

***CS442/542b: Artificial Intelligence II***  
***Prof. Olga Veksler***

***Lecture 4: Machine Learning***  
***Linear Classifier***

***Outline***

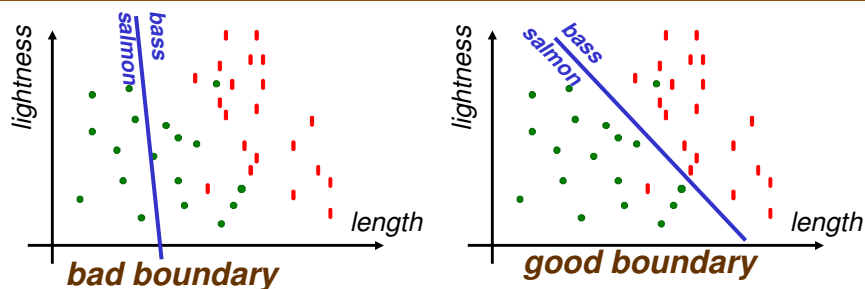
---

- Linear Classification or Linear Discriminant Functions
  - Introduction
  - 2 classes
  - Multiple classes
  - Optimization with gradient descent
  - Perceptron Criterion Function
    - Batch perceptron rule
    - Single sample perceptron rule
  - Minimum Squared Error (MSE) rule
    - Pseudoinverse
    - Gradient descent (Widrow-Hoff Procedure)

## ***Discriminant Function***

- Consider 2 class case, and let's denote them
  - class 1
  - class -1
- Want to build  $f(X,W)$  to give correct class of sample  $X$
- Let  $f(X,W) = \text{sign}(g(X,W))$ 
  - 1 if  $g(X,W)$  is positive
  - -1 if  $g(X,W)$  is negative
- $g(X,W)$  is called the ***discriminant function***

## ***Linear Discriminant Functions: Basic Idea***



- Have samples from 2 classes  $\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n$
- Assume 2 classes can be classified by a linear discriminant function  $g(X,W)$  with some unknown parameters  $W$
- Fit the “best”  $g(X,W)$  to data by optimizing over parameters  $W$
- What is best  $g(X,W)$ ?
  - simplest idea: minimize classification error on training data

## LDF: Introduction

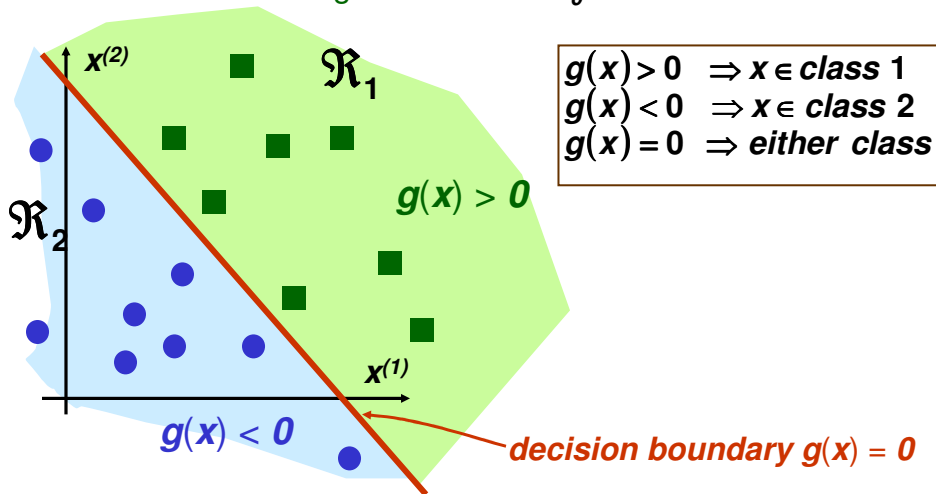
- Discriminant functions can be more general than linear
- For now, we will study linear discriminant functions
  - Simple model (should try simpler models first)
  - Analytically tractable
- Linear Discriminant functions are optimal for certain type of data
  - Gaussian distributions with equal covariance (don't worry if you don't know what a Gaussian is)
- May not be optimal for other data distributions, but they are very simple to use

## LDF: 2 Classes

- A discriminant function is linear if it can be written as

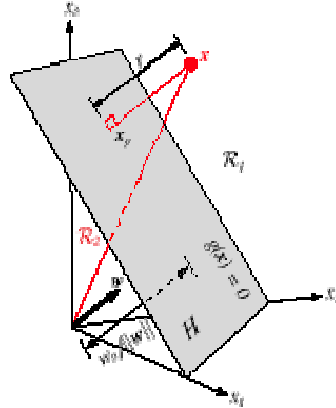
$$g(x) = w^t x + w_0$$

- $w$  is called the **weight vector** and  $w_0$  called **bias** or **threshold**



## LDF: 2 Classes

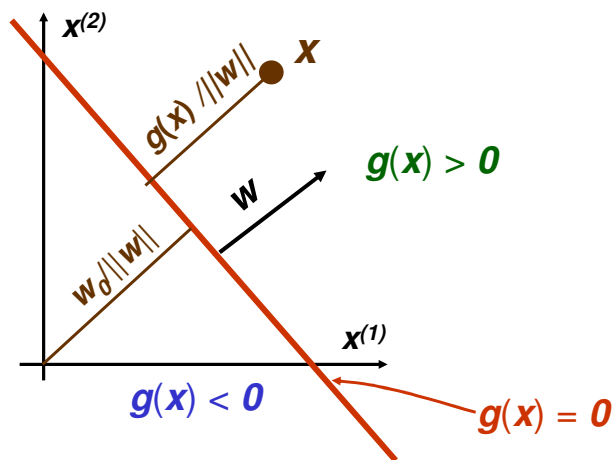
- Decision boundary  $g(\mathbf{x}) = \mathbf{w}^t \mathbf{x} + w_0 = 0$  is a **hyperplane**
  - set of vectors  $\mathbf{x}$  which for some scalars  $\alpha_0, \dots, \alpha_d$  satisfy  $\alpha_0 + \alpha_1 x^{(1)} + \dots + \alpha_d x^{(d)} = 0$
- A hyperplane is
  - a point in 1D
  - a line in 2D
  - a plane in 3D



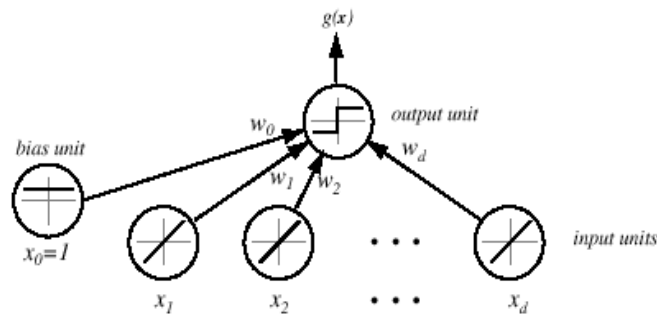
## LDF: 2 Classes

$$g(\mathbf{x}) = \mathbf{w}^t \mathbf{x} + w_0$$

- $\mathbf{w}$  determines orientation of the decision hyperplane
- $w_0$  determines location of the decision surface



## LDF: 2 Classes



**FIGURE 5.1.** A simple linear classifier having  $d$  input units, each corresponding to the values of the components of an input vector. Each input feature value  $x_i$  is multiplied by its corresponding weight  $w_i$ ; the effective input at the output unit is the sum all these products,  $\sum w_i x_i$ . We show in each unit its effective input-output function. Thus each of the  $d$  input units is linear, emitting exactly the value of its corresponding feature value. The single bias unit unit always emits the constant value 1.0. The single output unit emits a  $+1$  if  $\mathbf{w}^t \mathbf{x} + w_0 > 0$  or a  $-1$  otherwise. From: Richard O. Duda, Peter E. Hart, and David G. Stork, *Pattern Classification*. Copyright © 2001 by John Wiley & Sons, Inc.

## LDF: Many Classes

- Suppose we have  $m$  classes
- Define  $m$  linear discriminant functions

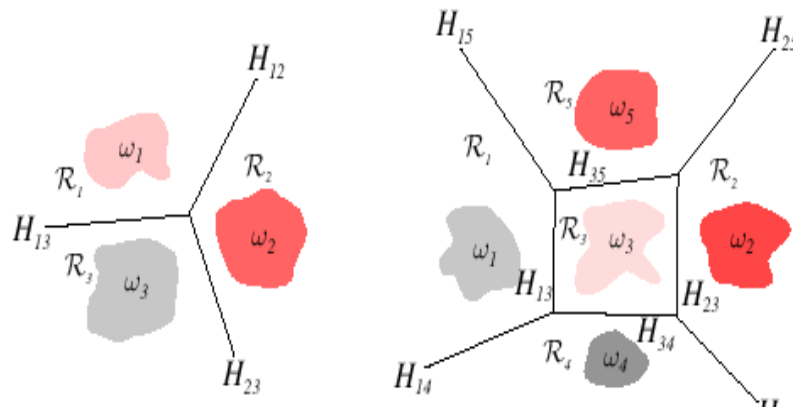
$$g_i(\mathbf{x}) = \mathbf{w}_i^t \mathbf{x} + w_{i0} \quad \mathbf{i} = 1, \dots, m$$

- Given  $\mathbf{x}$ , assign class  $c_i$  if

$$g_i(\mathbf{x}) \geq g_j(\mathbf{x}) \quad \forall j \neq i$$

- Such classifier is called a **linear machine**
- A linear machine divides the feature space into  $c$  decision regions, with  $g_i(\mathbf{x})$  being the largest discriminant if  $\mathbf{x}$  is in the region  $R_i$

## LDF: Many Classes

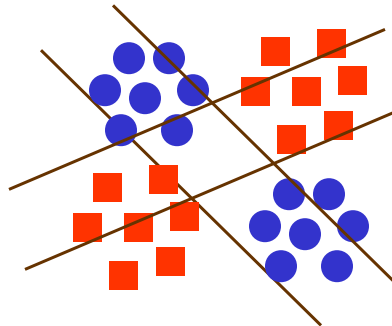


- Decision regions for a linear machine are **convex**
  - In particular, decision regions must be spatially contiguous

## LDF: Many Classes

- Thus applicability of linear machine to mostly limited to unimodal distributions
  - Most data is concentrated around one point in space

- Example:



- need non-contiguous decision regions
- thus linear machine will fail

## LDF: Augmented feature vector

- Linear discriminant function:  $g(\mathbf{x}) = \mathbf{w}^t \mathbf{x} + w_0$
- Can rewrite it:  $g(\mathbf{x}) = \underbrace{[w_0 \quad \mathbf{w}^t]}_{\substack{\text{new weight} \\ \text{vector } \mathbf{a}}} \underbrace{\begin{bmatrix} 1 \\ \mathbf{x} \end{bmatrix}}_{\substack{\text{new feature} \\ \text{vector } \mathbf{y}}} = \mathbf{a}^t \mathbf{y} = g(\mathbf{y})$
- $\mathbf{y}$  is called the **augmented feature vector**
- Added a dummy dimension to get a completely equivalent new **homogeneous** problem

*old problem*

$$g(\mathbf{x}) = \mathbf{w}^t \mathbf{x} + w_0$$

$$\begin{bmatrix} x_1 \\ \vdots \\ x_d \end{bmatrix}$$

*new problem*

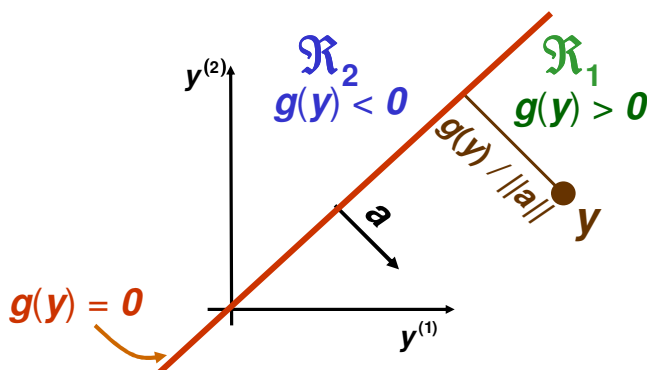
$$g(\mathbf{y}) = \mathbf{a}^t \mathbf{y}$$

$$\begin{bmatrix} 1 \\ x_1 \\ \vdots \\ x_d \end{bmatrix}$$

## LDF: Augmented feature vector

- Feature augmenting is done for simpler notation
- From now on we always assume that we have augmented feature vectors
  - Given samples  $\mathbf{x}_1, \dots, \mathbf{x}_n$  convert them to augmented samples  $\mathbf{y}_1, \dots, \mathbf{y}_n$  by adding a new dimension of value 1

$$\mathbf{y}_i = \begin{bmatrix} 1 \\ \mathbf{x}_i \end{bmatrix}$$



## LDF: Training Error

- For the rest of the lecture, assume we have 2 classes
- Samples  $\mathbf{y}_1, \dots, \mathbf{y}_n$  some in class 1, some in class 2
- Use these samples to determine weights  $\mathbf{a}$  in the discriminant function  $g(\mathbf{y}) = \mathbf{a}^t \mathbf{y}$
- What should be our criterion for determining  $\mathbf{a}$ ?
  - For now, suppose we want to minimize the training error (that is the number of misclassified samples  $\mathbf{y}_1, \dots, \mathbf{y}_n$ )
- Recall that  $g(\mathbf{y}_i) > 0 \Rightarrow \mathbf{y}_i$  classified  $\mathbf{c}_1$   
 $g(\mathbf{y}_i) < 0 \Rightarrow \mathbf{y}_i$  classified  $\mathbf{c}_2$
- Thus training error is 0 if 
$$\begin{cases} g(\mathbf{y}_i) > 0 & \forall \mathbf{y}_i \in \mathbf{c}_1 \\ g(\mathbf{y}_i) < 0 & \forall \mathbf{y}_i \in \mathbf{c}_2 \end{cases}$$

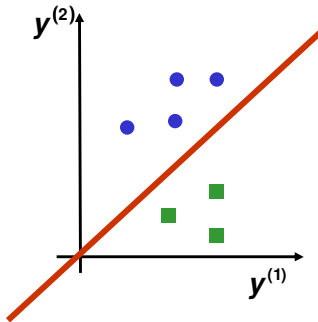
## LDF: Problem “Normalization”

- Thus training error is 0 if 
$$\begin{cases} \mathbf{a}^t \mathbf{y}_i > 0 & \forall \mathbf{y}_i \in \mathbf{c}_1 \\ \mathbf{a}^t \mathbf{y}_i < 0 & \forall \mathbf{y}_i \in \mathbf{c}_2 \end{cases}$$
- Equivalently, training error is 0 if 
$$\begin{cases} \mathbf{a}^t \mathbf{y}_i > 0 & \forall \mathbf{y}_i \in \mathbf{c}_1 \\ \mathbf{a}^t (-\mathbf{y}_i) > 0 & \forall \mathbf{y}_i \in \mathbf{c}_2 \end{cases}$$
- This suggest problem “normalization”:
  1. Replace all examples from class  $\mathbf{c}_2$  by their negative
$$\mathbf{y}_i \rightarrow -\mathbf{y}_i \quad \forall \mathbf{y}_i \in \mathbf{c}_2$$
  2. Seek weight vector  $\mathbf{a}$  s.t.
$$\mathbf{a}^t \mathbf{y}_i > 0 \quad \forall \mathbf{y}_i$$
    - If such  $\mathbf{a}$  exists, it is called a *separating* or *solution* vector
    - Original samples  $\mathbf{x}_1, \dots, \mathbf{x}_n$  can indeed be separated by a line then



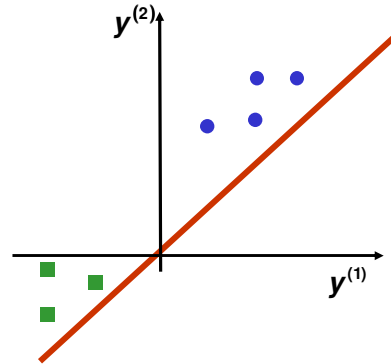
## LDF: Problem "Normalization"

before normalization



Seek a hyperplane that separates patterns from different categories

after "normalization"



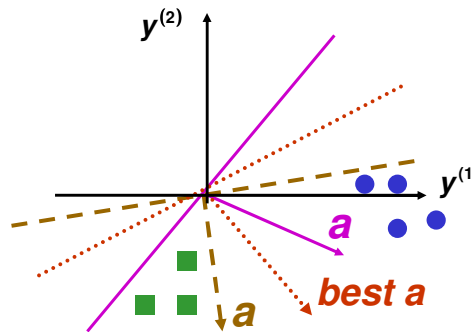
Seek hyperplane that puts *normalized* patterns on the same (positive) side



## LDF: Solution Region

- Find weight vector  $\mathbf{a}$  s.t. for all samples  $\mathbf{y}_1, \dots, \mathbf{y}_n$

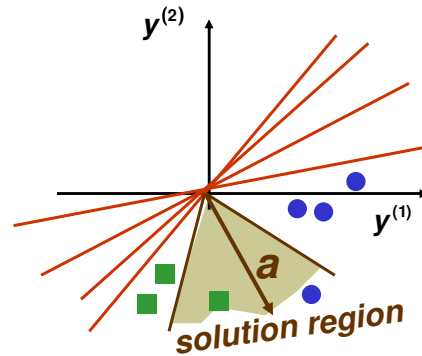
$$\mathbf{a}^t \mathbf{y}_i = \sum_{k=0}^d \mathbf{a}_k \mathbf{y}_i^{(k)} > 0$$



- In general, there are many such solutions  $\mathbf{a}$

## LDF: Solution Region

- **Solution region** for  $\mathbf{a}$ : set of all possible solutions
  - defined in terms of normal  $\mathbf{a}$  to the separating hyperplane



## Optimization

- Need to minimize a function of many variables

$$\mathbf{J}(\mathbf{x}) = \mathbf{J}(x_1, \dots, x_d)$$

- We know how to minimize  $\mathbf{J}(\mathbf{x})$ 
  - Take partial derivatives and set them to zero

$$\begin{bmatrix} \frac{\partial}{\partial x_1} \mathbf{J}(\mathbf{x}) \\ \vdots \\ \frac{\partial}{\partial x_d} \mathbf{J}(\mathbf{x}) \end{bmatrix} = \nabla \mathbf{J}(\mathbf{x}) = \mathbf{0}$$

*gradient*

- However solving analytically is not always easy

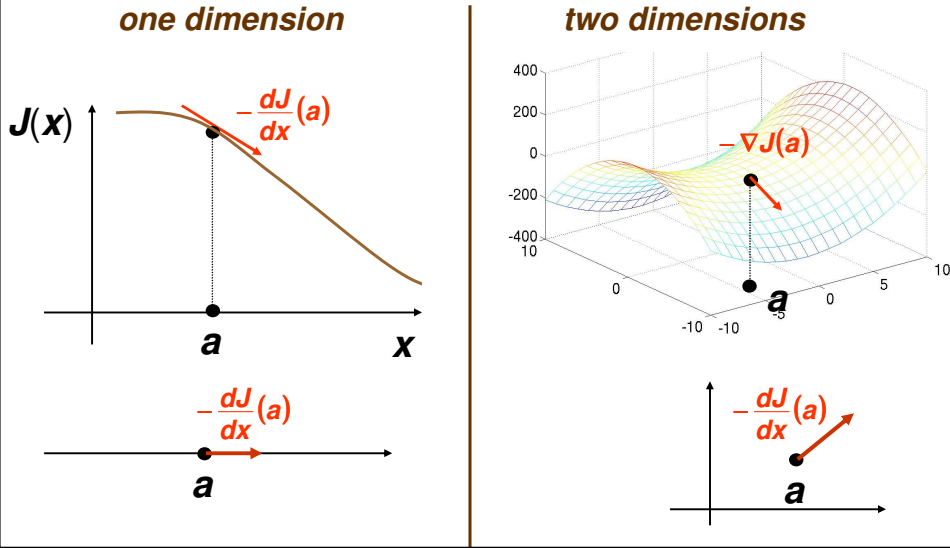
- Would you like to solve this system of nonlinear equations?

$$\begin{cases} \sin(x_1^2 + x_2^3) + e^{x_1^2} = 0 \\ \cos(x_1^2 + x_2^3) + \log(x_3^3)^{x_4^2} = 0 \end{cases}$$

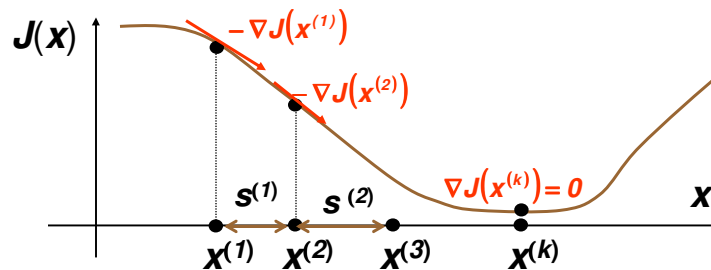
- Sometimes it is not even possible to write down an analytical expression for the derivative, we will see an example later today

## Optimization: Gradient Descent

- Gradient  $\nabla J(\mathbf{x})$  points in direction of steepest increase of  $J(\mathbf{x})$ , and  $-\nabla J(\mathbf{x})$  in direction of steepest decrease



## Optimization: Gradient Descent



**Gradient Descent** for minimizing any function  $J(\mathbf{x})$

set  $k = 1$  and  $\mathbf{x}^{(1)}$  to some initial guess for the weight vector

while  $\eta^{(k)} |\nabla J(\mathbf{x}^{(k)})| > \epsilon$

choose **learning rate**  $\eta^{(k)}$

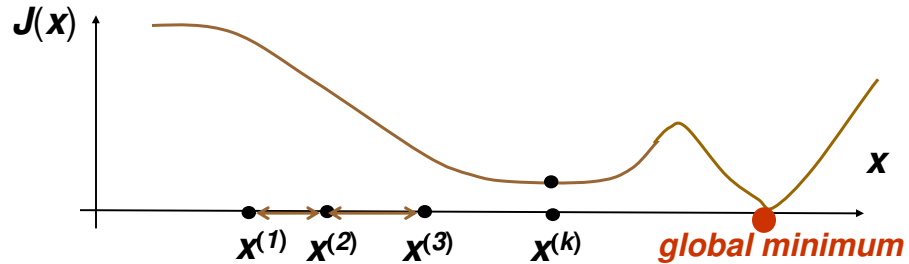
$$\mathbf{x}^{(k+1)} = \mathbf{x}^{(k)} - \eta^{(k)} \nabla J(\mathbf{x})$$

**(update rule)**

$$k = k + 1$$

## Optimization: Gradient Descent

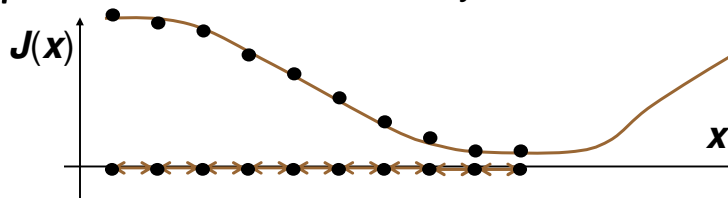
- Gradient descent is guaranteed to find only a local minimum



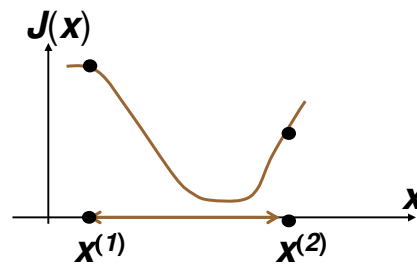
- Nevertheless gradient descent is very popular because it is simple and applicable to any function

## Optimization: Gradient Descent

- Main issue: how to set parameter  $\eta$  (**learning rate**)
- If  $\eta$  is too small, need too many iterations



- If  $\eta$  is too large may overshoot the minimum and possibly never find it (if we keep overshooting)



### LDF: Criterion Function

- Find weight vector  $\mathbf{a}$  s.t. for all samples  $\mathbf{y}_1, \dots, \mathbf{y}_n$

$$\mathbf{a}^t \mathbf{y}_i = \sum_{k=0}^d a_k y_i^{(k)} > 0$$

- Need criterion function  $\mathbf{J}(\mathbf{a})$  which is minimized when  $\mathbf{a}$  is a solution vector

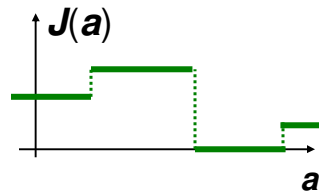
- Let  $Y_M$  be the set of examples misclassified by  $\mathbf{a}$

$$Y_M(\mathbf{a}) = \{\text{sample } \mathbf{y}_i \text{ s.t. } \mathbf{a}^t \mathbf{y}_i < 0\}$$

- First natural choice: number of misclassified examples

$$\mathbf{J}(\mathbf{a}) = |Y_M(\mathbf{a})|$$

- piecewise constant, gradient descent is useless



### LDF: Perceptron Criterion Function

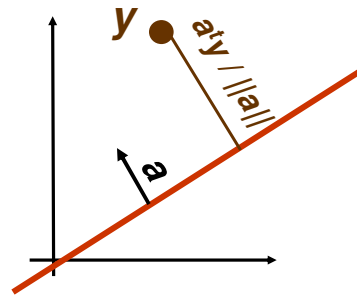
- Better choice: **Perceptron** criterion function

$$\mathbf{J}_p(\mathbf{a}) = \sum_{\mathbf{y} \in Y_M} (-\mathbf{a}^t \mathbf{y})$$

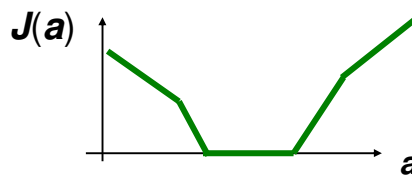
- If  $\mathbf{y}$  is misclassified,  $\mathbf{a}^t \mathbf{y} \leq 0$

- Thus  $\mathbf{J}_p(\mathbf{a}) \geq 0$

- $\mathbf{J}_p(\mathbf{a})$  is  $-\|\mathbf{a}\|$  times sum of distances of misclassified examples to decision boundary



- $\mathbf{J}_p(\mathbf{a})$  is piecewise linear and thus suitable for gradient descent



### LDF: Perceptron Batch Rule

$$J_p(\mathbf{a}) = \sum_{y \in Y_M} (-\mathbf{a}^t \mathbf{y})$$

- Gradient of  $J_p(\mathbf{a})$  is  $\nabla J_p(\mathbf{a}) = \sum_{y \in Y_M} (-\mathbf{y})$ 
  - $Y_M$  are samples misclassified by  $\mathbf{a}^{(k)}$
  - It is not possible to solve  $\nabla J_p(\mathbf{a}) = \mathbf{0}$  analytically because of  $Y_M$
- Update rule for gradient descent:  $\mathbf{x}^{(k+1)} = \mathbf{x}^{(k)} - \eta^{(k)} \nabla J(\mathbf{x})$

- Thus *gradient decent batch update rule* for  $J_p(\mathbf{a})$  is:

$$\mathbf{a}^{(k+1)} = \mathbf{a}^{(k)} + \eta^{(k)} \sum_{y \in Y_M} \mathbf{y}$$

- It is called *batch* rule because it is based on all misclassified examples

### LDF: Perceptron Single Sample Rule

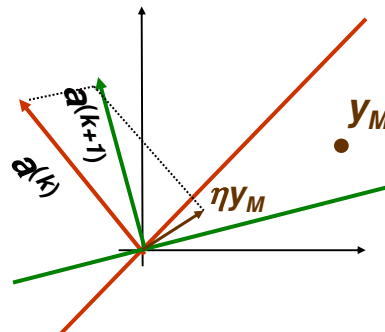
- Thus *gradient decent single sample rule* for  $J_p(\mathbf{a})$  is:

$$\mathbf{a}^{(k+1)} = \mathbf{a}^{(k)} + \eta^{(k)} \mathbf{y}_M$$

- note that  $\mathbf{y}_M$  is one sample misclassified by  $\mathbf{a}^{(k)}$
- must have a consistent way of visiting samples

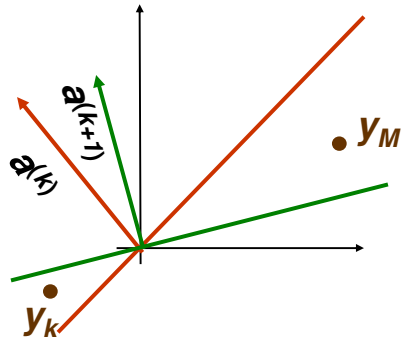
- Geometric Interpretation:

- $\mathbf{y}_M$  misclassified by  $\mathbf{a}^{(k)}$   
 $(\mathbf{a}^{(k)})^t \mathbf{y}_M \leq 0$
- $\mathbf{y}_M$  is on the wrong side of decision hyperplane
- adding  $\eta \mathbf{y}_M$  to  $\mathbf{a}$  moves new decision hyperplane in the right direction with respect to  $\mathbf{y}_M$

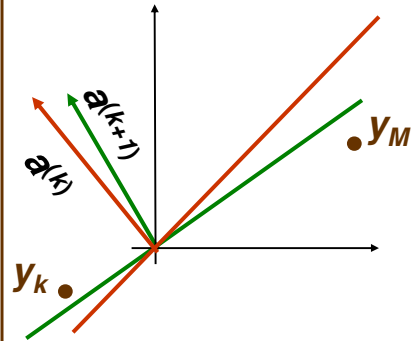


### LDF: Perceptron Single Sample Rule

$$\mathbf{a}^{(k+1)} = \mathbf{a}^{(k)} + \eta^{(k)} \mathbf{y}_M$$



$\eta$  is too large, previously correctly classified sample  $\mathbf{y}_k$  is now misclassified



$\eta$  is too small,  $\mathbf{y}_M$  is still misclassified

### LDF: Perceptron Example

name	features				grade
	good attendance?	tall?	sleeps in class?	chews gum?	
Jane	yes (1)	yes (1)	no (-1)	no (-1)	A
Steve	yes (1)	yes (1)	yes (1)	yes (1)	F
Mary	no (-1)	no (-1)	no (-1)	yes (1)	F
Peter	yes (1)	no (-1)	no (-1)	yes (1)	A

- **class 1:** students who get grade A
- **class 2:** students who get grade F

### LDF Example: Augment feature vector

	features					grade
<i>name</i>	<i>extra</i>	<i>good attendance?</i>	<i>tall?</i>	<i>sleeps in class?</i>	<i>chews gum?</i>	
Jane	1	yes (1)	yes (1)	no (-1)	no (-1)	A
Steve	1	yes (1)	yes (1)	yes (1)	yes (1)	F
Mary	1	no (-1)	no (-1)	no (-1)	yes (1)	F
Peter	1	yes (1)	no (-1)	no (-1)	yes (1)	A

- convert samples  $\mathbf{x}_1, \dots, \mathbf{x}_n$  to augmented samples  $\mathbf{y}_1, \dots, \mathbf{y}_n$  by adding a new dimension of value 1

### LDF: Perform “Normalization”

	features					grade
<i>name</i>	<i>extra</i>	<i>good attendance?</i>	<i>tall?</i>	<i>sleeps in class?</i>	<i>chews gum?</i>	
Jane	1	yes (1)	yes (1)	no (-1)	no (-1)	A
Steve	-1	yes (-1)	yes (-1)	yes (-1)	yes (-1)	F
Mary	-1	no (1)	no (1)	no (1)	yes (-1)	F
Peter	1	yes (1)	no (-1)	no (-1)	yes (1)	A

- Replace all examples from class  $\mathbf{c}_2$  by their negative

$$\mathbf{y}_i \rightarrow -\mathbf{y}_i \quad \forall \mathbf{y}_i \in \mathbf{c}_2$$

- Seek weight vector  $\mathbf{a}$  s.t.  $\mathbf{a}^t \mathbf{y}_i > 0 \quad \forall \mathbf{y}_i$



### LDF: Use Single Sample Rule

	features					grade
<i>name</i>	<i>extra</i>	<i>good attendance?</i>	<i>tall?</i>	<i>sleeps in class?</i>	<i>chews gum?</i>	
Jane	1	yes (1)	yes (1)	no (-1)	no (-1)	A
Steve	-1	yes (-1)	yes (-1)	yes (-1)	yes (-1)	F
Mary	-1	no (1)	no (1)	no (1)	yes (-1)	F
Peter	1	yes (1)	no (-1)	no (-1)	yes (1)	A

- Sample is misclassified if  $\mathbf{a}^t \mathbf{y}_i = \sum_{k=0}^4 \mathbf{a}_k \mathbf{y}_i^{(k)} < 0$
- gradient descent single sample rule:  $\mathbf{a}^{(k+1)} = \mathbf{a}^{(k)} + \eta^{(k)} \mathbf{y}_M$
- Set **fixed** learning rate to  $\eta^{(k)} = 1$ :  $\mathbf{a}^{(k+1)} = \mathbf{a}^{(k)} + \mathbf{y}_M$

### LDF: Gradient decent Example

- set equal initial weights  $\mathbf{a}^{(1)} = [0.25, 0.25, 0.25, 0.25]$
- visit all samples sequentially, modifying the weights for after finding a misclassified example

<i>name</i>	$\mathbf{a}^t \mathbf{y}$	<i>misclassified?</i>
Jane	$0.25*1+0.25*1+0.25*1+0.25*(-1)+0.25*(-1) > 0$	no
Steve	$0.25*(-1)+0.25*(-1)+0.25*(-1)+0.25*(-1)+0.25*(-1) < 0$	yes

- new weights

$$\begin{aligned} \mathbf{a}^{(2)} &= \mathbf{a}^{(1)} + \mathbf{y}_M = [0.25 \ 0.25 \ 0.25 \ 0.25 \ 0.25] + \\ &\quad + [-1 \ -1 \ -1 \ -1 \ -1] = \\ &= [-0.75 \ -0.75 \ -0.75 \ -0.75 \ -0.75] \end{aligned}$$

### LDF: Gradient decent Example

$$\mathbf{a}^{(2)} = [-0.75 \ -0.75 \ -0.75 \ -0.75 \ -0.75]$$

<i>name</i>	<i>a<sup>t</sup>y</i>	<i>misclassified?</i>
Mary	$-0.75*(-1)-0.75*1 -0.75 *1 -0.75 *1 -0.75*(-1) < 0$	yes

- new weights

$$\begin{aligned}\mathbf{a}^{(3)} &= \mathbf{a}^{(2)} + \mathbf{y}_M = [-0.75 \ -0.75 \ -0.75 \ -0.75 \ -0.75] + \\ &\quad + [-1 \ 1 \ 1 \ 1 \ -1] = \\ &= [-1.75 \ 0.25 \ 0.25 \ 0.25 \ -1.75]\end{aligned}$$

### LDF: Gradient decent Example

$$\mathbf{a}^{(3)} = [-1.75 \ 0.25 \ 0.25 \ 0.25 \ -1.75]$$

<i>name</i>	<i>a<sup>t</sup>y</i>	<i>misclassified?</i>
Peter	$-1.75 *1 +0.25* 1+0.25* (-1) +0.25 *(-1)-1.75*1 < 0$	yes

- new weights

$$\begin{aligned}\mathbf{a}^{(4)} &= \mathbf{a}^{(3)} + \mathbf{y}_M = [-1.75 \ 0.25 \ 0.25 \ 0.25 \ -1.75] + \\ &\quad + [1 \ 1 \ -1 \ -1 \ 1] = \\ &= [-0.75 \ 1.25 \ -0.75 \ -0.75 \ -0.75]\end{aligned}$$

### LDF: Gradient decent Example

$$\mathbf{a}^{(4)} = [-0.75 \quad 1.25 \quad -0.75 \quad -0.75 \quad -0.75]$$

name	$\mathbf{a}^t \mathbf{y}$	misclassified?
Jane	$-0.75 * 1 + 1.25 * 1 - 0.75 * 1 - 0.75 * (-1) - 0.75 * (-1) + 0$	no
Steve	$-0.75 * (-1) + 1.25 * (-1) - 0.75 * (-1) - 0.75 * (-1) - 0.75 * (-1) > 0$	no
Mary	$-0.75 * (-1) + 1.25 * 1 - 0.75 * 1 - 0.75 * 1 - 0.75 * (-1) > 0$	no
Peter	$-0.75 * 1 + 1.25 * 1 - 0.75 * (-1) - 0.75 * (-1) - 0.75 * 1 > 0$	no

- Thus the discriminant function is

$$g(\mathbf{y}) = -0.75 * y^{(0)} + 1.25 * y^{(1)} - 0.75 * y^{(2)} - 0.75 * y^{(3)} - 0.75 * y^{(4)}$$

- Converting back to the original features  $\mathbf{x}$ :

$$g(\mathbf{x}) = 1.25 * x^{(1)} - 0.75 * x^{(2)} - 0.75 * x^{(3)} - 0.75 * x^{(4)} - 0.75$$

### LDF: Gradient decent Example

- Converting back to the original features  $\mathbf{x}$ :

$$1.25 * x^{(1)} - 0.75 * x^{(2)} - 0.75 * x^{(3)} - 0.75 * x^{(4)} > 0.75 \Rightarrow \text{grade A}$$

$$1.25 * x^{(1)} - 0.75 * x^{(2)} - 0.75 * x^{(3)} - 0.75 * x^{(4)} < 0.75 \Rightarrow \text{grade F}$$



- This is just one possible solution vector

- If we started with weights  $\mathbf{a}^{(1)} = [0, 0.5, 0.5, 0, 0]$ , solution would be  $[-1, 1.5, -0.5, -1, -1]$

$$1.5 * x^{(1)} - 0.5 * x^{(2)} - x^{(3)} - x^{(4)} > 1 \Rightarrow \text{grade A}$$

$$1.5 * x^{(1)} - 0.5 * x^{(2)} - x^{(3)} - x^{(4)} < 1 \Rightarrow \text{grade F}$$

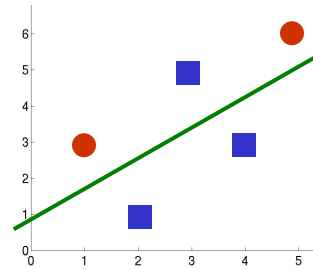
- In this solution, being tall is the least important feature

## LDF: Nonseparable Example

- Suppose we have 2 features and samples are:

- Class 1: [2,1], [4,3], [3,5]
- Class 2: [1,3] and [5,6]

- These samples are not separable by a line



- Still would like to get approximate separation by a line, good choice is shown in green

- some samples may be “noisy”, and it’s ok if they are on the wrong side of the line

- Get  $\mathbf{y}_1, \mathbf{y}_2, \mathbf{y}_3, \mathbf{y}_4$  by adding extra feature and “normalizing”

$$\mathbf{y}_1 = \begin{bmatrix} 1 \\ 2 \\ 1 \end{bmatrix} \quad \mathbf{y}_2 = \begin{bmatrix} 1 \\ 4 \\ 3 \end{bmatrix} \quad \mathbf{y}_3 = \begin{bmatrix} 1 \\ 3 \\ 5 \end{bmatrix} \quad \mathbf{y}_4 = \begin{bmatrix} -1 \\ -1 \\ -3 \end{bmatrix} \quad \mathbf{y}_5 = \begin{bmatrix} -1 \\ -5 \\ -6 \end{bmatrix}$$

## LDF: Nonseparable Example

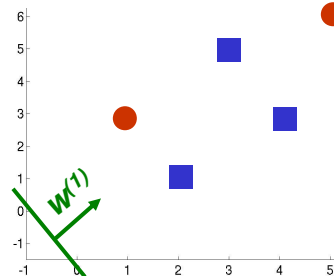
- Let’s apply Perceptron single sample algorithm

- initial equal weights  $\mathbf{a}^{(1)} = [1 \ 1 \ 1]$

- this is line  $\mathbf{x}^{(1)} + \mathbf{x}^{(2)} + 1 = 0$

- fixed learning rate  $\eta = 1$

$$\mathbf{a}^{(k+1)} = \mathbf{a}^{(k)} + \mathbf{y}_M$$



$$\mathbf{y}_1 = \begin{bmatrix} 1 \\ 2 \\ 1 \end{bmatrix} \quad \mathbf{y}_2 = \begin{bmatrix} 1 \\ 4 \\ 3 \end{bmatrix} \quad \mathbf{y}_3 = \begin{bmatrix} 1 \\ 3 \\ 5 \end{bmatrix} \quad \mathbf{y}_4 = \begin{bmatrix} -1 \\ -1 \\ -3 \end{bmatrix} \quad \mathbf{y}_5 = \begin{bmatrix} -1 \\ -5 \\ -6 \end{bmatrix}$$

- $\mathbf{y}_1^t \mathbf{a}^{(1)} = [1 \ 1 \ 1] * [1 \ 2 \ 1]^t > 0$  ✓

- $\mathbf{y}_2^t \mathbf{a}^{(1)} = [1 \ 1 \ 1] * [1 \ 4 \ 3]^t > 0$  ✓

- $\mathbf{y}_3^t \mathbf{a}^{(1)} = [1 \ 1 \ 1] * [1 \ 3 \ 5]^t > 0$  ✓

### LDF: Nonseparable Example

$$\mathbf{a}^{(1)} = [1 \ 1 \ 1] \quad \mathbf{a}^{(k+1)} = \mathbf{a}^{(k)} + \mathbf{y}_M$$

$$\mathbf{y}_1 = \begin{bmatrix} 1 \\ 2 \\ 1 \end{bmatrix} \quad \mathbf{y}_2 = \begin{bmatrix} 1 \\ 4 \\ 3 \end{bmatrix} \quad \mathbf{y}_3 = \begin{bmatrix} 1 \\ 3 \\ 5 \end{bmatrix} \quad \mathbf{y}_4 = \begin{bmatrix} -1 \\ -1 \\ -3 \end{bmatrix} \quad \mathbf{y}_5 = \begin{bmatrix} -1 \\ -5 \\ -6 \end{bmatrix}$$

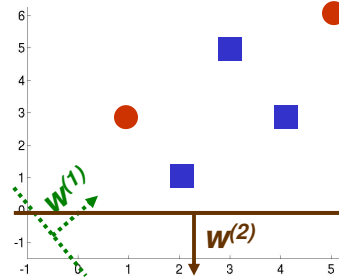
- $\mathbf{y}_4^t \mathbf{a}^{(1)} = [1 \ 1 \ 1]^* [-1 \ -1 \ -3]^t = -5 < 0$

$$\mathbf{a}^{(2)} = \mathbf{a}^{(1)} + \mathbf{y}_M = [1 \ 1 \ 1] + [-1 \ -1 \ -3] = [0 \ 0 \ -2]$$

- $\mathbf{y}_5^t \mathbf{a}^{(2)} = [0 \ 0 \ -2]^* [-1 \ -5 \ -6]^t = 12 > 0 \quad \checkmark$

- $\mathbf{y}_1^t \mathbf{a}^{(2)} = [0 \ 0 \ -2]^* [1 \ 2 \ 1]^t < 0$

$$\mathbf{a}^{(3)} = \mathbf{a}^{(2)} + \mathbf{y}_M = [0 \ 0 \ -2] + [1 \ 2 \ 1] = [1 \ 2 \ -1]$$



### LDF: Nonseparable Example

$$\mathbf{a}^{(3)} = [1 \ 2 \ -1] \quad \mathbf{a}^{(k+1)} = \mathbf{a}^{(k)} + \mathbf{y}_M$$

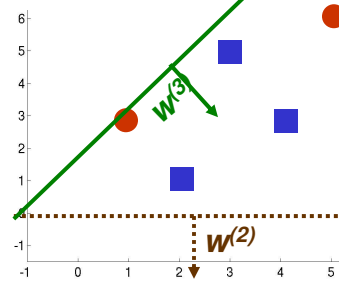
$$\mathbf{y}_1 = \begin{bmatrix} 1 \\ 2 \\ 1 \end{bmatrix} \quad \mathbf{y}_2 = \begin{bmatrix} 1 \\ 4 \\ 3 \end{bmatrix} \quad \mathbf{y}_3 = \begin{bmatrix} 1 \\ 3 \\ 5 \end{bmatrix} \quad \mathbf{y}_4 = \begin{bmatrix} -1 \\ -1 \\ -3 \end{bmatrix} \quad \mathbf{y}_5 = \begin{bmatrix} -1 \\ -5 \\ -6 \end{bmatrix}$$

- $\mathbf{y}_2^t \mathbf{a}^{(3)} = [1 \ 4 \ 3]^* [1 \ 2 \ -1]^t = 6 > 0 \quad \checkmark$

- $\mathbf{y}_3^t \mathbf{a}^{(3)} = [1 \ 3 \ 5]^* [1 \ 2 \ -1]^t > 0 \quad \checkmark$

- $\mathbf{y}_4^t \mathbf{a}^{(3)} = [-1 \ -1 \ -3]^* [1 \ 2 \ -1]^t = 0$

$$\mathbf{a}^{(4)} = \mathbf{a}^{(3)} + \mathbf{y}_M = [1 \ 2 \ -1] + [-1 \ -1 \ -3] = [0 \ 1 \ -4]$$



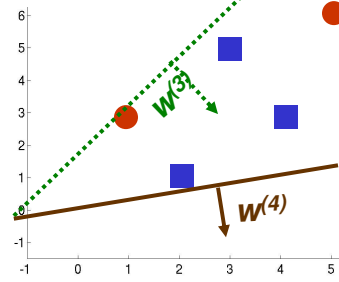
## LDF: Nonseparable Example

$$\mathbf{a}^{(4)} = [0 \ 1 \ -4] \quad \mathbf{a}^{(k+1)} = \mathbf{a}^{(k)} + \mathbf{y}_M$$

$$\mathbf{y}_1 = \begin{bmatrix} 1 \\ 2 \\ 1 \end{bmatrix} \quad \mathbf{y}_2 = \begin{bmatrix} 1 \\ 4 \\ 3 \end{bmatrix} \quad \mathbf{y}_3 = \begin{bmatrix} 1 \\ 3 \\ 5 \end{bmatrix} \quad \mathbf{y}_4 = \begin{bmatrix} -1 \\ -1 \\ -3 \end{bmatrix} \quad \mathbf{y}_5 = \begin{bmatrix} -1 \\ -5 \\ -6 \end{bmatrix}$$

- $\mathbf{y}_2^t \mathbf{a}^{(3)} = [1 \ 4 \ 3]^t [1 \ 2 \ -1]^t = 6 > 0 \checkmark$
- $\mathbf{y}_3^t \mathbf{a}^{(3)} = [1 \ 3 \ 5]^t [1 \ 2 \ -1]^t > 0 \checkmark$
- $\mathbf{y}_4^t \mathbf{a}^{(3)} = [-1 \ -1 \ -3]^t [1 \ 2 \ -1]^t = 0$

$$\mathbf{a}^{(4)} = \mathbf{a}^{(3)} + \mathbf{y}_M = [1 \ 2 \ -1]^t + [-1 \ -1 \ -3]^t = [0 \ 1 \ -4]^t$$



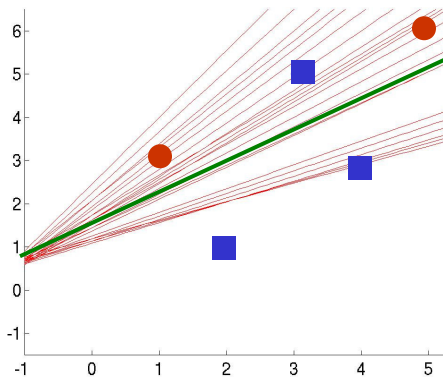
## LDF: Nonseparable Example

- we can continue this forever
  - there is no solution vector  $\mathbf{a}$  satisfying for all  $i$

$$\mathbf{a}^t \mathbf{y}_i = \sum_{k=0}^5 \mathbf{a}_k \mathbf{y}_i^{(k)} > 0$$

- need to stop but at a good point:

- solutions at iterations 900 through 915. Some are good some are not.
- How do we stop at a good solution?



### ***LDF: Convergence of Perceptron rules***

- If classes are linearly separable, and use fixed learning rate, that is for some constant  $\mathbf{c}$ ,  $\eta^{(k)} = \mathbf{c}$ 
  - *both single sample and batch perceptron rules converge to a correct solution* (could be any  $\mathbf{a}$  in the solution space)
- If classes are not linearly separable:
  - algorithm does not stop, it keeps looking for solution which does not exist
  - by choosing appropriate learning rate, can always ensure convergence:  $\eta^{(k)} \rightarrow \mathbf{0}$  as  $k \rightarrow \infty$
  - for example inverse linear learning rate:  $\eta^{(k)} = \frac{\eta^{(1)}}{k}$
  - for inverse linear learning rate convergence in the linearly separable case can also be proven
  - no guarantee that we stopped at a good point, but there are good reasons to choose inverse linear learning rate

### ***LDF: Perceptron Rule and Gradient decent***

- Linearly separable data
  - perceptron rule with gradient decent works well
- Linearly non-separable data
  - need to stop perceptron rule algorithm at a good point, this maybe tricky

#### ***Batch Rule***

- Smoother gradient because all samples are used

#### ***Single Sample Rule***

- easier to analyze
- Concentrates more than necessary on any isolated “noisy” training examples

## LDF: Minimum Squared-Error Procedures

- Idea: convert to easier and better understood problem

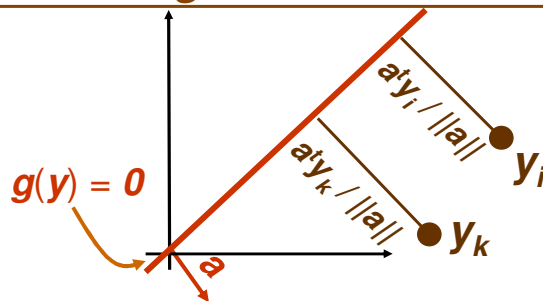
$\mathbf{a}^t \mathbf{y}_i > 0$  for all samples  $\mathbf{y}_i$   
solve system of linear inequalities



$\mathbf{a}^t \mathbf{y}_i = \mathbf{b}_i$  for all samples  $\mathbf{y}_i$   
solve system of linear equations

- MSE procedure
  - Choose **positive** constants  $\mathbf{b}_1, \mathbf{b}_2, \dots, \mathbf{b}_n$
  - try to find weight vector  $\mathbf{a}$  s.t.  $\mathbf{a}^t \mathbf{y}_i = \mathbf{b}_i$  for all samples  $\mathbf{y}_i$
  - If we can find weight vector  $\mathbf{a}$  such that  $\mathbf{a}^t \mathbf{y}_i = \mathbf{b}_i$  for all samples  $\mathbf{y}_i$ , then  $\mathbf{a}$  is a solution because  $\mathbf{b}_i$ 's are positive
  - consider all the samples (not just the misclassified ones)

## LDF: MSE Margins



- Since we want  $\mathbf{a}^t \mathbf{y}_i = \mathbf{b}_i$ , we expect sample  $\mathbf{y}_i$  to be at distance  $\mathbf{b}_i$  from the separating hyperplane (normalized by  $\|\mathbf{a}\|$ )
- Thus  $\mathbf{b}_1, \mathbf{b}_2, \dots, \mathbf{b}_n$  give relative expected distances or “margins” of samples from the hyperplane
- Should make  $\mathbf{b}_i$  small if sample  $i$  is expected to be near separating hyperplane, and make  $\mathbf{b}_i$  larger otherwise
- In the absence of any additional information, there are good reasons to set  $\mathbf{b}_1 = \mathbf{b}_2 = \dots = \mathbf{b}_n = 1$



## LDF: MSE Matrix Notation

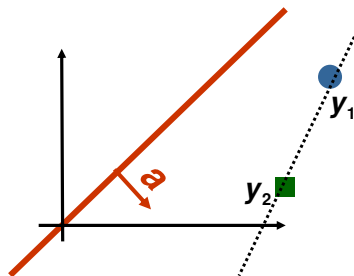
- Need to solve  $n$  equations 
$$\begin{cases} \mathbf{a}^t \mathbf{y}_1 = b_1 \\ \vdots \\ \mathbf{a}^t \mathbf{y}_n = b_n \end{cases}$$
- Using matrix notation:

$$\underbrace{\begin{bmatrix} y_1^{(0)} & y_1^{(1)} & \dots & y_1^{(d)} \\ y_2^{(0)} & y_2^{(1)} & \dots & y_2^{(d)} \\ \vdots & \vdots & \ddots & \vdots \\ y_n^{(0)} & y_n^{(1)} & \dots & y_n^{(d)} \end{bmatrix}}_{\mathbf{Y}} \underbrace{\begin{bmatrix} a_0 \\ a_1 \\ \vdots \\ a_d \end{bmatrix}}_{\mathbf{a}} = \underbrace{\begin{bmatrix} b_1 \\ b_2 \\ \vdots \\ b_n \end{bmatrix}}_{\mathbf{b}}$$

- Thus need to solve a linear system  $\mathbf{Y}\mathbf{a} = \mathbf{b}$

## LDF: Exact Solution is Rare

- Thus need to solve a linear system  $\mathbf{Y}\mathbf{a} = \mathbf{b}$ 
  - $\mathbf{Y}$  is an  $n$  by  $(d+1)$  matrix
- Exact solution can be found only if  $\mathbf{Y}$  is nonsingular and square, in which case the inverse  $\mathbf{Y}^{-1}$  exists
  - $\mathbf{a} = \mathbf{Y}^{-1}\mathbf{b}$
  - (number of samples) = (number of features + 1)
  - almost never happens in practice
  - in this case, guaranteed to find the separating hyperplane



## LDF: Approximate Solution

- Typically  $\mathbf{Y}$  is overdetermined, that is it has more rows (examples) than columns (features)
  - If it has more features than examples, should reduce dimensionality

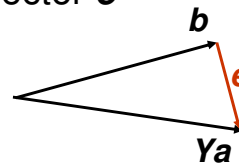
$$\mathbf{Y} \mathbf{a} = \mathbf{b}$$

- Need  $\mathbf{Y}\mathbf{a} = \mathbf{b}$ , but no exact solution exists for an overdetermined system of equation
  - More equations than unknowns
- Find an approximate solution  $\mathbf{a}$ , that is  $\mathbf{Y}\mathbf{a} \approx \mathbf{b}$ 
  - Note that approximate solution  $\mathbf{a}$  *does not* necessarily give the separating hyperplane in the separable case
  - But hyperplane corresponding to  $\mathbf{a}$  may still be a good solution, especially if there is no separating hyperplane

## LDF: MSE Criterion Function

- Minimum squared error approach: find  $\mathbf{a}$  which minimizes the length of the error vector  $\mathbf{e}$

$$\mathbf{e} = \mathbf{Y}\mathbf{a} - \mathbf{b}$$



- Thus minimize the *minimum squared error* criterion function:

$$\mathbf{J}_s(\mathbf{a}) = \|\mathbf{Y}\mathbf{a} - \mathbf{b}\|^2 = \sum_{i=1}^n (\mathbf{a}^t \mathbf{y}_i - b_i)^2$$

- Unlike the perceptron criterion function, we can optimize the minimum squared error criterion function analytically by setting the gradient to  $\mathbf{0}$

### LDF: Optimizing $J_s(\mathbf{a})$

$$J_s(\mathbf{a}) = \|\mathbf{Y}\mathbf{a} - \mathbf{b}\|^2 = \sum_{i=1}^n (\mathbf{a}^t \mathbf{y}_i - b_i)^2$$

- Let's compute the gradient:

$$\begin{aligned} \nabla J_s(\mathbf{a}) &= \begin{bmatrix} \frac{\partial J_s}{\partial a_0} \\ \vdots \\ \frac{\partial J_s}{\partial a_d} \end{bmatrix} = \frac{dJ_s}{d\mathbf{a}} = \sum_{i=1}^n \frac{d}{d\mathbf{a}} (\mathbf{a}^t \mathbf{y}_i - b_i)^2 \\ &= \sum_{i=1}^n 2(\mathbf{a}^t \mathbf{y}_i - b_i) \frac{d}{d\mathbf{a}} (\mathbf{a}^t \mathbf{y}_i - b_i) \\ &= \sum_{i=1}^n 2(\mathbf{a}^t \mathbf{y}_i - b_i) \mathbf{y}_i \\ &= 2\mathbf{Y}^t (\mathbf{Y}\mathbf{a} - \mathbf{b}) \end{aligned}$$

### LDF: Pseudo Inverse Solution

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

- Setting the gradient to 0:

$$2\mathbf{Y}^t (\mathbf{Y}\mathbf{a} - \mathbf{b}) = \mathbf{0} \Rightarrow \mathbf{Y}^t \mathbf{Y}\mathbf{a} = \mathbf{Y}^t \mathbf{b}$$

- Matrix  $\mathbf{Y}^t \mathbf{Y}$  is square (it has  $d+1$  rows and columns) and it is often non-singular
- If  $\mathbf{Y}^t \mathbf{Y}$  is non-singular, its inverse exists and we can solve for  $\mathbf{a}$  uniquely:

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

*pseudo inverse of  $\mathbf{Y}$*

$$((\mathbf{Y}^t \mathbf{Y})^{-1} \mathbf{Y}^t) \mathbf{Y} = (\mathbf{Y}^t \mathbf{Y})^{-1} (\mathbf{Y}^t \mathbf{Y}) = \mathbf{I}$$

- In matlab, simply  $\mathbf{a} = \mathbf{Y} \backslash \mathbf{b}$

### LDF: Minimum Squared-Error Procedures

- Only guaranteed the separating hyperplane if  $\mathbf{Y}\mathbf{a} > \mathbf{0}$ 
  - that is if all elements of vector  $\mathbf{Y}\mathbf{a} = \begin{bmatrix} \mathbf{a}^t \mathbf{y}_1 \\ \vdots \\ \mathbf{a}^t \mathbf{y}_n \end{bmatrix}$  are positive
- We have  $\mathbf{Y}\mathbf{a} \approx \mathbf{b}$
- That is  $\mathbf{Y}\mathbf{a} = \begin{bmatrix} \mathbf{b}_1 + \varepsilon_1 \\ \vdots \\ \mathbf{b}_n + \varepsilon_n \end{bmatrix}$  where  $\varepsilon$  may be negative
  - If  $\varepsilon_1, \dots, \varepsilon_n$  are small relative to  $\mathbf{b}_1, \dots, \mathbf{b}_n$ , then each element of  $\mathbf{Y}\mathbf{a}$  is positive, and  $\mathbf{a}$  gives a separating hyperplane
  - If approximation is not good,  $\varepsilon_i$  may be large and negative, for some  $i$ , thus  $\mathbf{b}_i + \varepsilon_i$  will be negative and  $\mathbf{a}$  is not a separating hyperplane
- Thus in linearly separable case, least squares solution  $\mathbf{a}$  does *not necessarily* give separating hyperplane
- But it will give a “reasonable” hyperplane

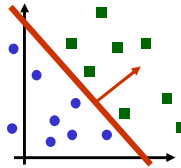
### LDF: Minimum Squared-Error Procedures

- We are free to choose  $\mathbf{b}$ . May be tempted to make  $\mathbf{b}$  large as a way to insure  $\mathbf{Y}\mathbf{a} \approx \mathbf{b} > \mathbf{0}$
- Does not work
  - Let  $\beta$  be a scalar, let's try  $\beta\mathbf{b}$  instead of  $\mathbf{b}$
  - if  $\mathbf{a}^*$  is a least squares solution to  $\mathbf{Y}\mathbf{a} = \mathbf{b}$ , then for any scalar  $\beta$ , least squares solution to  $\mathbf{Y}\mathbf{a} = \beta\mathbf{b}$  is  $\beta\mathbf{a}^*$ 

$$\begin{aligned} \arg \min_a \|\mathbf{Y}\mathbf{a} - \beta\mathbf{b}\|^2 &= \arg \min_a \beta^2 \|\mathbf{Y}(\mathbf{a} / \beta) - \mathbf{b}\|^2 \\ &= \arg \min_a \|\mathbf{Y}(\mathbf{a} / \beta) - \mathbf{b}\|^2 = \beta\mathbf{a}^* \end{aligned}$$
  - thus if for some  $i$ th element of  $\mathbf{Y}\mathbf{a}$  is less than 0, that is  $\mathbf{y}_i^t \mathbf{a} < 0$ , then  $\mathbf{y}_i^t (\beta\mathbf{a}) < 0$ ,
- Relative difference between components of  $\mathbf{b}$  matters, but not the size of each individual component

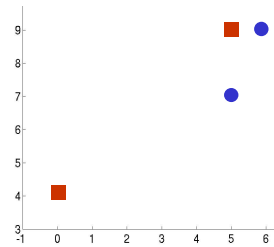
## LDF: How to choose $b$ in MSE Procedure?

- So far we assumed that constants  $b_1, b_2, \dots, b_n$  are positive but otherwise arbitrary
- Good choice is  $b_1 = b_2 = \dots = b_n = 1$ .
  - Interesting theoretical properties that we are not going to look at



## LDF: Example

- Class 1: (6 9), (5 7)
- Class 2: (5 9), (0 4)
- Set vectors  $y_1, y_2, y_3, y_4$  by adding extra feature and “normalizing”



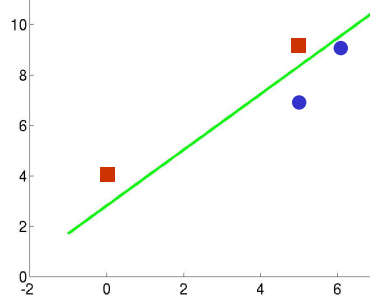
$$y_1 = \begin{bmatrix} 1 \\ 6 \\ 9 \end{bmatrix} \quad y_2 = \begin{bmatrix} 1 \\ 5 \\ 7 \end{bmatrix} \quad y_3 = \begin{bmatrix} -1 \\ -5 \\ -9 \end{bmatrix} \quad y_4 = \begin{bmatrix} -1 \\ 0 \\ -4 \end{bmatrix}$$

- Matrix  $Y$  is then 
$$Y = \begin{bmatrix} 1 & 6 & 9 \\ 1 & 5 & 7 \\ -1 & -5 & -9 \\ -1 & 0 & -4 \end{bmatrix}$$

### LDF: Example

- Choose  $\mathbf{b} = \begin{bmatrix} 1 \\ 1 \\ 1 \\ 1 \end{bmatrix}$
- In matlab,  $\mathbf{a} = \mathbf{Y} \backslash \mathbf{b}$  solves the least squares problem

$$\mathbf{a} = \begin{bmatrix} 2.7 \\ 1.0 \\ -0.9 \end{bmatrix}$$



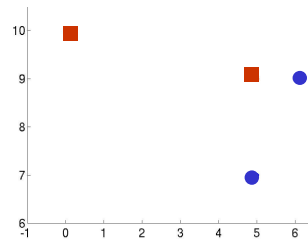
- Note  $\mathbf{a}$  is an approximation to  $\mathbf{Y}\mathbf{a} = \mathbf{b}$ , since no exact solution exists

$$\mathbf{Y}\mathbf{a} = \begin{bmatrix} 0.4 \\ 1.3 \\ 0.6 \\ 1.1 \end{bmatrix} \neq \begin{bmatrix} 1 \\ 1 \\ 1 \\ 1 \end{bmatrix}$$

- This solution does give a separating hyperplane since  $\mathbf{Y}\mathbf{a} > \mathbf{0}$

### LDF: Example

- Class 1: (6 9), (5 7)
- Class 2: (5 9), (0 10)
- The last sample is very far compared to others from the separating hyperplane



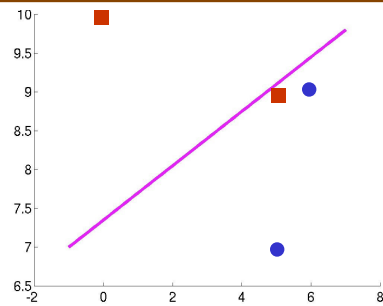
$$\mathbf{y}_1 = \begin{bmatrix} 1 \\ 6 \\ 9 \end{bmatrix} \quad \mathbf{y}_2 = \begin{bmatrix} 1 \\ 5 \\ 7 \end{bmatrix} \quad \mathbf{y}_3 = \begin{bmatrix} -1 \\ -5 \\ -9 \end{bmatrix} \quad \mathbf{y}_4 = \begin{bmatrix} -1 \\ 0 \\ -10 \end{bmatrix}$$

- Matrix  $\mathbf{Y} = \begin{bmatrix} 1 & 6 & 9 \\ 1 & 5 & 7 \\ -1 & -5 & -9 \\ -1 & 0 & -10 \end{bmatrix}$

### LDF: Example

- Choose  $\mathbf{b} = \begin{bmatrix} 1 \\ 1 \\ 1 \\ 1 \end{bmatrix}$
- In matlab,  $\mathbf{a} = \mathbf{Y} \backslash \mathbf{b}$  solves the least squares problem

$$\mathbf{a} = \begin{bmatrix} 3.2 \\ 0.2 \\ -0.4 \end{bmatrix}$$



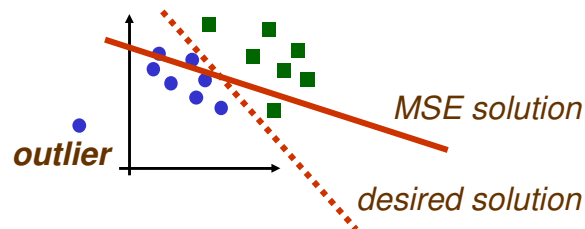
- Note  $\mathbf{a}$  is an approximation to  $\mathbf{Y}\mathbf{a} = \mathbf{b}$ , since no exact solution exists

$$\mathbf{Y}\mathbf{a} = \begin{bmatrix} 0.2 \\ 0.9 \\ -0.04 \\ 1.16 \end{bmatrix} \neq \begin{bmatrix} 1 \\ 1 \\ 1 \\ 1 \end{bmatrix}$$

- This solution does not give a separating hyperplane since  $\mathbf{a}^t \mathbf{y}_3 < 0$

### LDF: Example

- MSE pays too much attention to isolated “noisy” examples (such examples are called outliers)



- No problems with convergence though, and solution it gives ranges from reasonable to good

## LDF: Example

- we know that 4<sup>th</sup> point is far far from separating hyperplane
  - In practice look at points which are furthest from the decision boundary

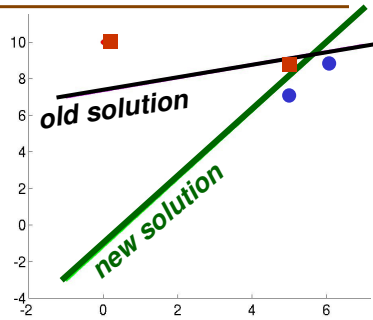
$$Ya = \begin{bmatrix} 0.2 \\ 0.9 \\ -0.04 \\ 1.16 \end{bmatrix}$$

- Set  $b_i$  larger for such points:  $b = \begin{bmatrix} 1 \\ 1 \\ 1 \\ 10 \end{bmatrix}$

- In Matlab, solve  $a = Yb$

$$a = \begin{bmatrix} -1.1 \\ 1.7 \\ -0.9 \end{bmatrix}$$

- Note  $a$  is an approximation to  $Ya = b$ , since  $Ya = \begin{bmatrix} 0.9 \\ 1.0 \\ 0.8 \\ 10.0 \end{bmatrix} \neq \begin{bmatrix} 1 \\ 1 \\ 1 \\ 10 \end{bmatrix}$
- This solution does give the separating hyperplane since  $Ya > 0$



## LDF: Gradient Descent for MSE solution

$$J_s(a) = \|Ya - b\|^2$$

- May wish to find MSE solution by gradient descent:
  - Computing the inverse of  $Y^tY$  may be too costly
  - $Y^tY$  may be close to singular if samples are highly correlated (rows of  $Y$  are almost linear combinations of each other)
    - computing the inverse of  $Y^tY$  is not numerically stable
- In the beginning of the lecture, computed the gradient:

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



### ***LDF: Widrow-Hoff Procedure***

---

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

- Thus the update rule for gradient descent:

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

- If  $\eta^{(k)} = \eta^{(1)} / k$  weight vector  $\mathbf{a}^{(k)}$  converges to the MSE solution  $\mathbf{a}$ , that is  $\mathbf{Y}^t(\mathbf{Y}\mathbf{a} - \mathbf{b}) = 0$
- *Widrow-Hoff procedure* reduces storage requirements by considering single samples sequentially:

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

### ***LDF: MSE for Multiple Classes***

---

- Suppose we have  $m$  classes
- Define  $m$  linear discriminant functions

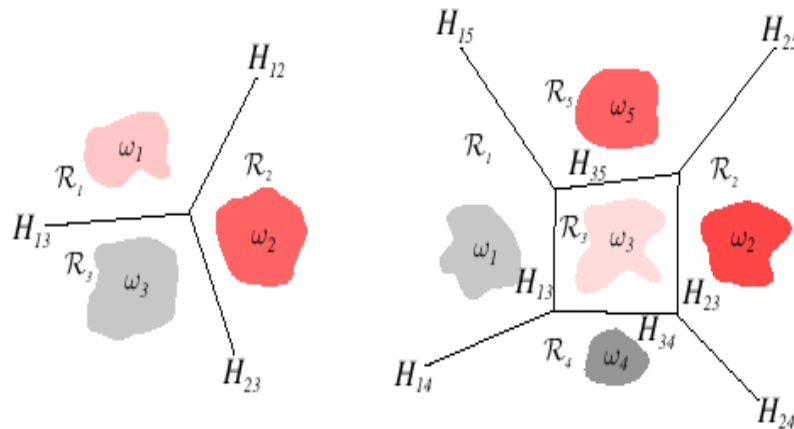
$$\mathbf{g}_i(\mathbf{x}) = \mathbf{w}_i^t \mathbf{x} + w_{i0} \quad i = 1, \dots, m$$

- Given  $\mathbf{x}$ , assign class  $c_i$  if

$$\mathbf{g}_i(\mathbf{x}) \geq \mathbf{g}_j(\mathbf{x}) \quad \forall j \neq i$$

- Such classifier is called a *linear machine*
- A linear machine divides the feature space into  $c$  decision regions, with  $\mathbf{g}_i(\mathbf{x})$  being the largest discriminant if  $\mathbf{x}$  is in the region  $R_i$

## LDF: Many Classes



## LDF: MSE for Multiple Classes

- We still use augmented feature vectors  $\mathbf{y}_1, \dots, \mathbf{y}_n$ 
  - but do not multiply by -1
- Define  $m$  linear discriminant functions

$$\mathbf{g}_i(\mathbf{y}) = \mathbf{a}_i^t \mathbf{y} \quad i = 1, \dots, m$$

- Given  $\mathbf{y}$ , assign class  $c_i$  if

$$\mathbf{a}_i^t \mathbf{y} \geq \mathbf{a}_j^t \mathbf{y} \quad \forall j \neq i$$

- For each class  $i$ , makes sense to seek weight vector  $\mathbf{a}_i$ , s.t.

$$\begin{cases} \mathbf{a}_i^t \mathbf{y} = 1 & \forall \mathbf{y} \in \text{class } i \\ \mathbf{a}_i^t \mathbf{y} = 0 & \forall \mathbf{y} \notin \text{class } i \end{cases}$$

- If we find such  $\mathbf{a}_1, \dots, \mathbf{a}_m$  the training error will be  $0$

### LDF: MSE for Multiple Classes

- For each class  $i$ , find weight vector  $\mathbf{a}_i$ , s.t.

$$\begin{cases} \mathbf{a}_i^t \mathbf{y} = 1 & \forall \mathbf{y} \in \text{class } i \\ \mathbf{a}_i^t \mathbf{y} = 0 & \forall \mathbf{y} \notin \text{class } i \end{cases}$$

- We can solve for each  $\mathbf{a}_i$  independently
- Let  $n_i$  be the number of samples in class  $i$
- Let  $\mathbf{Y}_i$  be matrix whose rows are samples from class  $i$ , so it has  $d+1$  columns and  $n_i$  rows
- Let's pile all samples in  $n$  by  $d+1$  matrix  $\mathbf{Y}$ :

$$\mathbf{Y} = \begin{bmatrix} \mathbf{Y}_1 \\ \mathbf{Y}_2 \\ \vdots \\ \mathbf{Y}_m \end{bmatrix} = \begin{bmatrix} \text{sample from class 1} \\ \text{sample from class 1} \\ \vdots \\ \text{sample from class } m \\ \text{sample from class } m \end{bmatrix}$$

### LDF: MSE for Multiple Classes

- Let  $\mathbf{b}_i$  be a column vector of length  $n$  which is  $\mathbf{0}$  everywhere except rows corresponding to samples from class  $i$ , where it is  $\mathbf{1}$ :

$$\mathbf{b}_i = \begin{bmatrix} 0 \\ \vdots \\ 1 \\ \vdots \\ 1 \\ \vdots \\ 0 \end{bmatrix} \left. \vphantom{\begin{bmatrix} 0 \\ \vdots \\ 1 \\ \vdots \\ 1 \\ \vdots \\ 0 \end{bmatrix}} \right\} \text{rows corresponding to samples from class } i$$

- We need to solve:  $\mathbf{Y}\mathbf{a}_i = \mathbf{b}_i$

$$\begin{bmatrix} \text{sample from class 1} \\ \text{sample from class 1} \\ \vdots \\ \text{sample from class } m \\ \text{sample from class } m \end{bmatrix} \begin{bmatrix} \text{weights } \mathbf{a}_i \end{bmatrix} = \begin{bmatrix} 0 \\ \vdots \\ 1 \\ \vdots \\ 1 \\ \vdots \\ 0 \end{bmatrix}$$

### **LDF: MSE for Multiple Classes**

- We need to solve  $Y\mathbf{a}_i = \mathbf{b}_i$
- Usually no exact solution since  $Y$  is overdetermined
- Use least squares to minimize norm of the error vector  $\|Y\mathbf{a}_i - \mathbf{b}_i\|$
- LSE solution with pseudoinverse:
 
$$\mathbf{a}_i = (Y^t Y)^{-1} Y^t \mathbf{b}_i$$
- Thus we need to solve  $m$  LSE problems, one for each class
- Can write these  $m$  LSE problems in one matrix

### **LDF: MSE for Multiple Classes**

- Let's pile all  $\mathbf{b}_i$  as columns in  $n$  by  $c$  matrix  $B$
- $$B = [\mathbf{b}_1 \ \dots \ \mathbf{b}_n]$$
- Let's pile all  $\mathbf{a}_i$  as columns in  $d + 1$  by  $m$  matrix  $A$

$$A = [\mathbf{a}_1 \ \dots \ \mathbf{a}_m] = \begin{bmatrix} \text{weights } a_1 \\ \text{weights } a_2 \\ \vdots \\ \text{weights } a_m \end{bmatrix}$$

- $m$  LSE problems can be represented in  $YA = B$ :

$$\begin{bmatrix} \text{sample from class1} \\ \text{sample from class1} \\ \text{sample from class2} \\ \text{sample from class3} \\ \text{sample from class3} \\ \text{sample from class3} \end{bmatrix} \begin{bmatrix} \text{weights for c1} \\ \text{weights for c2} \\ \text{weights for c3} \end{bmatrix} = \begin{bmatrix} 1 & 0 & 0 \\ 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \\ 0 & 0 & 1 \\ 0 & 0 & 1 \end{bmatrix}$$

**Y**                      **A**                      **B**

### ***LDF: MSE for Multiple Classes***

---

- Our objective function is:

$$\mathbf{J}(\mathbf{A}) = \sum_{i=1}^m \|\mathbf{Y}\mathbf{a}_i - \mathbf{b}_i\|^2$$

- $\mathbf{J}(\mathbf{A})$  is minimized with the use of pseudoinverse

$$\mathbf{A} = (\mathbf{Y}^t \mathbf{Y})^{-1} \mathbf{Y}\mathbf{B}$$

### ***LDF: Summary***

---

- ***Perceptron*** procedures
  - find a separating hyperplane in the linearly separable case,
  - do not converge in the non-separable case
  - can force convergence by using a decreasing learning rate, but are not guaranteed a reasonable stopping point
- ***MSE*** procedures
  - converge in separable and not separable case
  - may not find separating hyperplane if classes are linearly separable
  - use pseudoinverse if  $\mathbf{Y}^t \mathbf{Y}$  is not singular and not too large
  - use gradient descent (Widrow-Hoff procedure) otherwise