

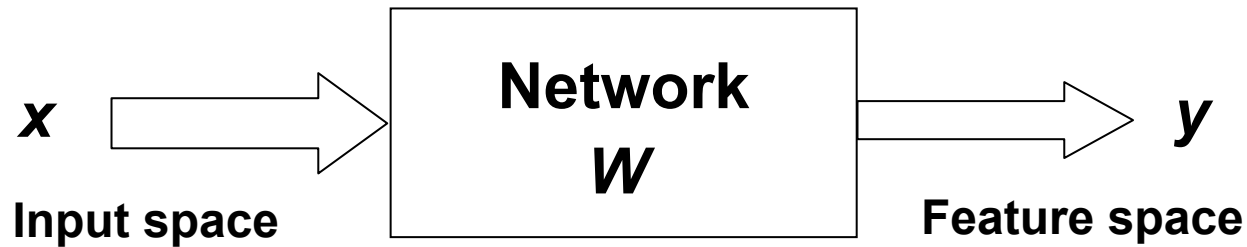
Discriminative Feature Extraction and Dimension Reduction

- PCA & LDA

Berlin Chen, 2004

Introduction

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



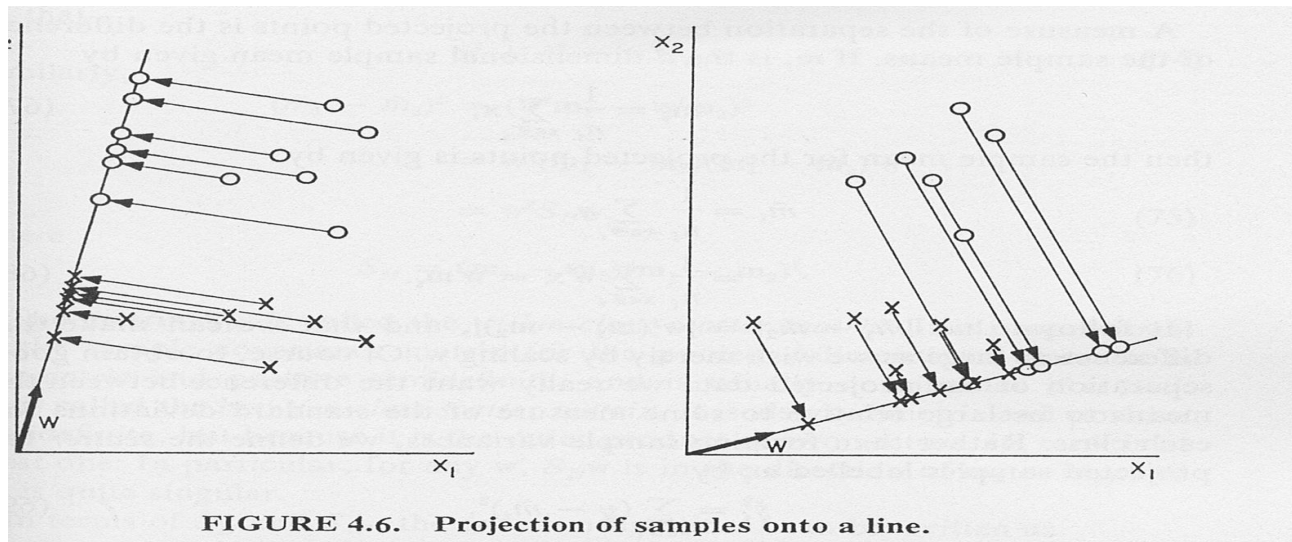
- Compute an input-output mapping based on some desirable properties

Introduction (cont.)

- Principal Component Analysis (PCA)
- Linear Discriminant Analysis (LDA)
- Heteroscedastic Discriminant Analysis (HDA)

Introduction (cont.)

- Formulation for discriminative feature extraction
 - Model-free (nonparametric)
 - Without prior information: e.g., PCA
 - With prior information: e.g., LDA
 - Model-dependent (parametric)
 - E.g., EM (Expectation-Maximization), MCE (Minimum Classification Error) Training



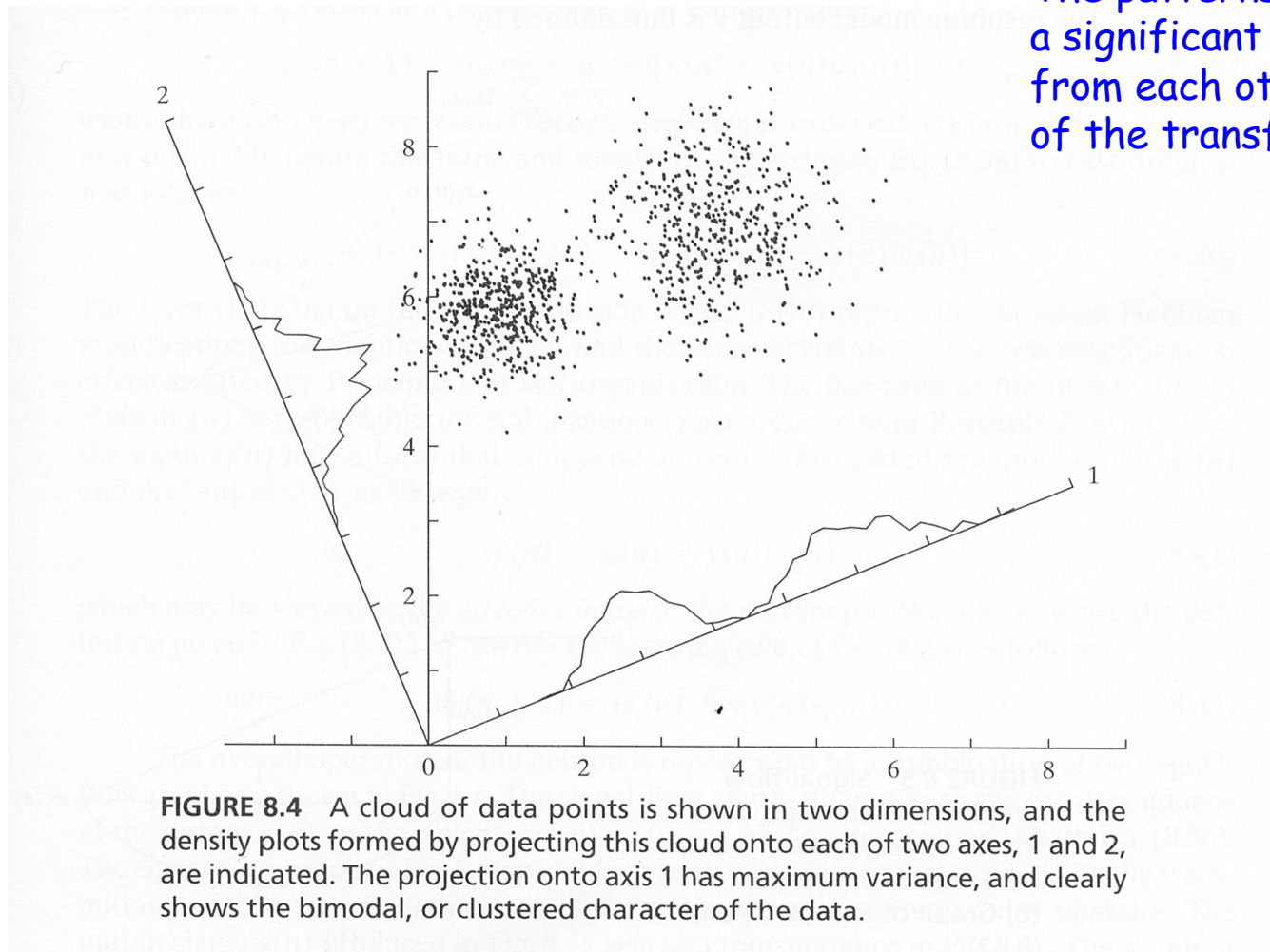
Principle Component Analysis (PCA)

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”

Principle Component Analysis (cont.)

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



Principle Component Analysis (cont.)

- Suppose \mathbf{x} is an n -dimensional zero mean random vector, $E_x \{ \mathbf{x} \} = \mathbf{0}$
 - 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

$$\mathbf{x} = \sum_{i=1}^n \underbrace{y_i}_{\text{The } i\text{-th component in the feature (mapped) space}} \boldsymbol{\varphi}_i = \boldsymbol{\Phi} \mathbf{y} \quad \text{where} \quad \mathbf{y} = [y_1 \quad \cdot \quad y_i \quad \cdot \quad y_n]^T$$
$$\boldsymbol{\Phi} = [\underbrace{\boldsymbol{\varphi}_1 \quad \cdot \quad \boldsymbol{\varphi}_i \quad \cdot \quad \boldsymbol{\varphi}_n}_{\text{The basis vectors}}]$$

Principle Component Analysis (cont.)

- 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_i also has the following properties

– Its mean is zero, too

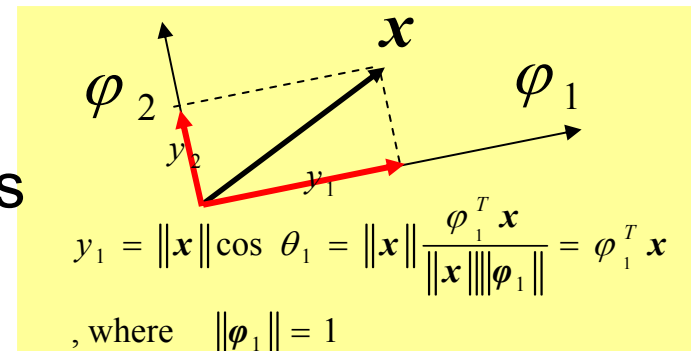
$$E\{y_i\} = E\{\varphi_i^T \mathbf{x}\} = \varphi_i^T E\{\mathbf{x}\} = \varphi_i^T \mathbf{0} = 0$$

– Its variance is

$$\sigma_i^2 = E\{y_i^2\} = E\{\varphi_i^T \mathbf{x} \mathbf{x}^T \varphi_i\} = \varphi_i^T E\{\mathbf{x} \mathbf{x}^T\} \varphi_i$$

$$\mathbf{R} = E\{\mathbf{x} \mathbf{x}^T\} = \frac{1}{N} \sum_i \mathbf{x}_i \mathbf{x}_i^T$$

$$= \varphi_i^T \mathbf{R} \varphi_i \quad \left[\mathbf{R} \text{ is the (auto-)correlation matrix of } \mathbf{x} \right]$$



Principle Component Analysis (cont.)

– 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

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

– Its variance is

$$\sigma_i^2 = E\{y_i^2\} = E\{\varphi_i^T \mathbf{x} \mathbf{x}^T \varphi_i\} = \varphi_i^T E\{\mathbf{x} \mathbf{x}^T\} \varphi_i$$

$$\mathbf{R} = E\{\mathbf{x} \mathbf{x}^T\} = \frac{1}{N} \sum_i \mathbf{x}_i \mathbf{x}_i^T$$

$$= \varphi_i^T \mathbf{R} \varphi_i \quad [\mathbf{R} \text{ is the (auto-)correlation matrix of } \mathbf{x}]$$

- The correlation between two projections y_i and y_j is

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

Principle Component Analysis (cont.)

- 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**

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

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

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

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

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

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

$$E \{ y_j \} = 0$$

$$\sigma_j^2 = E \{ y_j^2 \} - (E \{ y_j \})^2$$

$$= E \{ y_j^2 \}$$

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

We should discard the bases where the projections have lower variances

Principle Component Analysis (cont.)

- Minimum Mean-Squared Error Criterion

- If the orthonormal (basis) set $\boldsymbol{\varphi}_i$'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:

$$\mathbf{R} \boldsymbol{\varphi}_j = \lambda_j \boldsymbol{\varphi}_j$$

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

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

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

Principle Component Analysis (cont.)

- Minimum Mean-Squared Error Criterion
 - If the eigenvectors are retained associated with the m largest eigenvalues, the mean-squared error will be

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

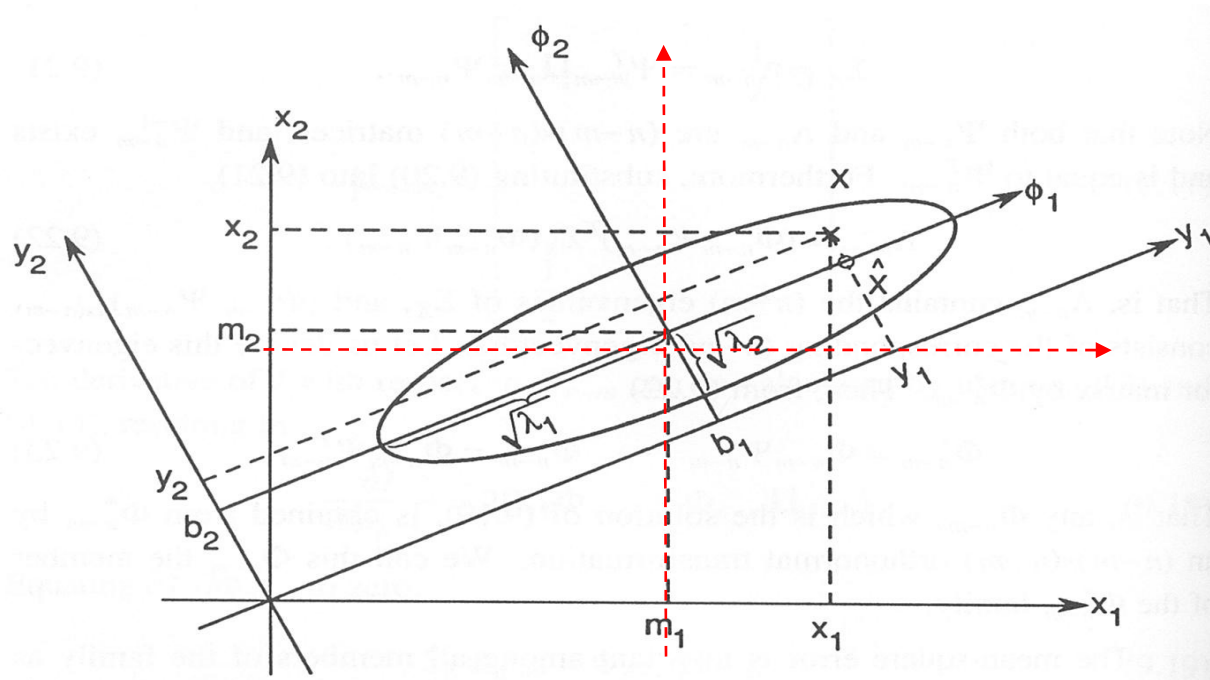
- Any two projections y_i and y_j will be mutually uncorrelated

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

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

Principle Component Analysis (cont.)

- An two-dimensional example of Principle Component Analysis



Principle Component Analysis (cont.)

- Minimum Mean-Squared Error Criterion

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

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

$$\frac{\partial \mathbf{x}^T \mathbf{y}}{\partial \mathbf{x}} = \mathbf{y}$$

To be minimized

constraints

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

Take derivation

$$\Rightarrow \forall_{m+1 \leq j \leq n} \frac{\partial J}{\partial \varphi_j} = 2 \mathbf{R} \varphi_j - \sum_{k=m+1}^n \mu_{jk} \varphi_k = 0 \quad \left(\text{where } \boldsymbol{\mu}_j^T = \left[\frac{1}{2} \mu_{j, m+1}, \dots, \frac{1}{2} \mu_{j, n} \right] \right)$$

$$\Rightarrow \forall_{m+1 \leq j \leq n} \mathbf{R} \varphi_j = \boldsymbol{\Phi}_{n-m} \boldsymbol{\mu}_j \quad \left(\text{where } \boldsymbol{\Phi}_{n-m} = [\varphi_{m+1} \dots \varphi_n] \right)$$

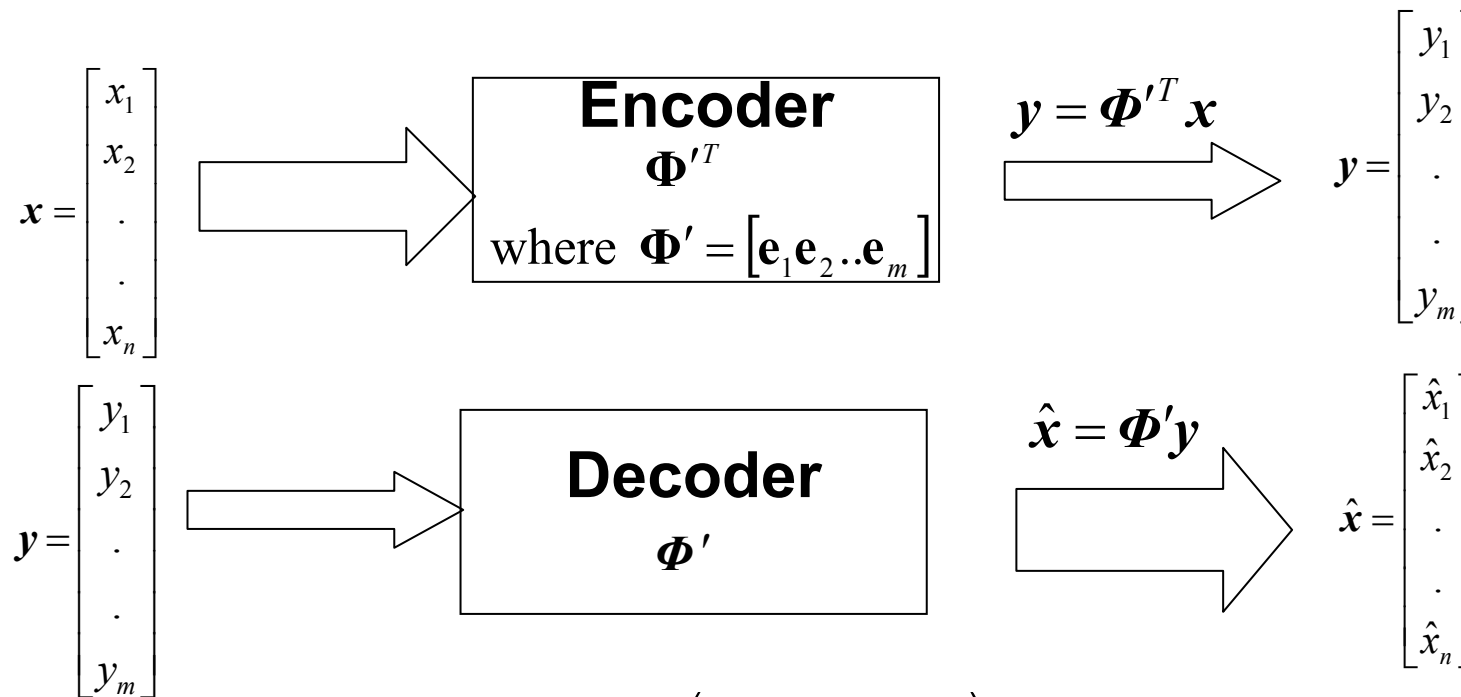
$$\Rightarrow \mathbf{R} [\varphi_{m+1} \dots \varphi_n] = \boldsymbol{\Phi}_{n-m} [\boldsymbol{\mu}_{m+1} \dots \boldsymbol{\mu}_n]$$

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

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

Principle Component Analysis (cont.)

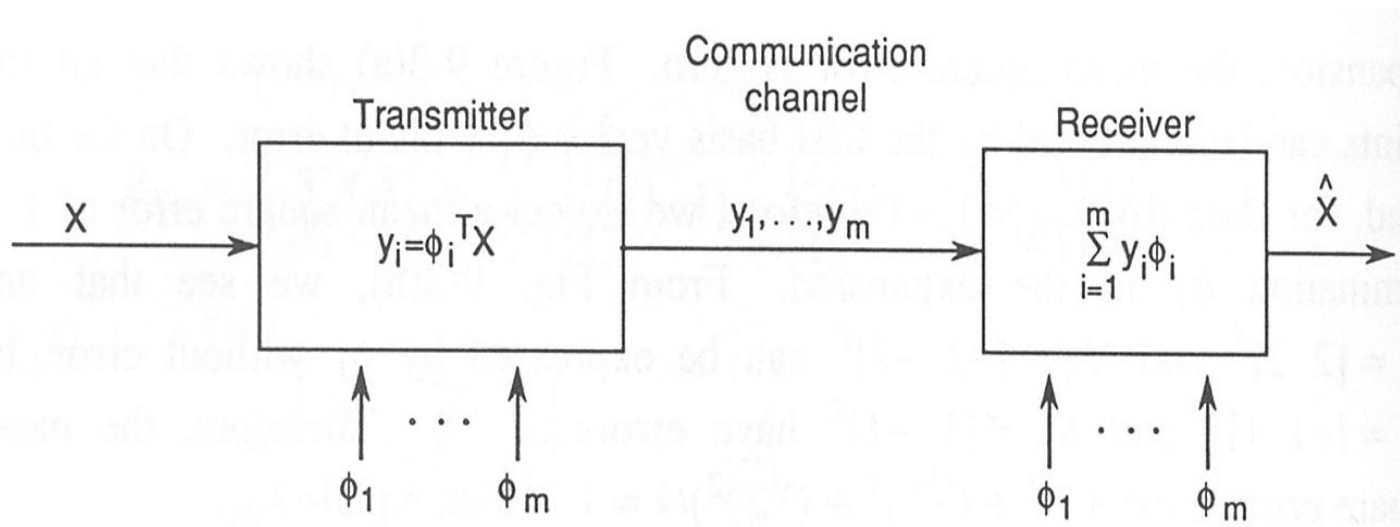
- Given an input vector \mathbf{x} with dimensional m
 - Try to construct a linear transform Φ' (Φ' is an $n \times m$ matrix $m < n$) such that the truncation result, $\Phi'^T \mathbf{x}$, is optimal in mean-squared error criterion



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

Principle Component Analysis (cont.)

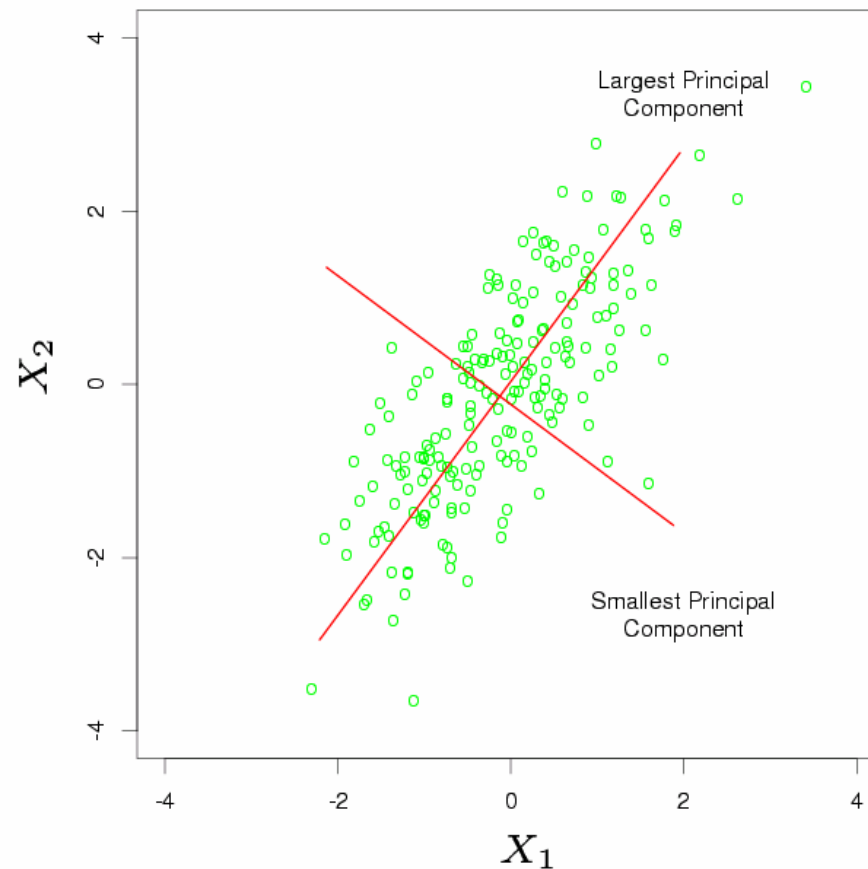
- Data compression in communication



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

Principle Component Analysis (cont.)

- Example 1: principal components of some data points



Principle Component Analysis (cont.)

- 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

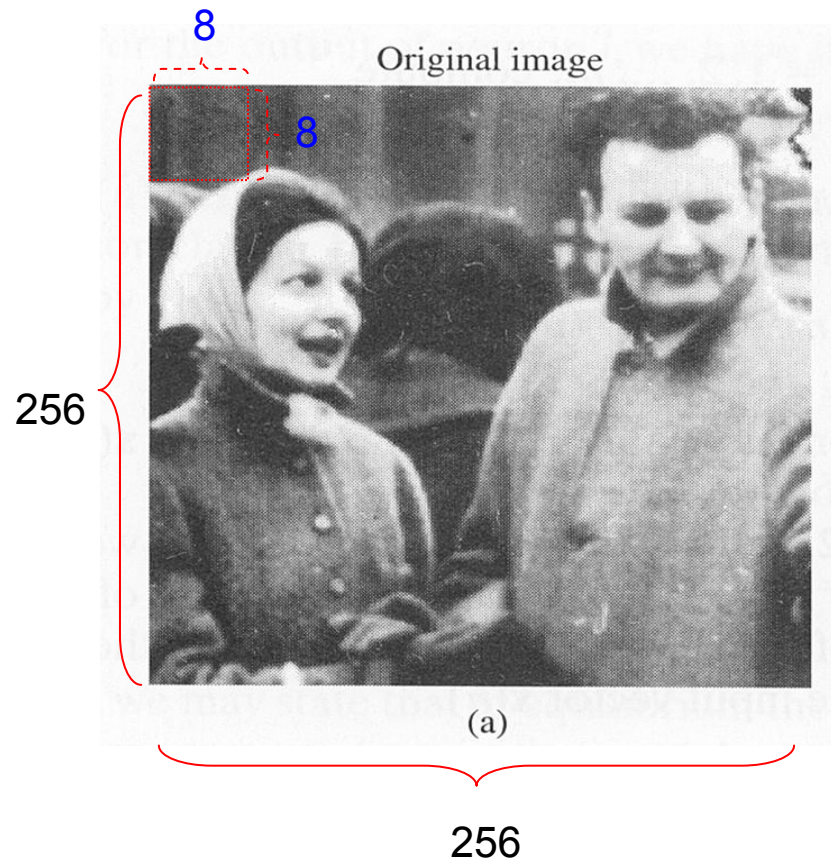
Feature	Eigenvalue
Feature 1	2.91082
Feature 2	0.92122
Feature 3	0.14735
Feature 4	0.02061

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

threshold for information content reserved

Principle Component Analysis (cont.)

- Example 3: Image Coding



Principle Component Analysis (cont.)

- Example 3: Image Coding (cont.)

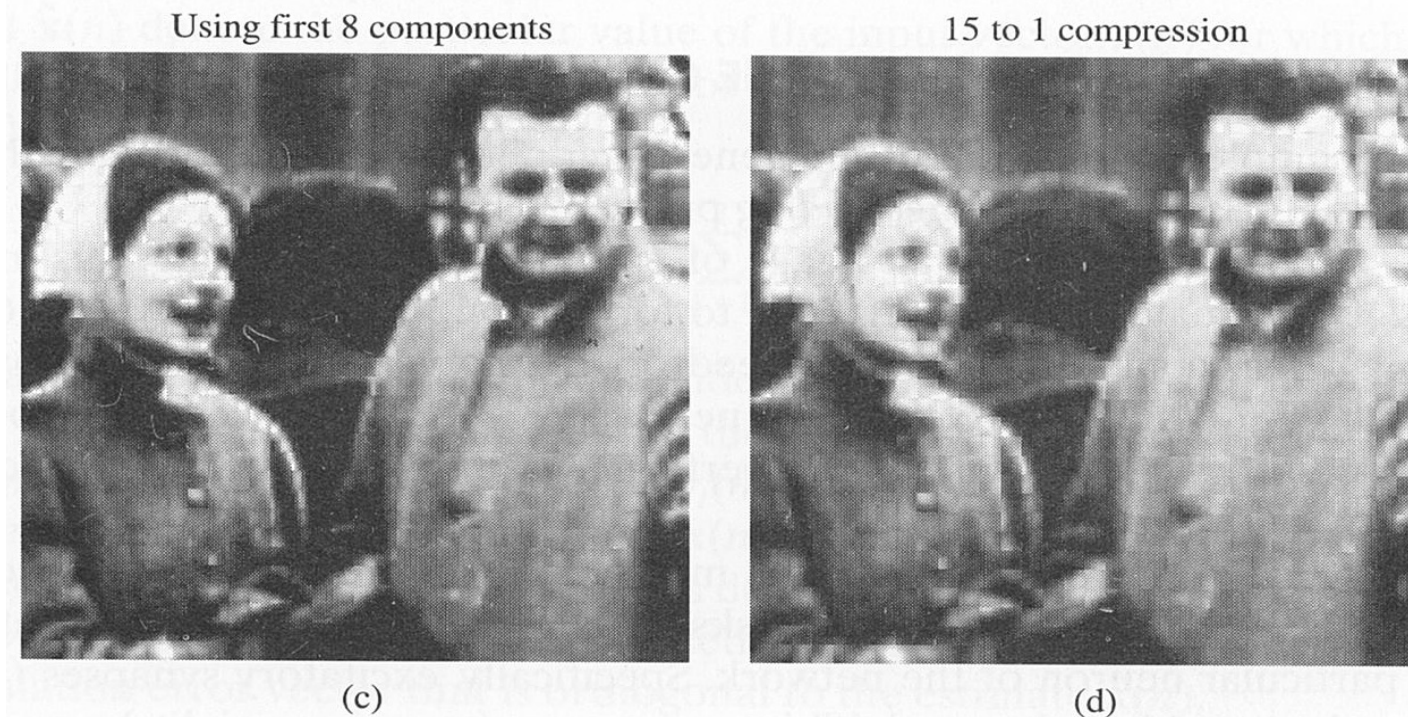


FIGURE 8.9 (a) An image of parents used in the image coding experiment. (b) 8×8 masks representing the synaptic weights learned by the GHA. (c) Reconstructed image of parents obtained using the dominant 8 principal components without quantization. (d) Reconstructed image of parents with 15 to 1 compression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

Principle Component Analysis

Eigenface and Eigenvoice

- Example 4: Eigenface in face recognition (1990)
 - Consider an individual image to be a linear combination of a small number of face components or “eigenface” 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

Principle Component Analysis

Eigenface and Eigenvoice (cont.)

- 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 themselves usually plausible faces, only directions of variations between faces

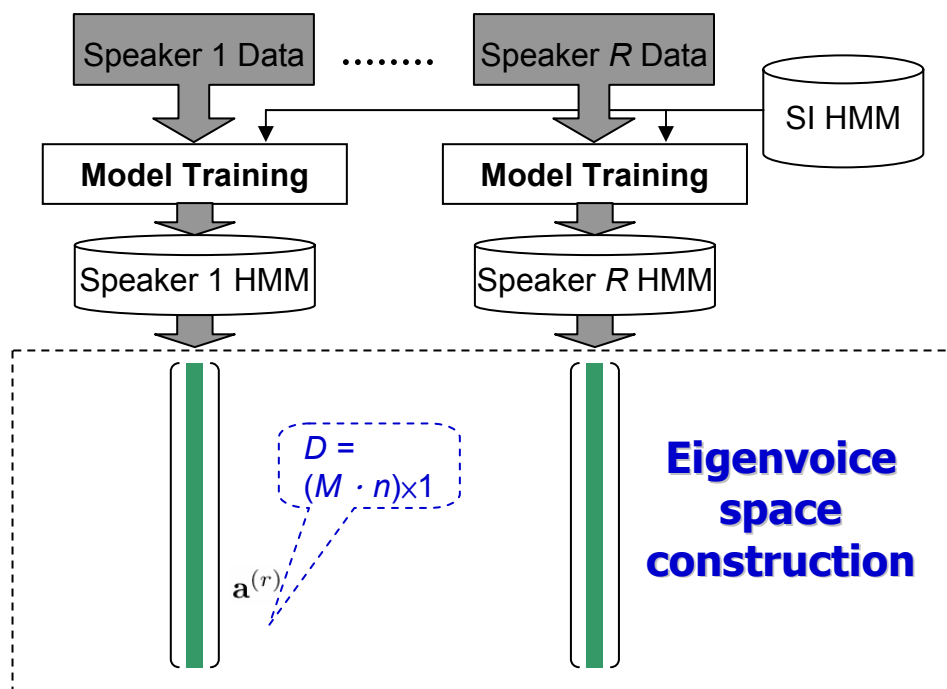
$$\hat{\mathbf{x}}_i = \bar{\mathbf{x}} + w_{i,1} \mathbf{e}(1) + w_{i,2} \mathbf{e}(2) + \dots + w_{i,K} \mathbf{e}(K)$$
$$\Rightarrow \mathbf{y}_i = [1, w_{i,1}, w_{i,2}, \dots, w_{i,K}]$$

Feature vector of a person i

Principle Component Analysis

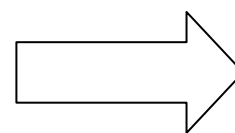
Eigenface and Eigenvoice (cont.)

- 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 model mean parameters (μ)



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

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



Principal Component Analysis

Principle Component Analysis

Eigenface and Eigenvoice (cont.)

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

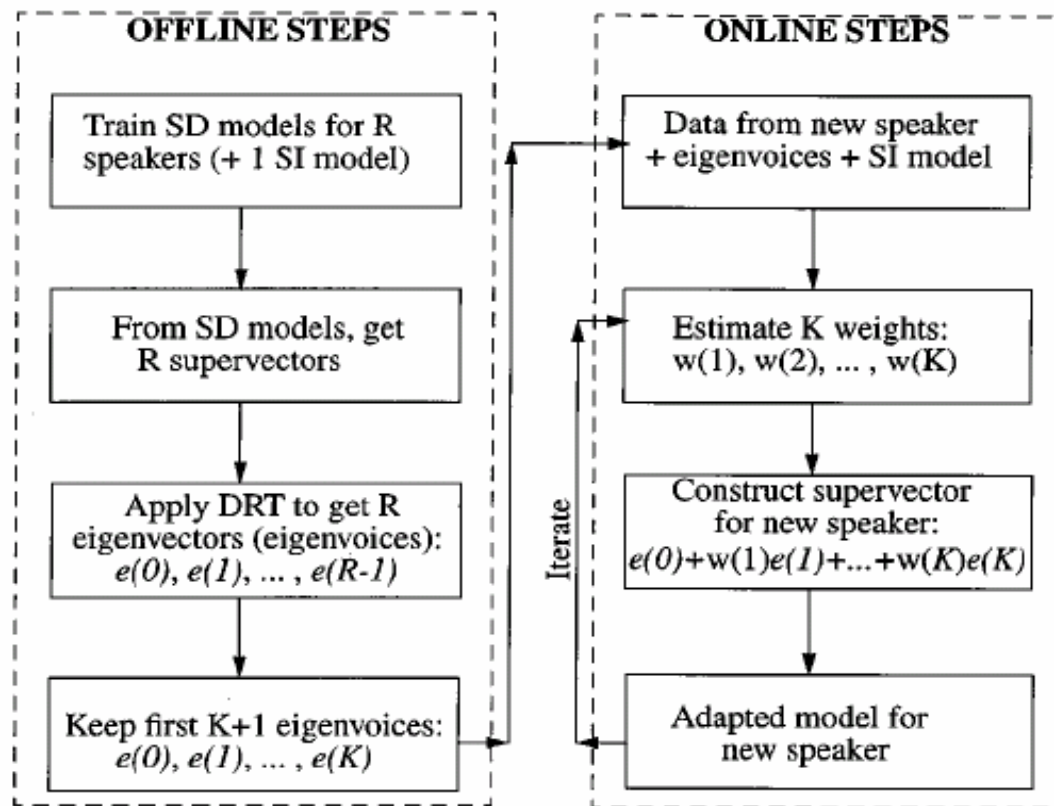


Fig. 1. Block diagram for eigenvoice speaker adaptation

Principle Component Analysis

Eigenface and Eigenvoice (cont.)

- 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

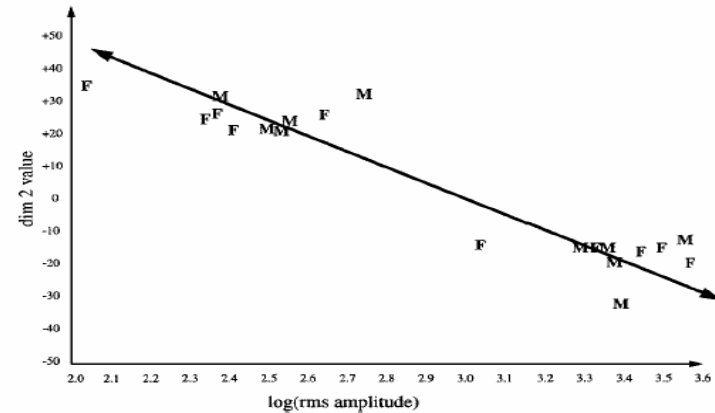


Fig. 3. Dimension 2 versus log(rms amplitude) for "A," extreme *M* and *F* in each speaker set

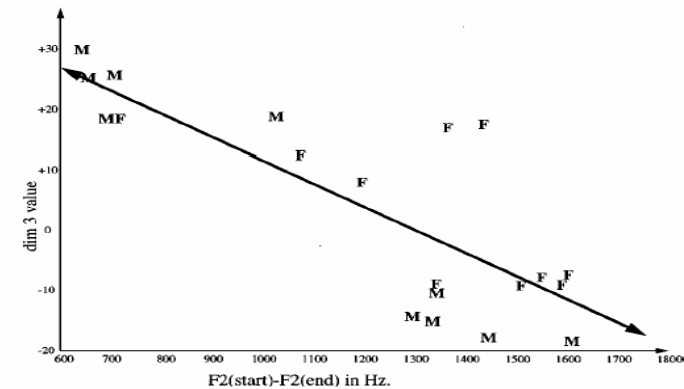


Fig. 4. Dimension 3 versus F2(start)-F2(end) for "U," extreme *M* and *F* in each speaker set

Linear Discriminant Analysis (LDA)

- 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 (cont.)

- 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

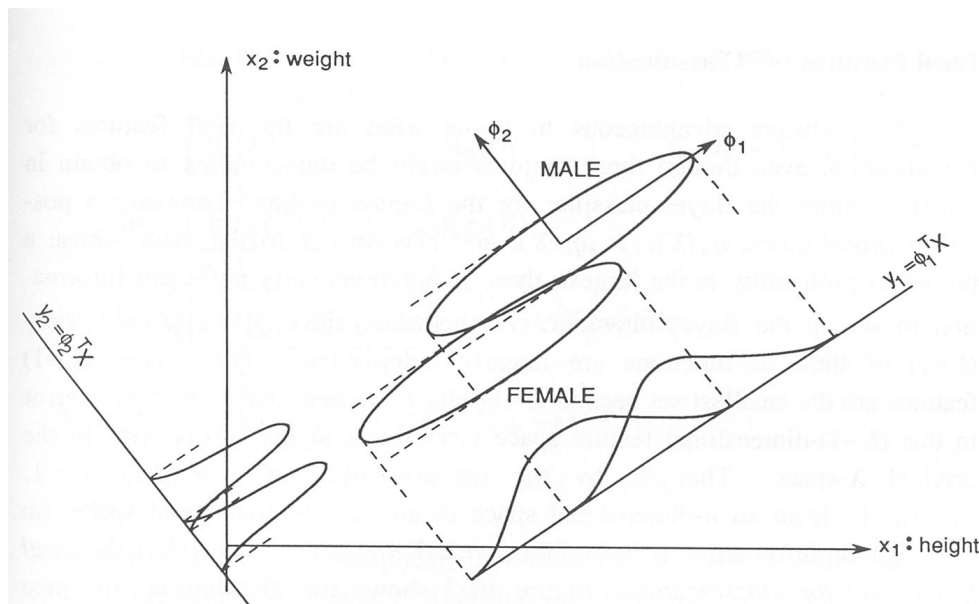


Fig. 10-1 An example of feature extraction for classification.

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

Linear Discriminant Analysis (cont.)

- 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

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

Linear Discriminant Analysis (cont.)

- 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

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

Linear Discriminant Analysis (cont.)

- 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 criterion 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 close 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 $\mathbf{S}_w^{-1} \mathbf{S}_b$

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

Linear Discriminant Analysis (cont.)

- Proof:

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

Or, for each column vector \mathbf{w}_i of \mathbf{W} , we want to find that :

The quadratic form has optimal solution : $\lambda_i = \frac{\mathbf{w}_i^T \mathbf{S}_b \mathbf{w}_i}{\mathbf{w}_i^T \mathbf{S}_w \mathbf{w}_i}$

$$\Rightarrow \frac{\partial \lambda_i}{\partial \mathbf{w}_i} = \frac{2\mathbf{S}_b \mathbf{w}_i (\mathbf{w}_i^T \mathbf{S}_w \mathbf{w}_i) - 2\mathbf{S}_w \mathbf{w}_i (\mathbf{w}_i^T \mathbf{S}_b \mathbf{w}_i)}{(\mathbf{w}_i^T \mathbf{S}_w \mathbf{w}_i)^2} = 0$$

$$\Rightarrow \frac{\mathbf{S}_b \mathbf{w}_i (\mathbf{w}_i^T \mathbf{S}_w \mathbf{w}_i)}{(\mathbf{w}_i^T \mathbf{S}_w \mathbf{w}_i)^2} - \frac{\mathbf{S}_w \mathbf{w}_i (\mathbf{w}_i^T \mathbf{S}_b \mathbf{w}_i)}{(\mathbf{w}_i^T \mathbf{S}_w \mathbf{w}_i)^2} = 0$$

$$\frac{\mathbf{S}_b \mathbf{w}_i}{\mathbf{w}_i^T \mathbf{S}_w \mathbf{w}_i} - \frac{\mathbf{S}_w \mathbf{w}_i}{\mathbf{w}_i^T \mathbf{S}_w \mathbf{w}_i} \lambda_i = 0 \quad \left(\because \lambda_i = \frac{\mathbf{w}_i^T \mathbf{S}_b \mathbf{w}_i}{\mathbf{w}_i^T \mathbf{S}_w \mathbf{w}_i} \right)$$

$$\Rightarrow \mathbf{S}_b \mathbf{w}_i - \lambda_i \mathbf{S}_w \mathbf{w}_i = 0 \Rightarrow \mathbf{S}_b \mathbf{w}_i = \lambda_i \mathbf{S}_w \mathbf{w}_i$$

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

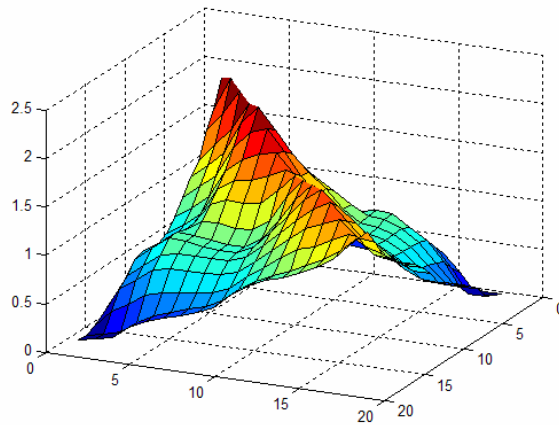
$$\left(\frac{\mathbf{F}}{\mathbf{G}} \right)' = \frac{\mathbf{F}'\mathbf{G} - \mathbf{G}'\mathbf{F}}{\mathbf{G}^2}$$

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

Linear Discriminant Analysis (cont.)

- Example 1: Experiments on Speech Signal Processing

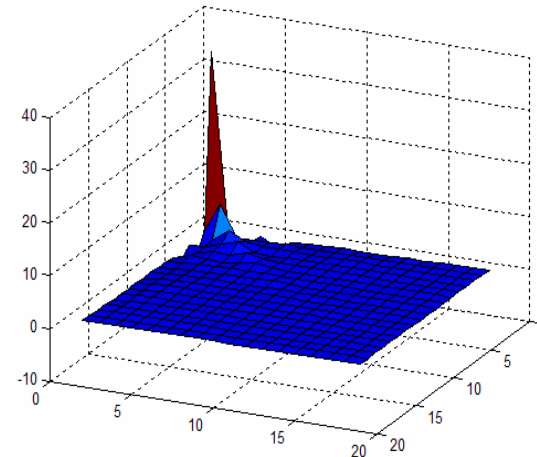
Covariance Matrix of the 18-Mel-filter-bank vectors



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

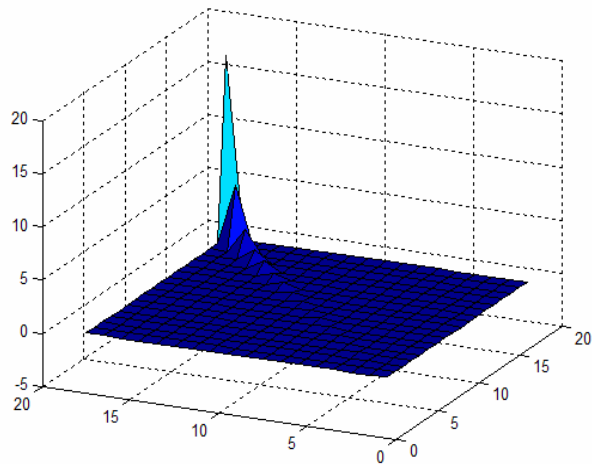
$$\Sigma' = \frac{1}{N} \sum_{\mathbf{y}_i} (\mathbf{y}_i - \bar{\mathbf{y}})(\mathbf{y}_i - \bar{\mathbf{y}})^T$$

After Cosine Transform

Experiments on Speech Signal Processing

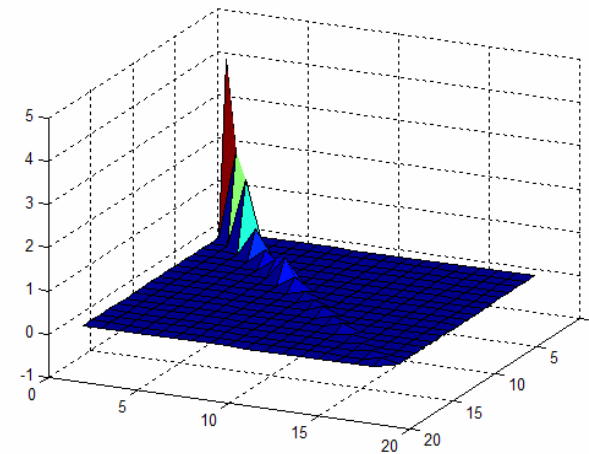
- Example1: Experiments on Speech Signal Processing (cont.)

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

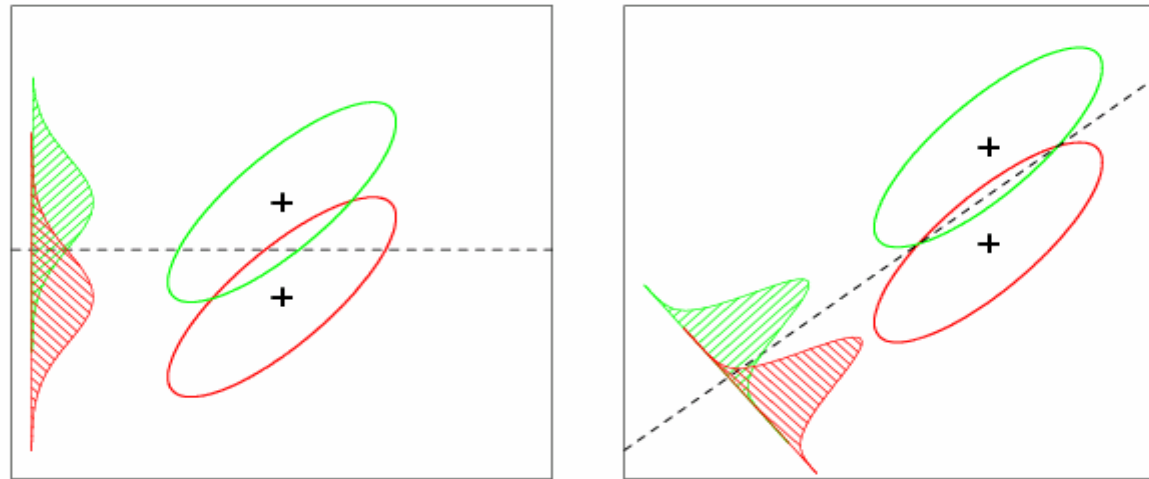


Figure 4.9: *Although the line joining the centroids defines the direction of greatest centroid spread, the projected data overlap because of the covariance (left panel). The discriminant direction minimizes this overlap for Gaussian data (right panel).*

LDA vs. HDA

- HDA: Heteroscedastic Discriminant Analysis

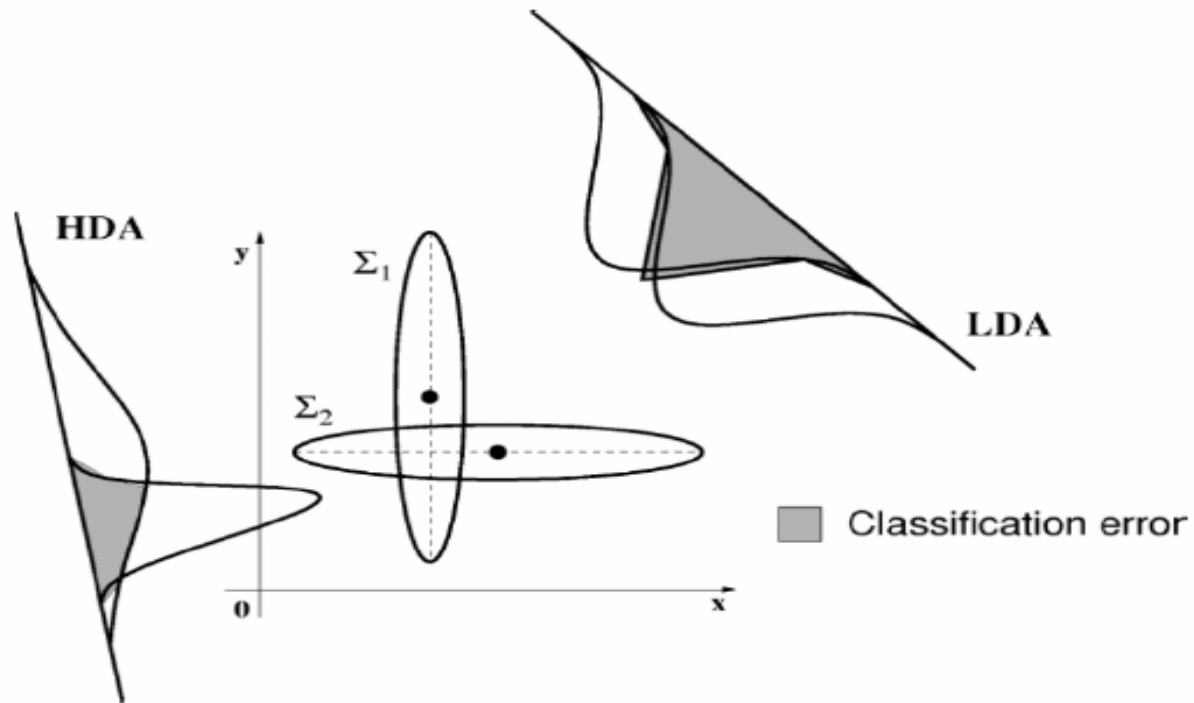


Fig. 1. Difference between LDA and HDA.

Heteroscedastic Discriminant Analysis (HDA) IBM, 2000

- Heteroscedastic : A set of statistical distributions having different variances
- LDA does not consider individual class covariances and may therefore generate suboptimal results
 - Modified the LDA objective function

$$H(\mathbf{W}) = \prod_{j=1}^J \left(\frac{|\mathbf{W}^T \mathbf{S}_b \mathbf{W}|}{|\mathbf{W}^T \mathbf{\Sigma}_j \mathbf{W}|} \right)^{N_j} = \frac{|\mathbf{W}^T \mathbf{S}_b \mathbf{W}|}{\prod_{j=1}^J |\mathbf{W}^T \mathbf{\Sigma}_j \mathbf{W}|^{N_j}}$$

- Take the log and rearrange terms

$$\log H(\mathbf{W}) = - \left(\sum_{j=1}^J N_j \log |\mathbf{W}^T \mathbf{\Sigma}_j \mathbf{W}| \right) + N \log |\mathbf{W}^T \mathbf{S}_b \mathbf{W}|$$

- However the dimensions of the HDA projection can often be highly correlated
 - An other transform can be further composed into HDA

Heteroscedastic Discriminant Analysis (cont.)

- The difference in the projections obtained from LDA and HDA for 2-class case

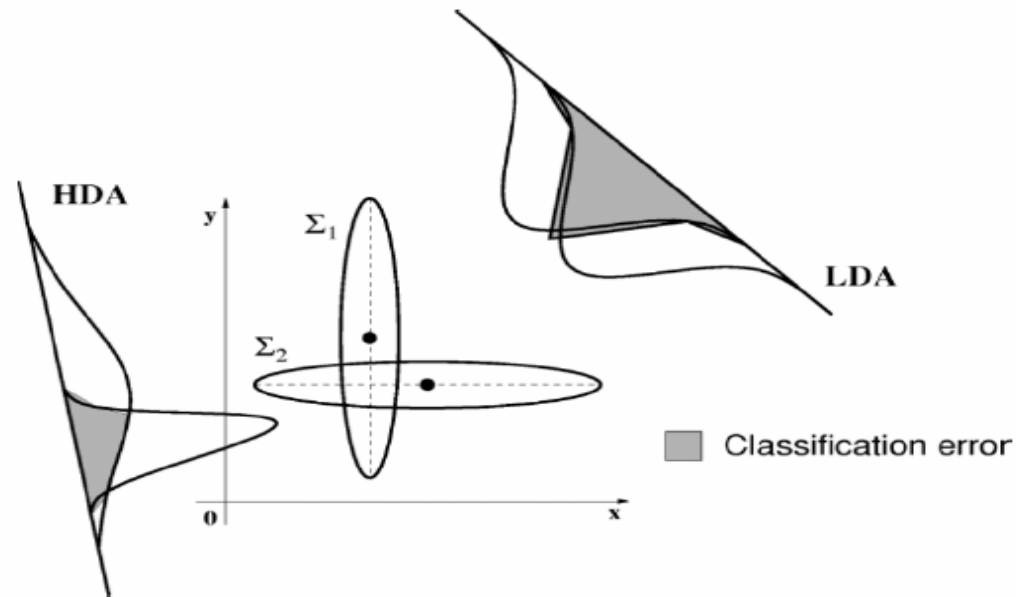


Fig. 1. Difference between LDA and HDA.

- 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