

# Linear Discriminant Functions

Minimum Squared Error Procedures:  
Pseudoinverse and Widrow-Hoff  
Ho-Kashyap Procedures

# Non-separable Behavior

- Perceptron is an error-correcting procedure
  - Modify weight only when an error is encountered
  - Success is due to relentless search for error-free solution
- Even if separating vector is found does not follow it will perform well on independent test data
  - Any set of fewer than  $2d$  samples is likely to be linearly separable

# Non-separable Case

- When design sets are large, points are most likely not linearly separable
- How does error correction procedures behave when points are not linearly separable?
  - Corrections will never cease in an error correction procedure
  - Infinite sequence of weight vectors
- Weight vectors are bounded
  - Empirical rule based on weight vector fluctuating near a terminal value
  - Averaging weight vectors can reduce risk of obtaining a bad solution

# Minimum Squared Error (MSE) Procedures

- Criterion function considered so far involve only misclassified samples
- Now consider a criterion that involves *all* of the samples, not just misclassified ones
- Previously we were interested in making all of the inner products  $\mathbf{a}^t \mathbf{y}_i$  positive
- Now try to make  $\mathbf{a}^t \mathbf{y}_i = b_i$  where  $b_i$  are some arbitrarily specified positive constants
- Thus replace the problem of solving a set of linear inequalities with more stringent but *better understood problem* of finding a solution to a set of linear equations

# Minimum Squared Error (MSE) Formulation

For all the samples  $y_1, y_2, \dots, y_n$  we want a weight vector  $a$  so that  $a^t y_i = b_i$  for some arbitrarily specified positive numbers.

**The matrix notation:**

Data Matrix  $\rightarrow$

$$\begin{pmatrix} y_{10} & y_{11} & \cdots & y_{1d} \\ y_{20} & y_{21} & \cdots & y_{2d} \\ \vdots & \vdots & \vdots & \vdots \\ y_{n0} & y_{n1} & \cdots & y_{nd} \end{pmatrix} \begin{pmatrix} a_0 \\ a_1 \\ \vdots \\ a_d \end{pmatrix} = \begin{pmatrix} b_0 \\ b_1 \\ \vdots \\ b_n \end{pmatrix} \text{ or}$$
$$Y a = b$$

Note that dimensionality is  $d+1$  by introducing 0 sub-script

**Error vector:**

$$e = Y a - b$$

**sum-of-squared-error criterion function:**

$$J_s(a) = \|Y a - b\|^2 = \sum^n (a^t y_i - b_i)^2$$

which we wish to minimize.

Can be solved either by a gradient descent procedure or a closed form solution based on gradient

Weight vector to be determined

# MSE Solution based on Pseudoinverse

The gradient:

$$\nabla J_s = \sum_{i=1}^n 2(a^t y_i - b_i) y_i = 2Y^t(Ya - b)$$

Set it to zero, we get

$$Y^t Y a = Y^t b$$

IF  $Y^t Y$  is *nonsingular*,

$$\begin{aligned} a &= (Y^t Y)^{-1} Y^t b \\ &= Y^\dagger b \end{aligned}$$

Closed form solution

The  $d$ -by- $n$  matrix  $Y^\dagger = (Y^t Y)^{-1} Y^t$  is called pseudoinverse of  $Y$ .

## Remarks:

1.  $Y^t Y$  is not always *nonsingular*. If we define  $Y^\dagger$  more generally by

$$Y^\dagger \equiv \lim_{\epsilon \rightarrow 0} (Y^t Y + \epsilon I)^{-1} Y^t$$

then the limit always exists and  $a = Y^\dagger b$  is a MSE solution to  $Y a = b$ .

2. For an arbitrarily fixed  $b$ , MSE solution may not be a separating vector.

# Example of Linear Classifier by Pseudoinverse

- $\omega_1$ :  $(1,2)^t$  and  $(2,0)^t$
- $\omega_2$ :  $(3,1)^t$  and  $(2,3)^t$

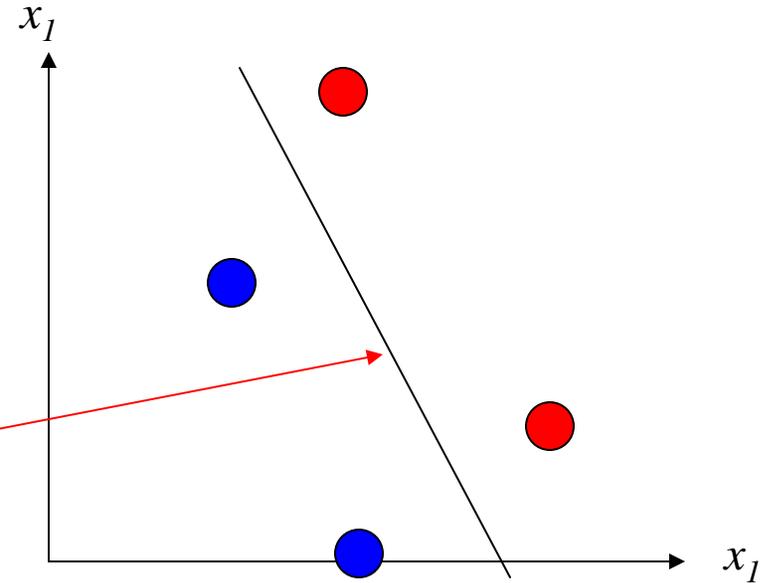
Sample Matrix ( $d = 1+2, n = 4$ )

$$Y = \begin{bmatrix} 1 & 1 & 2 \\ 1 & 2 & 0 \\ -1 & -3 & -1 \\ -1 & -2 & -3 \end{bmatrix}$$

Pseudo-inverse

$$Y^* = (Y^t Y)^{-1} Y^t = \begin{bmatrix} 5/4 & 13/12 & 3/4 & 7/12 \\ -1/2 & -1/6 & -1/2 & -1/6 \\ 0 & -1/3 & 0 & -1/3 \end{bmatrix}$$

$$a^t \begin{pmatrix} 1 \\ x_1 \\ x_2 \end{pmatrix} = 0$$



Assuming  $b = \begin{bmatrix} 1 \\ 1 \\ 1 \\ 1 \end{bmatrix}$

our solution is  $a = Y^t b = \begin{bmatrix} 11/3 \\ -4/3 \\ -2/3 \end{bmatrix}$

# Properties of MSE Procedure

1. Related to Fisher's Linear Discriminant
2. Asymptotic approximation to Bayes discriminant function
3. Can be formulated as a gradient descent procedure

# 1. MSE Relationship to Fisher's Linear Discriminant

- Show that with proper choice of the vector  $b$  the MSE discriminant function  $a^t y$  is directly related to Fisher's linear discriminant
- Assume first  $n_1$  samples are labelled  $\omega_1$  and second  $n_2$  samples are labelled  $\omega_2$

$$Y = \begin{bmatrix} 1_1 & X_1 \\ -1_2 & X_2 \end{bmatrix} \text{ where } 1_1 \text{ is a column vector } n_1 \text{ ones}$$

and  $X_i$  is a matrix whose rows are the samples labeled  $\omega_i$

Partition  $a$  and  $b$  correspondingly

$$a = \begin{bmatrix} w_0 \\ w \end{bmatrix} \text{ and with } b = \begin{bmatrix} \frac{n}{n_1} 1_1 \\ \frac{n}{n_2} 1_2 \end{bmatrix}$$

This special choice of  $b$  links the MSE solution to Fisher's Linear Discriminant

# MSE and Fisher's Linear Discriminant

- Define sample means  $m_i$  and pooled sample scatter matrix  $S_w$

$$m_i = \frac{1}{n_i} \sum_{x \in D_i} x, \quad i = 1, 2$$

$$S_w = \sum_i \sum_{x \in D_i} (x - m_i)(x - m_i)^t$$

- *and plug into MSE formulation yields*

$$w = \alpha n S_w^{-1} (m_1 - m_2)$$

*where  $\alpha$  is a scalar*

- *which is identical to the solution to the Fisher's linear discriminant except for a scale factor*
- *Decision rule: Decide  $\omega_1$  if  $w^t(x-m) > 0$ ; otherwise decide  $\omega_2$*

# MSE Relationship to Bayes

- If  $b=1_n$  MSE approaches a minimum squared error approximation to Bayes discriminant function

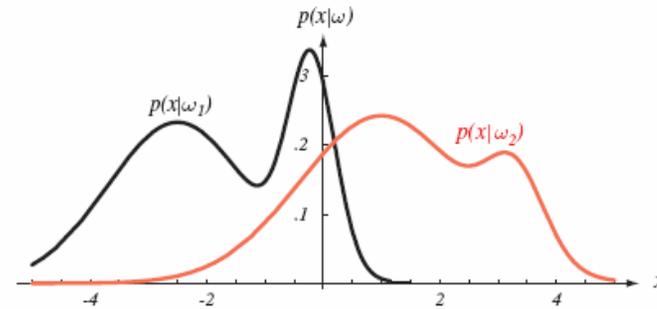
$$g_0(x) = P(\omega_1|x) - P(\omega_2|x)$$

in the limit as number of samples approaches infinity

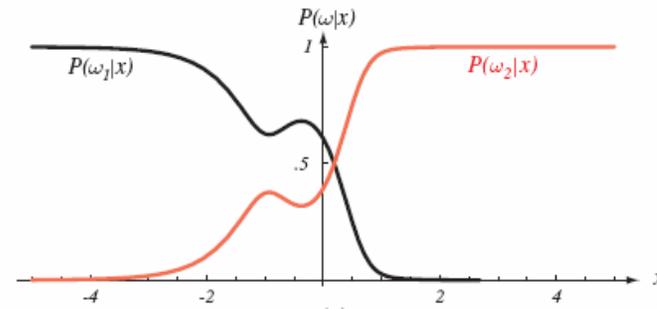
# 2. MSE Approximation to Bayes

- If  $b=1_n$  MSE solution approaches a minimum mean squared approximation to Bayes discriminant function

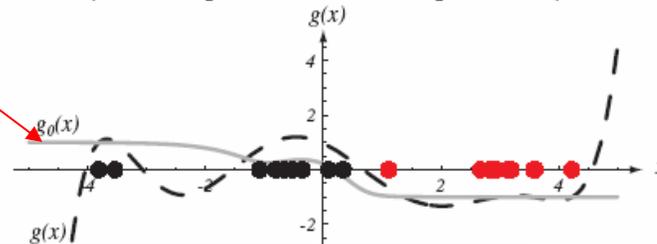
Class conditional densities



Posteriors



Bayes discriminant function



However MSE does not necessarily minimize probability of error

MSE solution (best approximation in region of data points)

### 3. MSE Solution using Gradient Descent:

- Criterion function  $J_s(\mathbf{a}) = \|\mathbf{Y}\mathbf{a}-\mathbf{b}\|^2$  could be minimized by gradient descent
- Advantage over pseudo-inverse:
  - Problem when  $\mathbf{Y}^t\mathbf{Y}$  is singular
  - Avoids need for working with large matrices
  - Computation involved is a feedback scheme that copes with roundoff or truncation

# Widrow-Hoff or Least Mean Squared (LMS) Rule

Since  $\Delta J_s = 2\mathbf{Y}^t(\mathbf{Y}\mathbf{a} - \mathbf{b})$

The obvious update rule is

$\mathbf{a}(1)$  arbitrary

$$\mathbf{a}(k+1) = \mathbf{a}(k) + \eta(k)\mathbf{Y}^t(\mathbf{Y}\mathbf{a}(k) - \mathbf{b})$$

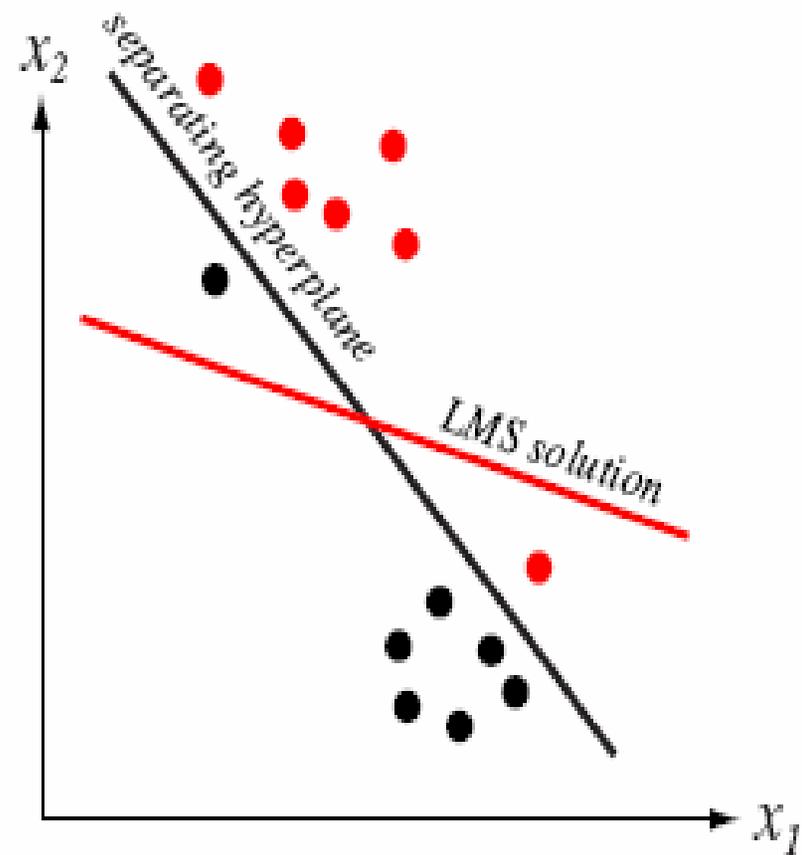
Can be reduced for storage requirement to the rule where samples are considered sequentially:

$\mathbf{a}(1)$  arbitrary

$$\mathbf{a}(k+1) = \mathbf{a}(k) + \eta(k)(\mathbf{b}(k) - \mathbf{a}^t(k)\mathbf{y}^k)\mathbf{y}^k$$

# Widrow-Hoff or LMS Property

- Relaxation procedures need not converge to a separating hyperplane even if there exists one
- Widrow-Hoff does converge
- However solution need not give a separating vector even if one exists



# Ho-Kashyap Procedures

- Procedures considered so far:
  - Perceptron and relaxation procedures finding separating vectors if samples are linearly separable, but do not converge on non-separable problems
  - MSE procedures yield a weight vector whether samples are linearly separable or not, but no guarantee that vector is a separating vector
- Criterion function is  $J_s(\mathbf{a}) = \|\mathbf{Y}\mathbf{a} - \mathbf{b}\|^2$ 
  - Subject to constraint  $b > 0$
  - How to obtain both a separating vector  $\mathbf{a}$  and a margin  $b$ ?

# Ho-Kashyap Procedures

## 5.9 The Ho-Kashyap Procedures

---

**The Descent Procedure on the two variable function:**

The gradient of  $J_s$  with respect to  $a$ :

$$\nabla_a J_s = 2Y^t(Ya - b)$$

and the gradient of  $J_s$  with respect to  $b$ :

$$\nabla_b J_s = -2(Ya - b)$$

For any value of  $b$ , we can always have

$$a = Y^\dagger b$$

which makes  $\nabla_a J_s = 0$  and minimizing  $J_s$  with respect to  $a$ .

# Ho-Kashyap

Constraints on  $b$ :

1.  $b > 0$
2. we don't want  $b(k) \rightarrow 0$ .

We can get a sequence of  $b(k)$  by starting  $b > 0$  and refusing to reduce any of its components.

In the updating rule for  $b$ :

$$b(k+1) = b(k) - \eta(k) \nabla_b J_s = b(k) + 2\eta(k)(Ya - b)$$

we take out all the positive components from  $\nabla_b J_s$  or equivalently we take only positive components in the second term on the right side.

# Ho-Kashyap Rule

Then we get the *Ho-Kashyap rule* for minimizing  $J_s(a, b)$ :

$$\begin{cases} b(1) > 0 \\ b(k+1) = b(k) + 2\eta(k)e^+(k) \end{cases}$$

where

$$e(k) = Ya(k) - b(k) \quad \leftarrow \text{Error vector}$$

$$e^+(k) = \frac{1}{2}(e(k) + |e(k)|) \quad \leftarrow \begin{array}{l} \text{Positive} \\ \text{Part of} \\ \text{Error vector} \end{array}$$

and

$$a(k) = Y^\dagger b(k), \quad 1, 2, \dots$$

# Ho-Kashyap Algorithm

## Ho-Kashyap Algorithm:

```
1 begin initialize  $a, b, \eta < 1$ , threshold  $b_{min}, k_{max}$ 
2   do  $k \leftarrow (k + 1) \bmod n$ 
3      $e \leftarrow Ya - b$ 
4      $e^+ \leftarrow 1/2(e + |e|)$ 
5      $b \leftarrow a + 2\eta(k)e^+$ 
6      $a \leftarrow Y^\dagger b(k)$ 
7     if  $|e| \leq b_{min}$  then return  $a, b$  and exit
8   util  $k = k_{max}$ 
9   print "NO SOLUTION FOUND"
10 end
```

# Ho-Kashyap Conclusion

## Remarks:

1. If the training samples are linearly separable then Ho-Kashyap Algorithm will always have a solution.
2. The pseudoinverse  $Y^\dagger$  depends only on  $Y$  and needs to be computed once. In the case  $Y^tY$  is nonsingular,  $Y^\dagger = (Y^tY)^{-1}$  and otherwise it is in a different form.

**In the nonseparable case there will be an indication of nonseparability. However there is no bound on the number of steps needed to disclose nonseparability**

# Linear Programming Algorithms

- Perceptron, relaxation, Ho-Kashyap are gradient descent procedures for solving simultaneous inequalities.
- Task can also be approached with LP
- Minimize objective function  $z = \alpha^t u$  subject to constraints  $Au \geq \beta$
- Simplex algorithm

Surfaces of constant  $z$  are shown in gray  
While constraints are shown in red  
Simplex finds the extremum of  $z$

