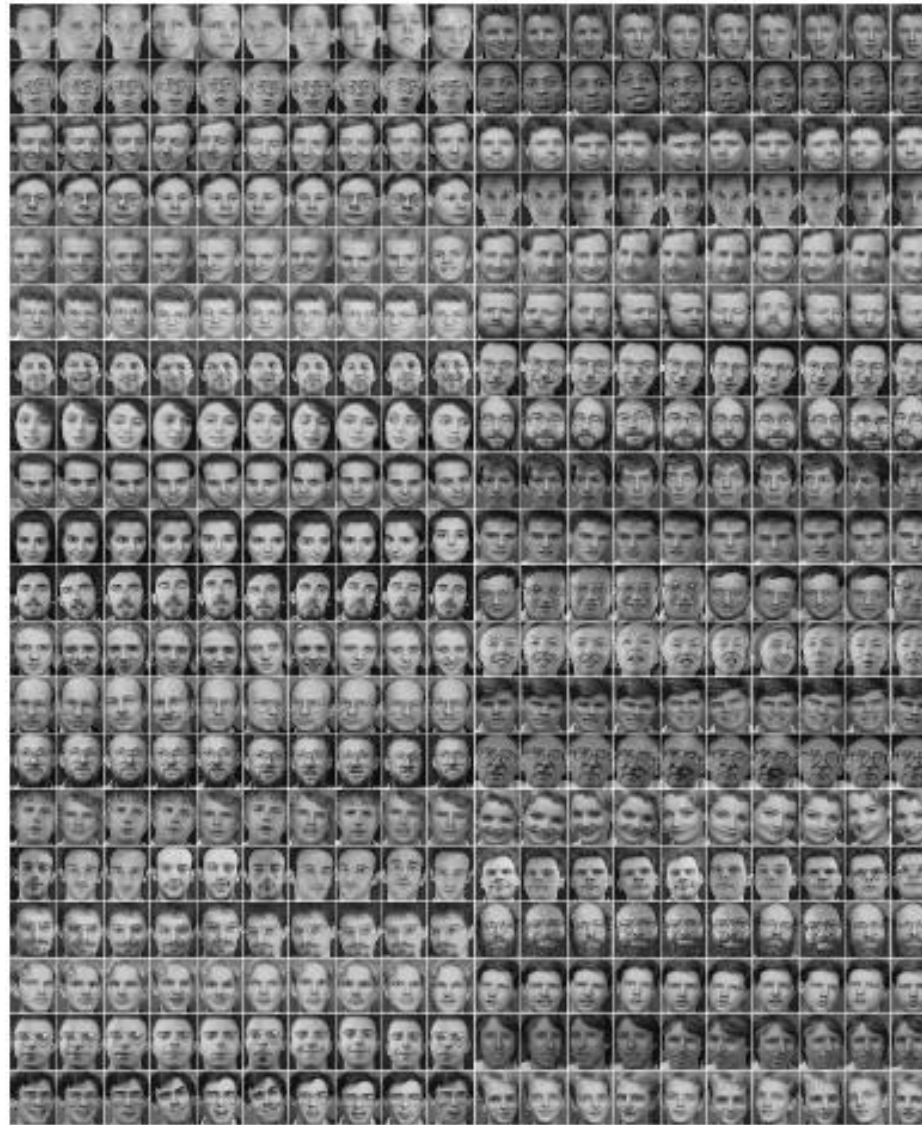# Face Recognition


## Arnaldo J. Abrantes

- **Goal**: Given a set of faces with known identity (training set) and a set of faces not identified, belonging to the same group of people (test set), we intend to identify each person in the test set

- **Difficulty**: The faces are immersed in a high dimensional space

- **Strategy**: Reduce the dimensionality of space through linear combination of features – Projection in sub-spaces

- **Two methods:**
  – *Eigenfaces* – PCA (*Principal Component Analysis*)
  – *Fisherfaces* – MDA (*Multiple Discriminant Analysis*)

- Let $\mathbf{a}, \mathbf{b} \in \mathfrak{R}^n$ be two vectors

$$\mathbf{a} = \left(a_1, \ldots, a_n\right)^T \qquad \mathbf{b} = \left(b_1, \ldots, b_n\right)^T$$



  - Vector norm

$$\|\mathbf{a}\| = \left(\mathbf{a}^T \mathbf{a}\right)^{1/2}$$

  - Inner product

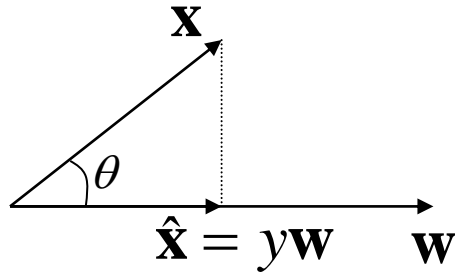$$< \mathbf{a}, \mathbf{b} > = \|\mathbf{a}\| \|\mathbf{b}\| \cos \theta = \mathbf{a}^T \mathbf{b}$$

  - Angle between two vectors

$$\theta = \cos^{-1} \frac{\mathbf{a}^T \mathbf{b}}{\left(\mathbf{a}^T \mathbf{a} \mathbf{b}^T \mathbf{b}\right)^{1/2}}$$

- Consider the orthogonal projection

$$\mathbf{w}^T\mathbf{x} = \|\mathbf{x}\|\|\mathbf{w}\|\cos\theta$$

$$\cos\theta = \frac{y\|\mathbf{w}\|}{\|\mathbf{x}\|}$$

$$y = \frac{\mathbf{w}^T\mathbf{x}}{\|\mathbf{w}\|^2}$$

$\hat{\mathbf{x}} = y\mathbf{w}$

- Hypothesis: orthonormal basis

$$\left(\mathbf{w}_1, \ldots, \mathbf{w}_m\right) \quad m \le n$$

$$\mathbf{w}_i^T\mathbf{w}_j = \begin{cases} 1 & se \quad i = j \\ 0 & se \quad i \ne j \end{cases}$$
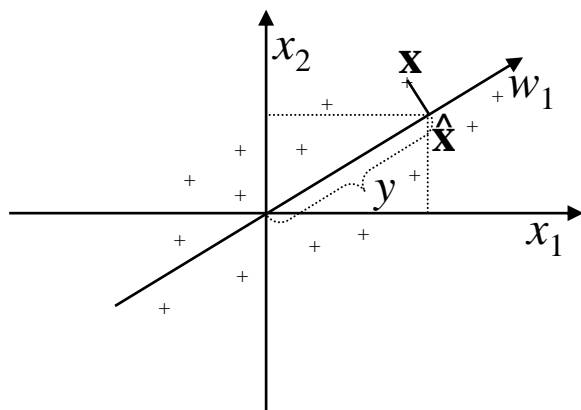
$$\hat{\mathbf{x}} = \sum_{i=1}^{m} y_i\mathbf{w}_i$$

$$y_i = \mathbf{w}_i^T\mathbf{x}$$

- Interpretation 1: $\mathbf{x} \in \mathfrak{R}^n \xrightarrow{\ W\ } \mathbf{y} \in \mathfrak{R}^m$ $\qquad W = \left(\mathbf{w}_1 \middle| \mathbf{w}_2 \middle| \dots \middle| \mathbf{w}_m\right)$

$$\mathbf{y} = \left(y_1, \dots, y_m\right)^T$$

$$\mathbf{y} = W^T \mathbf{x} \qquad\qquad n \; \mathrm{x} \; m$$

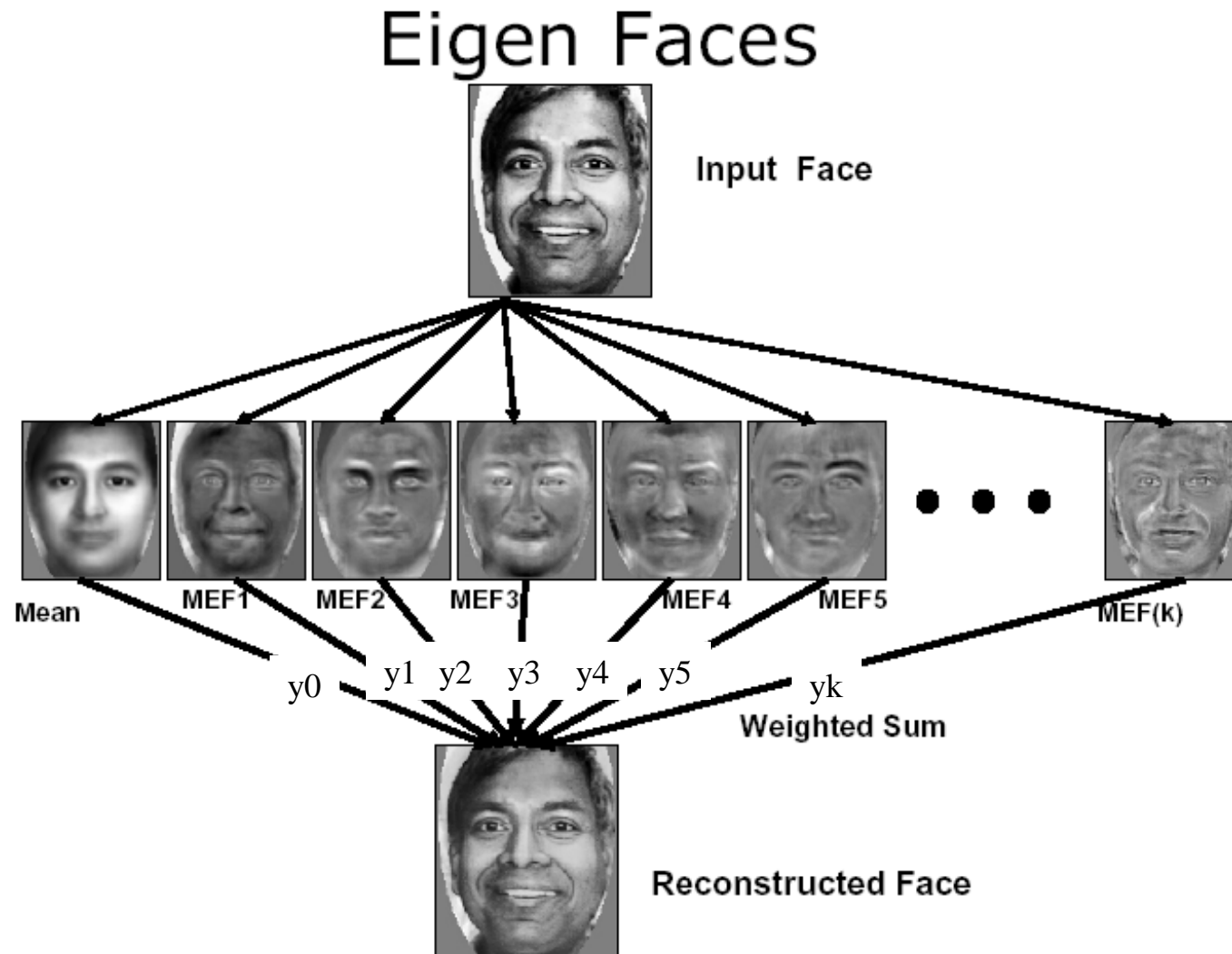- Interpretations 2: $\mathbf{x} \in \mathfrak{R}^n \xrightarrow{\ P\ } \hat{\mathbf{x}} \in \mathfrak{R}^n$ $\qquad\qquad P = WW^T$

$$\hat{\mathbf{x}} = y_1 \mathbf{w}_1 + \dots + y_m \mathbf{w}_m$$

$$= \mathbf{w}_1^T \mathbf{x} \mathbf{w}_1 + \dots + \mathbf{w}_m^T \mathbf{x} \mathbf{w}_m$$

$$= \left(\mathbf{w}_1 \mathbf{w}_1^T\right)\mathbf{x} + \dots + \left(\mathbf{w}_m \mathbf{w}_m^T\right)\mathbf{x}$$

$$= P\mathbf{x}$$

Projection operator:
- Idempotent $- P^2 = I$

$$\hat{\mathbf{x}} = WW^T \mathbf{x} \qquad\qquad \hat{\mathbf{x}} = W\mathbf{y}$$

Eigen Faces

Input Face

Mean  MEF1  MEF2  MEF3  MEF4  MEF5  MEF(k)

y0  y1  y2  y3  y4  y5  yk

Weighted Sum

Reconstructed Face

- How to choose the basis of representation?

$$W = (\mathbf{w}_1, \mathbf{w}_2, \ldots, \mathbf{w}_m)$$

- Consider having a training set with N vectors (faces)

$$\mathbf{x}_1, \ldots, \mathbf{x}_N$$

- Minimizing least squares criterion (PCA)

$$J_m = \sum_{k=1}^{N} \left\| \mathbf{x}_k - \hat{\mathbf{x}}_k \right\|^2$$

$$\hat{\mathbf{x}}_k = \boldsymbol{\mu} + \sum_{i=1}^{m} y_{ki} \mathbf{w}_i$$

$$\boldsymbol{\mu} = \frac{1}{N} \sum_{k=1}^{N} \mathbf{x}_k$$

Mean Face

- What are the "optimal" values for $y_{ki}$ and $\mathbf{w}_i$?

- Orthogonal projection

$$\frac{dJ_1}{dy_{k1}} = 2y_{k1} - 2\mathbf{w}_1^T(\mathbf{x}_k - \boldsymbol{\mu}) = 0$$

$$y_{k1} = \mathbf{w}_1^T(\mathbf{x}_k - \boldsymbol{\mu})$$

$$J_1 = \sum_{k=1}^{N} \left\| \mathbf{x}_k - (\boldsymbol{\mu} + y_{k1}\mathbf{w}_1) \right\|^2$$

$$= \sum_{k=1}^{N} \left\| y_{k1}\mathbf{w}_1 - (\mathbf{x}_k - \boldsymbol{\mu}) \right\|^2$$

$$= \sum_{k=1}^{N} y_{k1}^2 - 2\sum_{k=1}^{N} y_{k1}\mathbf{w}_1^T(\mathbf{x}_k - \boldsymbol{\mu}) + \sum_{k=1}^{N} \left\| \mathbf{x}_k - \boldsymbol{\mu} \right\|^2$$

- In what direction?

$$J_1 = \sum_{k=1}^{N} y_{k1}^2 - 2\sum_{k=1}^{N} y_{k1}^2 + \sum_{k=1}^{N} \left\| \mathbf{x}_k - \boldsymbol{\mu} \right\|^2$$

$$= -\sum_{k=1}^{N} \left[ \mathbf{w}_1^T(\mathbf{x}_k - \boldsymbol{\mu}) \right]^2 + \sum_{k=1}^{N} \left\| \mathbf{x}_k - \boldsymbol{\mu} \right\|^2$$

$$= -\sum_{k=1}^{N} \mathbf{w}_1^T(\mathbf{x}_k - \boldsymbol{\mu})(\mathbf{x}_k - \boldsymbol{\mu})^T \mathbf{w}_1 + \sum_{k=1}^{N} \left\| \mathbf{x}_k - \boldsymbol{\mu} \right\|^2$$

$$= -\mathbf{w}_1^T S \mathbf{w}_1 + \sum_{k=1}^{N} \left\| \mathbf{x}_k - \boldsymbol{\mu} \right\|^2 \quad \longleftarrow \quad \textbf{Minimize}$$

Scatter matrix:

$$S = \sum_{k=1}^{N} (\mathbf{x}_k - \boldsymbol{\mu})(\mathbf{x}_k - \boldsymbol{\mu})^T$$

$\uparrow$

$(n \times n)$

9

# Optimization problem with restrictions

- Minimize $J_1$ is the same as:
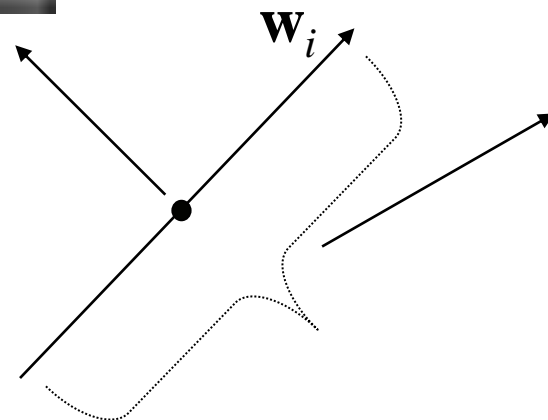
  - Maximize $$\mathbf{w}_1^T S \mathbf{w}_1$$

  - Subject to the restriction $\|\mathbf{w}_1\| = 1$

- Method of Lagrange multipliers $u = \mathbf{w}_1^T S \mathbf{w}_1 - \lambda\left(\mathbf{w}_1^T \mathbf{w}_1 - 1\right)$

$$\frac{\partial u}{\partial w_1} = 2S\mathbf{w}_1 - 2\lambda\mathbf{w}_1 = 0 \longrightarrow S\mathbf{w}_1 = \lambda\mathbf{w}_1$$

$$\mathbf{w}_1^T S \mathbf{w}_1 = \lambda \longleftarrow$$ Maximizing $u$ corresponds to choose the eigenvector whose eigenvalue is maximum

- How to generalize this reasoning to $m > 1$?

$$\mathbf{w}_i$$

- In general, $S$ matrix has a dimension of $n$, very high. However, $\text{rank}(S) \leq N - 1$ ⬋ (Usually much less than $n$)

- Solution

  – Define matrix $A$ as:
  $$A = \left(\mathbf{x}_1 - \boldsymbol{\mu} \middle| \mathbf{x}_2 - \boldsymbol{\mu} \middle| \cdots \middle| \mathbf{x}_N - \boldsymbol{\mu}\right)$$

  $$S = AA^T$$
  $$S\mathbf{w}_i = \lambda_i \mathbf{w}_i \qquad\qquad (n \text{ x } N)$$

  – Compute the eigenvector/eigenvalues of lowest dimension matrix $R$

  $$R = A^T A$$
  $$R\mathbf{v}_i = \varepsilon_i \mathbf{v}_i$$

  – Relate the two sets

  $$A^T A\mathbf{v}_i = \varepsilon_i \mathbf{v}_i$$
  $$AA^T A\mathbf{v}_i = \varepsilon_i A\mathbf{v}_i$$
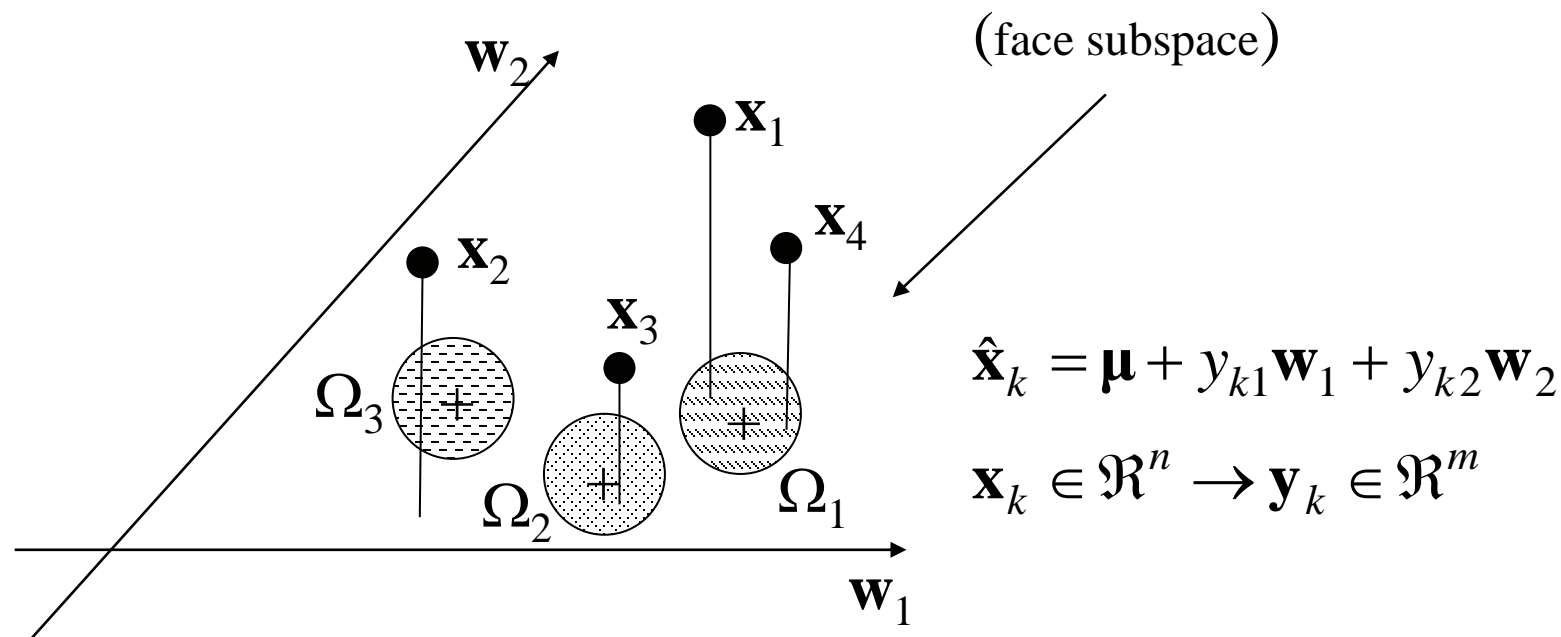  $$S(A\mathbf{v}_i) = \varepsilon_i (A\mathbf{v}_i)$$

  $$\begin{cases} \mathbf{w}_i = A\mathbf{v}_i \\ \lambda_i = \varepsilon_i \end{cases}$$

- Prepare the training set consisting of $N$ faces, $\mathbf{x}_1 \ldots \mathbf{x}_N$, properly aligned

- Determine the mean face $\boldsymbol{\mu}$

- Define matrix $A$ (size of $n \times N$) whose columns contain the "AC components" of the faces from the training set

- Compute the eigenvectors and eigenvalues of matrix $R = A^T A$ (size of $N \times N$)

- Select $m$ (maximum of $N$-1) eigenvectors from $R$, associated to the highest eigenvalues. Define matrix $V$ ($N \times m$), formed by the $m$ eigenvectors of $R$

- Get matrix $W$ ($n \times m$) by the relation $W = AV$ (not forgetting that the $W$ columns form an orthonormal basis)

- Develop a face classifier (e.g., nearest neighbor) based on models trained with the observations (feature vectors) of the projections in the face subspace, $\mathbf{y} = W^T(\mathbf{x}\text{-}\boldsymbol{\mu})$

- 2 *eigenfaces* – $\mathbf{w}_1, \mathbf{w}_2$
- 3 classes (known persons, $\Omega_1, .. \Omega_3$)
- 4 faces to classify $(\mathbf{x}_1, \ldots, \mathbf{x}_4)$



$$\hat{\mathbf{x}}_k = \boldsymbol{\mu} + y_{k1}\mathbf{w}_1 + y_{k2}\mathbf{w}_2$$

$$\mathbf{x}_k \in \Re^n \rightarrow \mathbf{y}_k \in \Re^m$$

# Application example: Face retrieval from a database



Query face

The 15 most similar faces in a universe of 7562 faces

- The PCA method is used to find a set of "optimal" directions in order to perform an efficient data **representation**;

- However, the directions that are useful to **represent** may not be the best for **discriminating**

<div align="right">

O Q O
Q O Q
Q O O

</div>

- The methods referred as Multi-Discriminant Analysis (MDA) try to find directions useful to perform an efficient discrimination

- Consider again the training set consisting of $N$ vectors (faces)

$$\mathbf{x}_1,\ldots,\mathbf{x}_N \qquad \mathbf{x}_i \in \mathfrak{R}^n$$

- However, consider now that the set is classified in 2 classes,
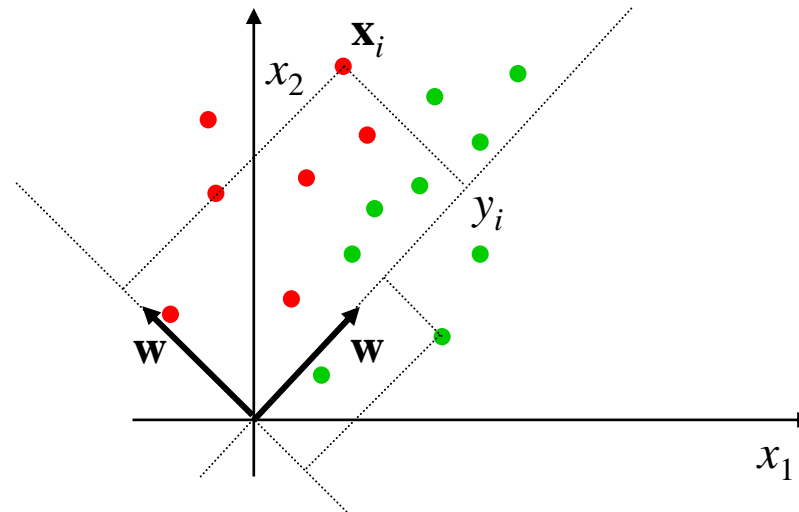  - $n_1$ samples belong to class $\Omega_1$
  - $n_2$ samples belong to class $\Omega_2$ 
  
  where $(n_1 + n_2 = N)$

- It is intended to project each of the $N$ vectors into a given straight line (direction), obtaining the corresponding samples1d, $(y_1,\ldots,y_N)$

$$y_i = \mathbf{w}^T \mathbf{x}_i$$

$$\mathfrak{R}^n \to \mathfrak{R}$$



- **Question**: What is the "best" direction that separate (discriminate) samples from the two classes?

- **First idea**: Maximize the difference between the sample means after the projection

    - Sample means of the two classes $\quad \mu_i = \dfrac{1}{n_i} \displaystyle\sum_{\mathbf{x} \in \Omega_i} \mathbf{x} \qquad i = 1,2$

    - Sample means after the projection $\quad \tilde{\mu}_i = \dfrac{1}{n_i} \displaystyle\sum_{y \in \Psi_i} y$

      $$= \dfrac{1}{n_i} \sum_{x \in \Omega_i} \mathbf{w}^T \mathbf{x} = \mathbf{w}^T \mu_i$$

- Criteria to optimize

    $$J(\mathbf{w}) = \left| \tilde{\mu}_1 - \tilde{\mu}_2 \right|$$
    $$= \left| \mathbf{w}^T \left( \mu_1 - \mu_2 \right) \right|$$

    - Does not work!

- **Idea**: It is not enough just to separate the means of the two clusters; it is also necessary that the scatter measure round the mean should be reduced

  - *Scatter measure* (variance)

  $$\tilde{s}_i^2 = \sum_{y \in \Psi_i} (y - \tilde{\mu}_i)^2 \qquad i = 1,2$$

- Fisher Criterion

  $$J(\mathbf{w}) = \frac{|\tilde{\mu}_1 - \tilde{\mu}_2|^2}{\tilde{s}_1^2 + \tilde{s}_2^2}$$

  - Expressed explicitly as a function of $\mathbf{w}$

  $$J(\mathbf{w}) = \frac{\mathbf{w}^T S_B \mathbf{w}}{\mathbf{w}^T S_w \mathbf{w}}$$

  - Scatter Matrices

  $$S_B = (\boldsymbol{\mu}_1 - \boldsymbol{\mu}_2)(\boldsymbol{\mu}_1 - \boldsymbol{\mu}_2)^T \quad \longleftarrow \quad \text{Inter-classes}$$

  $$S_i = \sum_{\mathbf{x} \in \Omega_i} (\mathbf{x} - \boldsymbol{\mu}_i)(\mathbf{x} - \boldsymbol{\mu}_i)^T$$

  $$S_w = S_1 + S_2 \qquad \qquad \longleftarrow \quad \text{Intra-classes}$$

- Cost function to maximize

$$J(\mathbf{w}) = \frac{\mathbf{w}^T S_b \mathbf{w}}{\mathbf{w}^T S_w \mathbf{w}}$$

scalar values

- Necessary condition

$$\frac{dJ}{d\mathbf{w}} = \frac{2 S_b \mathbf{w}(\mathbf{w}^T S_w \mathbf{w}) - 2 S_w \mathbf{w}(\mathbf{w}^T S_b \mathbf{w})}{(\mathbf{w}^T S_w \mathbf{w})^2} = 0$$
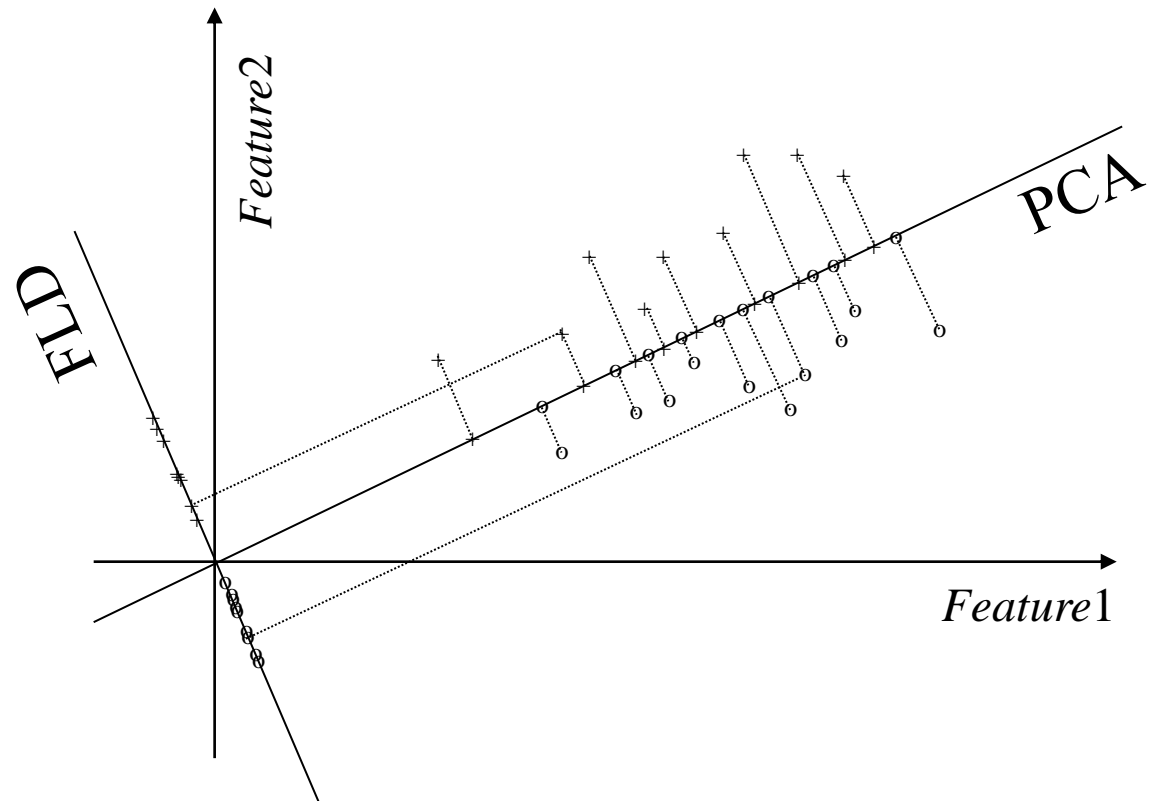
- That is, $\quad S_b \mathbf{w} = \lambda S_w \mathbf{w}$     Generalized eigenvalues/eigenvectors

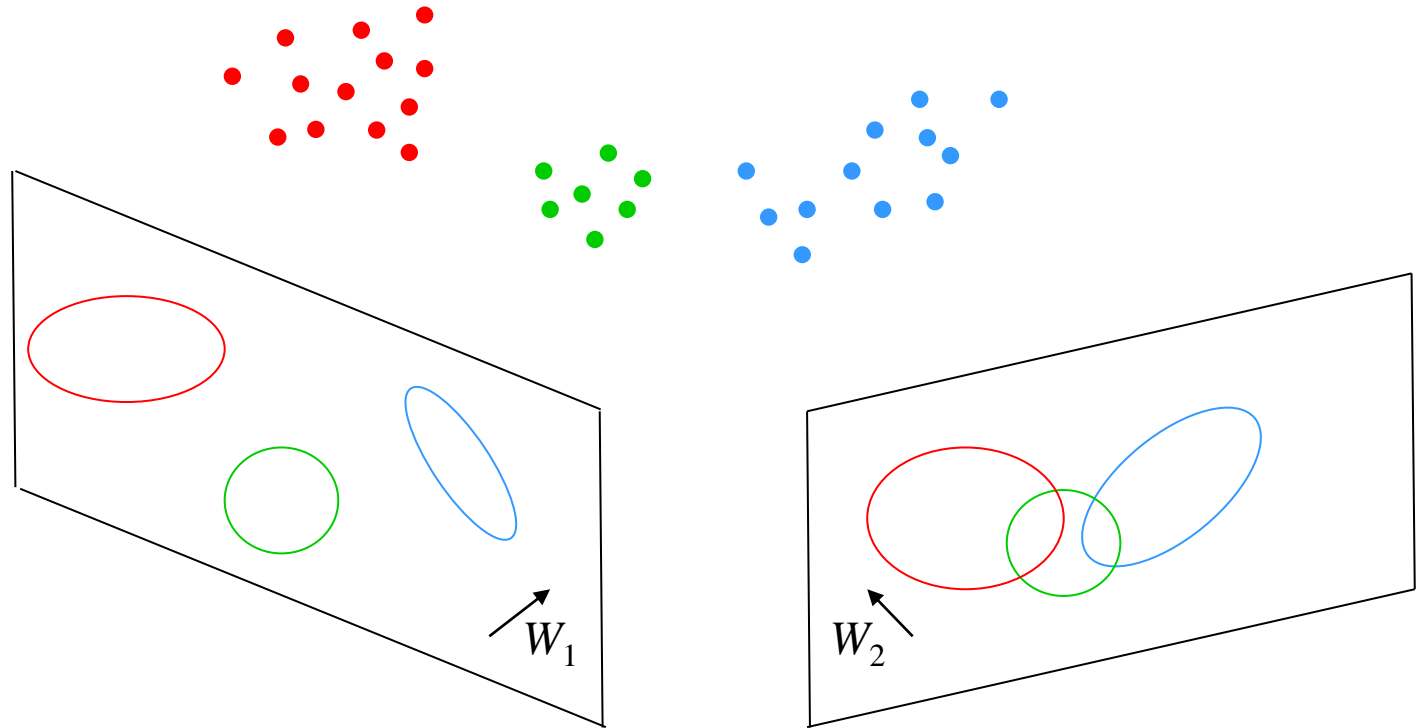- Note that $\quad S_b \mathbf{w} = \alpha(\boldsymbol{\mu}_1 - \boldsymbol{\mu}_2)$

- **Solution** $\quad\longrightarrow\quad \mathbf{w} = S_w^{-1}(\boldsymbol{\mu}_1 - \boldsymbol{\mu}_2)$     $S_w$ can not be singular!

- What is the best way to combine features?
- Example: Combine two features in one.
  - Criteria of minimizing the square error (PCA)
  - Criteria of the scatter ratio (inter and intra classes – FLD).

$W_1$

$W_2$

# Multi-Discriminant Analysis (MDA)

- Number of classes, $c$, is greater than 2
  - The aim is to obtain $c$-1 discriminant functions
  $$\Re^n \rightarrow \Re^{c-1}$$

- The total scatter is given by $\quad S_T = \sum_{\mathbf{x}} (\mathbf{x} - \mathbf{\mu})(\mathbf{x} - \mathbf{\mu})^T$

- $S_T$ can be decomposed into intra ($S_w$) and inter ($S_b$) scatter matrices

$$S_T = \sum_{i=1}^{c} \sum_{x \in \Omega_i} (\mathbf{x} - \mathbf{\mu}_i + \mathbf{\mu}_i - \mathbf{\mu})(\mathbf{x} - \mathbf{\mu}_i + \mathbf{\mu}_i - \mathbf{\mu})^T$$

$$= \sum_{i=1}^{c} \sum_{x \in \Omega_i} (\mathbf{x} - \mathbf{\mu}_i)(\mathbf{x} - \mathbf{\mu}_i)^T + \sum_{i=1}^{c} \sum_{x \in \Omega_i} (\mathbf{\mu}_i - \mathbf{\mu})(\mathbf{\mu}_i - \mathbf{\mu})^T$$

$$= S_w + \underbrace{\sum_{i=1}^{c} n_i (\mathbf{\mu}_i - \mathbf{\mu})(\mathbf{\mu}_i - \mathbf{\mu})^T}_{S_b}$$

$$S_w = \sum_{i=1}^{c} S_i$$

$$S_i = \sum_{\mathbf{x} \in \Omega_i} (\mathbf{x} - \mathbf{\mu}_i)(\mathbf{x} - \mathbf{\mu}_i)^T$$

- **Goal**: Determine the scatter matrix after the projection

$$\left(\Re^n \rightarrow \Re^{c-1}\right) \qquad\qquad W = \left(\mathbf{w}_1 \big| \mathbf{w}_2 \big| \dots \big| \mathbf{w}_{c-1}\right)$$

$$\mathbf{y} = W^T \mathbf{x} \qquad\qquad\qquad \uparrow$$

$$n \text{ x } c\text{-}1$$

- Define

$$\tilde{S}_w = \sum_{i=1}^{c} \sum_{\mathbf{y} \in \Psi_i} (\mathbf{y} - \tilde{\mathbf{\mu}}_i)(\mathbf{y} - \tilde{\mathbf{\mu}}_i)^T \qquad \tilde{\mathbf{\mu}}_i = \frac{1}{n_i} \sum_{\mathbf{y} \in \Psi_i} \mathbf{y}$$

$$\tilde{S}_b = \sum_{i=1}^{c} n_i (\tilde{\mathbf{\mu}}_i - \tilde{\mathbf{\mu}})(\tilde{\mathbf{\mu}}_i - \tilde{\mathbf{\mu}})^T \qquad \tilde{\mathbf{\mu}} = \frac{1}{N} \sum_{i=1}^{c} n_i \tilde{\mathbf{\mu}}_i$$

- The following relations are obtained,

$$\tilde{S}_w = W^T S_w W \qquad\qquad \tilde{S}_b = W^T S_b W$$

- Generalizing the above procedure (case $c = 2$), for the case of multiple classes, the following optimization criterion is obtained

$$J(W) = \frac{|W^T S_b W|}{|W^T S_w W|} \quad \longleftarrow \quad \text{Determinants ratio}$$

- Necessary condition, which must comply with the solution

$$S_b \mathbf{w}_i = \lambda S_w \mathbf{w}_i$$

- **Difficulty**
  - In the problem of the faces, the scatter matrix, $S_w$ is singular
  - At best, there are $N\text{-}c$ non-zero eigenvalues

$$\text{rank}(S_w) \leq N - c \ll n$$

- The solution called *fisherfaces* involves 2 steps:

1. Using the PCA method to obtain an intermediate subspace

$$\Re^n \rightarrow \Re^{N-c}$$

$$\mathbf{y}_{pca} = W_{pca}^T \mathbf{x} \qquad W_{pca} = \arg \max_W \left| W^T S_T W \right|$$

2. Then, uses the MDA method to obtain the $c$-1 discriminant directions

$$\Re^{N-c} \rightarrow \Re^{c-1}$$

$$\mathbf{y}_{fld} = W_{fld}^T \mathbf{y}_{pca} \qquad W_{fld} = \arg \max_W \frac{\left| W^T W_{pca}^T S_b W_{pca} W \right|}{\left| W^T W_{pca}^T S_w W_{pca} W \right|}$$

$$= W_{fld}^T W_{pca}^T \mathbf{x}$$

- Given the faces $\mathbf{x}_1 \ldots \mathbf{x}_N$, properly aligned and classified into class $c$ (i = 1,...,c), each with $n_i$ elements
- Determine the mean face $\boldsymbol{\mu}$, and the mean face of each class $\boldsymbol{\mu}_i$
- Determine the scatter matrix $S_T$ ($n$ x $n$)
  - Determine $m$ (maximum of $N$-$c$) non-zero eigenvectors, $\mathbf{W}_{pca}$ (see algorithm *eigenfaces*)
- Determine the $S_b$ ($n$ x $n$) and $S_w$ ($n$ x $n$) matrices
- Compute

$$\breve{S}_b = W_{pca}^T S_b W_{pca}$$
$$\breve{S}_w = W_{pca}^T S_w W_{pca}$$

- Determine the $c$-1 "larger eigenvectors" from the matrix

$$\breve{S}_w^{-1} \breve{S}_b \qquad \longleftarrow \qquad m \text{ x } m$$