



Rensselaer

Lecture 10: Linear Discriminant Functions (2)

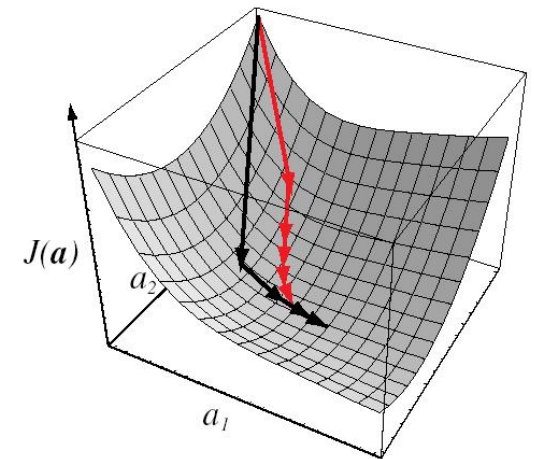
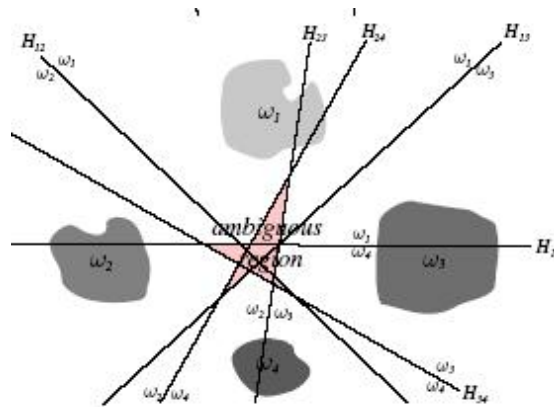
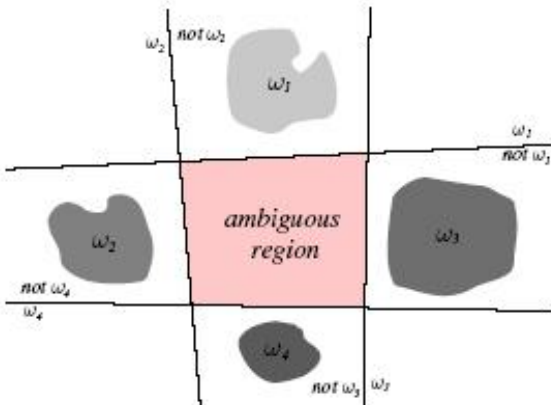
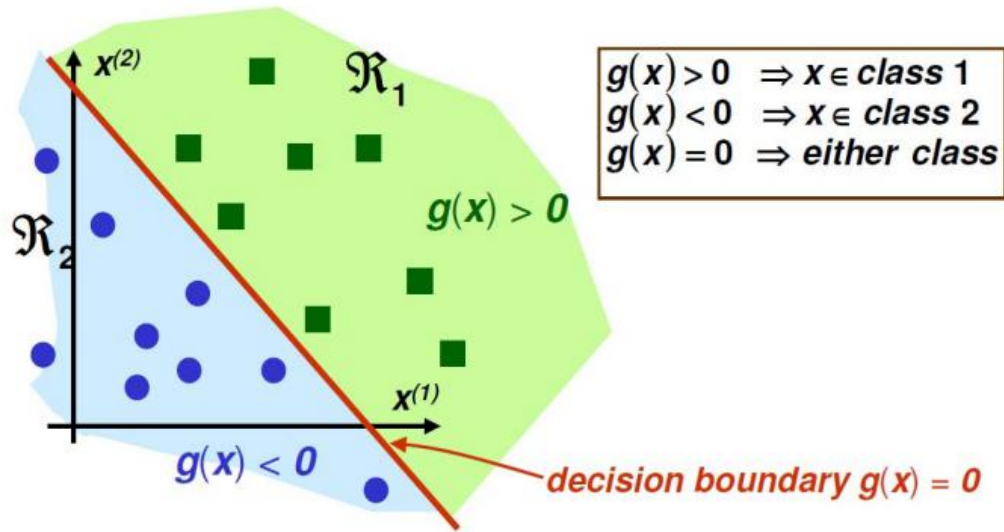
Dr. Chengjiang Long

Computer Vision Researcher at Kitware Inc.

Adjunct Professor at RPI.

Email: longc3@rpi.edu

Recap Previous Lecture



Outline

- Perceptron Rule
- Minimum Squared-Error Procedure
- Ho–Kashyap Procedure

Outline

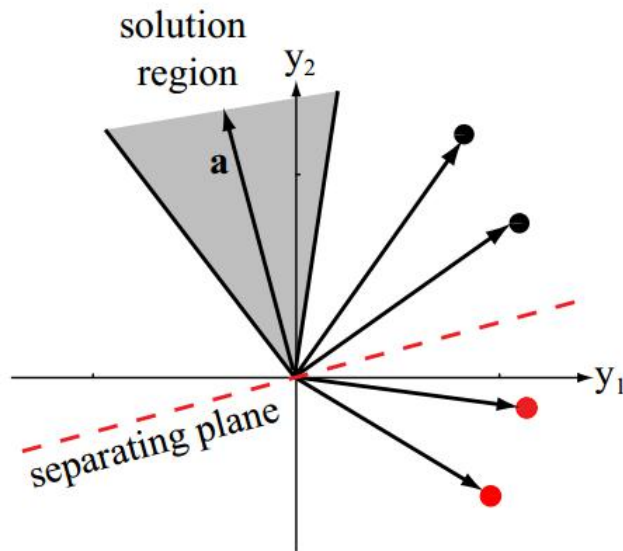
- **Perceptron Rule**
- Minimum Squared-Error Procedure
- Ho–Kashyap Procedure

"Dual" Problem

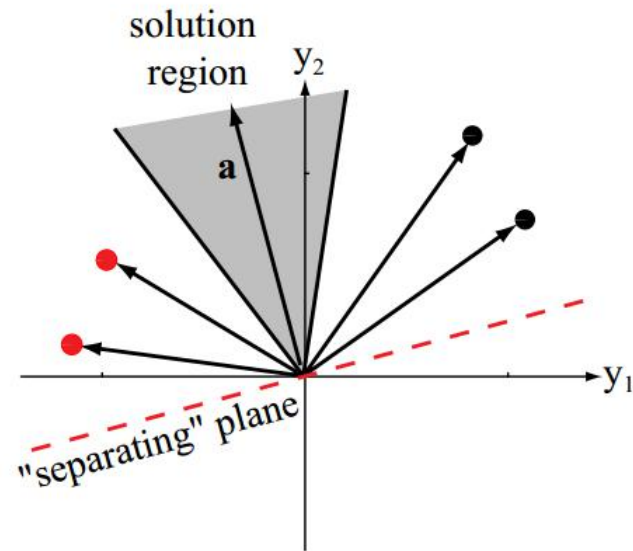
Classification rule:

If $\alpha^t \mathbf{y}_i > 0$ assign \mathbf{y}_i to ω_1
else if $\alpha^t \mathbf{y}_i < 0$ assign \mathbf{y}_i to ω_2

- If \mathbf{y}_i in ω_2 , replace \mathbf{y}_i by $-\mathbf{y}_i$
- Find α such that: $\alpha^t \mathbf{y}_i > 0$



Seek a hyperplane that **separates** patterns from different categories



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

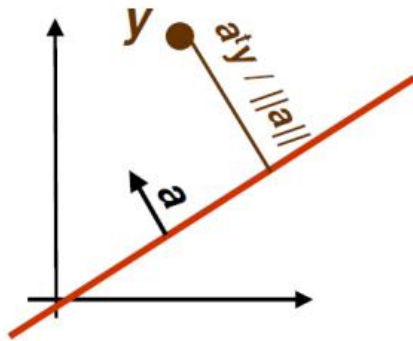
Perceptron rule

- Use Gradient Descent assuming that the error function to be minimized is:

$$J_p(\boldsymbol{\alpha}) = \sum_{\mathbf{y} \in Y(\boldsymbol{\alpha})} (-\boldsymbol{\alpha}^t \mathbf{y})$$

the set of samples
misclassified by $\boldsymbol{\alpha}$.

- If $Y(\boldsymbol{\alpha})$ is empty, $J_p(\boldsymbol{\alpha})=0$; otherwise, $J_p(\boldsymbol{\alpha}) \geq 0$.
- $J_p(\boldsymbol{\alpha})$ is $\|\boldsymbol{\alpha}\|$ times the sum of distances of misclassified.
- $J_p(\boldsymbol{\alpha})$ is piecewise linear and thus suitable for gradient descent.



Perceptron Batch Rule

- The gradient of $J_p(\mathbf{\alpha})$ is:

$$J_p(\mathbf{\alpha}) = \sum_{\mathbf{y} \in Y(\mathbf{\alpha})} (-\mathbf{\alpha}^t \mathbf{y}) \quad \Rightarrow \quad \nabla J_p = \sum_{\mathbf{y} \in Y(\mathbf{\alpha})} (-\mathbf{y})$$

- It is not possible to solve analytically $\nabla J_p = 0$.
- The perceptron update rule is obtained using gradient descent:

$$\mathbf{\alpha}(k+1) = \mathbf{\alpha}(k) + \eta(k) \sum_{\mathbf{y} \in Y(\mathbf{\alpha})} \mathbf{y}$$

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

Perceptron Single Sample Rule

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

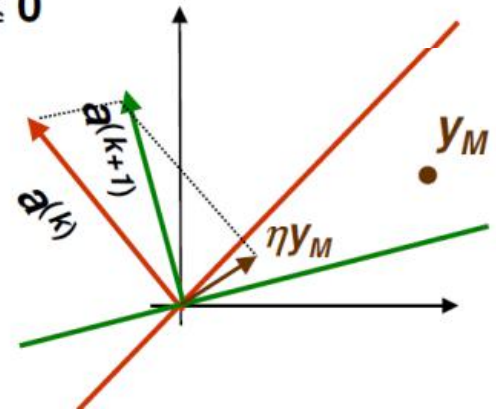
$$\mathbf{a}^{(k+1)} = \mathbf{a}^{(k)} + \eta^{(k)} \mathbf{y}_M$$

- Note that y_M is one sample misclassified by $\mathbf{a}^{(k)}$
- Must have a consistent way of visiting samples

- Geometric Interpretation:

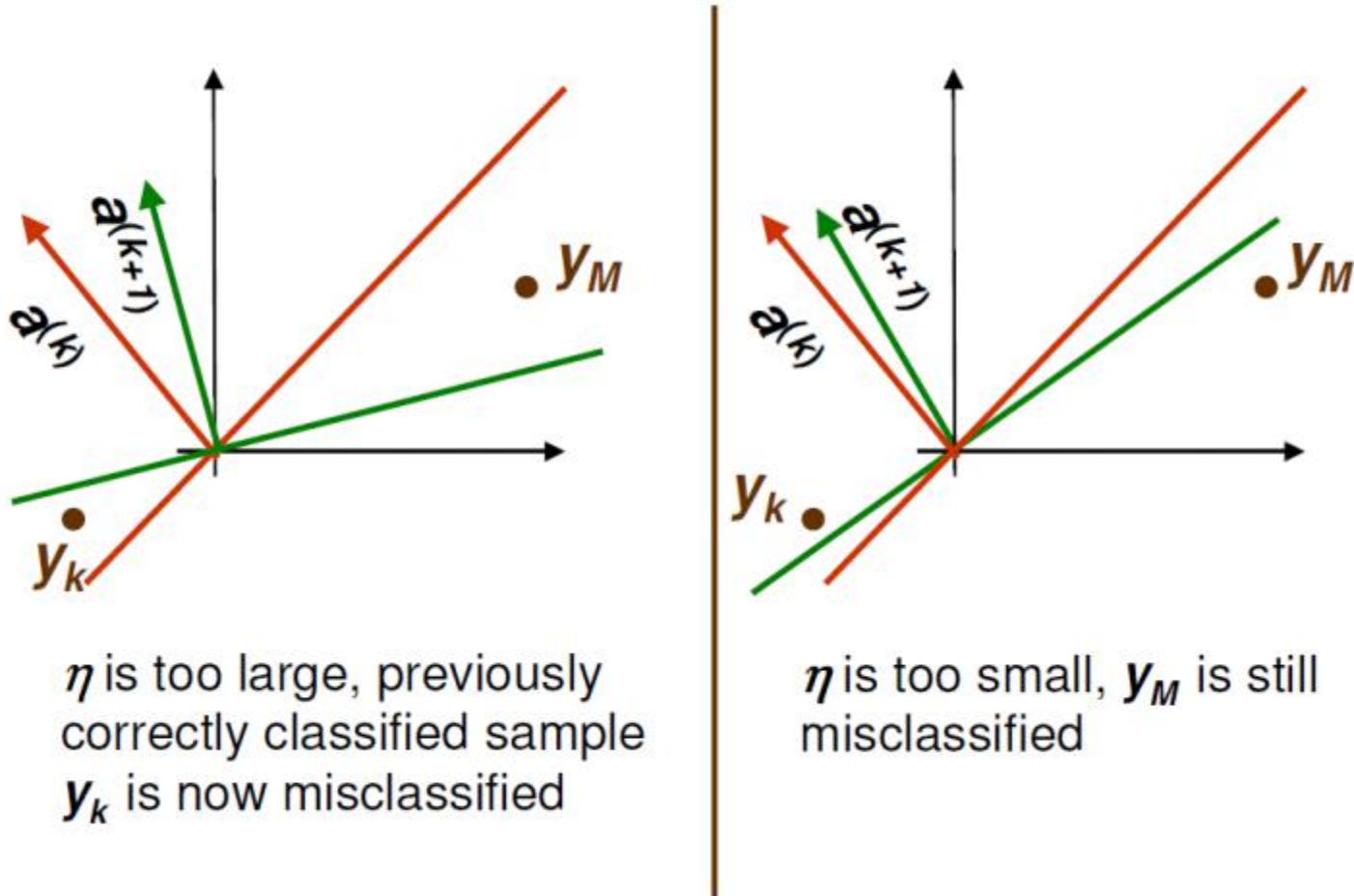
- Note that y_M is one sample misclassified by $(\mathbf{a}^{(k)})^t \mathbf{y}_M \leq 0$
- y_M is on the wrong side of decision hyperplane
- Adding ηy_M to \mathbf{a} moves the new decision hyperplane in the right direction with respect to y_M

$$(\mathbf{a}^{(k)})^t \mathbf{y}_M \leq 0$$



Perceptron Single Sample Rule

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



Perceptron Example

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

- Class 1: students who get A
- Class 2: students who get F

Perceptron Example

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

- **Augment samples** by adding an extra feature (dimension) equal to 1

Perceptron Example

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

- **Normalize:**

- Replace all examples from class 2 by their negative values $\mathbf{y}_i \rightarrow -\mathbf{y}_i \quad \forall \mathbf{y}_i \in \mathbf{c}_2$
- Seek \mathbf{a} such that: $\mathbf{a}^t \mathbf{y}_i > 0 \quad \forall \mathbf{y}_i$

Perceptron Example

	features				grade
<i>name</i>	<i>good attendance?</i>	<i>tall?</i>	<i>sleeps in class?</i>	<i>chews gum?</i>	
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

- **Single Sample Rule:**

- 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: $\mathbf{a}^{(k+1)} = \mathbf{a}^{(k)} + \eta^{(k)} \sum_{\mathbf{y} \in Y_M} \mathbf{y}$
- Set η fixed learning rate to $\eta^{(k)} = 1$: $\mathbf{a}^{(k+1)} = \mathbf{a}^{(k)} + \mathbf{y}_M$

Perceptron Example

- Set equal initial weights

$$\mathbf{a}^{(1)} = [0.25, 0.25, 0.25, 0.25, 0.25]$$

- Visit all samples sequentially, modifying the weights after each misclassified example

<i>name</i>	<i>a'y</i>	<i>misclassified?</i>
Jane	$0.25*1+0.25*1+0.25*1+0.25*(-1)+0.25*(-1) > 0$	<i>no</i>
Steve	$0.25*(-1)+0.25*(-1)+0.25*(-1)+0.25*(-1)+0.25*(-1) < 0$	<i>yes</i>

- New weights

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

Perceptron Example

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

<i>name</i>	$\mathbf{a}^t \mathbf{y}$	<i>misclassified?</i>
Mary	$-0.75*(-1) - 0.75*1 - 0.75*1 - 0.75*1 - 0.75*(-1) < 0$	yes

- New weights

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

Perceptron Example

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

<i>name</i>	$\mathbf{a}^t \mathbf{y}$	<i>misclassified?</i>
Peter	$-1.75 * 1 + 0.25 * 1 + 0.25 * (-1) + 0.25 * (-1) - 1.75 * 1 < 0$	yes

- New weights

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

Perceptron Example

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

<i>name</i>	<i>a^ty</i>	<i>misclassified?</i>
Jane	$-0.75 * 1 + 1.25 * 1 - 0.75 * 1 - 0.75 * (-1) - 0.75 * (-1) + 0$	<i>no</i>
Steve	$-0.75 * (-1) + 1.25 * (-1) - 0.75 * (-1) - 0.75 * (-1) - 0.75 * (-1) > 0$	<i>no</i>
Mary	$-0.75 * (-1) + 1.25 * 1 - 0.75 * 1 - 0.75 * 1 - 0.75 * (-1) > 0$	<i>no</i>
Peter	$-0.75 * 1 + 1.25 * 1 - 0.75 * (-1) - 0.75 * (-1) - 0.75 * 1 > 0$	<i>no</i>

- 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} :

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

Perceptron Example

- Converting back to the original features x :

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

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



- This is just one possible solution vector.
- If we started with weights $\mathbf{a}^{(1)} = [0, 0.5, 0.5, 0, 0]$, the solution would be $[-1, 1.5, -0.5, -1, -1]$

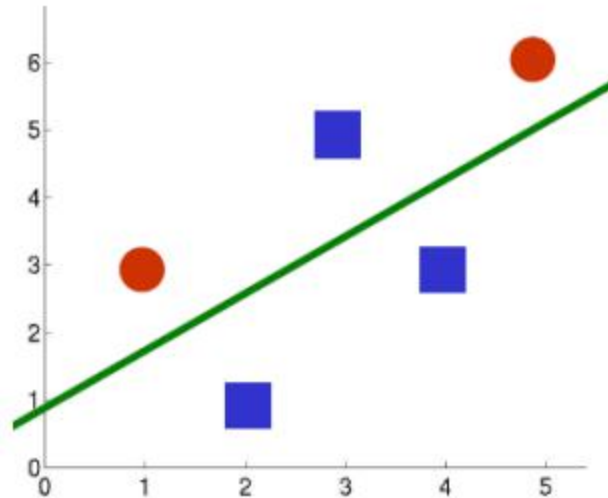
$$1.5 * x^{(1)} - 0.5 * x^{(2)} - x^{(3)} - x^{(4)} > 1 \Rightarrow \textit{grade A}$$

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

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

LDF: Non-separable Example

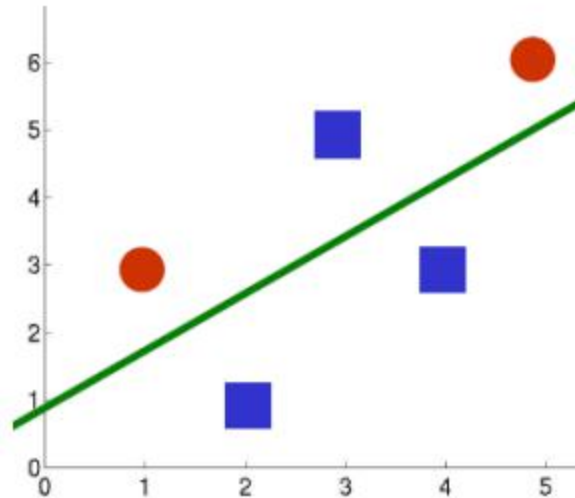
- Suppose we have 2 features and the samples are:
 - Class 1: [2,1], [4,3], [3,5]
 - Class 2: [1,3] and [5,6]
- These samples are not separable by a line
- Still would like to get approximate separation by a line
 - A good choice is shown in green
 - Some samples may be “noisy”, and we could accept them being misclassified



LDF: Non-separable Example

- Obtain y_1, y_2, y_3, y_4 by adding extra feature and “normalizing”

$$y_1 = \begin{bmatrix} 1 \\ 2 \\ 1 \end{bmatrix} \quad y_2 = \begin{bmatrix} 1 \\ 4 \\ 3 \end{bmatrix} \quad y_3 = \begin{bmatrix} 1 \\ 3 \\ 5 \end{bmatrix} \quad y_4 = \begin{bmatrix} -1 \\ -1 \\ -3 \end{bmatrix} \quad y_5 = \begin{bmatrix} -1 \\ -5 \\ -6 \end{bmatrix}$$



LDF: Non-separable Example

- Apply Perceptron single sample algorithm
- Initial equal weights

$$\mathbf{a}^{(1)} = [1 \ 1 \ 1]$$

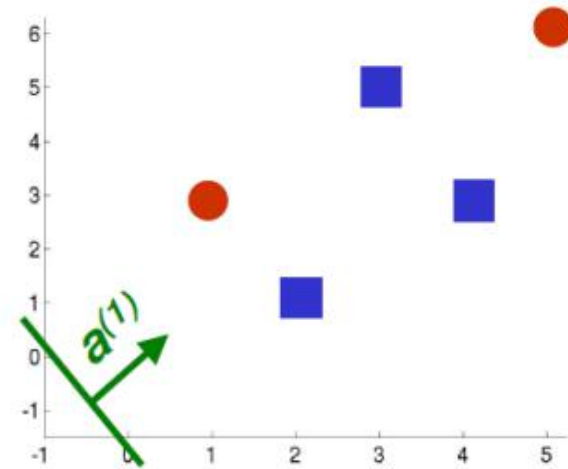
– Line equation $x^{(1)} + x^{(2)} + 1 = 0$

- Fixed learning rate $\eta = 1$

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

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

- $\mathbf{y}_1^t \mathbf{a}^{(1)} = [1 \ 1 \ 1] * [1 \ 2 \ 1]^t > 0 \quad \checkmark$
- $\mathbf{y}_2^t \mathbf{a}^{(1)} = [1 \ 1 \ 1] * [1 \ 4 \ 3]^t > 0 \quad \checkmark$
- $\mathbf{y}_3^t \mathbf{a}^{(1)} = [1 \ 1 \ 1] * [1 \ 3 \ 5]^t > 0 \quad \checkmark$



LDF: Non-separable Example

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

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

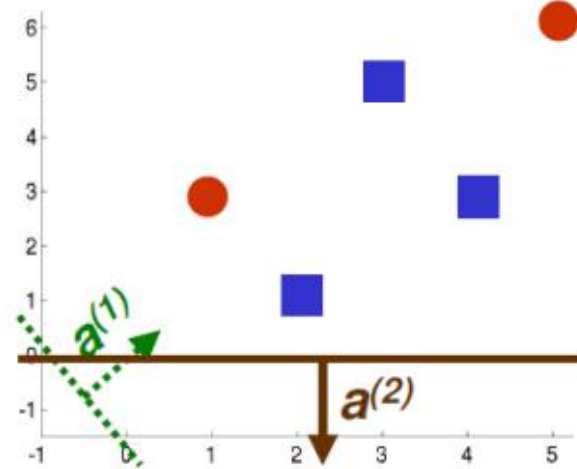
- $\mathbf{y}_4^t \mathbf{a}^{(1)} = [1 \ 1 \ 1]^* [-1 \ -1 \ -3]^t = -5 < 0$

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

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

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

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



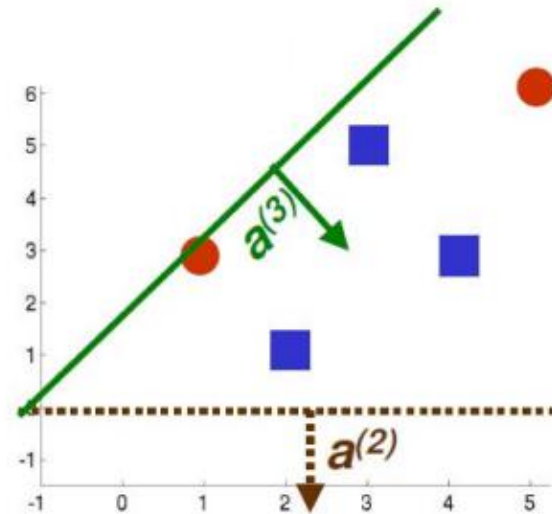
LDF: Non-separable Example

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

$$y_1 = \begin{bmatrix} 1 \\ 2 \\ 1 \end{bmatrix} \quad y_2 = \begin{bmatrix} 1 \\ 4 \\ 3 \end{bmatrix} \quad y_3 = \begin{bmatrix} 1 \\ 3 \\ 5 \end{bmatrix} \quad y_4 = \begin{bmatrix} -1 \\ -1 \\ -3 \end{bmatrix} \quad y_5 = \begin{bmatrix} -1 \\ -5 \\ -6 \end{bmatrix}$$

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

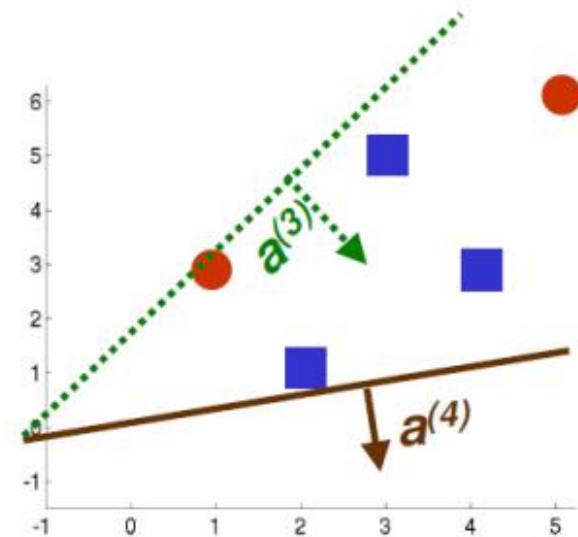
$$\mathbf{a}^{(4)} = \mathbf{a}^{(3)} + \mathbf{y}_M = [1 \ 2 \ -1] + [-1 \ -1 \ -3] = [0 \ 1 \ -4]$$



LDF: Non-separable Example

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

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

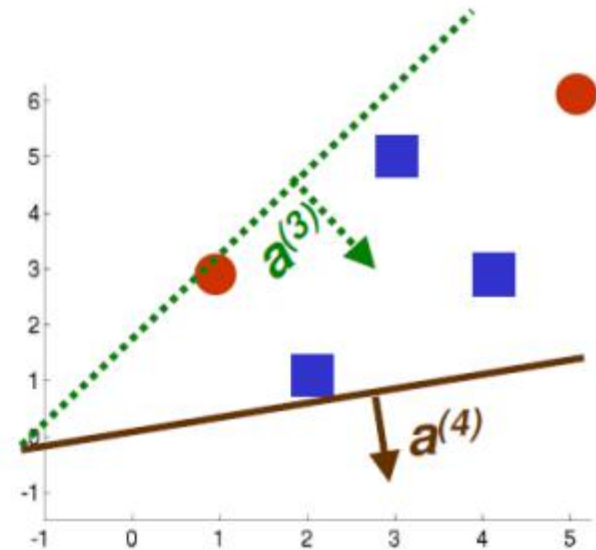


- $\mathbf{y}_5^t \mathbf{a}^{(4)} = [-1 \ -5 \ -6] \cdot [0 \ 1 \ -4] = 19 > 0$
- $\mathbf{y}_1^t \mathbf{a}^{(4)} = [1 \ 2 \ 1] \cdot [0 \ 1 \ -4] = -2 < 0$
-

LDF: Non-separable Example

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

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



- $\mathbf{y}_5^t \mathbf{a}^{(4)} = [-1 \ -5 \ -6] * [0 \ 1 \ -4] = 19 > 0$
- $\mathbf{y}_1^t \mathbf{a}^{(4)} = [1 \ 2 \ 1] * [0 \ 1 \ -4] = -2 < 0$
-

LDF: Non-separable Example

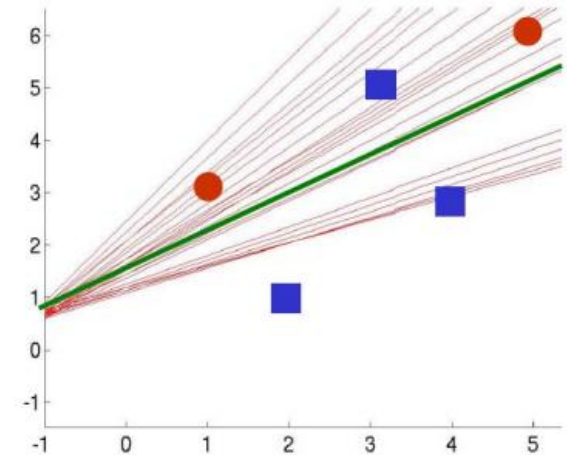
- We can continue this forever.
- There is no solution vector \mathbf{a} satisfying for all \mathbf{x}_i

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

- Need to stop but at a good point
- Will not converge in the nonseparable case
- To ensure convergence can set

$$\eta^{(k)} = \frac{\eta^{(1)}}{k}$$

- However we are not guaranteed that we will stop at a good point



Convergence of Perceptron Rules

- If classes are linearly separable and we use fixed learning rate, that is for $\eta(k) = \text{const}$
- Then, **both the single sample and batch perceptron rules** converge to a correct solution (could be any a in the solution space)
- If classes are not linearly separable:
 - The algorithm does not stop, it keeps looking for a solution which does not exist
 - By choosing appropriate learning rate, we can always ensure convergence:
 - For example inverse linear learning rate:
 - 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

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

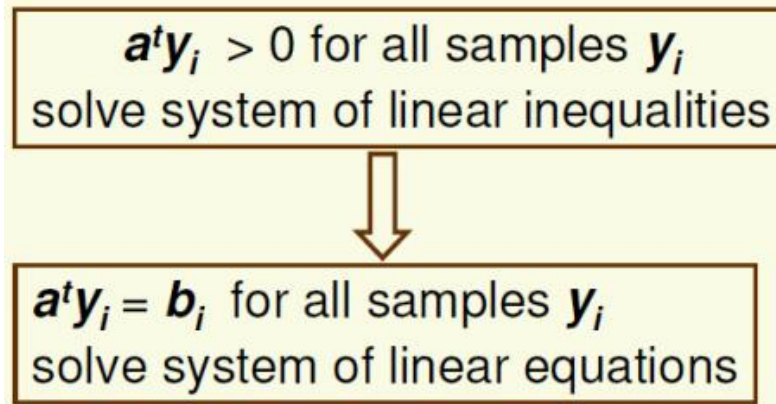
- easier to analyze
- Concentrates more than necessary on any isolated “noisy” training examples

Outline

- Perceptron Rule
- **Minimum Squared-Error Procedure**
- Ho–Kashyap Procedure

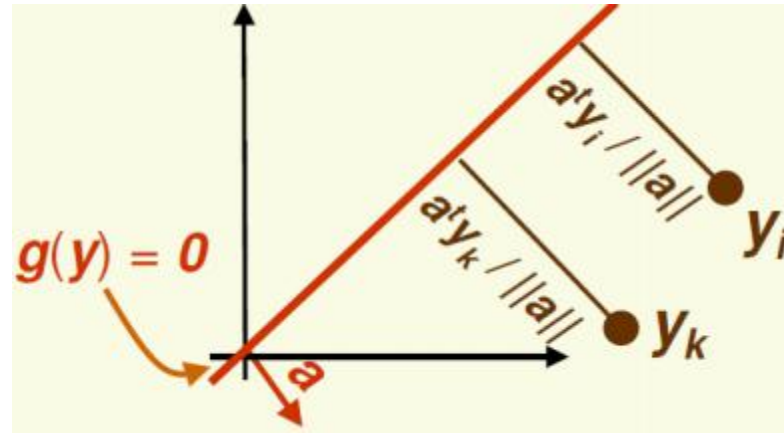
Minimum Squared-Error Procedures

- Idea: convert to easier and better understood problem



- MSE procedure
 - Choose positive constants b_1, b_2, \dots, b_n
 - Try to find weight vector \mathbf{a} such that $\mathbf{a}^t \mathbf{y}_i = b_i$ for all samples \mathbf{y}_i
 - If we can find such a vector, then \mathbf{a} is a solution because the b_i 's are positive
 - Consider all the samples (not just the misclassified ones)

MSE Margins



- If $a'y_i = b_i$, y_i must be at distance b_i from the separating hyperplane (normalized by $\|a\|$)
- Thus b_1, b_2, \dots, b_n give relative expected distances or “margins” of samples from the hyperplane
- Should make b_i small if sample i is expected to be near separating hyperplane, and large otherwise
- In the absence of any additional information, set $b_1 = b_2 = \dots = b_n = 1$

MSE Matrix Notation

- Need to solve n equations
- In matrix form $Y\mathbf{a}=\mathbf{b}$

$$\begin{cases} \mathbf{a}^t \mathbf{y}_1 = b_1 \\ \vdots \\ \mathbf{a}^t \mathbf{y}_n = b_n \end{cases}$$

$$\underbrace{\begin{bmatrix} \mathbf{y}_1^{(0)} & \mathbf{y}_1^{(1)} & \dots & \mathbf{y}_1^{(d)} \\ \mathbf{y}_2^{(0)} & \mathbf{y}_2^{(1)} & \dots & \mathbf{y}_2^{(d)} \\ \vdots & \vdots & \ddots & \vdots \\ \mathbf{y}_n^{(0)} & \mathbf{y}_n^{(1)} & \dots & \mathbf{y}_n^{(d)} \end{bmatrix}}_{\mathbf{Y}} \underbrace{\begin{bmatrix} a_0 \\ a_1 \\ \vdots \\ a_d \end{bmatrix}}_{\mathbf{a}} = \underbrace{\begin{bmatrix} b_1 \\ b_2 \\ \vdots \\ b_n \end{bmatrix}}_{\mathbf{b}}$$

Exact Solution is Rare

- Need to solve a linear system $Ya = b$
 - Y is an $n \times (d + 1)$ matrix
- Exact solution only if Y is non-singular and square (the inverse Y^{-1} exists)
 - $a = Y^{-1}b$
 - (number of samples) = (number of features + 1)
 - Almost never happens in practice
 - Guaranteed to find the separating hyperplane

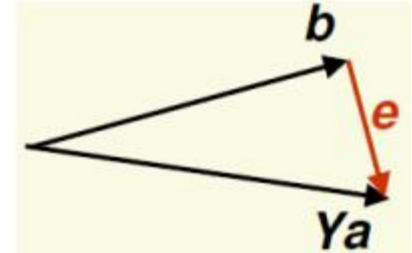
Approximate Solution

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

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 0

Computing the Gradient

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

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

Pseudo-Inverse Solution

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

- Setting the gradient to 0:

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

- The 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:

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

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

MSE Procedures

- Only guaranteed separating hyperplane if $Y\mathbf{a} \geq 0$
 - That is if all elements of vector $Y\mathbf{a}$ are positive

$$Y\mathbf{a} = \begin{bmatrix} b_1 + \varepsilon_1 \\ \vdots \\ b_n + \varepsilon_n \end{bmatrix}$$

- where ε may be negative
- If $\varepsilon_1, \dots, \varepsilon_n$ are small relative to b_1, \dots, b_n , then each element of $Y\mathbf{a}$ is positive, and \mathbf{a} gives a separating hyperplane
 - If the approximation is not good, ε_i may be large and negative, for some i , thus $b_i + \varepsilon_i$ will be negative and \mathbf{a} is not a separating hyperplane
- In linearly separable case, least squares solution \mathbf{a} does not necessarily give separating hyperplane

MSE Procedures

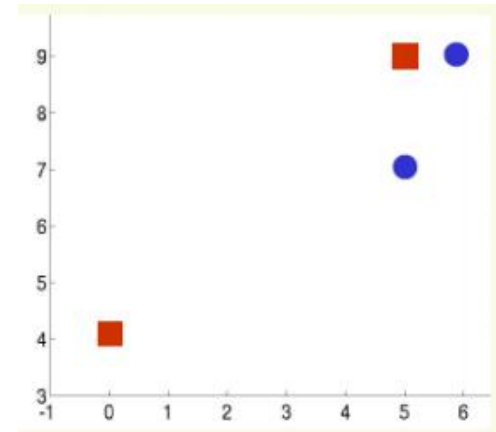
- We are free to choose b . We may be tempted to make b large as a way to ensure $Y\mathbf{a} = b > 0$
 - Does not work
 - Let β be a scalar, let's try βb instead of b
- If \mathbf{a}^* is a least squares solution to $Y\mathbf{a} = b$, then for any scalar β , the least squares solution to $Y\mathbf{a}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

$$\arg \min_a \|Y\mathbf{a} - \beta b\|^2 = \arg \min_a \beta^2 \|Y(\mathbf{a} / \beta) - b\|^2 = \beta \mathbf{a}^*$$

- Thus if the i -th element of $Y\mathbf{a}$ is less than 0, that is $\mathbf{y}_i^t \mathbf{a} < 0$, then $\mathbf{y}_i^t (\beta \mathbf{a}) < 0$
 - The relative difference between components of b matters, but not the size of each individual component

LDF using MSE: Example 1

- Class 1: (6 9), (5 7)
- Class 2: (5 9), (0 4)
- Add extra feature and “normalize”



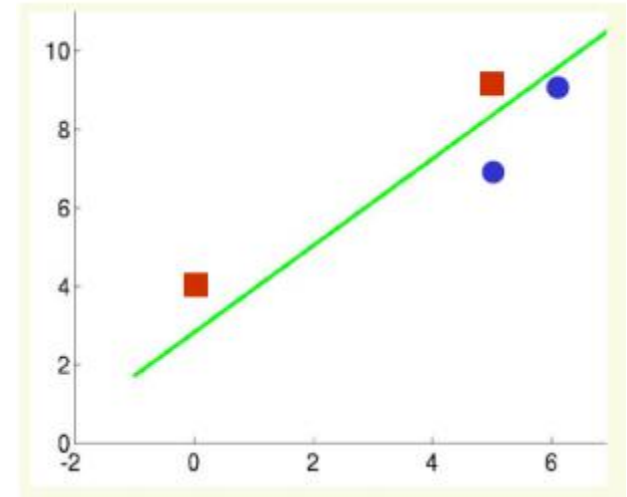
$$y_1 = \begin{bmatrix} 1 \\ 6 \\ 9 \end{bmatrix} \quad y_2 = \begin{bmatrix} 1 \\ 5 \\ 7 \end{bmatrix} \quad y_3 = \begin{bmatrix} -1 \\ -5 \\ -9 \end{bmatrix} \quad y_4 = \begin{bmatrix} -1 \\ 0 \\ -4 \end{bmatrix}$$

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

LDF using MSE: Example 1

- Choose $\mathbf{b}=[1 \ 1 \ 1 \ 1]^T$
- In Matlab, $\mathbf{a}=\mathbf{Y}\backslash\mathbf{b}$ solves the least squares problem

$$\mathbf{a} = \begin{bmatrix} 2.66 \\ 1.045 \\ -0.944 \end{bmatrix}$$

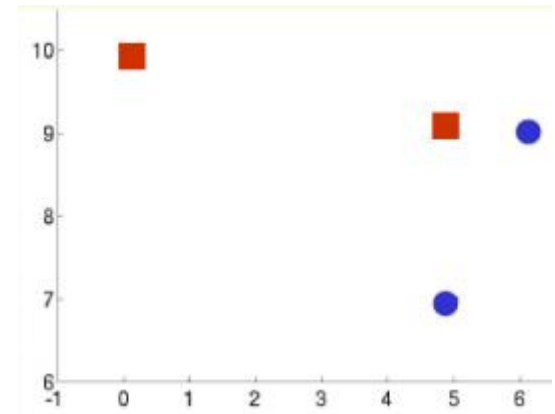


- Note \mathbf{a} is an approximation to $\mathbf{Y}\mathbf{a} = \mathbf{b}$, since no exact solution exists
- This solution gives a separating hyperplane since $\mathbf{Y}\mathbf{a} > 0$

$$\mathbf{Y}\mathbf{a} = \begin{bmatrix} 0.44 \\ 1.28 \\ 0.61 \\ 1.11 \end{bmatrix}$$

LDF using MSE: Example 2

- 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



$$y_1 = \begin{bmatrix} 1 \\ 6 \\ 9 \end{bmatrix} \quad y_2 = \begin{bmatrix} 1 \\ 5 \\ 7 \end{bmatrix} \quad y_3 = \begin{bmatrix} -1 \\ -5 \\ -9 \end{bmatrix} \quad y_4 = \begin{bmatrix} -1 \\ 0 \\ -10 \end{bmatrix}$$

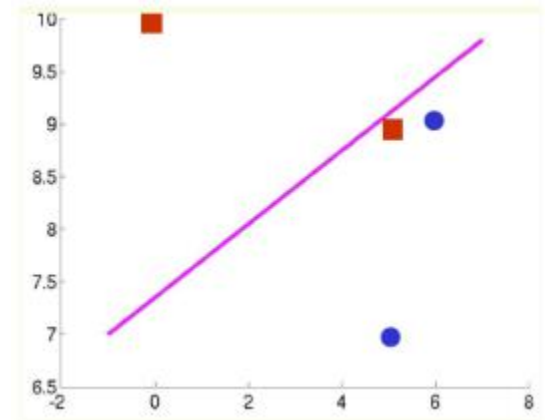
$$Y = \begin{bmatrix} 1 & 6 & 9 \\ 1 & 5 & 7 \\ -1 & -5 & -9 \\ -1 & 0 & -10 \end{bmatrix}$$

LDF using MSE: Example 2

- Choose $\mathbf{b}=[1 \ 1 \ 1 \ 1]^T$
- In Matlab, $\mathbf{a}=\mathbf{Y}\backslash\mathbf{b}$ solves the least squares problem

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

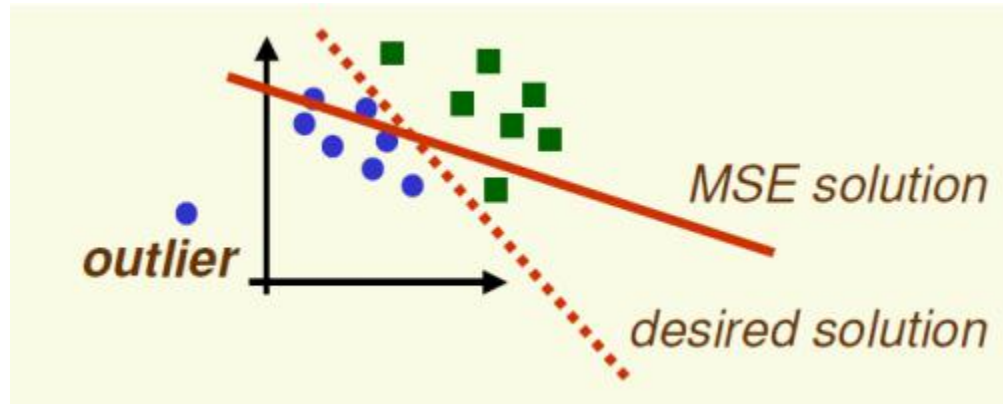
$$\mathbf{Y}\mathbf{a} = \begin{bmatrix} 0.2 \\ 0.9 \\ -0.04 \\ 1.16 \end{bmatrix} \neq \begin{bmatrix} 1 \\ 1 \\ 1 \\ 1 \end{bmatrix}$$



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

LDF using MSE: Example 2

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



- No problems with convergence
- Solution ranges from reasonable to good

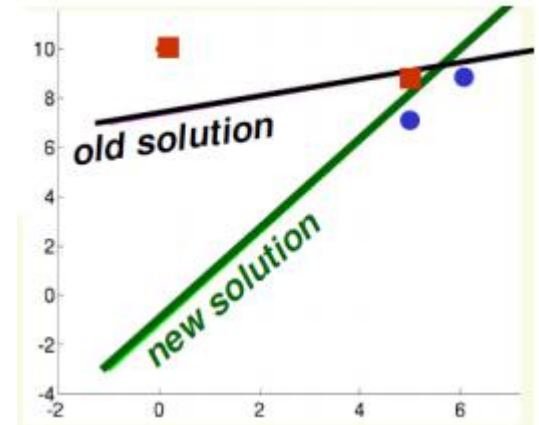
LDF using MSE: Example 2

- We can see that the 4-th point is vary far from separating hyperplane
 - In practice we don't know this
- A more appropriate b could be
- In Matlab, $\mathbf{a}=\mathbf{Y}\backslash\mathbf{b}$ solves the least squares problem

$$\mathbf{b} = \begin{bmatrix} 1 \\ 1 \\ 1 \\ 10 \end{bmatrix}$$

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

$$\mathbf{Y}\mathbf{a} = \begin{bmatrix} 0.9 \\ 1.0 \\ 0.8 \\ 10.0 \end{bmatrix} \neq \begin{bmatrix} 1 \\ 1 \\ 1 \\ 10 \end{bmatrix}$$



- This solution gives the separating hyperplane since $\mathbf{Y}\mathbf{a} > 0$

Gradient Descent for MSE

$$\mathbf{J}_s(\mathbf{a}) = \|\mathbf{Y}\mathbf{a} - \mathbf{b}\|^2$$

- May wish to find MSE solution by gradient descent:
 1. Computing the inverse of $\mathbf{Y}^t\mathbf{Y}$ may be too costly
 2. $\mathbf{Y}^t\mathbf{Y}$ may be close to singular if samples are highly correlated (rows of \mathbf{Y} are almost linear combinations of each other) computing the inverse of $\mathbf{Y}^t\mathbf{Y}$ is not numerically stable
- As shown before, the gradient is:

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

Widrow-Hoff Procedure

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

- Thus the update rule for gradient descent is:

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

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

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

Outline

- Perceptron Rule
- Minimum Squared-Error Procedure
- **Ho-Kashyap Procedure**

Ho-Kashyap Procedure

- In the MSE procedure, if \mathbf{b} is chosen arbitrarily, finding separating hyperplane is not guaranteed.
- Suppose training samples are linearly separable. Then there is \mathbf{a}^s and positive \mathbf{b}^s s.t.

$$Y\mathbf{a}^s = \mathbf{b}^s > 0$$

- If we knew \mathbf{b}^s could apply MSE procedure to find the separating hyperplane
- Idea: find both \mathbf{a}^s and \mathbf{b}^s
- Minimize the following criterion function, restricting to positive \mathbf{b} :

$$J_{HK}(\mathbf{a}, \mathbf{b}) = \|\mathbf{Y}\mathbf{a} - \mathbf{b}\|^2$$

$$J_{HK}(\mathbf{a}^s, \mathbf{b}^s) = 0$$

Ho-Kashyap Procedure

$$\mathbf{J}_{HK}(\mathbf{a}, \mathbf{b}) = \|\mathbf{Y}\mathbf{a} - \mathbf{b}\|^2$$

- As usual, take partial derivatives w.r.t. \mathbf{a} and \mathbf{b}

$$\nabla_{\mathbf{a}} \mathbf{J}_{HK} = 2\mathbf{Y}^t(\mathbf{Y}\mathbf{a} - \mathbf{b}) = \mathbf{0}$$

$$\nabla_{\mathbf{b}} \mathbf{J}_{HK} = -2(\mathbf{Y}\mathbf{a} - \mathbf{b}) = \mathbf{0}$$

- Use modified gradient descent procedure to find a minimum of $\mathbf{J}_{HK}(\mathbf{a}, \mathbf{b})$
- Alternate the two steps below until convergence:
 - ① Fix \mathbf{b} and minimize $\mathbf{J}_{HK}(\mathbf{a}, \mathbf{b})$ with respect to \mathbf{a}
 - ② Fix \mathbf{a} and minimize $\mathbf{J}_{HK}(\mathbf{a}, \mathbf{b})$ with respect to \mathbf{b}

Ho-Kashyap Procedure

$$\nabla_a \mathbf{J}_{HK} = 2\mathbf{Y}^t(\mathbf{Y}\mathbf{a} - \mathbf{b}) = \mathbf{0} \quad \nabla_b \mathbf{J}_{HK} = -2(\mathbf{Y}\mathbf{a} - \mathbf{b}) = \mathbf{0}$$

- Alternate the two steps below until convergence:
 - ① Fix \mathbf{b} and minimize $\mathbf{J}_{HK}(\mathbf{a}, \mathbf{b})$ with respect to \mathbf{a}
 - ② Fix \mathbf{a} and minimize $\mathbf{J}_{HK}(\mathbf{a}, \mathbf{b})$ with respect to \mathbf{b}
- Step (1) can be performed with pseudoinverse
 - For fixed \mathbf{b} minimum of $\mathbf{J}_{HK}(\mathbf{a}, \mathbf{b})$ with respect to \mathbf{a} is found by solving

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

– Thus

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

Ho-Kashyap Procedure

- Step 2: fix \mathbf{a} and minimize $J_{HK}(\mathbf{a}, \mathbf{b})$ with respect to \mathbf{b}
- We can't use $\mathbf{b} = \mathbf{Y}\mathbf{a}$ because \mathbf{b} has to be positive
- Solution: use modified gradient descent
- Regular gradient descent rule:

$$\mathbf{b}^{(k+1)} = \mathbf{b}^{(k)} - \eta^{(k)} \nabla_{\mathbf{b}} \mathbf{J}(\mathbf{a}^{(k)}, \mathbf{b}^{(k)})$$

- If any components of $\nabla_{\mathbf{b}} \mathbf{J}$ are positive, \mathbf{b} will decrease and can possibly become negative

$$\mathbf{b}^{(k+1)} = \begin{bmatrix} 1 \\ 1 \\ 1 \end{bmatrix} - 2 * \begin{bmatrix} 2 \\ -3 \\ -2 \end{bmatrix} = \begin{bmatrix} -3 \\ 7 \\ 5 \end{bmatrix}$$

Ho-Kashyap Procedure

- Start with positive \mathbf{b} , follow negative gradient but refuse to decrease any components of \mathbf{b}
- This can be achieved by setting all the positive components of $\nabla_{\mathbf{b}}\mathbf{J}$ to 0

$$\mathbf{b}^{(k+1)} = \mathbf{b}^{(k)} - \eta \frac{1}{2} [\nabla_{\mathbf{b}}\mathbf{J}(\mathbf{a}^{(k)}, \mathbf{b}^{(k)}) - |\nabla_{\mathbf{b}}\mathbf{J}(\mathbf{a}^{(k)}, \mathbf{b}^{(k)})|]$$

here $|v|$ denotes vector we get after applying absolute value to all elements of v

$$\mathbf{b}^{(k+1)} = \begin{bmatrix} 1 \\ 1 \\ 1 \end{bmatrix} - 2 * \frac{1}{2} \left[\begin{bmatrix} -2 \\ -3 \\ -2 \end{bmatrix} - \begin{bmatrix} 2 \\ 3 \\ 2 \end{bmatrix} \right] = \begin{bmatrix} 1 \\ 1 \\ 1 \end{bmatrix} - \begin{bmatrix} 0 \\ -6 \\ -4 \end{bmatrix} = \begin{bmatrix} 1 \\ 7 \\ 5 \end{bmatrix}$$

- Not doing steepest descent anymore, but we are still doing descent **and** ensure that \mathbf{b} is positive

Ho-Kashyap Procedure

$$\mathbf{b}^{(k+1)} = \mathbf{b}^{(k)} - \eta \frac{1}{2} [\nabla_{\mathbf{b}} \mathbf{J}(\mathbf{a}^{(k)}, \mathbf{b}^{(k)}) - |\nabla_{\mathbf{b}} \mathbf{J}(\mathbf{a}^{(k)}, \mathbf{b}^{(k)})|]$$

$$\nabla_{\mathbf{b}} \mathbf{J} = -2(\mathbf{Y}\mathbf{a} - \mathbf{b}) = \mathbf{0}$$

Let $\mathbf{e}^{(k)} = \mathbf{Y}\mathbf{a}^{(k)} - \mathbf{b}^{(k)} = -\frac{1}{2} \nabla_{\mathbf{b}} \mathbf{J}(\mathbf{a}^{(k)}, \mathbf{b}^{(k)})$

Then

$$\begin{aligned} \mathbf{b}^{(k+1)} &= \mathbf{b}^{(k)} - \eta \frac{1}{2} [-2\mathbf{e}^{(k)} - |2\mathbf{e}^{(k)}|] \\ &= \mathbf{b}^{(k)} + \eta [\mathbf{e}^{(k)} + |\mathbf{e}^{(k)}|] \end{aligned}$$

Ho-Kashyap Procedure

- The final Ho–Kashyap procedure:

0) Start with arbitrary $\mathbf{a}^{(1)}$ and $\mathbf{b}^{(1)} > 0$, let $k = 1$

repeat steps (1) through (4)

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

2) Solve for $\mathbf{b}^{(k+1)}$ using $\mathbf{a}^{(k)}$ and $\mathbf{b}^{(k)}$

$$\mathbf{b}^{(k+1)} = \mathbf{b}^{(k)} + \eta[\mathbf{e}^{(k)} + \|\mathbf{e}^{(k)}\|]$$

3) Solve for $\mathbf{a}^{(k+1)}$ using $\mathbf{b}^{(k+1)}$

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

4) $k = k + 1$

until $\mathbf{e}^{(k)} \geq 0$ or $k > k_{max}$ or $\mathbf{b}^{(k+1)} = \mathbf{b}^{(k)}$

- For convergence, learning rate should be fixed between $0 < \eta < 1$.

Ho-Kashyap Procedure

$$\mathbf{b}^{(k+1)} = \mathbf{b}^{(k)} + \eta [\mathbf{e}^{(k)} + |\mathbf{e}^{(k)}|]$$

- What if $\mathbf{e}^{(k)}$ is negative for all components?

$$\mathbf{b}^{(k+1)} = \mathbf{b}^{(k)} \text{ and corrections stop}$$

- Write $\mathbf{e}^{(k)}$ out:

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

- Multiply by \mathbf{Y}^t :

$$\mathbf{Y}^t\mathbf{e}^{(k)} = \mathbf{Y}^t(\mathbf{Y}(\mathbf{Y}^t\mathbf{Y})^{-1}\mathbf{Y}^t\mathbf{b}^{(k)} - \mathbf{b}^{(k)}) = \mathbf{Y}^t\mathbf{b}^{(k)} - \mathbf{Y}^t\mathbf{b}^{(k)} = \mathbf{0}$$

- Thus

$$\mathbf{Y}^t\mathbf{e}^{(k)} = \mathbf{0}$$

Ho-Kashyap Procedure

- Suppose training samples are linearly separable. Then there is \mathbf{a}^s and positive \mathbf{b}^s s.t

$$\mathbf{Y}\mathbf{a}^s = \mathbf{b}^s > \mathbf{0}$$

- Multiply both sides by $(\mathbf{e}^{(k)})^t$

$$\mathbf{0} = (\mathbf{e}^{(k)})^t \mathbf{Y}\mathbf{a}^s = (\mathbf{e}^{(k)})^t \mathbf{b}^s$$

- Either by $\mathbf{e}^{(k)} = \mathbf{0}$ or one of its components is positive

Ho-Kashyap Procedure

- In the linearly separable case,
 - $\mathbf{e}^{(k)} = 0$, found solution, stop
 - one of components of $\mathbf{e}^{(k)}$ is positive, algorithm continues
- In non separable case,
 - $\mathbf{e}^{(k)}$ will have only negative components eventually, thus found proof of nonseparability
 - No bound on how many iteration need for the proof of nonseparability

Example

- Class 1: (6,9), (5,7)
- Class 2: (5,9), (0, 10)

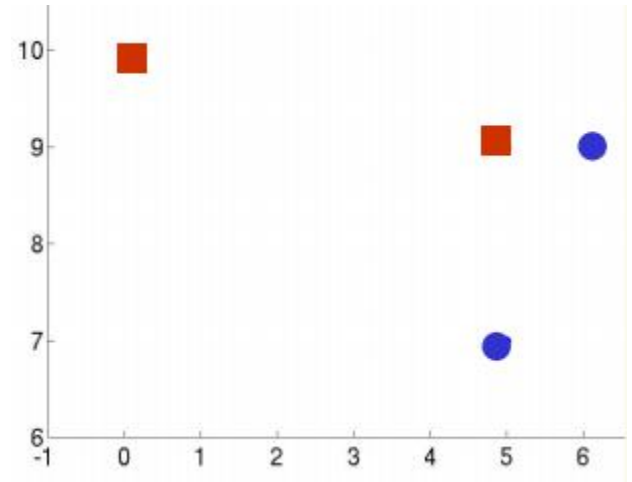
- Matrix $\mathbf{Y} = \begin{bmatrix} 1 & 6 & 9 \\ 1 & 5 & 7 \\ -1 & -5 & -9 \\ -1 & 0 & -10 \end{bmatrix}$

- Start with $\mathbf{a}^{(1)} = \begin{bmatrix} 1 \\ 1 \\ 1 \\ 1 \end{bmatrix}$ and $\mathbf{b}^{(1)} = \begin{bmatrix} 1 \\ 1 \\ 1 \\ 1 \end{bmatrix}$

- Use fixed learning $\eta = 0.9$

- At the start

$$\mathbf{Y}\mathbf{a}^{(1)} = \begin{bmatrix} 16 \\ 13 \\ -15 \\ -11 \end{bmatrix}$$



Example

- Iteration 1:

$$\mathbf{e}^{(1)} = \mathbf{Y}\mathbf{a}^{(1)} - \mathbf{b}^{(1)} = \begin{bmatrix} 16 \\ 13 \\ -15 \\ -11 \end{bmatrix} - \begin{bmatrix} 1 \\ 1 \\ 1 \\ 1 \end{bmatrix} = \begin{bmatrix} 15 \\ 12 \\ -16 \\ -12 \end{bmatrix}$$

- solve for $\mathbf{b}^{(2)}$ using $\mathbf{a}^{(1)}$ and $\mathbf{b}^{(1)}$

$$\mathbf{b}^{(2)} = \mathbf{b}^{(1)} + 0.9[\mathbf{e}^{(1)} + \|\mathbf{e}^{(1)}\|] = \begin{bmatrix} 1 \\ 1 \\ 1 \\ 1 \end{bmatrix} + 0.9 \left[\begin{bmatrix} 15 \\ 12 \\ -16 \\ -12 \end{bmatrix} + \begin{bmatrix} 15 \\ 12 \\ 16 \\ 12 \end{bmatrix} \right] = \begin{bmatrix} 28 \\ 22.6 \\ 1 \\ 1 \end{bmatrix}$$

- solve for $\mathbf{a}^{(2)}$ using $\mathbf{b}^{(2)}$

$$\mathbf{a}^{(2)} = (\mathbf{Y}^t \mathbf{Y})^{-1} \mathbf{Y}^t \mathbf{b}^{(2)} = \begin{bmatrix} -2.6 & 4.7 & 1.6 & -0.5 \\ 0.16 & -0.1 & -0.1 & 0.2 \\ 0.26 & -0.5 & -0.2 & -0.1 \end{bmatrix} * \begin{bmatrix} 28 \\ 22.6 \\ 1 \\ 1 \end{bmatrix} = \begin{bmatrix} 34.6 \\ 2.7 \\ -3.8 \end{bmatrix}$$

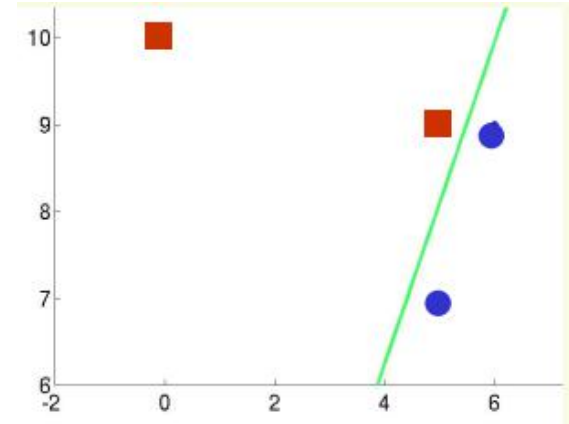
Example

- Continue iterations until $\mathbf{Y}\mathbf{a} > 0$
 - In practice, continue until minimum component of $\mathbf{Y}\mathbf{a}$ is less than 0.01
- After 104 iterations converged to solution

$$\mathbf{a} = \begin{bmatrix} -34.9 \\ 27.3 \\ -11.3 \end{bmatrix} \quad \mathbf{b} = \begin{bmatrix} 28 \\ 23 \\ 1 \\ 147 \end{bmatrix}$$

- \mathbf{a} does gives a separating hyperplane

$$\mathbf{Y}\mathbf{a} = \begin{bmatrix} 27.2 \\ 22.5 \\ 0.14 \\ 1.48 \end{bmatrix}$$



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 even 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
- Ho–Kashyap procedures
 - always converge
 - find separating hyperplane in the linearly separable case
 - more costly

Q & A