#### Linear classifiers

Machine Learning

Hamid R Rabiee – Zahra Dehghanian Spring 2025



Sharif University of Technology

## Classification problem

- Given: Training set
  - labeled set of N input-output pairs  $D = \{(\mathbf{x}^{(i)}, y^{(i)})\}_{i=1}^{N}$
  - ▶  $y \in \{1, \dots, K\}$
- Goal: Given an input x, assign it to one of K classes
- Examples:
  - Spam filter
  - Handwritten digit recognition



- Decision boundaries are linear in x, or linear in some given set of functions of x
- Linearly separable data: data points that can be exactly classified by a linear decision surface.
- Why linear classifier?
  - Even when they are not optimal, we can use their simplicity
    - are relatively easy to compute
    - In the absence of information suggesting otherwise, linear classifiers are an attractive candidates for initial, trial classifiers.



## Two Category

$$f(x; w) = w^T x + w_0 = w_0 + w_1 x_1 + \dots + w_d x_d$$
  

$$x = [x_1 x_2 \dots x_d]$$

- $\mathbf{w} = [w_1 \, w_2 \, \dots \, w_d]$
- $w_0$ : bias
- if  $w^T x + w_0 \ge 0$  then  $\mathcal{C}_1$
- else  $C_2$

Decision surface (boundary):  $w^T x + w_0 = 0$ 

 $\boldsymbol{w}$  is orthogonal to every vector lying within the decision surface



## Example





## Linear classifier: Two Category

- P Decision boundary is a (d-1)-dimensional hyperplane H in the d-dimensional feature space
  - The orientation of H is determined by the normal vector  $[w_1, ..., w_d]$
  - $w_0$  determines the location of the surface.
    - The normal distance from the origin to the decision surface is  $\frac{w_0}{\|w\|}$

$$\boldsymbol{x} = \boldsymbol{x}_{\perp} + r \frac{\boldsymbol{w}}{\|\boldsymbol{w}\|}$$
$$\boldsymbol{w}^{T} \boldsymbol{x} + \boldsymbol{w}_{0} = r \|\boldsymbol{w}\| \Rightarrow r = \frac{\boldsymbol{w}^{T} \boldsymbol{x} + \boldsymbol{w}_{0}}{\|\boldsymbol{w}\|}$$

gives a signed measure of the perpendicular distance r of the point x from the decision surface





Sharif University of Technology

## Linear boundary: geometry





## Non-linear decision boundary

- Choose non-linear features
- Classifier still linear in parameters w







## Cost Function for linear classification

- Finding linear classifiers can be formulated as an optimization problem:
  - Select how to measure the prediction loss
    - Based on the training set  $D = \{(x^{(i)}, y^{(i)})\}_{i=1}^n$ , a cost function J(w) is defined
  - Solve the resulting optimization problem to find parameters:
    - Find optimal  $\hat{f}(\mathbf{x}) = f(\mathbf{x}; \hat{\mathbf{w}})$  where  $\hat{\mathbf{w}} = \underset{\mathbf{w}}{\operatorname{argmin}} J(\mathbf{w})$
- Criterion or cost functions for classification:
  - We will investigate several cost functions for the classification problem



## SSE cost function for classification

#### SSE cost function is not suitable for classification:

- Least square loss penalizes 'too correct' predictions (that they lie a long way on the correct side of the decision)
- Least square loss also lack robustness to noise

K = 2





Sharif University of Technology

## SSE cost function for classification







Sharif University of Technology



## SSE cost function for classification

Is it more suitable if we set 
$$f(x; w) = g(w^T x)$$
?  

$$J(w) = \sum_{i=1}^{N} (\operatorname{sign}(w^T x^{(i)}) - y^{(i)})^2 \qquad (\operatorname{sign}(w^T x) - y)^2$$

$$\operatorname{sign}(z) = \begin{cases} -1, & z < 0 \\ 1, & z \ge 0 \end{cases}$$

$$w^T x$$

• J(w) is a piecewise constant function shows the number of misclassifications





Sharif University of Technology

Training error incurred in classifying training samples

## Perceptron algorithm

- Linear classifier
- Two-class:  $y \in \{-1,1\}$

• 
$$y = -1$$
 for  $C_2$ ,  $y = 1$  for  $C_1$ 

Goal: ∀i, 
$$x^{(i)} \in C_1 \Rightarrow w^T x^{(i)} > 0$$
∀i,  $x^{(i)} \in C_2 \Rightarrow w^T x^{(i)} < 0$ 

• 
$$f(\mathbf{x}; \mathbf{w}) = \operatorname{sign}(\mathbf{w}^T \mathbf{x})$$

## Perceptron criterion

$$J_P(\boldsymbol{w}) = -\sum_{i \in \mathcal{M}} \boldsymbol{w}^T \boldsymbol{x}^{(i)} \boldsymbol{y}^{(i)}$$

 $\mathcal{M}$ : subset of training data that are misclassified

Many solutions? Which solution among them?



Sharif University of Technology

?

## Cost function



# of misclassifications as a cost function

Perceptron's cost function

There may be many solutions in these cost functions

[Duda, Hart, and Stork, 2002]

Sharif University of Technology

## Batch Perceptron

"Gradient Descent" to solve the optimization problem:

$$\boldsymbol{w}^{t+1} = \boldsymbol{w}^t - \eta \nabla_{\!\!\boldsymbol{w}} J_P(\boldsymbol{w}^t)$$
$$\nabla_{\!\!\boldsymbol{w}} J_P(\boldsymbol{w}) = -\sum_{i \in \mathcal{M}} \boldsymbol{x}^{(i)} \boldsymbol{y}^{(i)}$$

Batch Perceptron converges in finite number of steps for linearly separable data:

```
Initialize w
Repeat
w = w + \eta \sum_{i \in \mathcal{M}} x^{(i)} y^{(i)}
Until \eta \sum_{i \in \mathcal{M}} x^{(i)} y^{(i)} < \theta
```



### Stochastic gradient descent for Perceptron

- Single-sample perceptron:
  - If  $x^{(i)}$  is misclassified:

$$\boldsymbol{w}^{t+1} = \boldsymbol{w}^t + \eta \boldsymbol{x}^{(i)} \boldsymbol{y}^{(i)}$$

- Perceptron convergence theorem: for linearly separable data
  - If training data are linearly separable, the single-sample perceptron is also guaranteed to find a solution in a finite number of steps

Fixed-Increment single sample Perceptron

 $\eta$  can be set to I and proof still works  $\xrightarrow{}$ 

```
w \leftarrow 0

t \leftarrow 0

repeat

t \leftarrow t + 1

i \leftarrow t \mod N

if x^{(i)} is misclassified then

w = w + x^{(i)}y^{(i)}

Until all patterns properly classified
```









Sharif University of Technology

## Perceptron: Example





Change *w* in a direction that corrects the error

[Bishop]



Sharif University of Technology

## Convergence of Perceptron



[Duda, Hart & Stork, 2002]

? For data sets that are not linearly separable, the single-sample perceptron learning algorithm will never converge



## Pocket algorithm

- For the data that are not linearly separable due to noise:
  - Keeps in its pocket the best *w* encountered up to now.

```
Initialize w
for t = 1, ..., T
i \leftarrow t \mod N
if x^{(i)} is misclassified then
w^{new} = w + x^{(i)}y^{(i)}
if E_{train}(w^{new}) < E_{train}(w) then
w = w^{new}
```

end

$$E_{train}(\boldsymbol{w}) = \frac{1}{N} \sum_{n=1}^{N} \left[ sign(\boldsymbol{w}^{T} \boldsymbol{x}^{(n)}) \neq y^{(n)} \right]$$



## Multi-class classification

- Solutions to multi-category problems:
  - Converting the problem to a set of two-class problems
  - Extend the learning algorithm to support multi-class:
    - A function  $f_i(\mathbf{x})$  for each class *i* is found

$$\square \ \hat{y} = \underset{i=1,\dots,c}{\operatorname{argmax}} f_i(\boldsymbol{x})$$







? <a href="https://forms.gle/vKRbyVVsWRKcZuqr8">https://forms.gle/vKRbyVVsWRKcZuqr8</a>



**Regression: Probabilistic perspective** 

## Multi-class classification

- Solutions to multi-category problems:
  - Converting the problem to a set of two-class problems
  - Extend the learning algorithm to support multi-class:
    - A function  $f_i(\mathbf{x})$  for each class *i* is found

$$\Box \ \hat{y} = \operatorname*{argmax}_{i=1,\ldots,c} f_i(x)$$





## Converting multi-class problem to a set of two-class problems

#### "one versus rest" or "one against all"

- For each class  $C_i$ , a linear discriminant function that separates samples of  $C_i$  from all the other samples is found.
  - Totally linearly separable

#### "one versus one"

- c(c-1)/2 linear discriminant functions are used, one to separate samples of a pair of classes.
  - Pairwise linearly separable



## Multi-class classification



## Multi-class classification



## Multi-class classification: ambiguity

Converting the multi-class problem to a set of two-class problems can lead to **regions in which the classification is undefined** 



[Duda, Hart & Stork, 2002]



## Multi-class classification

- Solutions to multi-category problems:
  - Converting the problem to a set of two-class problems
  - Extend the learning algorithm to support multi-class:
    - A function  $f_i(\mathbf{x})$  for each class *i* is found

$$\square \hat{y} = \operatorname*{argmax}_{i=1,\dots,c} f_i(\boldsymbol{x})$$

 $\boldsymbol{x}$  is assigned to class  $C_i$  if  $f_i(\boldsymbol{x}) > f_j(\boldsymbol{x}) \quad \forall j \neq i$ 





## **Discriminant Functions**

- **P** Discriminant functions: A discriminant function  $f_i(x)$  for each class  $C_i$  (i = 1, ..., K):
  - x is assigned to class  $C_i$  if:

$$f_i(\boldsymbol{x}) > f_j(\boldsymbol{x}) \quad \forall j \neq i$$

Thus, we can easily divide the feature space into K decision regions

$$\forall \boldsymbol{x}, f_i(\boldsymbol{x}) > f_j(\boldsymbol{x}) \quad \forall j \neq i \Rightarrow \boldsymbol{x} \in \mathcal{R}_i$$

 $\mathcal{R}_i$ : Region of the *i*-th class

- Decision surfaces (or boundaries) can also be found using discriminant functions
  - Boundary of the  $\mathcal{R}_i$  and  $\mathcal{R}_j$  separating samples of these two categories:

 $\forall \boldsymbol{x}, f_i(\boldsymbol{x}) = f_j(\boldsymbol{x})$ 



## Discriminant functions

- **P** Discriminant function can directly assign each vector x to a specific class k
- A popular way of representing a classifier
  - Many classification methods are based on discriminant functions
- Assumption: the classes are taken to be disjoint
  - > The input space is thereby divided into **decision regions** 
    - boundaries are called **decision boundaries** or decision surfaces.



## Discriminant Functions: Two-Category

- ?
- Decision surface:  $f(\mathbf{x}) = 0$
- For two-category problem, we can only find a function  $f : \mathbb{R}^d \to \mathbb{R}$ 
  - $f_1(x) = f(x)$
  - $f_2(x) = -f(x)$



# Multi-class classification: linear machine

- A discriminant function  $f_i(\mathbf{x}) = \mathbf{w}_i^T \mathbf{x} + w_{i0}$  for each class  $C_i$  (*i* = 1, ..., *K*):
  - x is assigned to class  $C_i$  if:

 $f_i(\boldsymbol{x}) > f_j(\boldsymbol{x}) \quad \forall j \neq i$ 

- Decision surfaces (boundaries) can also be found using discriminant functions
  - Boundary of the contiguous  $\mathcal{R}_i$  and  $\mathcal{R}_j$ :  $\forall x, f_i(x) = f_j(x)$

• 
$$(w_i - w_j)^T x + (w_{i0} - w_{j0}) = 0a$$

Decision regions are convex



# Multi-class classification: linear machine

#### Decision regions are convex

• Linear machines are most suitable for problems where  $p(\mathbf{x}|C_i)$  are unimodal.

$$\begin{aligned} x_1, x_2 \in \mathcal{R}_i \Rightarrow \forall j \neq i, f_i(x_1) \geq f_j(x_1) \\ f_i(x_2) \geq f_j(x_2) \end{aligned}$$
  
$$\Rightarrow \alpha f_i(x_1) + (1 - \alpha) f_i(x_2) \geq \alpha f_j(x_1) + (1 - \alpha) f_j(x_2) \\ f_i \text{ is linear } \Rightarrow f_i(\alpha x_1 + (1 - \alpha) x_2) \geq f_j(\alpha x_1 + (1 - \alpha) x_2) \\ \Rightarrow \alpha x_1 + (1 - \alpha) x_2 \in \mathcal{R}_i \end{aligned}$$

Convex region definition:  $\forall x_1, x_2 \in \mathcal{R}, 0 \le \alpha \le 1 \Rightarrow \alpha x_1 + (1 - \alpha) x_2 \in \mathcal{R}$ 



### Multi-class classification: linear machine



[Duda, Hart & Stork, 2002]



Sharif University of Technology

## Perceptron: multi-class

$$\hat{y} = \operatorname*{argmax}_{i=1,...,c} \boldsymbol{w}_i^T \boldsymbol{x}$$
$$J_P(\boldsymbol{W}) = -\sum_{i \in \mathcal{M}} \left( \boldsymbol{w}_{y^{(i)}} - \boldsymbol{w}_{\hat{y}^{(i)}} \right)^T \boldsymbol{x}^{(i)}$$

 $\mathcal{M}$ : subset of training data that are misclassified  $\mathcal{M} = \{i | \hat{y}^{(i)} \neq y^{(i)}\}$ 

> Initialize  $W = [w_1, ..., w_c], k \leftarrow 0$ repeat  $k \leftarrow (k + 1) \mod N$ if  $x^{(i)}$  is misclassified then  $w_{\hat{y}^{(i)}} = w_{\hat{y}^{(i)}} - x^{(i)}$  $w_{y^{(i)}} = w_{y^{(i)}} + x^{(i)}$ Until all patterns properly classified



Sharif University of Technology

?

## Linear Discriminant Analysis (LDA)

- Fisher's Linear Discriminant Analysis :
  - Dimensionality reduction
    - Finds linear combinations of features with large ratios of betweengroups scatters to within-groups scatters (as discriminant new variables)
  - Classification
    - Predicts the class of an observation x by first projecting it to the space of discriminant variables and then classifying it in this space



- ? What is a good criterion?
  - ? Separating different classes in the projected space







- ? What is a good criterion?
  - ? Separating different classes in the projected space







- ? What is a good criterion?
  - ? Separating different classes in the projected space







- ? Between-class distance
  - ? Squared distance between the centroids of different classes





Sharif University of Technology



## LDA Problem

- Problem definition:
  - K = 2 classes
  - $\{(x^{(i)}, y^{(i)})\}_{i=1}^{N}$  training samples with  $N_1$  samples from the first class ( $C_1$ ) and  $N_2$  samples from the second class ( $C_2$ )
  - Goal: finding the best direction w that we hope to enable accurate classification
- The projection of sample x onto a line in direction w is  $w^T x$
- What is the measure of the separation between the projected points of different classes?



## Measure of Separation in the Projected Direction

- The direction of the line jointing the class means is the solution of the following problem:
  - Maximizes the separation of the projected class means

$$\max_{w} J(w) = (\mu'_{1} - \mu'_{2})^{2}$$
  
s.t.  $||w|| = 1$   
$$\mu'_{1} = w^{T} \mu_{1} \qquad \mu_{1} = \frac{\sum_{x^{(i)} \in \mathcal{C}_{1}} x^{(i)}}{N_{1}}$$
  
$$\mu'_{2} = w^{T} \mu_{2} \qquad \mu_{2} = \frac{\sum_{x^{(i)} \in \mathcal{C}_{2}} x^{(i)}}{N_{2}}$$

- What is the problem with the criteria considering only  $|\mu_1' \mu_2'|$ ?
  - It does not consider the variances of the classes in the projected direction



### Measure of Separation in the Projected Direction

Is the direction of the line jointing the class means a good candidate for w?







- ? Between-class distance
  - ? Squared distance between the centroids/means of different classes
- ? Within-class distance
  - ? Accumulated squared distance of an instance to the centroid/ mean of its class





- Fisher idea: maximize a function that will give
  - Iarge separation between the projected class means
  - while also achieving a small variance within each class, thereby minimizing the class overlap.

$$J(\boldsymbol{w}) = \frac{|\mu_1' - \mu_2'|^2}{s_1'^2 + s_2'^2}$$



♥ The scatters of projected data are:

$$s_1^{\prime 2} = \sum_{\boldsymbol{x}^{(i)} \in \mathcal{C}_1} (\boldsymbol{w}^T \boldsymbol{x}^{(i)} - \boldsymbol{w}^T \boldsymbol{\mu}_1)^2$$
$$s_2^{\prime 2} = \sum_{\boldsymbol{x}^{(i)} \in \mathcal{C}_2} (\boldsymbol{w}^T \boldsymbol{x}^{(i)} - \boldsymbol{w}^T \boldsymbol{\mu}_1)^2$$



Sharif University of Technology

$$J(\boldsymbol{w}) = \frac{|\mu_1' - \mu_2'|^2}{s_1'^2 + s_2'^2}$$

$$|\mu'_1 - \mu'_2|^2 = |\mathbf{w}^T \boldsymbol{\mu}_1 - \mathbf{w}^T \boldsymbol{\mu}_2|^2$$
  
=  $\mathbf{w}^T (\boldsymbol{\mu}_1 - \boldsymbol{\mu}_2) (\boldsymbol{\mu}_1 - \boldsymbol{\mu}_2)^T \mathbf{w}$ 

$$s_1'^2 = \sum_{\boldsymbol{x}^{(i)} \in \mathcal{C}_1} (\boldsymbol{w}^T \boldsymbol{x}^{(i)} - \boldsymbol{w}^T \boldsymbol{\mu}_1)^2$$
$$= \boldsymbol{w}^T \left( \sum_{\boldsymbol{x}^{(i)} \in \mathcal{C}_1} (\boldsymbol{x}^{(i)} - \boldsymbol{\mu}_1) (\boldsymbol{x}^{(i)} - \boldsymbol{\mu}_1)^T \right) \boldsymbol{w}$$



Sharif University of Technology

$$J(\boldsymbol{w}) = \frac{\boldsymbol{w}^T \boldsymbol{S}_B \boldsymbol{w}}{\boldsymbol{w}^T \boldsymbol{S}_W \boldsymbol{w}}$$

Between-class  
scatter matrix 
$$\qquad S_B = (\boldsymbol{\mu}_1 - \boldsymbol{\mu}_2)(\boldsymbol{\mu}_1 - \boldsymbol{\mu}_2)^T$$

Within-class scatter matrix

$$\boldsymbol{S}_W = \boldsymbol{S}_1 + \boldsymbol{S}_2$$

$$S_{1} = \sum_{\boldsymbol{x}^{(i)} \in \mathcal{C}_{1}} (\boldsymbol{x}^{(i)} - \boldsymbol{\mu}_{1}) (\boldsymbol{x}^{(i)} - \boldsymbol{\mu}_{1})^{T}$$
$$S_{2} = \sum_{\boldsymbol{x}^{(i)} \in \mathcal{C}_{2}} (\boldsymbol{x}^{(i)} - \boldsymbol{\mu}_{2}) (\boldsymbol{x}^{(i)} - \boldsymbol{\mu}_{2})^{T}$$
scatter matrix=N×covariance matrix

Sharif University of Technology

## LDA Derivation

$$J(\mathbf{w}) = \frac{\mathbf{w}^{T} \mathbf{S}_{B} \mathbf{w}}{\mathbf{w}^{T} \mathbf{S}_{W} \mathbf{w}}$$
$$\frac{\partial \mathbf{w}^{T} \mathbf{S}_{B} \mathbf{w}}{\partial \mathbf{w}} \times \mathbf{w}^{T} \mathbf{S}_{W} \mathbf{w} - \frac{\partial \mathbf{w}^{T} \mathbf{S}_{W} \mathbf{w}}{\partial \mathbf{w}} \times \mathbf{w}^{T} \mathbf{S}_{B} \mathbf{w}}{\left(\mathbf{w}^{T} \mathbf{S}_{W} \mathbf{w}\right)^{2}} = \frac{\left(2\mathbf{S}_{B} \mathbf{w}\right) \mathbf{w}^{T} \mathbf{S}_{W} \mathbf{w} - \left(2\mathbf{S}_{W} \mathbf{w}\right) \mathbf{w}^{T} \mathbf{S}_{B} \mathbf{w}}{\left(\mathbf{w}^{T} \mathbf{S}_{W} \mathbf{w}\right)^{2}}$$

 $\mathbf{w}^T \mathbf{S} \mathbf{w}$ 

$$\frac{\partial J(\mathbf{w})}{\partial \mathbf{w}} = 0 \Longrightarrow \mathbf{S}_{B} \mathbf{w} = \lambda \mathbf{S}_{W} \mathbf{w}$$



## LDA Derivation

 $S_B w$  (for any vector w) points in the same direction as  $\mu_1 - \mu_2$ :

$$S_B w = (\boldsymbol{\mu}_1 - \boldsymbol{\mu}_2)(\boldsymbol{\mu}_1 - \boldsymbol{\mu}_2)^T w \propto (\boldsymbol{\mu}_1 - \boldsymbol{\mu}_2)$$
$$w \propto S_W^{-1}(\boldsymbol{\mu}_1 - \boldsymbol{\mu}_2)$$

Thus, we can solve the eigenvalue problem immediately



## LDA Algorithm

- ${f P}$  Find  ${m \mu}_1$  and  ${m \mu}_2$  as the mean of class 1 and 2 respectively
- Find  $S_1$  and  $S_2$  as scatter matrix of class 1 and 2 respectively
- $\flat S_W = S_1 + S_2$
- $S_B = (\mu_1 \mu_2)(\mu_1 \mu_2)^T$
- Feature Extraction
  - $w = S_w^{-1}(\mu_1 \mu_2)$  as the eigenvector corresponding to the largest eigenvalue of  $S_w^{-1}S_b$
- Classification
  - $w = S_w^{-1}(\mu_1 \mu_2)$
  - Using a threshold on  $w^T x$ , we can classify x



# Multi-class classification: target coding scheme

- Target values:
  - ▶ Binary classification: a target variable  $y \in \{0,1\}$
  - Multiple classes (K > 2):
    - Target class  $C_j$ :  $y_j = 1$  $\forall i \neq j \ y_i = 0$  One of K coding





?

 $J(\boldsymbol{W}) = \sum_{k=1}^{K} \|\boldsymbol{X}\boldsymbol{w}_k - \boldsymbol{y}_k\|_2^2$ 

SE cost function for manti-elass-y)

$$\boldsymbol{X} = \begin{bmatrix} 1 & x_1^{(1)} & \cdots & x_d^{(1)} \\ 1 & x_1^{(2)} & \cdots & x_d^{(2)} \\ \vdots & \vdots & \ddots & \vdots \\ 1 & x_1^{(n)} & \cdots & x_d^{(n)} \end{bmatrix} \quad \boldsymbol{y}_k = \begin{bmatrix} y_k^{(1)} \\ \vdots \\ y_k^{(n)} \end{bmatrix}$$
$$\boldsymbol{W} = \begin{bmatrix} \boldsymbol{W}_1 & \cdots & \boldsymbol{W}_K \end{bmatrix}$$
$$\boldsymbol{W} = \begin{bmatrix} \boldsymbol{W}_1 & \cdots & \boldsymbol{W}_K \end{bmatrix}$$
$$\boldsymbol{Y} = \begin{bmatrix} \boldsymbol{y}_1 & \cdots & \boldsymbol{y}_K \end{bmatrix}$$
$$\boldsymbol{\nabla}_{\boldsymbol{W}} / (\boldsymbol{W}) = \boldsymbol{0} \Rightarrow \widehat{\boldsymbol{W}} = (\boldsymbol{X}^T \boldsymbol{X})^{-1} \boldsymbol{X} \boldsymbol{Y}$$



Sharif University of Technology

### SSE for multi-class: example

Low performance of the SSE cost function for the classification problem







#### Think about multi-class perceptron, LDA



#### Multi-class LDA

$$W = [w_1 w_2 ... w_{K-1}]$$
  
 $x' = W^T x$ 

- Means and scatters after transform  $x' = W^T x$ :
  - $\flat \ \mathbf{S}_B' = \mathbf{W}^T \mathbf{S}_B \mathbf{W}$
  - $\flat \ \mathbf{S}'_W = \mathbf{W}^T \mathbf{S}_W \mathbf{W}$



## Multi-Class LDA: Objective Function

- We seek a transformation matrix W that in some sense "<u>maximizes the ratio of the between-class scatter to the</u> <u>within-class scatter</u>".
- A simple scalar measure of scatter is the determinant of the scatter matrix.



## Multi-class LDA (MDA)

- - The projection from a d-dimensional space to a (C-1)-dimensional space (tacitly assumed that  $d \ge C$ ).

$$\boldsymbol{S}_{W} = \sum_{j=1}^{K} \boldsymbol{S}_{j}$$
$$\boldsymbol{S}_{B} = \sum_{j=1}^{K} N_{j} (\boldsymbol{\mu}_{j} - \boldsymbol{\mu}) (\boldsymbol{\mu}_{j} - \boldsymbol{\mu})^{T}$$

$$\mu_{j} = \frac{\sum_{x^{(i)} \in \mathcal{C}_{j}} x^{(i)}}{N_{j}} \quad j = 1, ..., K$$
  
$$\mu = \frac{\sum_{i=1}^{N} x^{(i)}}{N}$$
  
$$S_{j} = \sum_{x^{(i)} \in \mathcal{C}_{j}} (x^{(i)} - \mu_{j}) (x^{(i)} - \mu_{j})^{T} \qquad j = 1, ..., K$$



## Multi-class LDA: Objective function

$$J(\boldsymbol{W}) = \frac{|\boldsymbol{W}^T \boldsymbol{S}_B \boldsymbol{W}|}{|\boldsymbol{W}^T \boldsymbol{S}_W \boldsymbol{W}|} \quad \text{determinant}$$

- The solution of the problem where  $W = [w_1 w_2 ... w_{C-1}]$ :  $S_B w_i = \lambda_i S_W w_i$
- It is a generalized eigenvectors problem.



?

## **Objective Function: Trace Ratio**

maximizing between-class squared distances and minimizing within-class squared distances:

$$J(\boldsymbol{W}) = \frac{tr(\boldsymbol{W}^T \boldsymbol{S}_B \boldsymbol{W})}{tr(\boldsymbol{W}^T \boldsymbol{S}_W \boldsymbol{W})}$$

Between-class distance = trace of between-class scatter Within-class distance = trace of within-class scatter



#### **Objective Function: Trace Ratio**

$$S_W = \sum_{j=1}^C \sum_{\mathbf{x}^{(i)} \in \mathcal{C}_j} (\mathbf{x}^{(i)} - \boldsymbol{\mu}_j) (\mathbf{x}^{(i)} - \boldsymbol{\mu}_j)^T$$
  
$$S_B = \sum_{j=1}^C N_j (\boldsymbol{\mu}_j - \boldsymbol{\mu}) (\boldsymbol{\mu}_j - \boldsymbol{\mu})^T$$

$$tr(\mathbf{S}_{W}) = \sum_{j=1}^{C} \sum_{\mathbf{x}^{(i)} \in \mathcal{C}_{j}} (\mathbf{x}^{(i)} - \boldsymbol{\mu}_{j})^{T} (\mathbf{x}^{(i)} - \boldsymbol{\mu}_{j})$$
$$= \sum_{j=1}^{C} \sum_{\mathbf{x}^{(i)} \in \mathcal{C}_{j}} \|\mathbf{x}^{(i)} - \boldsymbol{\mu}_{j}\|^{2}$$
$$tr(\mathbf{S}_{B}) = \sum_{j=1}^{C} N_{j} (\boldsymbol{\mu}_{j} - \boldsymbol{\mu})^{T} (\boldsymbol{\mu}_{j} - \boldsymbol{\mu}) = \sum_{j=1}^{C} N_{j} \|\boldsymbol{\mu}_{j} - \boldsymbol{\mu}\|^{2}$$



Sharif University of Technology

## Objective Function: Trace Ratio vs. Ratio Trace

$$tr(S_w) = \sum_{i=1}^{c} \sum_{\substack{j \in class \ i}} \sum_{\substack{k \in class \ i}} \|\mathbf{x}_j - \mathbf{x}_k\|^2$$
$$tr(S_b) = \sum_{\substack{i=1 \ j \neq i}}^{c} \sum_{\substack{j \in class \ i}}^{c} \sum_{\substack{k \in class \ j}}^{c} \|\mathbf{x}_j - \mathbf{x}_k\|^2$$



## Multi-class LDA: Other objective function

There are many possible choices of criterion for multiclass LDA, e.g.:

$$J(\boldsymbol{W}) = tr(\boldsymbol{S}_{W}^{\prime-1}\boldsymbol{S}_{B}^{\prime}) = tr((\boldsymbol{W}^{T}\boldsymbol{S}_{W}\boldsymbol{W})^{-1}(\boldsymbol{W}^{T}\boldsymbol{S}_{B}\boldsymbol{W}))$$

- The solution is given by solving a generalized eigenvalue problem  $S_w^{-1}S_b$ 
  - Solution: eigen vectors corresponding to the largest eigen values constitute the new variables



## LDA Criterion Limitation

- ♥ When  $\mu_1 = \mu_2$ , LDA criterion can not lead to a proper projection (J(w) = 0)
  - However, discriminatory information in the scatter of the data may be helpful



- If classes are non-linearly separable they may have large overlap when projected to any line
- LDA implicitly assumes Gaussian distribution of samples of each class





## **Issues in LDA**

- ♥ Singularity or undersampled problem (when N < d)
  - Example: gene expression data, images, text documents
- Can reduces dimension only to  $d' \leq C 1$ 
  - $rank(\boldsymbol{S}_B) \leq C 1$ 
    - $S_B$  is the sum of *C* matrices  $(\mu_j \mu)(\mu_j \mu)^T$  of rank (at most) one and only C 1 of these are independent,
    - ▶  $\Rightarrow$  atmost C 1 nonzero eigenvalues and the desired weight vectors correspond to these nonzero eigenvalues.





- 2 Although LDA often provide more suitable features for classification tasks, LDA may fails in some situations such as:
  - when the number of samples per class is small (overfitting problem of LDA)
  - when the training data non-uniformly sample the underlying distribution
  - $\Box$  when the number of the desired features is more than C-1

- Advances in the recent decade:
  - Semi-supervised feature extraction
  - Nonlinear dimensionality reduction



## Applications

- ? Face recognition
  - ? Belhumeour et al., PAMI'97
- ? Image retrieval
  - ? Swets and Weng, PAMI'96
- ? Gene expression data analysis
  - ? Dudoit et al., JASA'02; Ye et al., TCBB'04
- ? Protein expression data analysis
  - ? Lilien et al., Comp. Bio.'03
- ? Text mining
  - Park et al., SIMAX'03; Ye et al., PAMI'04
- ? Medical image analysis
  - ? Dundar, SDM'05





- ? C. Bishop, "Pattern Recognition and Machine Learning", Chapter 4.1.
- ? Course CE-717, Dr. M.Soleymani

