4: *System and human system engineering testing of EPS*
- c5: *Relation of user-perceived response time to error measurement*
- m1: *The generation of random, binary, unordered trees*
- m2: *The intersection graph of paths in trees*
- m3: *Graph minors IV: Widths of trees and well-quasi-ordering*
- m4: *Graph minors: A survey*

Terms

Documents

	c1	c2	c3	c4	c5	m1	m2	m3	m4
<i>human</i>	1	0	0	1	0	0	0	0	0
<i>interface</i>	1	0	1	0	0	0	0	0	0
<i>computer</i>	1	1	0	0	0	0	0	0	0
<i>user</i>	0	1	1	0	1	0	0	0	0
<i>system</i>	0	1	1	2	0	0	0	0	0
<i>response</i>	0	1	0	0	1	0	0	0	0
<i>time</i>	0	1	0	0	1	0	0	0	0
<i>EPS</i>	0	0	1	1	0	0	0	0	0
<i>survey</i>	0	1	0	0	0	0	0	0	1
<i>trees</i>	0	0	0	0	0	1	1	1	0
<i>graph</i>	0	0	0	0	0	0	1	1	1
<i>minors</i>	0	0	0	0	0	0	0	1	1

LSA (6/7)

2-D Plot of Terms and Docs from Example

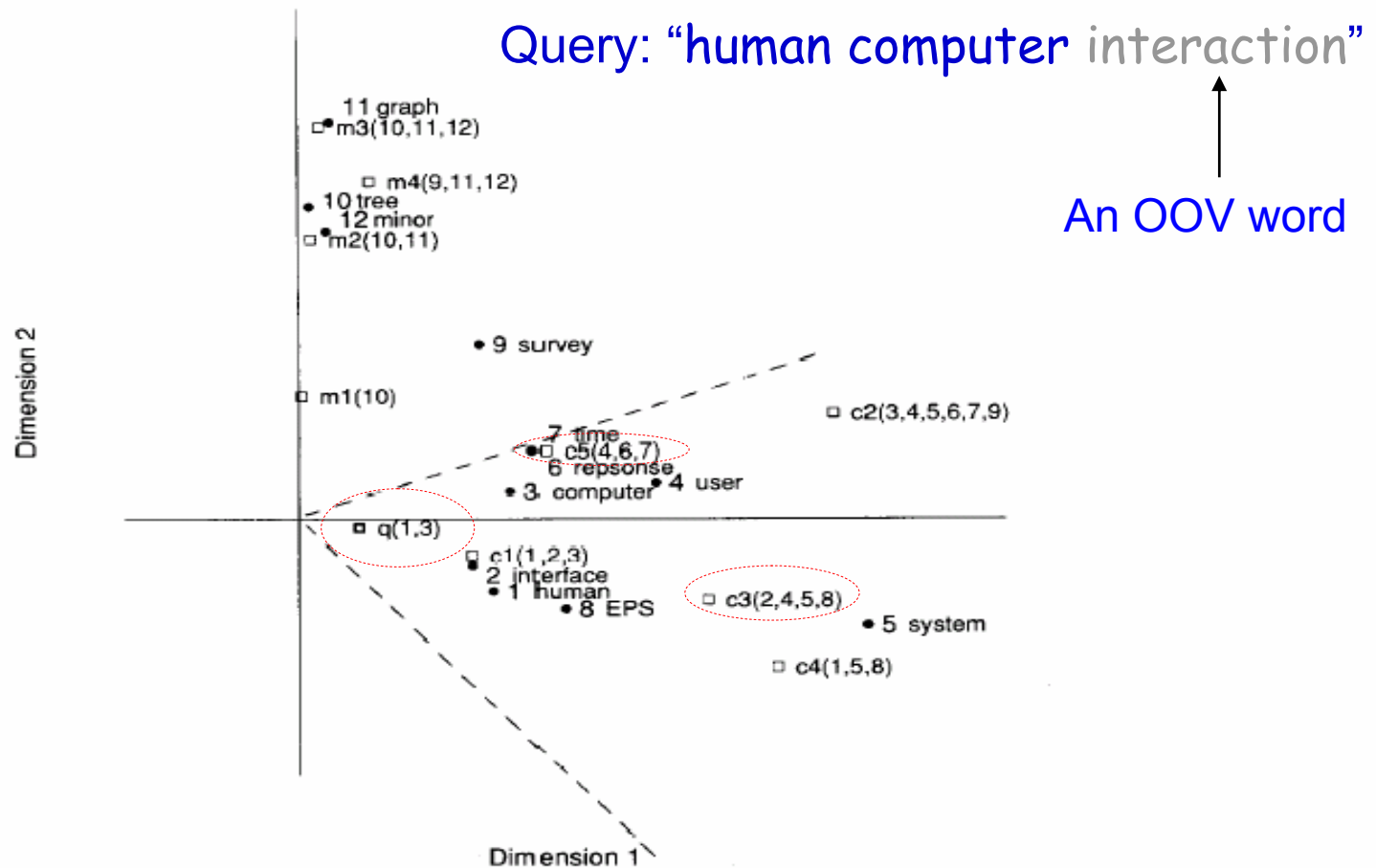
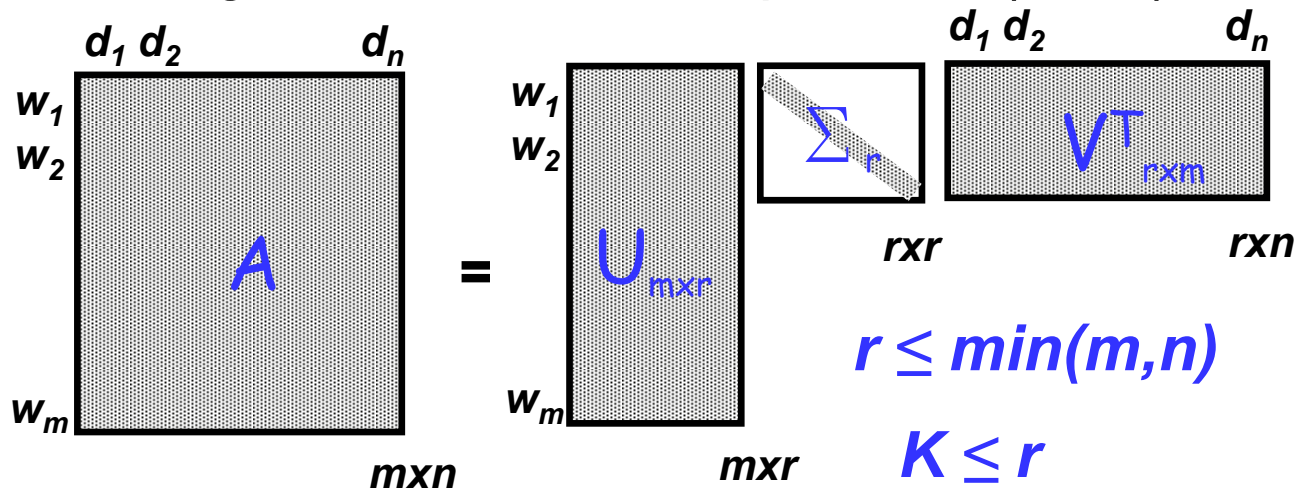


FIG. 1. A two-dimensional plot of 12 Terms and 9 Documents from the same TM set. Terms are represented by filled circles. Documents are shown as open squares, and component terms are indicated parenthetically. The query ("human computer interaction") is represented as a pseudo-document at point q . Axes are scaled for Document-Document or Term-Term comparisons. The dotted cone represents the region whose points are within a cosine of .9 from the query q . All documents about human-computer (c1-c5) are "near" the query (i.e., within this cone), but none of the graph theory documents (m1-m4) are nearby. In this reduced space, even documents c3 and c5 which share no terms with the query are near it.



LSA (7/7)

- Singular Value Decomposition (SVD)



Row $A \in R^n$
 Col $A \in R^m$
 Both U and V has orthonormal column vectors

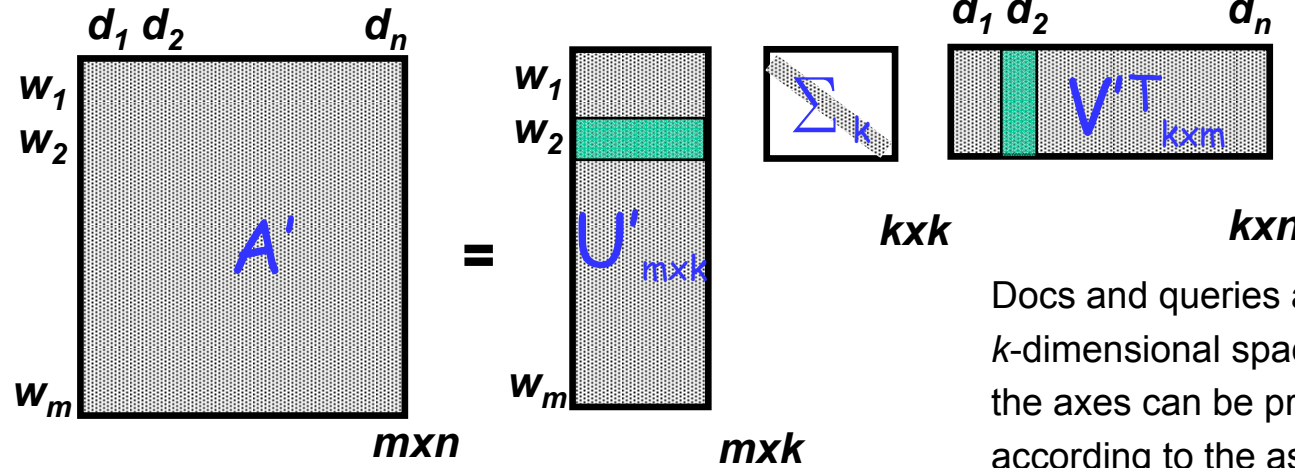
$$U^T U = I_{r \times r}$$

$$V^T V = I_{r \times r}$$

$$r \leq \min(m, n)$$

$$K \leq r$$

$$\|A\|_F^2 \geq \|A'\|_F^2$$



$$\|A\|_F^2 = \sum_{i=1}^m \sum_{j=1}^n a_{ij}^2$$

Docs and queries are represented in a k -dimensional space. The quantities of the axes can be properly weighted according to the associated diagonal values of Σ_k

LSA Derivations (1/7)

- Singular Value Decomposition (SVD)

- $A^T A$ is symmetric $n \times n$ matrix

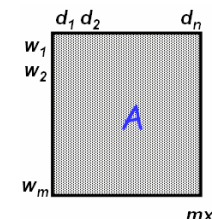
- All eigenvalues λ_j are nonnegative real numbers

$$\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_n \geq 0 \quad \Sigma^2 = \text{diag}(\lambda_1, \lambda_1, \dots, \lambda_n)$$

- All eigenvectors v_j are orthonormal ($\in R^n$)

$$V_{n \times n} = [v_1 v_2 \dots v_n] \quad v_j^T v_j = 1 \quad (V^T V = I_{n \times n})$$

sigma $\sigma_j = \sqrt{\lambda_j}, j = 1, \dots, n$



- Define **singular values**:

- As the square roots of the eigenvalues of $A^T A$

- As the lengths of the vectors Av_1, Av_2, \dots, Av_n

For $\lambda_i \neq 0, i=1, \dots, r,$
 $\{Av_1, Av_2, \dots, Av_r\}$ is an
 orthogonal basis of Col A

$$\sigma_1 = \|Av_1\|$$

$$\sigma_2 = \|Av_2\|$$

.....

$$\|Av_i\|^2 = v_i^T A^T A v_i = v_i^T \lambda_i v_i = \lambda_i$$

$$\Rightarrow \|Av_i\| = \sigma_i$$

LSA Derivations (2/7)

- $\{Av_1, Av_2, \dots, Av_r\}$ is an **orthogonal** basis of **Col A**

$$Av_i \bullet Av_j = (Av_i)^T Av_j = v_i^T A^T Av_j = \lambda_j v_i^T v_j = 0$$

- Suppose that A (or $A^T A$) has rank $r \leq n$

$$\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_r > 0, \quad \lambda_{r+1} = \lambda_{r+2} = \dots = \lambda_n = 0$$

- Define an **orthonormal** basis $\{u_1, u_2, \dots, u_r\}$ for Col A

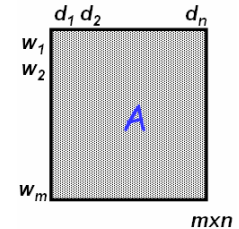
$$u_i = \frac{1}{\|Av_i\|} Av_i = \frac{1}{\sigma_i} Av_i \Rightarrow \sigma_i u_i = Av_i$$

u_i also an
orthonormal matrix
($m \times r$)

$$\Rightarrow [u_1 \ u_2 \ \dots \ u_r] \Sigma_r = A [v_1 \ v_2 \ \dots \ v_r]$$

V : an orthonormal matrix ($n \times r$)

Known in advance



- Extend to an orthonormal basis $\{u_1, u_2, \dots, u_m\}$ of R^m

$$\Rightarrow [u_1 \ u_2 \ \dots \ u_r \ \dots \ u_m] \Sigma = A [v_1 \ v_2 \ \dots \ v_r \ \dots \ v_n]$$

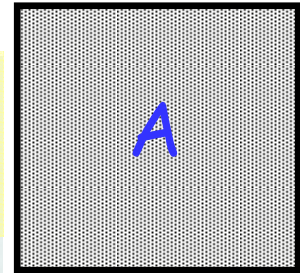
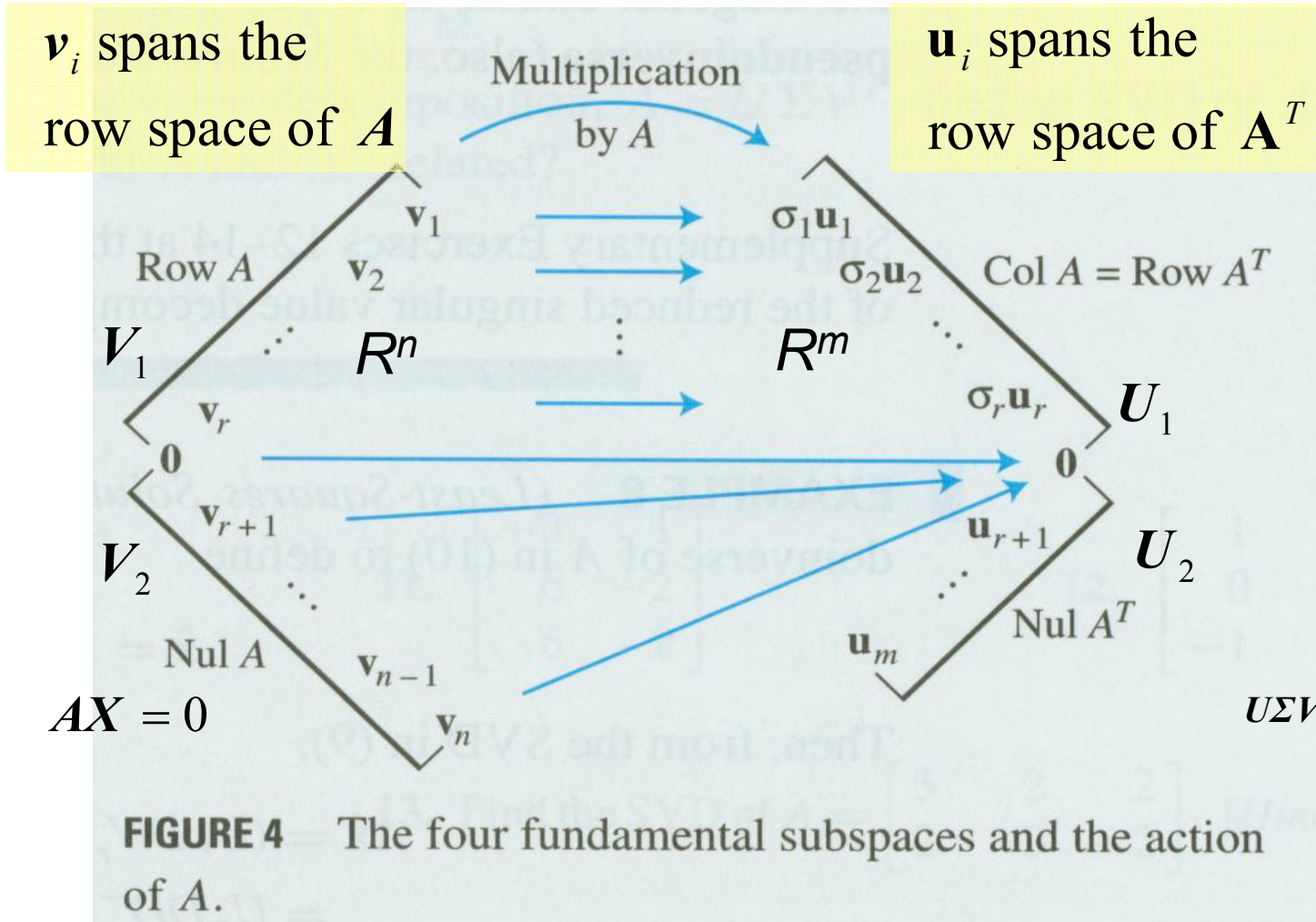
$$\Rightarrow U \Sigma = AV \Rightarrow U \Sigma V^T = A \underbrace{V V^T}$$

$$\Rightarrow A = U \Sigma V^T \quad \Sigma_{m \times n} = \begin{pmatrix} \Sigma_r & \mathbf{0}_{r \times (n-r)} \\ \mathbf{0}_{(m-r) \times r} & \mathbf{0}_{(m-r) \times (n-r)} \end{pmatrix} \quad I_{n \times n} \quad ?$$

$$\|A\|_F^2 = \sum_{i=1}^m \sum_{j=1}^n a_{ij}^2$$

$$\|A\|_F^2 = \sigma_1^2 + \sigma_2^2 + \dots + \sigma_r^2 \quad ?$$

LSA Derivations (3/7)



$m \times n$

$$\begin{aligned}
 U \Sigma V^T &= \begin{pmatrix} U_1 & U_2 \end{pmatrix} \begin{pmatrix} \Sigma_1 & 0 \\ 0 & 0 \end{pmatrix} \begin{pmatrix} V_1^T \\ V_2^T \end{pmatrix} \\
 &= U_1 \Sigma_1 V_1^T \\
 &= A V_1 V_1^T \quad \boxed{U \Sigma = A V} \\
 &= A
 \end{aligned}$$

LSA Derivations (4/7)

- Additional Explanations

- Each row of U is related to the projection of a corresponding row of A onto the basis formed by columns of V

$$A = U\Sigma V^T$$

$$\Rightarrow AV = U\Sigma V^T V = U\Sigma \Rightarrow U\Sigma = AV$$

- the i -th entry of a row of U is related to the projection of a corresponding row of A onto the i -th column of V

- Each row of V is related to the projection of a corresponding row of A^T onto the basis formed by U

$$A = U\Sigma V^T$$

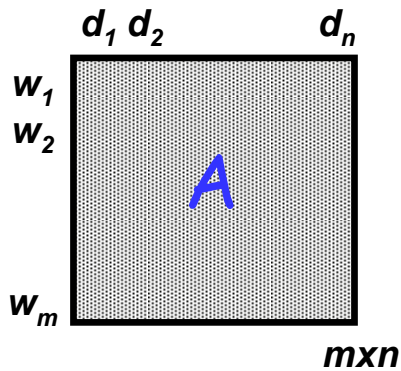
$$\Rightarrow A^T U = (U\Sigma V^T)^T U = V\Sigma U^T U = V\Sigma$$

$$\Rightarrow V\Sigma = A^T U$$

- the i -th entry of a row of V is related to the projection of a corresponding row of A^T onto the i -th column of U

LSA Derivations (5/7)

- Fundamental comparisons based on SVD
 - The original word-document matrix (A)



- compare two terms \rightarrow dot product of two rows of A
 - or an entry in AA^T
- compare two docs \rightarrow dot product of two columns of A
 - or an entry in $A^T A$
- compare a term and a doc \rightarrow each individual entry of A

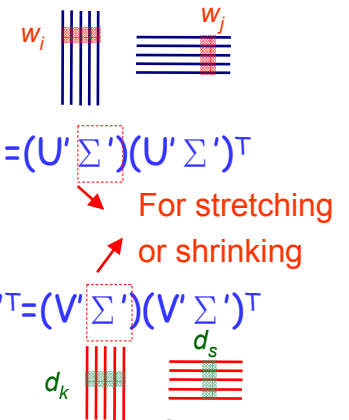
- The new word-document matrix (A')

$$U' = U_{m \times k}$$

$$\Sigma' = \Sigma_k$$

$$V' = V_{n \times k}$$

- compare two terms $A'A^T = (U' \Sigma' V'^T) (U' \Sigma' V'^T)^T = U' \Sigma' V'^T V' \Sigma'^T U'^T = (U' \Sigma') (U' \Sigma')^T$
 - \rightarrow dot product of two rows of $U' \Sigma'$
- compare two docs $A^T A = (U' \Sigma' V'^T)^T (U' \Sigma' V'^T) = V' \Sigma'^T U'^T U' \Sigma' V'^T = (V' \Sigma') (V' \Sigma')^T$
 - \rightarrow dot product of two rows of $V' \Sigma'$
- compare a query word and a doc \rightarrow each individual entry of A'



LSA Derivations (6/7)

- **Fold-in:** find representations for pseudo-docs q
 - For objects (new queries or docs) that did not appear in the original analysis
 - Fold-in a new $m \times 1$ query (or doc) vector

$$\hat{q}_{1 \times k} = \left(q^T \right)_{1 \times m} U_{m \times k} \Sigma^{-1}_{k \times k}$$

Just like a row of V Query represented by the weighted sum of its constituent term vectors The separate dimensions are differentially weighted

- Cosine measure between the query and doc vectors in the latent semantic space

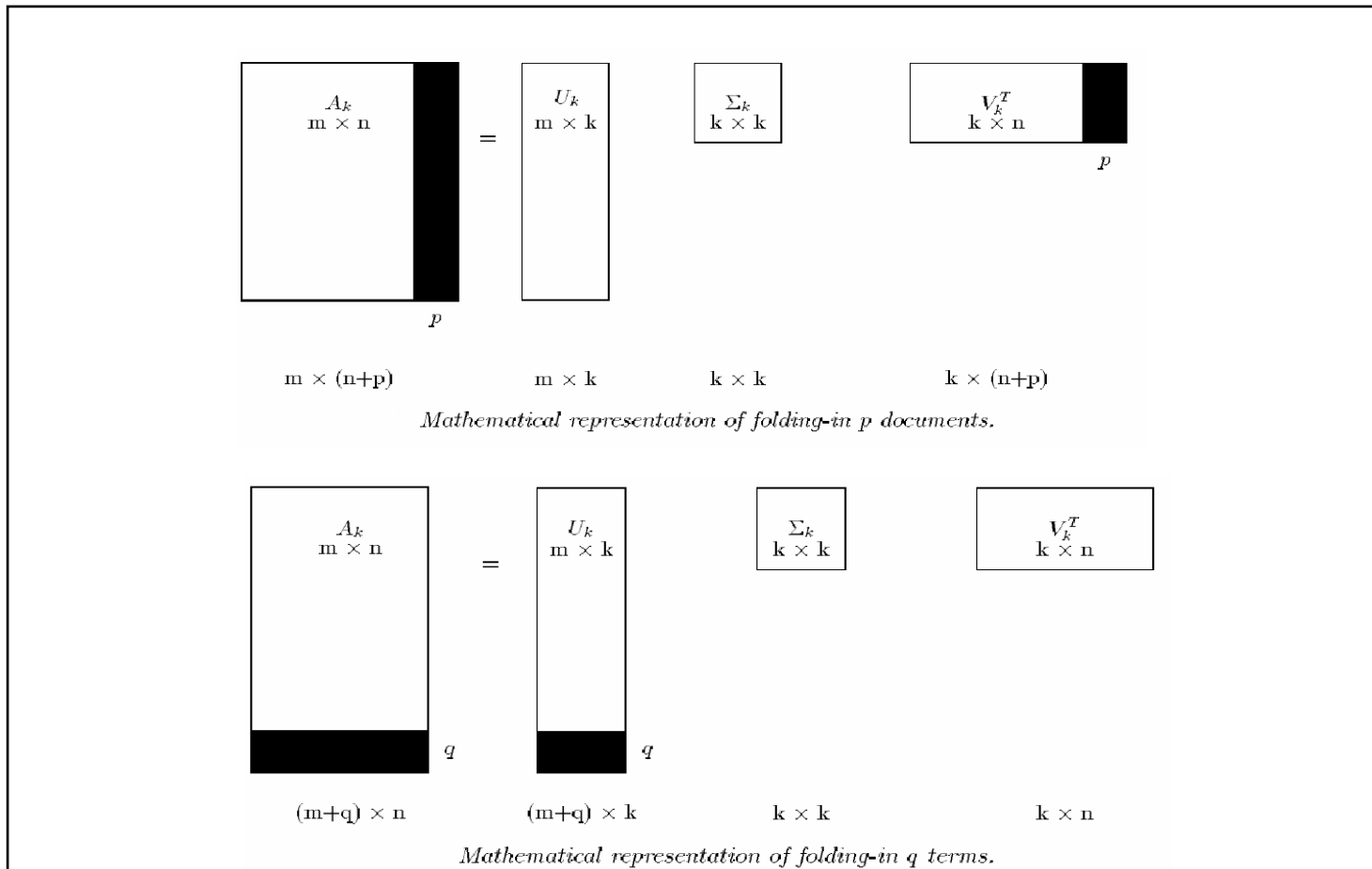
$$\text{sim} \left(\hat{q}, \hat{d} \right) = \text{coine} \left(\hat{q} \Sigma, \hat{d} \Sigma \right) = \frac{\hat{q} \Sigma^2 \hat{d}^T}{\left| \hat{q} \Sigma \right| \left| \hat{d} \Sigma \right|}$$

row vectors

LSA Derivations (7/7)

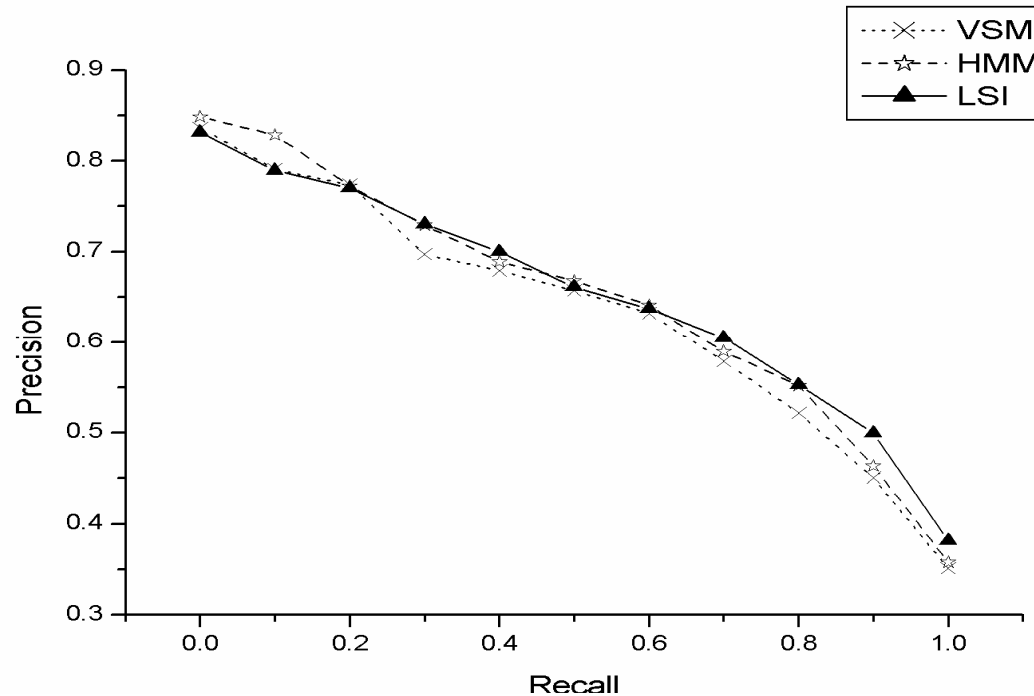
- Fold-in a new $1 \times n$ term vector

$$\hat{t}_{1 \times k} = t_{1 \times n} V_{n \times k} \Sigma_{k \times k}^{-1}$$



LSA Example

- Experimental results
 - HMM is consistently better than VSM at all recall levels
 - LSA is better than VSM at higher recall levels



Recall-Precision curve at 11 standard recall levels evaluated on TDT-3 SD collection. (Using word-level indexing terms)

LSA: Conclusions

- Advantages
 - A clean formal framework and a clearly defined optimization criterion (least-squares)
 - Conceptual simplicity and clarity
 - Handle synonymy problems (“heterogeneous vocabulary”)
 - Good results for high-recall search
 - Take term co-occurrence into account
- Disadvantages
 - High computational complexity
 - LSA offers only a partial solution to polysemy
 - E.g. bank, bass,...

LSA Toolkit: SVDLIBC (1/5)

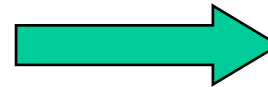
- Doug Rohde's SVD C Library version 1.3 is based on the [SVDPACKC](#) library
- Download it at <http://tedlab.mit.edu/~dr/>

LSA Toolkit: SVDLIBC (2/5)

- Given a sparse term-doc matrix

- E.g., 4 terms and 3 docs

	Doc		
Term	2.3	0.0	4.2
	0.0	1.3	2.2
	3.8	0.0	0.5
	0.0	0.0	0.0



Row #Tem	Col. # Doc	Nonzero entries
4	3	6
2		2 nonzero entries at Col 0
0	2.3	Col 0, Row 0
2	3.8	Col 0, Row 2
1		1 nonzero entry at Col 1
1	1.3	Col 1, Row 1
3		3 nonzero entries at Col 2
0	4.2	Col 2, Row 0
1	2.2	Col 2, Row 1
2	0.5	Col 2, Row 2

- Each entry is weighted by *TFxIDF* score

- Perform SVD to obtain corresponding term and doc vectors represented in the latent semantic space
- Evaluate the information retrieval capability of the LSA approach by using varying sizes (e.g., 100, 200, .., 600 etc.) of LSA dimensionality

LSA Toolkit: SVDLIBC (3/5)

- Example: term-docmatrix

Indexing Term no.	Doc no.	Nonzero entries
51253	2265	218852
77		
508	7.725771	
596	16.213399	
612	13.080868	
709	7.725771	
713	7.725771	
744	7.725771	
1190	7.725771	
1200	16.213399	
1259	7.725771	

- SVD command (IR_svd.bat)

`svd -r st -o LSA100 -d 100 Term-Doc-Matrix` **output** → **LSA100-Ut**
LSA100-S
LSA100-Vt

sparse matrix input prefix of output files No. of reserved eigenvectors name of sparse matrix input

LSA Toolkit: SVDLIBC (4/5)

- **LSA100-Ut**

51253 words

100 51253

0.003 0.001

0.002 0.002

word vector (u^T): 1x100

- **LSA100-S**

100

2686.18

829.941

559.59

....

100 eigenvalues

- **LSA100-Vt** 2265 docs

100 2265

0.021 0.035

0.012 0.022

doc vector (v^T): 1x100

LSA Toolkit: SVDLIBC (5/5)

- Fold-in a new $m \times 1$ query vector

$$\hat{q}_{1 \times k} = \left(q^T \right)_{1 \times m} U_{m \times k} \Sigma_{k \times k}^{-1}$$

Just like a row of V

Query represented by the weighted sum of its constituent term vectors

The separate dimensions are differentially weighted

TFxIDF weighted beforehand

- Cosine measure between the query and doc vectors in the latent semantic space

$$\text{sim}(\hat{q}, \hat{d}) = \text{coine}(\hat{q}\Sigma, \hat{d}\Sigma) = \frac{\hat{q}\Sigma^2\hat{d}^T}{|\hat{q}\Sigma| |\hat{d}\Sigma|}$$