# Discriminative Feature Extraction and Dimension Reduction

Berlin Chen, 2002

# Introduction

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



$x$ → | **Network** $W$ | → $y$

**Input space** — **Feature space**

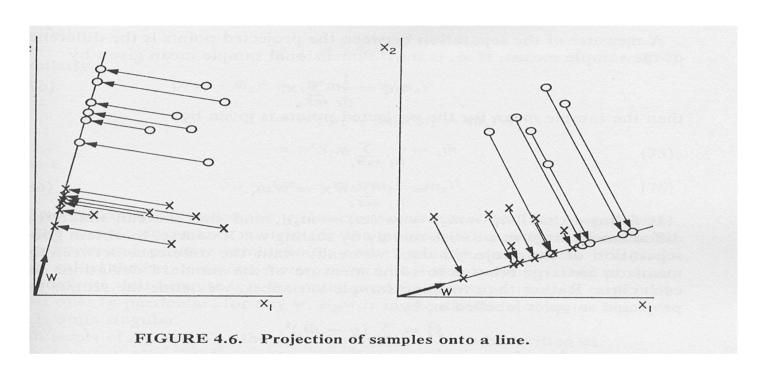  - Compute an input-output mapping based on some desirable properties

# Introduction

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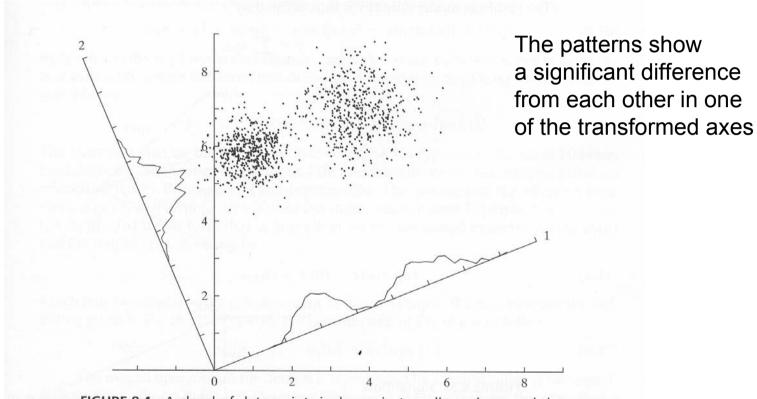- Heteroscedastic Discriminant Analysis (HDA)

# Introduction



FIGURE 4.6.   Projection of samples onto a line.

- Formulation
  - Model-free (nonparametric)
    - With/without prior information
  - Model-dependent (parametric)

# 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 (PCA)



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.

# Principle Component Analysis (PCA)

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

  - **x** can be represented without error by the summation of *n* linearly independent vectors

$$x = \sum_{i=i}^{n} \underbrace{y_i}\, \varphi_i = \Phi y \qquad \text{where} \quad y = \begin{bmatrix} y_1 & . & y_i & . & y_n \end{bmatrix}^T$$

$$\Phi = \underbrace{\begin{bmatrix} \varphi_1 & . & \varphi_i & . & \varphi_n \end{bmatrix}}$$

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

The basis vectors

# Principle Component Analysis (PCA)

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

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

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

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

- $y_i$ also has the following properties
  - Its mean is zero, too

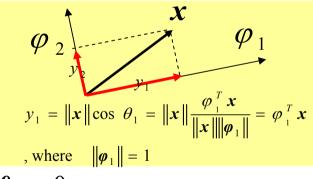$$E\{y_i\} = E\{\boldsymbol{\varphi}_i^T \boldsymbol{x}\} = \boldsymbol{\varphi}_i^T E\{\boldsymbol{x}\} = \boldsymbol{\varphi}_i^T \boldsymbol{0} = 0$$

  - Its variance is

$$\sigma_i^2 = E\{y_i^2\} = E\{\boldsymbol{\varphi}_i^T \boldsymbol{x}\boldsymbol{x}^T \boldsymbol{\varphi}_i\} = \boldsymbol{\varphi}_i^T E\{\boldsymbol{x}\boldsymbol{x}^T\}\boldsymbol{\varphi}_i$$

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

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

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

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

# Principle Component Analysis (PCA)

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

- $y_i$ also has the following properties

  – Its mean is zero, too

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

  – Its variance is

$$\sigma_i^2 = E\{y_i^2\} = E\{\varphi_i^T \boldsymbol{xx}^T \varphi_i\} = \varphi_i^T E\{\boldsymbol{xx}^T\}\varphi_i \qquad \boldsymbol{R} = E\{\boldsymbol{xx}^T\} = \frac{1}{N}\sum_i \boldsymbol{x}_i \boldsymbol{x}_i^T$$

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

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

$$E\{y_i y_j\} = E\left\{\left(\varphi_i^T \boldsymbol{x}\right)\left(\varphi_j^T \boldsymbol{x}\right)^T\right\} = E\{\varphi_i^T \boldsymbol{xx}^T \varphi_j\}$$

$$= \varphi_i^T E\{\boldsymbol{xx}^T\}\varphi_j = \varphi_i^T \boldsymbol{R} \varphi_j$$

# Principle Component Analysis (PCA)

- ## Minimum Mean-Squared Error Criterion
  - We want to choose only $m$ of $\varphi_i$'s that we still can approximate $x$ well in **mean-squared error criterion**

$$x = \sum_{i=1}^{n} y_i \varphi_i = \sum_{i=1}^{m} y_i \varphi_i + \sum_{j=m+1}^{n} y_j \varphi_j$$

$$\hat{x}(m) = \sum_{i=1}^{m} y_i \varphi_i$$

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

$$= E\left\{ \sum_{j=m+1}^{n} \sum_{k=m+1}^{n} y_j y_k \varphi_j^T \varphi_k \right\}$$

$$\boxed{\begin{array}{l} E\{y_j\} = 0 \\ \sigma_j^2 = E\{y_j^2\} - (E\{y_j\})^2 \\ \quad = E\{y_j^2\} \end{array}} \qquad = \sum_{j=m+1}^{n} E\{y_j^2\} \qquad \boxed{\because \varphi_j^T \varphi_k = \begin{cases} 1 & \text{if } j = k \\ 0 & \text{if } j \neq k \end{cases}}$$

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

We should discard the bases where the projections have lower variances

10

# Principle Component Analysis (PCA)

- ## Minimum Mean-Squared Error Criterion

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

    - They will have the property that:

$$\boldsymbol{R}\,\varphi_j = \lambda_j\,\varphi_j$$

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

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

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

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

# Principle Component Analysis (PCA)

- ## Minimum Mean-Squared Error Criterion

  - If the eigenvectors are retained associated with the *m* largest eigenvalues, the mean-squared error will be
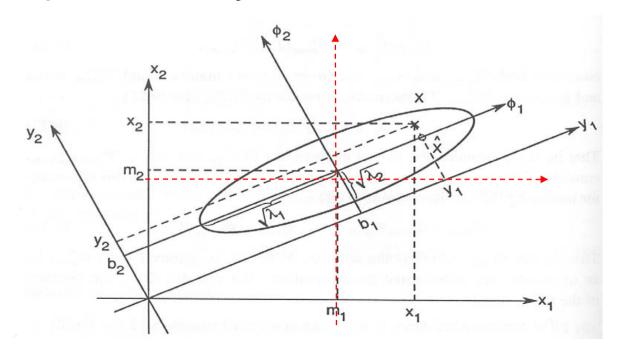
  $$\overline{\varepsilon}_{eigen}(m) = \sum_{j=m+1}^{n} \lambda_j \quad \left( \text{where } \lambda_1 \geq ... \geq \lambda_m \geq ... \geq \lambda_n \right)$$

  - 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
      - Gaussians and diagonal matrices

# Principle Component Analysis (PCA)

- An two-dimensional example of Principle Component Analysis

# Principle Component Analysis (PCA)

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

To be minimized     constraints

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

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

Take derivation

$$\Rightarrow \forall_{m+1 \le j \le n} \frac{\partial J}{\partial \varphi_j} = 2R\varphi_j - 2\sum\limits_{k=m+1}^{n} \mu_{jk}\varphi_k = 0 \quad \left(\text{where } \mu_j^T = \left[\mu_{j\ m+1}....\mu_{jn}\right]\right)$$

$$\Rightarrow \forall_{m+1 \le j \le n} \ R\varphi_j = \Phi_{n-m}\mu_j \quad \left(\text{where } \Phi_{n-m} = \left[\varphi_{m+1}....\varphi_n\right]\right)$$

$$\Rightarrow R\left[\varphi_{m+1}....\varphi_n\right] = \Phi_{n-m}\left[\mu_{m+1}.....\mu_n\right]$$

$$\Rightarrow R\Phi_{n-m} = \Phi_{n-m}U_{n-m} \quad \left(\text{where } U_{n-m} = \left[\mu_{m+1}....\mu_n\right]\right)$$

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

# Principle Component Analysis (PCA)

- Given an input vector $x$ with dimensional $m$
  - Try to construct a linear transform $\Phi'$ ($\Phi'$ is an nxm 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 $\Phi' = [e_1 e_1 .. e_l]$

$$y = \Phi'^T x$$

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

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

**Decoder**
$$\Phi'$$

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

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

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

# Principle Component Analysis (PCA)

- Data compression in communication



Communication channel

Transmitter     Receiver

$$X \quad \longrightarrow \quad y_i = \phi_i^T X \quad \longrightarrow \quad y_1, \ldots, y_m \quad \longrightarrow \quad \sum_{i=1}^{m} y_i \phi_i \quad \longrightarrow \quad \hat{X}$$

$$\phi_1 \quad \cdots \quad \phi_m \qquad\qquad \phi_1 \quad \cdots \quad \phi_m$$

- – 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 (PCA)
## Hebbian-based Maximum Eigenfilter

$$y(n) = \sum_{i=1}^{m} w_i(n) x_i(n)$$

$$w_i(n+1) = \frac{w_i(n) + \eta y(n) x_i(n)}{\left( \sum_{j=1}^{m} \left[ w_j(n) + \eta y(n) x_j(n) \right]^2 \right)^{\frac{1}{2}}}$$

$$w_i(n+1) \approx w_i(n) + \eta y(n) \underbrace{\left[ x_i(n) - y(n) w_i(n) \right]}_{x_i'(n)}$$

**It had been proved that**

$$\lim_{n \to \infty} w(n) \to \varphi_1 \text{ (the first principal component)}$$



(a)

(b)

**FIGURE 8.5** Signal-flow graph representation of maximum eigenfilter. (a) Graph of Eq. (8.36). (b) Graph of Eqs. (8.41) and (8.42).

# Principle Component Analysis (PCA)
## Hebbian-based Principal Analysis

- The Hebbian-based maximum eigenfilter can be expanded into a single layer feedforward network for principal component analysis (sanger, 1989)



$$y_j(n) = \sum_{i=1}^{m} w_{ji}(n) x_i(n), \quad j = 1, ..., J$$

$$\Delta w_{ji}(n) = \eta y_j(n) \underbrace{\left[ x_i(n) - \sum_{k=1}^{j} w_{ki}(n) y_k(n) \right]}_{x_i'(n)}$$

$$w_{ji}(n+1) = w_{ji}(n) + \Delta w_{ji}(n)$$

**It had been proved that**

$$\lim_{n \to \infty} \Delta w_j(n) \to 0$$

$$\lim_{n \to \infty} w_j(n) \to \varphi_j \text{ (the } j \text{ - th principal component)}$$

# Principle Component Analysis (PCA)
## Hebbian-based Principal Analysis

- Example: Image Coding



8x8
Non-overlapping
image block

8

256

256

Original image

(a)

Weights

(b)

# Principle Component Analysis (PCA)
## Hebbian-based Principal Analysis

- Example: Image Coding



**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 (PCA)
## Adaptive Principal Components Extraction

- Both feedward and lateral connections are used

$$y_j(n) = \boldsymbol{w}_j^T(n)\boldsymbol{x}(n) + \boldsymbol{a}_j^T \boldsymbol{y}_{j-1}(n)$$

$$\boldsymbol{w}_j(n+1) = \boldsymbol{w}_j(n) + \eta\left[y_j(n)\boldsymbol{x}(n) - y_j^2(n)\boldsymbol{w}_j(n)\right]$$

$$\boldsymbol{a}_j(n+1) = \boldsymbol{a}_j(n) - \eta\left[y_j(n)\boldsymbol{y}_{j-1}(n) + y_j^2(n)\boldsymbol{a}_j(n)\right]$$



**FIGURE 8.11** Network with feedforward and lateral connections for deriving the APEX algorithm.

21

# Principle Component Analysis (PCA)
## Eigenface and Eigenvoice

- 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
  - 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 eigenface from "eigenface 1" onwards model the variations from this average face

# Principle Component Analysis (PCA)
## Eigenface and Eigenvoice

- ## Eigenface in face recognition (1990)
  - ### 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

# Principle Component Analysis (PCA)
## Eigenface and Eigenvoice

- **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 ($\mu$)



Let each new speaker $S$ be represented by a point $P$ in $K$-space

$$P = e(0) + w(1) * e1 + \cdots + w(K) * e(K).$$

**Principal Component Analysis**

Speaker 1 Data ........ Speaker $R$ Data

SI HMM

**Model Training**   **Model Training**

Speaker 1 HMM   Speaker $R$ HMM

$D = (M \cdot n) \times 1$

$\mathbf{a}^{(r)}$

**Eigenvoice space construction**

# Principle Component Analysis (PCA)
## Eigenface and Eigenvoice

- Eigenvoice in speaker adaptation



Fig. 1. Block diagram for eigenvoice speaker adaptation

# Principle Component Analysis (PCA)
## Eigenface and Eigenvoice

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

# Linear Discriminant Analysis

- 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.

# Linear Discriminant Analysis (LDA)

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

  - The sample mean is: $\overline{\boldsymbol{x}} = \dfrac{1}{N} \sum\limits_{i=1}^{N} \boldsymbol{x}_i$

  - The class sample means are: $\overline{\boldsymbol{x}}_j = \dfrac{1}{N_j} \sum\limits_{g(\boldsymbol{x}_i)=j} \boldsymbol{x}_i$

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

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

  $$ \boldsymbol{S}_w = \dfrac{1}{N} \sum\limits_{j} N_j \boldsymbol{\Sigma}_j $$

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

  $$ \boldsymbol{S}_b = \dfrac{1}{N} \sum\limits_{j} N_j (\overline{\boldsymbol{x}}_j - \overline{\boldsymbol{x}})(\overline{\boldsymbol{x}}_j - \overline{\boldsymbol{x}})^T $$

# Linear Discriminant Analysis (LDA)

- 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

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

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

$$= W^T S_w W$$

# Linear Discriminant Analysis (LDA)

- If the transform $W = [w_1 \, w_2 \, .... \, w_m]$ is applied
  - The **average between-class variation** will be
  $$\widetilde{S}_b = W^T S_b W$$
  - Try to find optimal $W$ such that the following criterion function is maximized

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

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

    $$S_b w_i = \lambda_i S_w w_i$$

    - That is, $w_i$'s are the eigenvectors corresponding to the largest eigenvalues of

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

# Linear Discriminant Analysis (LDA)

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

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

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

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

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

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

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

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

# Heteroscedastic Discriminant Analysis (HDA)

- 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

$$H\left(\boldsymbol{W}\right) = \prod_{j=1}^{J}\left(\frac{\left|\boldsymbol{W}^{T}\boldsymbol{S}_{b}\boldsymbol{W}\right|}{\left|\boldsymbol{W}^{T}\boldsymbol{\Sigma}_{j}\boldsymbol{W}\right|}\right)^{Nj} = \frac{\left|\boldsymbol{W}^{T}\boldsymbol{S}_{b}\boldsymbol{W}\right|}{\prod_{j=1}^{J}\left|\boldsymbol{W}^{T}\boldsymbol{\Sigma}_{j}\boldsymbol{W}\right|^{Nj}}$$

  - Take the log and rearrange terms

$$\log H\left(\boldsymbol{W}\right) = -\left(\sum_{j=1}^{J}N_{j}\log\left|\boldsymbol{W}^{T}\boldsymbol{\Sigma}_{j}\boldsymbol{W}\right|\right) + N\log\left|\boldsymbol{W}^{T}\boldsymbol{S}_{b}\boldsymbol{W}\right|$$

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

# Heteroscedastic Discriminant Analysis (HDA)

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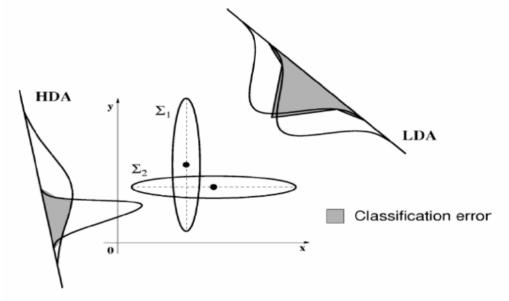


Fig. 1.  Difference between LDA and HDA.

- 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