# Data Analysis and Dimension Reduction

## - PCA, LDA and LSA

Berlin Chen
Department of Computer Science & Information Engineering
National Taiwan Normal University
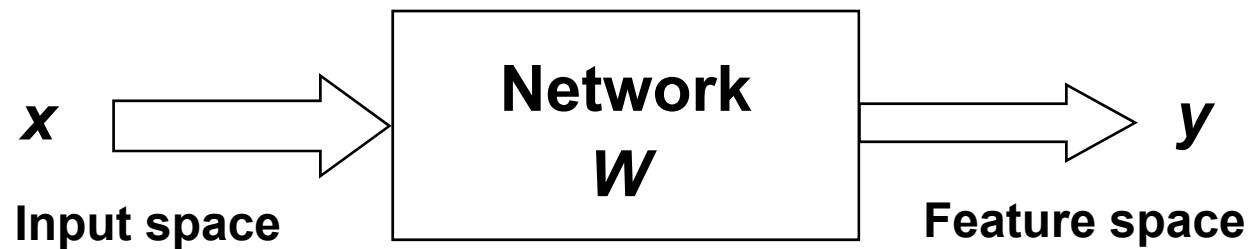
References:
1. *Introduction to Machine Learning* , Chapter 6
2. *Data Mining: Concepts, Models, Methods and Algorithms*, Chapter 3

# Introduction (1/3)

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

$x$ → [ **Network** $W$ ] → $y$

**Input space**          **Feature space**

  - 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)

# Bivariate Random Variables

- If the random variables *X* and *Y* have a certain joint distribution that describes a bivariate random variable

$$Z = \begin{bmatrix} X \\ Y \end{bmatrix}$$ bivariate random variable

$$\Rightarrow \mu_Z = \begin{bmatrix} \mu_X \\ \mu_Y \end{bmatrix}$$ mean vector

$$\Rightarrow \Sigma_Z = \begin{bmatrix} \sigma_{X,X} & \sigma_{X,Y} \\ \sigma_{X,Y} & \sigma_{Y,Y} \end{bmatrix}$$ covariance matrix

variance $\sigma_{X,X} = \sigma_X^2 = \mathbf{E}\left[(X - \mathbf{E}[X])^2\right] = \mathbf{E}\left[X^2\right] - (\mathbf{E}[X])^2$

covariance $\sigma_{X,Y} = \sigma_{Y,X} = \mathbf{E}\left[(X - \mathbf{E}[X])(Y - \mathbf{E}[Y])\right] = \mathbf{E}[XY] - \mathbf{E}[X]\mathbf{E}[Y]$

# Multivariate Random Variables

- If the random variables $X_1, X_2, \ldots, X_n$ have a certain joint distribution that describes a multivariate random variable

$$X = \begin{bmatrix} X_1 \\ \vdots \\ X_n \end{bmatrix} \qquad \text{multivariate random variable}$$

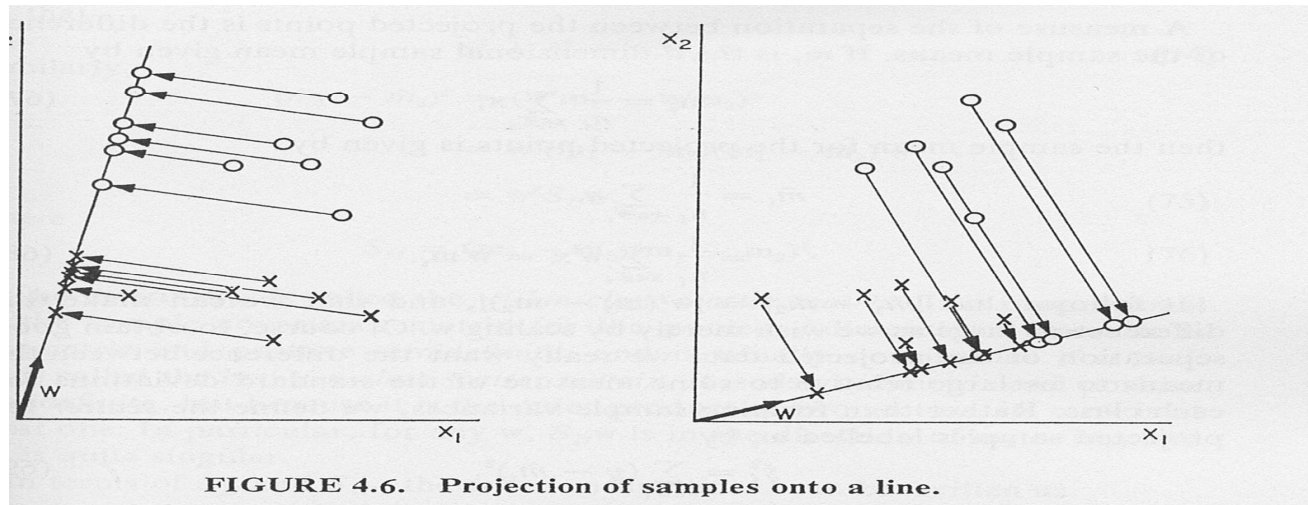$$\Rightarrow \mu_X = \begin{bmatrix} \mu_{X_1} \\ \vdots \\ \mu_{X_n} \end{bmatrix} \qquad \text{mean vector}$$

$$\Rightarrow \Sigma_Z = \begin{bmatrix} \sigma_{X_1,X_1} & \cdots & \sigma_{X_1,X_n} \\ \vdots & \vdots & \vdots \\ \sigma_{X_n,X_1} & \cdots & \sigma_{X_n,X_n} \end{bmatrix} \qquad \text{covariance matrix}$$

variance $\sigma_{X_i,X_i} = \sigma_{X_i}^2 = \mathbf{E}\left[(X_i - \mathbf{E}[X_i])^2\right] = \mathbf{E}\left[X_i^2\right] - (\mathbf{E}[X_i])^2$

covariance $\sigma_{X_i,X_j} = \sigma_{X_j,X_i} = \mathbf{E}\left[(X_i - \mathbf{E}[X_i])(X_j - \mathbf{E}[X_j])\right] = \mathbf{E}\left[X_i X_j\right] - \mathbf{E}[X_i]\mathbf{E}[X_j]$

# 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
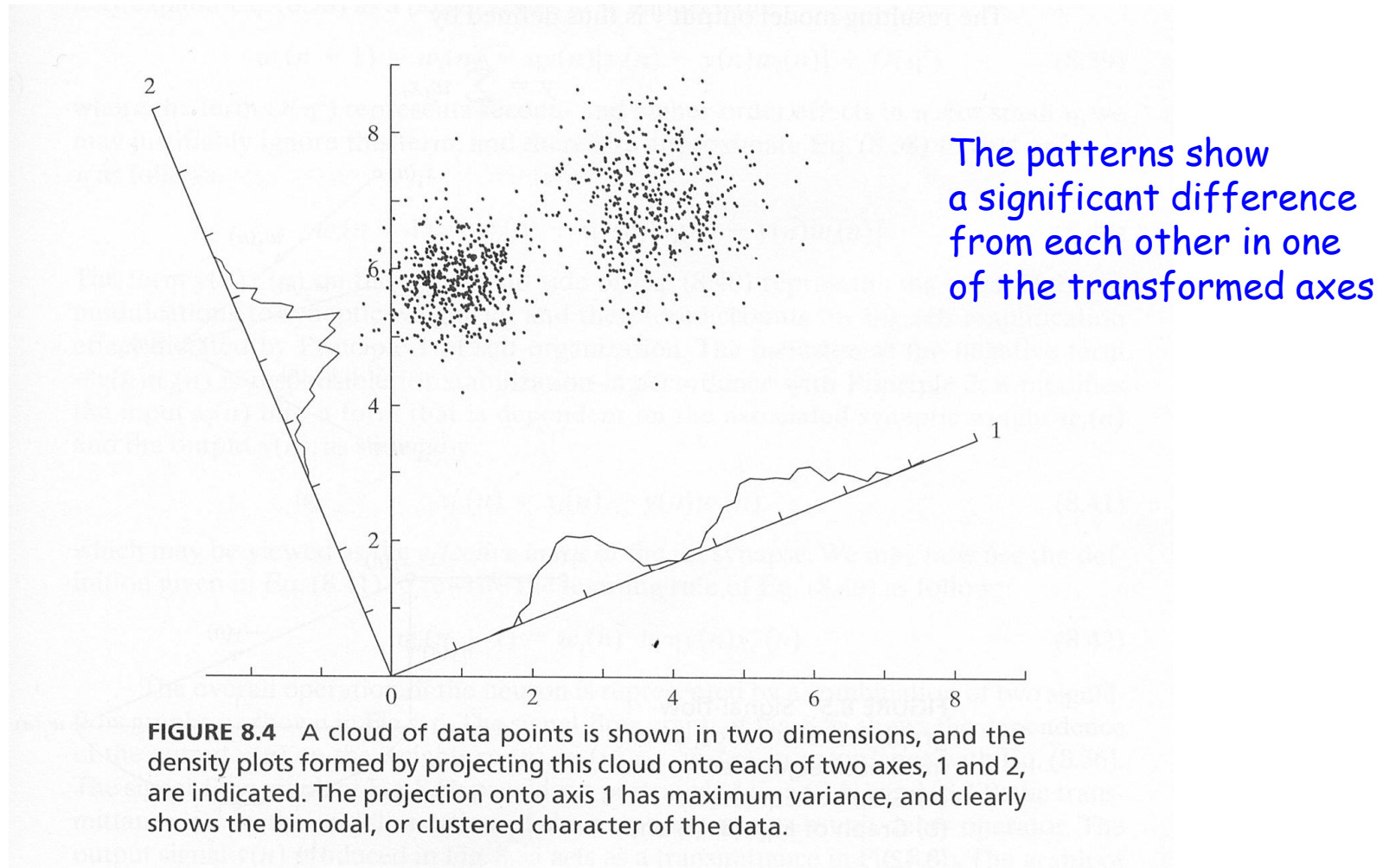    - PLSA with multinomial latent cluster distributions



FIGURE 4.6. Projection of samples onto a line.

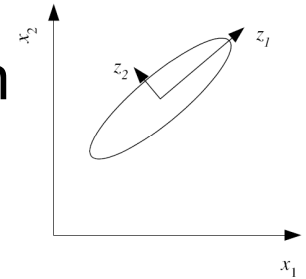# Principal Component Analysis (PCA) (1/2)

- ## Known as Karhunen-Loěve 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)



The patterns show a significant difference from each other in one of the transformed axes

**FIGURE 8.4** A cloud of data points is shown in two dimensions, and the density plots formed by projecting this cloud onto each of two axes, 1 and 2, are indicated. The projection onto axis 1 has maximum variance, and clearly shows the bimodal, or clustered character of the data.

# PCA Derivations (1/13)

- ## Suppose $\boldsymbol{x}$ is an *n*-dimensional zero mean random vector, $\boldsymbol{\mu} = \mathbf{E}\{\boldsymbol{x}\} = \boldsymbol{0}$

  - If $\boldsymbol{x}$ is not zero-mean, we can subtract the mean before processing the following analysis

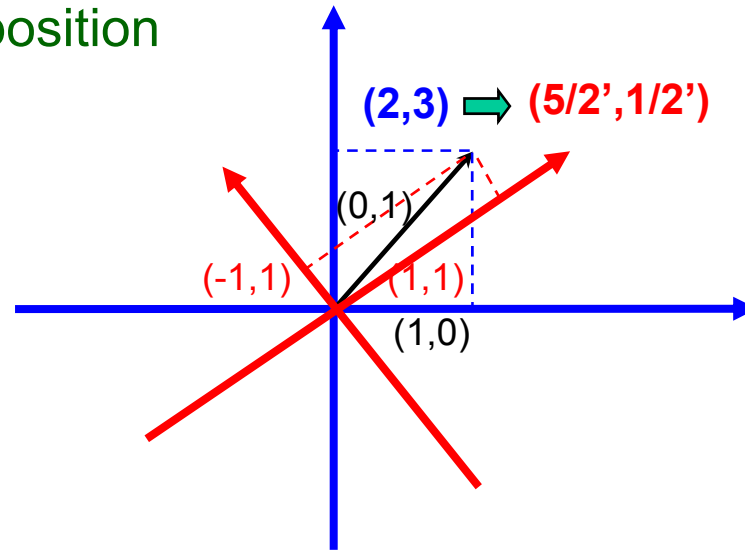  - $\boldsymbol{x}$ can be represented without error by the summation of *n* linearly independent vectors

$$\boldsymbol{x} = \sum_{i=i}^{n} y_i \boldsymbol{\varphi}_i = \boldsymbol{\Phi}\boldsymbol{y} \quad \text{where} \quad \mathbf{y} = \begin{bmatrix} y_1 & . & y_i & . & y_n \end{bmatrix}^T$$
$$\boldsymbol{\Phi} = \begin{bmatrix} \boldsymbol{\varphi}_1 & . & \boldsymbol{\varphi}_i & . & \boldsymbol{\varphi}_n \end{bmatrix}$$

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

The basis vectors

# PCA Derivations (2/13)

Subspace Decomposition



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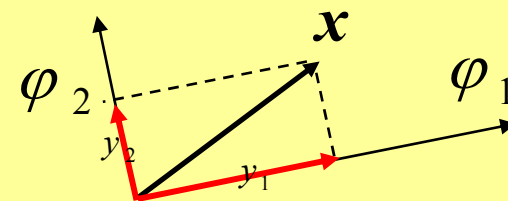
orthogonal basis sets

# PCA Derivations (3/13)

– Further assume the column (basis) vectors of the matrix $\boldsymbol{\Phi}$ form an orthonormal set

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

• Such that $y_i$ is equal to the projection of $\boldsymbol{x}$ on $\boldsymbol{\varphi}_i$

$$\forall_i \quad y_i = \boldsymbol{x}^T \boldsymbol{\varphi}_i = \boldsymbol{\varphi}_i^T \boldsymbol{x}$$

$$y_1 = \|\boldsymbol{x}\| \cos \theta_1 = \|\boldsymbol{x}\| \frac{\boldsymbol{\varphi}_1^T \boldsymbol{x}}{\|\boldsymbol{x}\| \|\boldsymbol{\varphi}_1\|} = \boldsymbol{\varphi}_1^T \boldsymbol{x}$$

$$, \text{where} \quad \|\boldsymbol{\varphi}_1\| = 1$$

# PCA Derivations (4/13)

– Further assume the column (basis) vectors of the matrix $\boldsymbol{\Phi}$ form an orthonormal set

$$\boldsymbol{\Sigma} = \mathbf{E}\left\{ (\boldsymbol{x} - \boldsymbol{\mu})(\boldsymbol{x} - \boldsymbol{\mu})^T \right\}$$

$$\approx \left( \frac{1}{N} \sum_{i=1}^{N} \boldsymbol{x}_i \boldsymbol{x}_i^T \right) \underbrace{- \boldsymbol{\mu\mu}^T}_{0}$$

- $y_i$ also has the following properties

  – Its mean is zero, too

$$\mathbf{E}\left\{ y_i \right\} = \mathbf{E}\left\{ \boldsymbol{\varphi}_i^T \boldsymbol{x} \right\} = \boldsymbol{\varphi}_i^T \mathbf{E}\left\{ \boldsymbol{x} \right\} = \boldsymbol{\varphi}_i^T \boldsymbol{0} = 0$$

$$\boldsymbol{R} = \mathbf{E}\left\{ \boldsymbol{x}\boldsymbol{x}^T \right\} \approx \frac{1}{N} \sum_{i} \boldsymbol{x}_i \boldsymbol{x}_i^T$$

  – Its variance is

$$\sigma_i^2 = \mathbf{E}\left\{ y_i^2 \right\} - \left[ \mathbf{E}\{ y \} \right]^2 = \mathbf{E}\left\{ y_i^2 \right\} = \mathbf{E}\left\{ \boldsymbol{\varphi}_i^T \boldsymbol{x}\boldsymbol{x}^T \boldsymbol{\varphi}_i \right\} = \boldsymbol{\varphi}_i^T \mathbf{E}\left\{ \boldsymbol{x}\boldsymbol{x}^T \right\} \boldsymbol{\varphi}_i$$

$$= \boldsymbol{\varphi}_i^T \boldsymbol{R} \boldsymbol{\varphi}_i \qquad \left[ \boldsymbol{R} \text{ is the (auto-)correlation matrix of } \boldsymbol{x} \right]$$

- The correlation between two projections $y_i$ and $y_j$ is

$$\mathbf{E}\left\{ y_i y_j \right\} = \mathbf{E}\left\{ \left( \boldsymbol{\varphi}_i^T \boldsymbol{x} \right)\left( \boldsymbol{\varphi}_j^T \boldsymbol{x} \right)^T \right\} = \mathbf{E}\left\{ \boldsymbol{\varphi}_i^T \boldsymbol{x}\boldsymbol{x}^T \boldsymbol{\varphi}_j \right\}$$

$$= \boldsymbol{\varphi}_i^T \mathbf{E}\left\{ \boldsymbol{x}\boldsymbol{x}^T \right\} \boldsymbol{\varphi}_j = \boldsymbol{\varphi}_i^T \boldsymbol{R} \boldsymbol{\varphi}_j$$

# 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 $\boldsymbol{x}$ well in **mean-squared error criterion**

original vector

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

reconstructed vector

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

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

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

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

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

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

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

We should discard the bases where the projections have lower variances

- ## Minimum Mean-Squared Error Criterion

  – If the orthonormal (basis) set $\varphi_i\text{'s}$ is selected to be the eigenvectors of the correlation matrix $R$, associated with eigenvalues $\lambda_i\text{'s}$

  - They will have the property that:

  $R$ is real and symmetric, therefore its eigenvectors $R$ form a orthonormal set

  $$R\,\varphi_j = \lambda_j\,\varphi_j$$

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

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

  $$\overline{\varepsilon}(m) = \sum_{j=m+1}^{n} \sigma_j^2$$

  $$= \sum_{j=m+1}^{n} \varphi_j^T R\,\varphi_j = \sum_{j=m+1}^{n} \varphi_j^T \lambda_j \varphi_j = \sum_{j=m+1}^{n} \lambda_j$$

- ## Minimum Mean-Squared Error Criterion

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

  $$\bar{\varepsilon}_{eigen}(m) = \sum_{j=m+1}^{n} \lambda_j \quad \left(\text{where} \, \lambda_1 \geq \ldots \geq \lambda_m \geq \ldots \geq \lambda_n \geq 0\right)$$

  - Any two projections $y_i$ and $y_j$ will be mutually uncorrelated

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

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

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

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

# PCA Derivations (8/13)

- A Two-dimensional Example of Principle Component Analysis

- ## 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

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

Define: $J = \sum_{j=m+1}^{n} \varphi_j^T R \varphi_j - \sum_{j=m+1}^{n} \sum_{k=m+1}^{n} u_{jk}\left(\varphi_j^T \varphi_k - \delta_{jk}\right)$

Partial Differentiation

$$\Rightarrow \forall_{m+1 \leq j \leq n} \frac{\partial J}{\partial \varphi_j} = 2R\varphi_j - 2\sum_{k=m+1}^{n} u_{jk}\varphi_k = 0 \quad \left(\text{where } u_j^T = \begin{bmatrix} u_{j\ m+1} & \ldots & u_{jn} \end{bmatrix}\right)$$

$$\Rightarrow \forall_{m+1 \leq j \leq n} \quad R\varphi_j = \Phi_{n-m} u_j \quad \left(\text{where } \Phi_{n-m} = \begin{bmatrix} \varphi_{m+1} & \ldots & \varphi_n \end{bmatrix}\right)$$

$$\Rightarrow R\begin{bmatrix} \varphi_{m+1} & \ldots & \varphi_n \end{bmatrix} = \Phi_{n-m}\begin{bmatrix} u_{m+1} & \ldots & u_n \end{bmatrix}$$

$$\Rightarrow R\Phi_{n-m} = \Phi_{n-m}U_{n-m} \quad \left(\text{where } U_{n-m} = \begin{bmatrix} u_{m+1} & \ldots & u_n \end{bmatrix}\right)$$

Have a particular solution if $U_{n-m}$ is a diagonal matrix and its diagonal elements is the eigenvalues $\lambda_{m+1}\ldots\lambda_n$ of $R$ and $\varphi_{m+1}\ldots\varphi_n$ is their corresponding eigenvectors

# PCA Derivations (10/13)

- ## Given an input vector $x$ with dimension $m$

  - Try to construct a linear transform $\Phi'$ ($\Phi'$ is an $n \times m$ matrix $m < n$) such that the truncation result, $\Phi'^T x$, is optimal in mean-squared error criterion

$$x = \begin{bmatrix} x_1 \\ x_2 \\ . \\ . \\ x_n \end{bmatrix}$$

**Encoder**

$\Phi'^T$

where $\Phi' = \begin{bmatrix} \varphi_1 \varphi_2 .. \varphi_m \end{bmatrix}$

$y = \Phi'^T x$

$$y = \begin{bmatrix} y_1 \\ y_2 \\ . \\ . \\ y_m \end{bmatrix}$$

$$y = \begin{bmatrix} y_1 \\ y_2 \\ . \\ . \\ y_m \end{bmatrix}$$

**Decoder**

$\Phi'$

$\hat{x} = \Phi' y$

$$\hat{x} = \begin{bmatrix} \hat{x}_1 \\ \hat{x}_2 \\ . \\ . \\ \hat{x}_n \end{bmatrix}$$

$$\text{minimize } E_x \left( (\hat{x}\text{-}x)^T (\hat{x}\text{-}x) \right)$$

# 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 $\dfrac{\lambda_1 + \lambda_2 + \cdots + \lambda_m}{\lambda_1 + \lambda_2 + \cdots + \lambda_m + \cdots + \lambda_n} \geq Threshold$



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

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

# PCA Derivations (13/13)

- PCA finds a linear transform **W** such that the <span style="color:blue">sum</span> of **average between-class variation** and **average within-class variation** is maximal

$$J(W) = \left|\widetilde{S}\right| \overset{?}{=} \left|\widetilde{S}_w + \widetilde{S}_b\right| = \left|W^T S_w W + W^T S_b W\right|$$

$$S = \frac{1}{N}\sum_i (x_i - \overline{x})(x_i - \overline{x})^T$$

<span style="color:blue">sample index</span>

$$S_w = \frac{1}{N}\sum_j N_j \Sigma_j$$

<span style="color:blue">class index</span>

$$S_b = \frac{1}{N}\sum_j N_j (\overline{x}_j - \overline{x})(\overline{x}_j - \overline{x})^T$$

$$\widetilde{S}_w = W^T S_w W$$

$$\widetilde{S}_b = W^T S_b W$$

Try to show that:
$$S = S_w + S_b$$

# 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

$$\mathbf{x}_1 = \begin{bmatrix} x_{1,1} \\ x_{1.2} \\ . \\ . \\ . \\ x_{1.n} \end{bmatrix}, \mathbf{x}_2 = \begin{bmatrix} x_{2,1} \\ x_{2.2} \\ . \\ . \\ . \\ x_{2.n} \end{bmatrix}, \dots\dots, \quad \mathbf{x}_L = \begin{bmatrix} x_{L,1} \\ x_{L.2} \\ . \\ . \\ . \\ x_{L.n} \end{bmatrix}$$

  - Steps
    - Convert each of the *L* reference images into a vector of floating point numbers representing light intensity in each pixel
    - Calculate the coverance/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 $K$ ($K \le L$) 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

$$\hat{\mathbf{x}}_i = \overline{\mathbf{x}} + w_{i,1}\mathbf{e}(1) + w_{i,2}\mathbf{e}(2) + \ldots + w_{i,K}\mathbf{e}(K)$$
$$\Rightarrow \mathbf{y}_i = \left[1, w_{i,1}, w_{i,2}, \ldots, w_{i,K}\right]$$

*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**

**A projected face image**



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



**?**

(Indicate directions of variations between faces )

# 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 ($\mu$)



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

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

SI HMM model

# 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 $W$ such that the ratio of **average between-class variation** over **average within-class variation** is maximal



Fig. 10-1 An example of feature extraction for classification.

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 $\boldsymbol{x}_i$ with dimensionality $n$, each of them is belongs to one of the $J$ classes $g(\boldsymbol{x}_i) = j, \quad j \in \{1, 2, ...., J\}, g(\cdot)$ is class index

  - The sample mean is: $\displaystyle \overline{\boldsymbol{x}} = \frac{1}{N} \sum_{i=1}^{N} \boldsymbol{x}_i$

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

  - The class sample covariances are: $\displaystyle \boldsymbol{\Sigma}_j = \frac{1}{N_j} \sum_{g(\boldsymbol{x}_i) = j} (\boldsymbol{x}_i - \overline{\boldsymbol{x}}_j)(\boldsymbol{x}_i - \overline{\boldsymbol{x}}_j)^T$

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

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

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

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

# LDA Derivations (2/4)

- If the transform $W = [w_1 \, w_2 \, .... \, w_m]$ is applied

  - The sample vectors will be $y_i = W^T x_i$

  - The sample mean will be $\bar{y} = \dfrac{1}{N} \sum\limits_{i=1}^{N} W^T x_i = W^T \left( \dfrac{1}{N} \sum\limits_{i=1}^{N} x_i \right) = W^T \bar{x}$

  - The class sample means will be $\bar{y}_j = \dfrac{1}{N_j} \sum\limits_{g(x_i)=j} W^T x_i = W^T \bar{x}_j$

  - The **average within-class variation** will be

$$\tilde{S}_w = \frac{1}{N} \sum_j N_j \left\{ \frac{1}{N_j} \cdot \sum_{g(x_i)=j} \left( W^T x_i - \frac{1}{N_j} \sum_{g(x_i)=j} \left( W^T x_i \right) \right) \left( W^T x_i - \frac{1}{N_j} \sum_{g(x_i)=j} \left( W^T x_i \right) \right)^T \right\}$$

$$= W^T \left\{ \frac{1}{N} \sum_j N_j \Sigma_j \right\} W$$

$$= W^T S_w W$$

$$\begin{array}{ccc} \bar{x} & \xrightarrow{W} & \bar{y} \\ S_w & & \tilde{S}_w \\ S_b & & \tilde{S}_b \end{array}$$

# LDA Derivations (3/4)

- If the transform $W = [w_1 w_2 .... w_m]$ is applied
  - Similarly, the **average between-class variation** will be

    $$\widetilde{S}_b = W^T S_b W$$

  - Try to find optimal $W$ such that the following objective function is maximized

    $$J(W) = \frac{|\widetilde{S}_b|}{|\widetilde{S}_w|} = \frac{|W^T S_b W|}{|W^T S_w W|}$$

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

      $$S_b w_i = \lambda_i S_w w_i$$

    - That is, $w_i$'s are the eigenvectors corresponding to the largest eigenvalues of $S_w^{-1} S_b$

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

# LDA Derivations (4/4)

- Proof:

$$\because \hat{W} = \underset{\hat{w}}{\arg\max}\, J(W) = \underset{\hat{w}}{\arg\max}\, \frac{\left|\tilde{S}_b\right|}{\left|\tilde{S}_w\right|} = \underset{\hat{w}}{\arg\max}\, \frac{\left|W^T S_b W\right|}{\left|W^T S_w W\right|}$$

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

The qradtic form has optimal solution : $\lambda_i = \dfrac{w^T_i S_b w_i}{w^T_i S_w w_i}$

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

$$\Rightarrow \frac{\partial \lambda_i}{\partial w_i} = \frac{2 S_b w_i \left(w^T_i S_w w_i\right) - 2 S_w w_i \left(w^T_i S_b w_i\right)}{\left(w^T_i S_w w_i\right)^2} = 0$$

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

$$\Rightarrow \frac{S_b w_i \left(w^T_i S_w w_i\right)}{\left(w^T_i S_w w_i\right)^2} - \frac{S_w w_i \left(w^T_i S_b w_i\right)}{\left(w^T_i S_w w_i\right)^2} = 0$$

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

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

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

# LDA Examples: Feature Transformation (1/2)

- Example1: Experiments on Speech Signal Processing

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



Calculated using Year-99's 5471 files

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

Covariance Matrix of the 18-cepstral vectors



Calculated using Year-99's 5471 files

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

**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 approaches 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 (2/4)

# 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



2000筆原始資料經PCA轉換後分布圖

2000筆原始資料經LDA轉換後分布圖

# 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)

- **Dimension Reduction and Feature Extraction**
  - **PCA**

<span style="color:blue">feature space</span>

$$X \xrightarrow{\quad n \quad} \boxed{y_i = \varphi_i^T X} \xrightarrow{\quad Y \quad}_{k} \boxed{\sum_{i=1}^{k} y_i \varphi_i} \xrightarrow{\quad n \quad} \hat{X}$$

$\varphi_1 \qquad \varphi_k \qquad\qquad \varphi_1 \qquad \varphi_k$

<span style="color:orange">orthonormal basis</span>

- **SVD (in LSA)**

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



$A \quad (mxn)$

$U' \quad (mxr)$, k — <span style="color:blue">latent semantic space</span>

$\Sigma' \quad (kxk), \ (rxr)$

$V'^T \quad (rxn)$, k — <span style="color:blue">latent semantic space</span>

$A' \quad (mxn)$

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

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

# 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 $x$ :



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

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

- Frameworks to circumvent vocabulary mismatch

**Doc** ➡ **terms** ➡ **structure model**

doc expansion

literal term matching

latent semantic
structure retrieval

query expansion

**Query** ➡ **terms** ➡ **structure model**

# 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 |
| *human* | 1 | 0 | 0 | 1 | 0 | 0 | 0 | 0 | 0 |
| *interface* | 1 | 0 | 1 | 0 | 0 | 0 | 0 | 0 | 0 |
| *computer* | 1 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| *user* | 0 | 1 | 1 | 0 | 1 | 0 | 0 | 0 | 0 |
| *system* | 0 | 1 | 1 | 2 | 0 | 0 | 0 | 0 | 0 |
| *response* | 0 | 1 | 0 | 0 | 1 | 0 | 0 | 0 | 0 |
| *time* | 0 | 1 | 0 | 0 | 1 | 0 | 0 | 0 | 0 |
| *EPS* | 0 | 0 | 1 | 1 | 0 | 0 | 0 | 0 | 0 |
| *survey* | 0 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 1 |
| *trees* | 0 | 0 | 0 | 0 | 0 | 1 | 1 | 1 | 0 |
| *graph* | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 1 | 1 |
| *minors* | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 1 |

## 2-D Plot of Terms and Docs from Example



Query: "human computer interaction"

An OOV word

FIG. 1. A two-dimensional plot of 12 Terms and 9 Documents from the sampe 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

$U^T U = I_{rXr}$

$V^T V = I_{rXr}$

$r \leq min(m,n)$

$K \leq r$

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

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

$A$ ($mxn$) = $U_{mxr}$ ($mxr$) $\Sigma_r$ ($rxr$) $V^T_{rxm}$ ($rxn$)

$A'$ ($mxn$) = $U'_{mxk}$ ($mxk$) $\Sigma_k$ ($kxk$) $V'^T_{kxm}$ ($kxn$)

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 $\Sigma_k$

# LSA Derivations (1/7)

- ## Singular Value Decomposition (SVD)
  - $A^TA$ is symmetric $n$x$n$ matrix
    - All eigenvalues $\lambda_j$ are nonnegative real numbers

    $$\lambda_1 \geq \lambda_2 \geq .... \geq \lambda_n \geq 0 \quad \Sigma^2 = diag(\lambda_1, \lambda_1, ..., \lambda_n)$$

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

    $$V_{n \times n} = \begin{bmatrix} v_1 v_2 ... v_n \end{bmatrix} \quad v_j^T v_j = 1 \quad (V^T V = I_{nxn})$$

    $$\text{sigma} \quad \sigma_j = \sqrt{\lambda_j}, \ j = 1, ..., n$$

    - Define **singular values:**
      - As the square roots of the eigenvalues of $A^TA$
      - As the lengths of the vectors $Av_1$, $Av_2$, ...., $Av_n$

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

$$\sigma_1 = \|Av_1\|$$
$$\sigma_2 = \|Av_2\|$$
.....

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

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

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

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

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

  - Define an **orthonormal** basis $\{u_1, u_2, \ldots, u_r\}$ 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 (mxr)

$V$: an orthonormal matrix (nxr)

  $$\Rightarrow \begin{bmatrix} u_1\ u_2 \ldots u_r \end{bmatrix} \Sigma_r = A \begin{bmatrix} v_1\ v_2 & v_r \end{bmatrix}$$

Known in advance

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

$$\Rightarrow \begin{bmatrix} u_1\ u_2 \ldots u_r \ldots u_m \end{bmatrix} \Sigma = A \begin{bmatrix} v_1\ v_2 \ldots v_r \ldots v_n \end{bmatrix}$$

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

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

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

$$\|A\|_F^2 = \sigma_1^2 + \sigma_2^2 + \ldots + \sigma_r^2 \quad \textcolor{red}{?}$$

# LSA Derivations (3/7)

$v_i$ spans the row space of $A$

$u_i$ spans the row space of $\mathbf{A}^T$

Multiplication by $A$

$A$

*mxn*

Row $A$

$v_1$ $\longrightarrow$ $\sigma_1 u_1$

$v_2$ $\longrightarrow$ $\sigma_2 u_2$ Col $A$ = Row $A^T$

$V_1$ $R^n$ $R^m$

$v_r$ $\sigma_r u_r$ $U_1$

$\mathbf{0}$ $\mathbf{0}$

$v_{r+1}$ $u_{r+1}$ $U_2$

$V_2$ Nul $A^T$

Nul $A$ $v_{n-1}$ $u_m$

$AX = 0$ $v_n$

**FIGURE 4** The four fundamental subspaces and the action of $A$.

$\mathbf{U}$ $\mathbf{V}^T$

$$U\Sigma V^T = (U_1 \quad U_2)\begin{pmatrix} \Sigma_1 & \mathbf{0} \\ \mathbf{0} & \mathbf{0} \end{pmatrix}\begin{pmatrix} V_1^T \\ V_2^T \end{pmatrix}$$

$$= U_1 \Sigma_1 V_1^T$$

$$= A V_1 V_1^T \qquad U\Sigma = AV$$

$$= A$$

# 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 $V$

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

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

    - 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 $U$

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

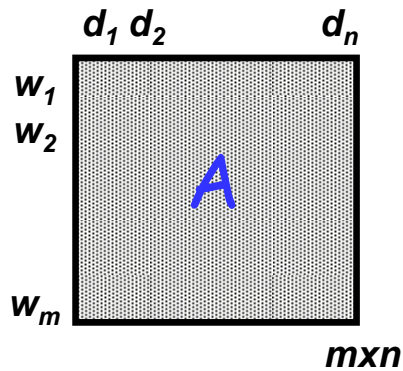    $$\Rightarrow A^T U = \left(U\Sigma V^T\right)^T U = V\Sigma U^T U = V\Sigma$$

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

    - 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 $U$

- ## Fundamental comparisons based on SVD

  - ### The original word-document matrix (A)

    $d_1$ $d_2$ ... $d_n$

    $w_1$
    $w_2$

    A

    $w_m$

    $mxn$

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

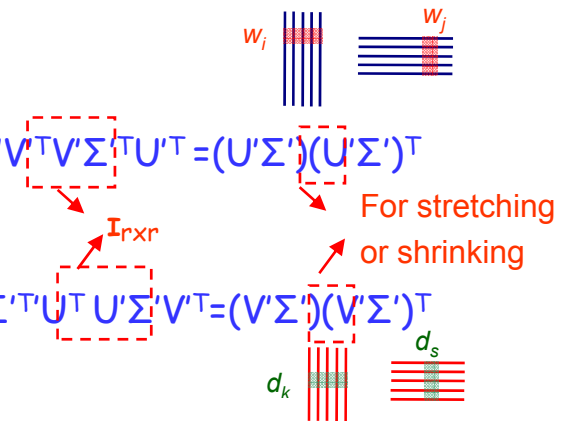  - ### The new word-document matrix (A')

    $U'=U_{m \times k}$
    $\Sigma'=\Sigma_k$
    $V'=V_{n \times k}$

    - compare two terms
      → dot product of two rows of $U'\Sigma'$

      $A'A'^\mathsf{T}=(U'\Sigma'V'^\mathsf{T})\,(U'\Sigma'V'^\mathsf{T})^\mathsf{T}=U'\Sigma'V'^\mathsf{T}V'\Sigma'^\mathsf{T}U'^\mathsf{T}=(U'\Sigma')(U'\Sigma')^\mathsf{T}$

      $I_{rxr}$

      For stretching or shrinking

    - compare two docs
      → dot product of two rows of $V'\Sigma'$

      $A'^\mathsf{T}A'=(U'\Sigma'V'^\mathsf{T})^\mathsf{T}\,(U'\Sigma'V'^\mathsf{T})=V'\Sigma'^\mathsf{T}U'^\mathsf{T}U'\Sigma'V'^\mathsf{T}=(V'\Sigma')(V'\Sigma')^\mathsf{T}$

    - compare a query word and a doc → each individual entry of A'

# LSA Derivations (6/7)

- **Fold-in**: find representations for pesudo-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

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

The separate dimensions are differentially weighted

Just like a row of V

Query represented by the weighted sum of it constituent term vectors

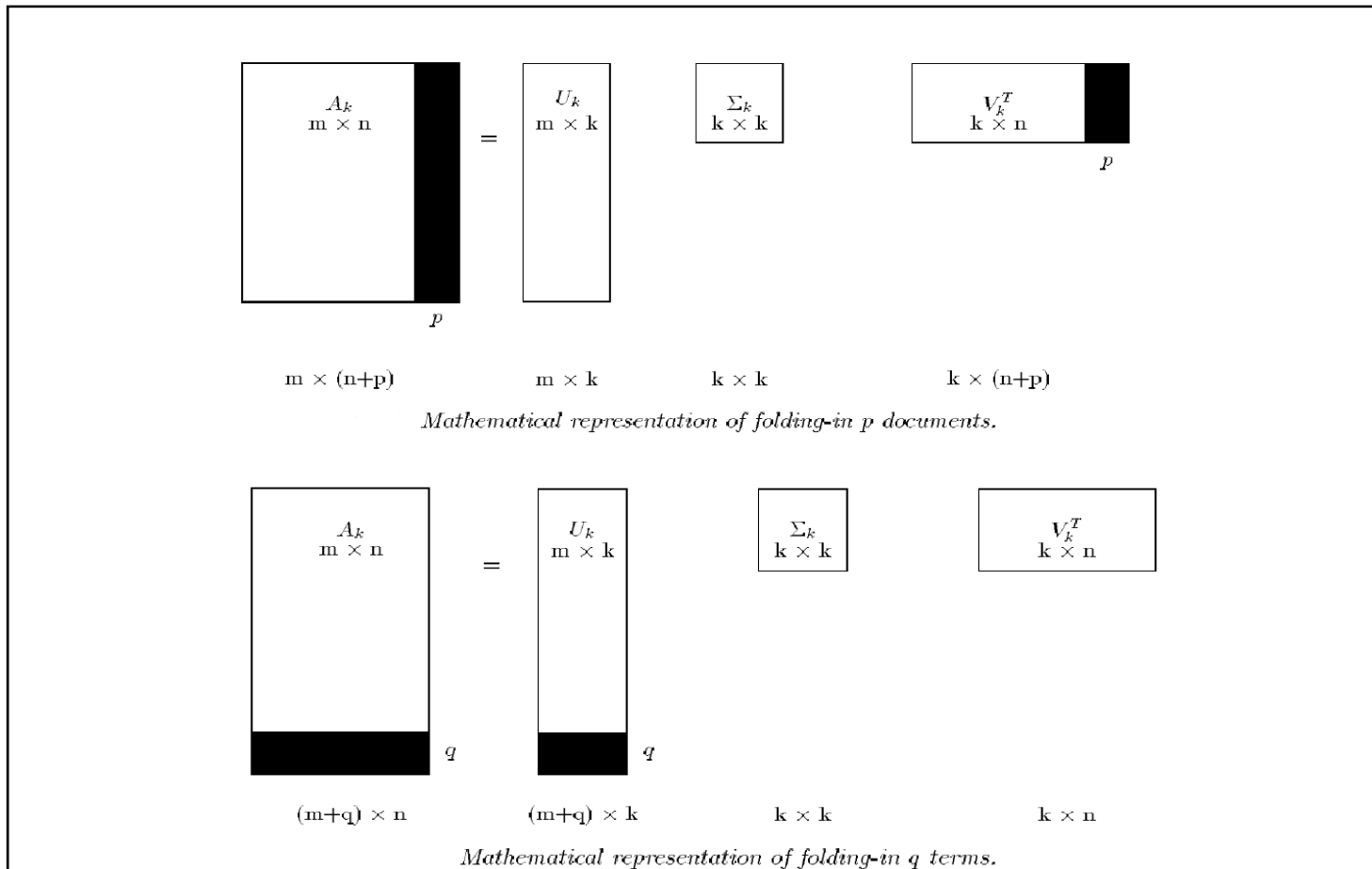    - Cosine measure between the query and doc vectors in the latent semantic space

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

row vectors

# LSA Derivations (7/7)

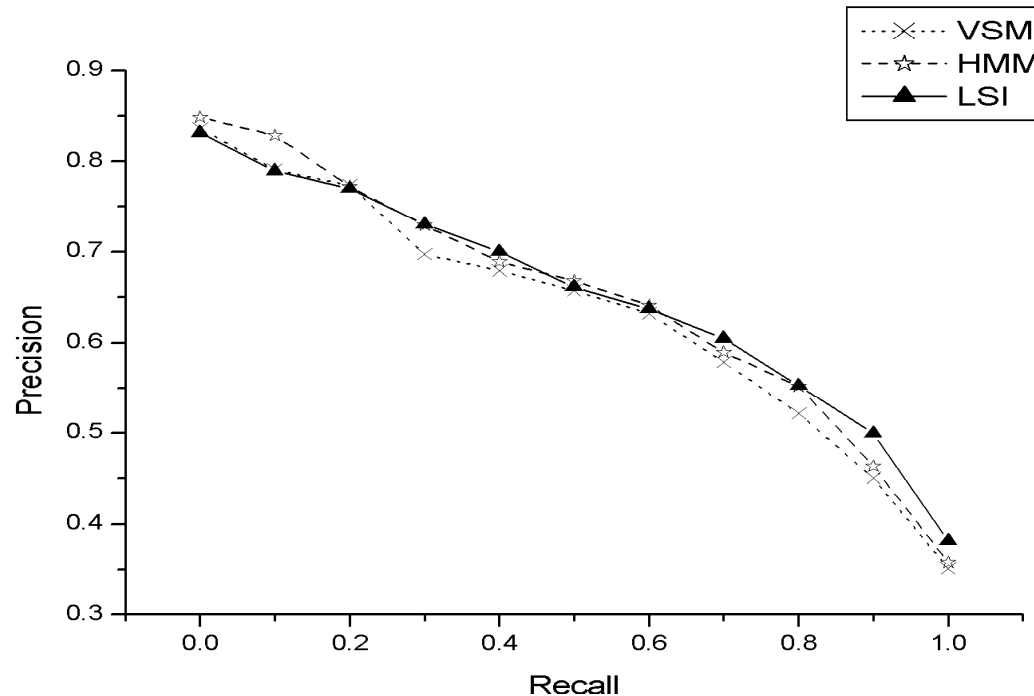- Fold-in a new 1 x n term vector

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



Mathematical representation of folding-in p documents.



Mathematical representation of folding-in q terms.

# 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,…
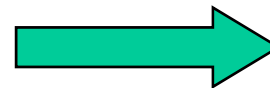
# LSA Toolkit: SVDLIBC (1/5)

- 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



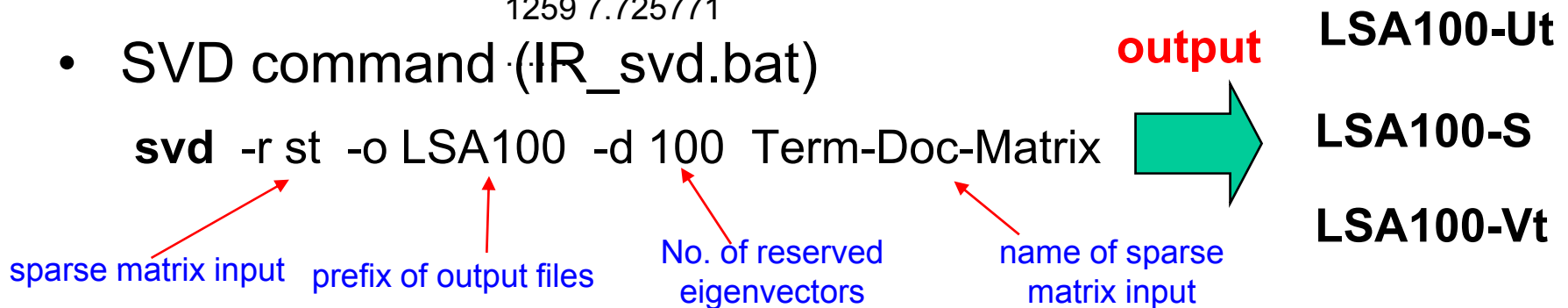  - Each entry is weighted by *TFxIDF* score

- **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

sparse matrix input

prefix of output files

No. of reserved eigenvectors

name of sparse matrix input

**output**

**LSA100-Ut**

**LSA100-S**

**LSA100-Vt**

# 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*x1 query vector

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

The separate dimensions are differentially weighted

Just like a row of V

Query represented by the weighted sum of it constituent term vectors

*TFxIDF* weighted beforehand

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

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