

# ***Minimum Squared Error***

# LDF: Minimum Squared-Error Procedures

- Idea: convert to easier and better understood problem

Perceptron

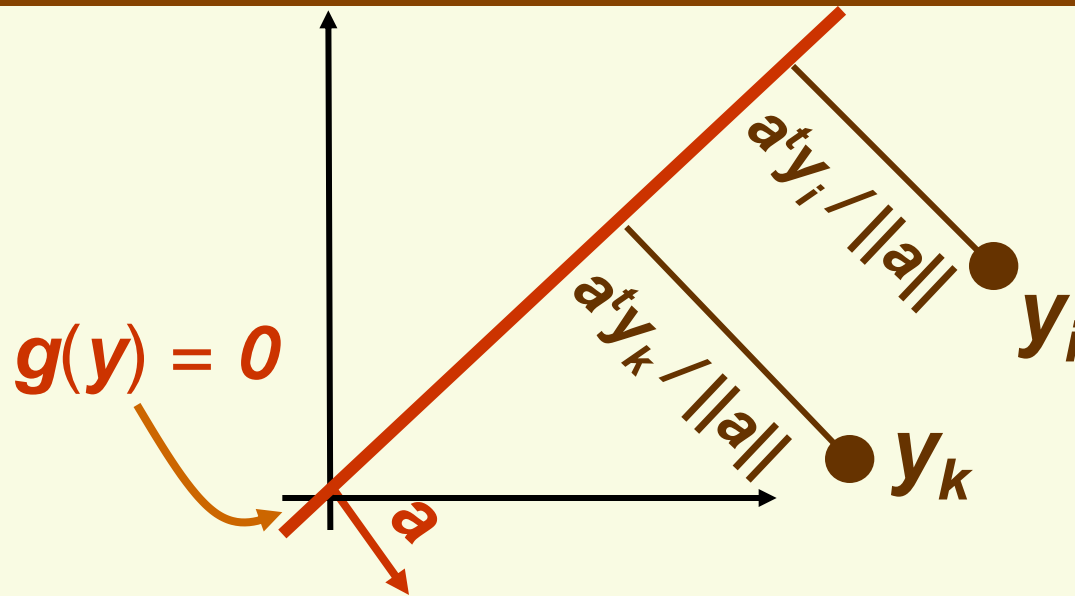
→  $\mathbf{a}^t \mathbf{y}_i > 0$  for all samples  $\mathbf{y}_i$   
solve system of linear inequalities



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

- MSE procedure
  - Choose **positive** constants  $b_1, b_2, \dots, b_n$
  - try to find weight vector  $\mathbf{a}$  s.t.  $\mathbf{a}^t \mathbf{y}_i = b_i$  for all samples  $\mathbf{y}_i$
  - If we can find weight vector  $\mathbf{a}$  such that  $\mathbf{a}^t \mathbf{y}_i = b_i$  for all samples  $\mathbf{y}_i$ , then  $\mathbf{a}$  is a solution because  $b_i$ '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  $\begin{cases} \mathbf{a}^t \mathbf{y}_1 = b_1 \\ \vdots \\ \mathbf{a}^t \mathbf{y}_n = b_n \end{cases}$
- Introduce matrix notation:

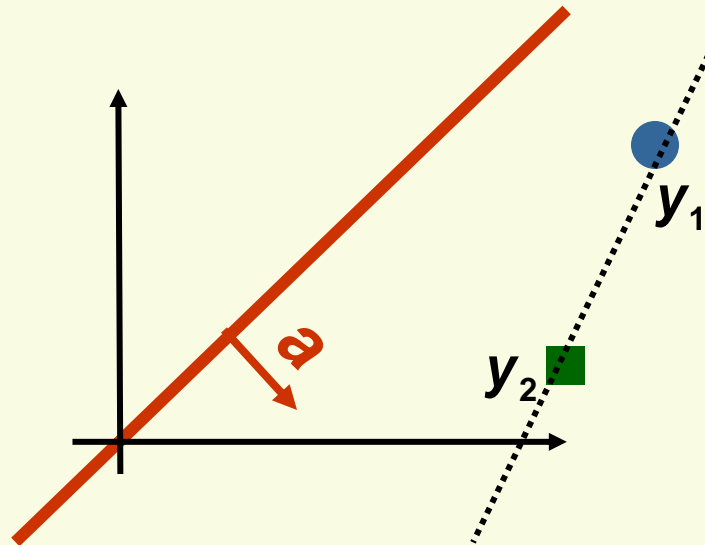
$$\underbrace{\begin{bmatrix} \mathbf{y}_1^{(0)} & \mathbf{y}_1^{(1)} & \dots & \mathbf{y}_1^{(d)} \\ \mathbf{y}_2^{(0)} & \mathbf{y}_2^{(1)} & \dots & \mathbf{y}_2^{(d)} \\ \vdots & \vdots & \ddots & \vdots \\ \mathbf{y}_n^{(0)} & \mathbf{y}_n^{(1)} & \dots & \mathbf{y}_n^{(d)} \end{bmatrix}}_{\mathbf{Y}} \underbrace{\begin{bmatrix} \mathbf{a}_0 \\ \mathbf{a}_1 \\ \vdots \\ \mathbf{a}_d \end{bmatrix}}_{\mathbf{a}} = \underbrace{\begin{bmatrix} b_1 \\ b_2 \\ \vdots \\ b_n \end{bmatrix}}_{\mathbf{b}}$$

- 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  $Y\mathbf{a} = \mathbf{b}$ 
  - $Y$  is an  $n$  by  $(d + 1)$  matrix
- Exact solution can be found only if  $Y$  is nonsingular and square, in which case the inverse  $Y^{-1}$  exists
  - $\mathbf{a} = Y^{-1}\mathbf{b}$
  - (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  $Y$  is overdetermined, that is it has more rows (examples) than columns (features)
  - If it has more features than examples, should reduce dimensionality

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

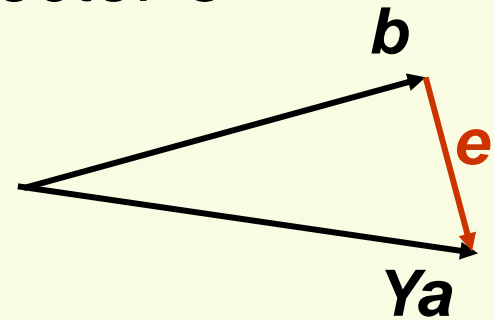
- Need  $Ya = b$ , but no exact solution exists for an overdetermined system of equation
  - More equations than unknowns
- Find an approximate solution  $a$ , that is  $Ya \approx b$ 
  - Note that approximate solution  $a$  *does not* necessarily give the separating hyperplane in the separable case
  - But hyperplane corresponding to  $a$  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:

$$\nabla J_s(\mathbf{a}) = \begin{bmatrix} \frac{\partial J_s}{\partial \mathbf{a}_0} \\ \vdots \\ \frac{\partial J_s}{\partial \mathbf{a}_d} \end{bmatrix} = 2\mathbf{Y}^t(\mathbf{Y}\mathbf{a} - \mathbf{b})$$

- Setting the gradient to 0:

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



# ***LDF: Pseudo Inverse Solution***

---

- 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  $\mathbf{Y}$*

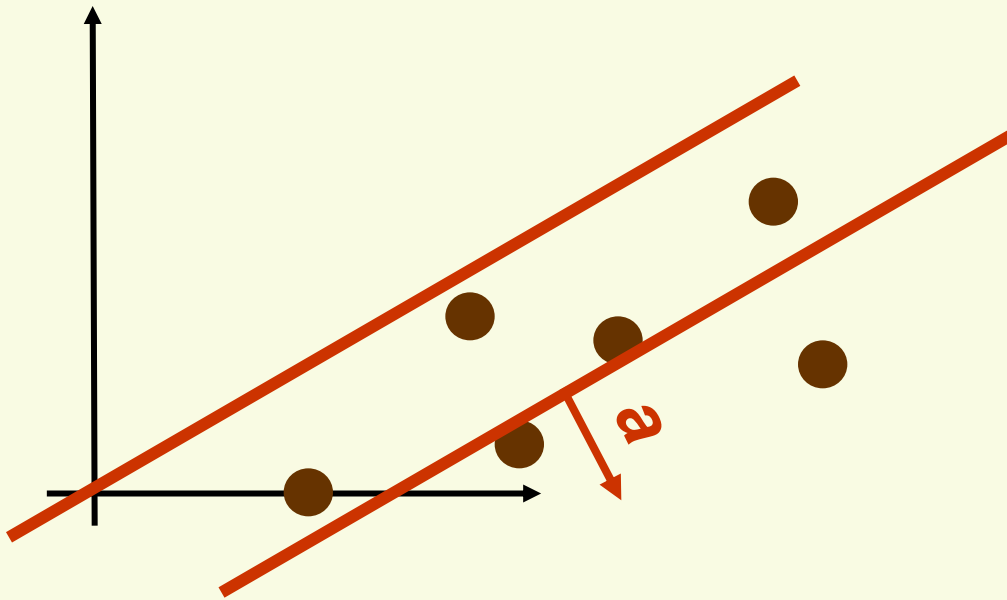
$$\left( (\mathbf{Y}^t\mathbf{Y})^{-1} \mathbf{Y}^t \right) \mathbf{Y} = (\mathbf{Y}^t\mathbf{Y})^{-1} (\mathbf{Y}^t\mathbf{Y}) = \mathbf{I}$$

# LDF: Minimum Squared-Error Procedures

- If  $b_1 = \dots = b_n = 1$ , MSE procedure is equivalent to finding a hyperplane of best fit through the samples  $\mathbf{y}_1, \dots, \mathbf{y}_n$

$$J_s(\mathbf{a}) = \|\mathbf{Y}\mathbf{a} - \mathbf{1}_n\|^2$$

$$\mathbf{1}_n = \begin{bmatrix} 1 \\ \vdots \\ 1 \end{bmatrix} \Bigg\} n$$



- Then we shift this line to the origin, if this line was a good fit, all samples will be classified correctly

# 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* gives 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  $\beta\mathbf{a}^*$

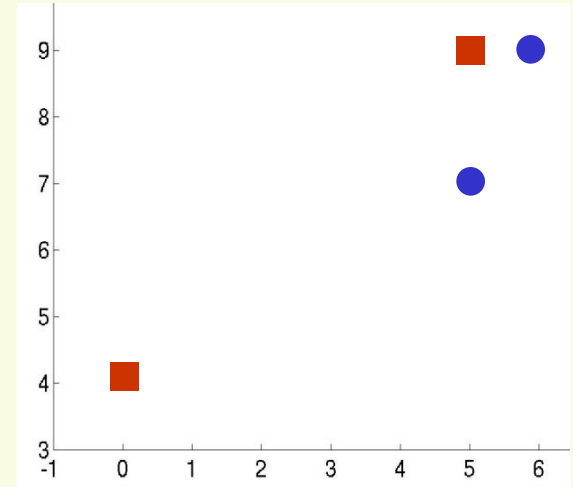
$$\begin{aligned} \underset{\mathbf{a}}{\operatorname{argmin}} \|\mathbf{Y}\mathbf{a} - \beta\mathbf{b}\|^2 &= \underset{\mathbf{a}}{\operatorname{argmin}} \beta^2 \|\mathbf{Y}(\mathbf{a}/\beta) - \mathbf{b}\|^2 \\ &= \underset{\mathbf{a}}{\operatorname{argmin}} \|\mathbf{Y}(\mathbf{a}/\beta) - \mathbf{b}\|^2 = \beta\mathbf{a}^* \end{aligned}$$

- 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: Example

- Class 1: (6 9), (5 7)
- Class 2: (5 9), (0 4)
- Set vectors  $\mathbf{y}_1, \mathbf{y}_2, \mathbf{y}_3, \mathbf{y}_4$  by adding extra feature and “normalizing”



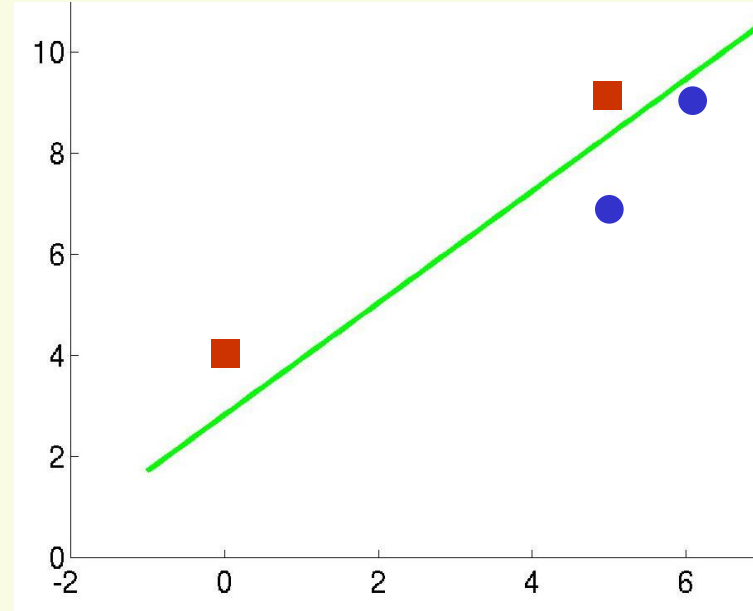
$$\mathbf{y}_1 = \begin{bmatrix} 1 \\ 6 \\ 9 \end{bmatrix} \quad \mathbf{y}_2 = \begin{bmatrix} 1 \\ 5 \\ 7 \end{bmatrix} \quad \mathbf{y}_3 = \begin{bmatrix} -1 \\ -5 \\ -9 \end{bmatrix} \quad \mathbf{y}_4 = \begin{bmatrix} -1 \\ 0 \\ -4 \end{bmatrix}$$

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

# 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

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



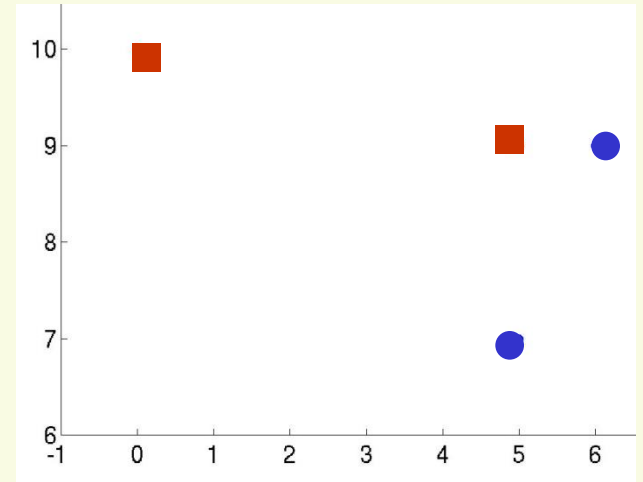
- Note  $\mathbf{a}$  is an approximation to  $\mathbf{Y}\mathbf{a} = \mathbf{b}$ , since no exact solution exists

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

- 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



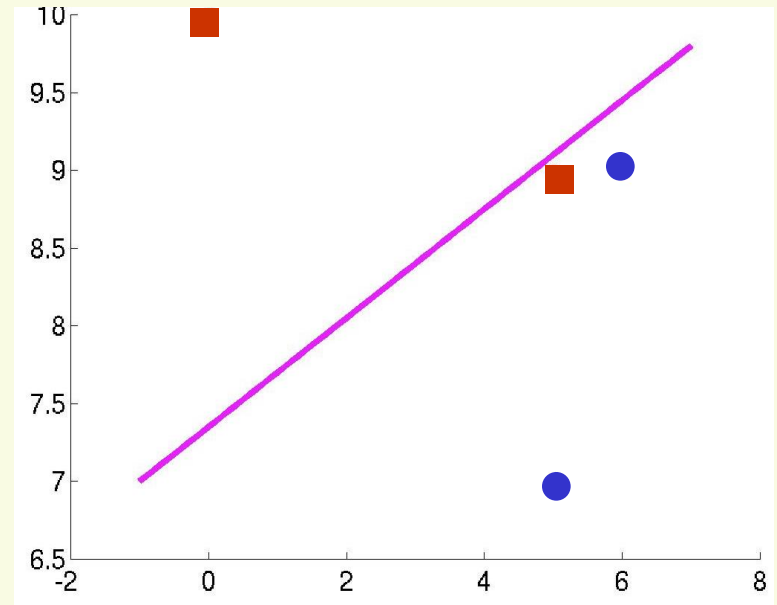
$$y_1 = \begin{bmatrix} 1 \\ 6 \\ 9 \end{bmatrix} \quad y_2 = \begin{bmatrix} 1 \\ 5 \\ 7 \end{bmatrix} \quad y_3 = \begin{bmatrix} -1 \\ -5 \\ -9 \end{bmatrix} \quad y_4 = \begin{bmatrix} -1 \\ 0 \\ -10 \end{bmatrix}$$

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

# 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

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

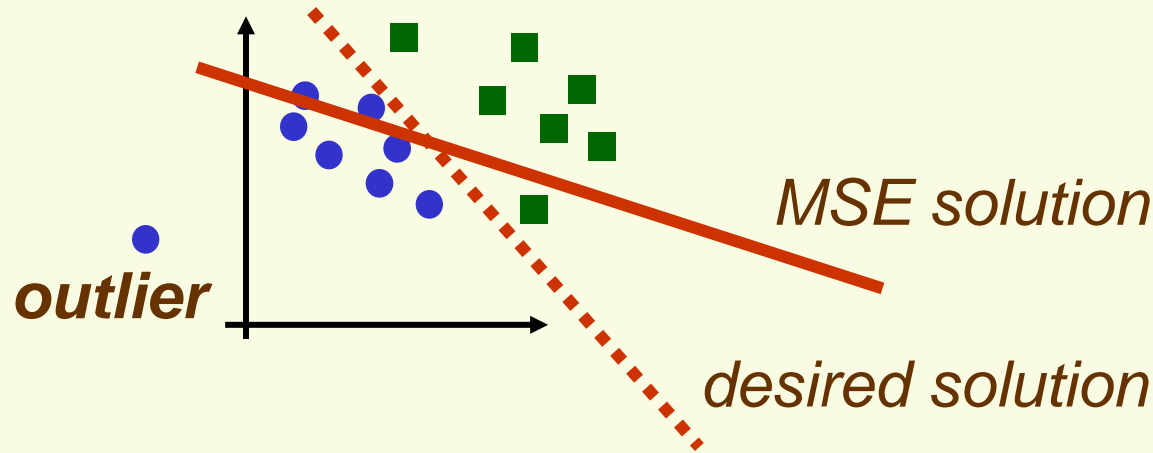
- This solution does not give a separating hyperplane since  $\mathbf{a}^t \mathbf{y}_3 < 0$



## *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 we don't know this

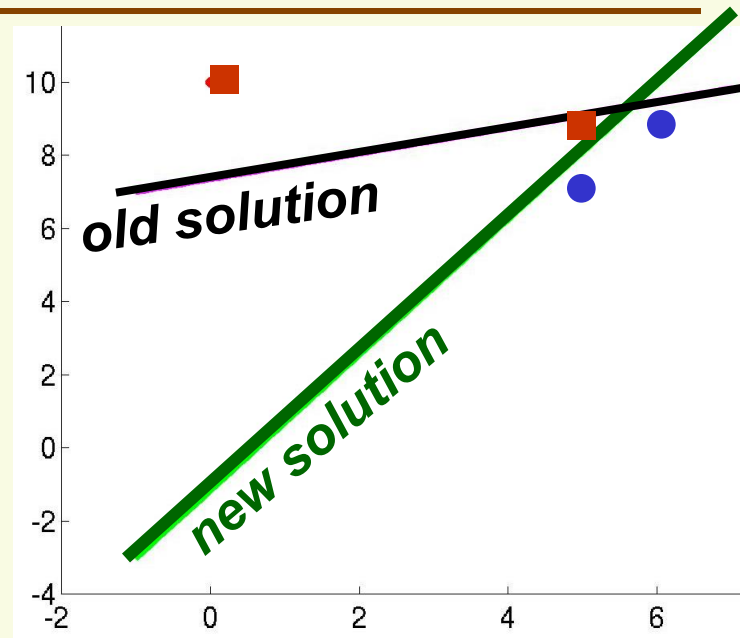
- Thus appropriate  $\mathbf{b} = \begin{bmatrix} 1 \\ 1 \\ 1 \\ 10 \end{bmatrix}$

- In Matlab, solve  $\mathbf{a} = \mathbf{Y} \backslash \mathbf{b}$

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

- Note  $\mathbf{a}$  is an approximation to  $\mathbf{Y}\mathbf{a} = \mathbf{b}$ ,  $\mathbf{Y}\mathbf{a} = \begin{bmatrix} 0.9 \\ 1.0 \\ 0.8 \\ 10.0 \end{bmatrix} \neq \begin{bmatrix} 1 \\ 1 \\ 1 \\ 10 \end{bmatrix}$

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



# ***LDF: Gradient Descent for MSE solution***

---

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

- 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
- In the beginning of the lecture, computed the gradient:

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

# ***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

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

- Given  $\mathbf{x}$ , assign class  $\mathbf{c}_i$  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  $\mathbf{c}$  decision regions, with  $\mathbf{g}_i(\mathbf{x})$  being the largest discriminant if  $\mathbf{x}$  is in the region  $R_i$

# LDF: MSE for Multiple Classes

---

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

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

- Let  $\mathbf{Y}_i$  be matrix whose rows are samples from class  $i$ , so it has  $\mathbf{d} + 1$  columns and  $n_i$  rows
- Let's pile all samples in  $n$  by  $\mathbf{d} + 1$  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  $\mathbf{1}$ :

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

# LDF: MSE for Multiple Classes

---

- Let's pile all  $\mathbf{b}_i$  as columns in  $n$  by  $c$  matrix  $\mathbf{B}$

$$\mathbf{B} = [\mathbf{b}_1 \cdots \mathbf{b}_n]$$

- Let's pile all  $\mathbf{a}_i$  as columns in  $d + 1$  by  $m$  matrix  $\mathbf{A}$

$$\mathbf{A} = [\mathbf{a}_1 \cdots \mathbf{a}_m] = \begin{bmatrix} \text{weights } \mathbf{a}_1 \\ \text{weights } \mathbf{a}_2 \\ \vdots \\ \text{weights } \mathbf{a}_m \end{bmatrix}$$

- $m$  LSE problems can be represented in  $\mathbf{YA} = \mathbf{B}$ :

$$\begin{array}{c} \begin{bmatrix} \text{sample from class1} \\ \text{sample from class1} \\ \text{sample from class2} \\ \text{sample from class3} \\ \text{sample from class3} \\ \text{sample from class3} \end{bmatrix} \\ \mathbf{Y} \end{array} \begin{array}{c} \begin{bmatrix} \text{weights for c1} \\ \text{weights for c2} \\ \text{weights for c3} \end{bmatrix} \\ \mathbf{A} \end{array} = \begin{array}{c} \begin{bmatrix} 1 & 0 & 0 \\ 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \\ 0 & 0 & 1 \\ 0 & 0 & 1 \end{bmatrix} \\ \mathbf{B} \end{array}$$



# ***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