

Pattern Recognition (Αναγνώριση Προτύπων)

Linear Discriminant Functions (Γραμμικές Συναρτήσεις Διάκρισης)

Panos Trahanias

UNIVERSITY OF CRETE
DEPARTMENT of COMPUTER
SCIENCE



Linear Discriminant Functions

- **Goal:** Formulate linear –with respect to feature vector \mathbf{x} - discriminant functions that define hyperplanes as decision surfaces.
- **Why?** Simple form, straightforward implementation, optimal for Gaussian pdfs.
- **How:** Formulate the parameter (weights) estimation problem as an optimization of a criterion (cost) function.
- **What is the criterion function?** A real valued function of the weights (parameters) that is amenable to minimization, e.g. probability of misclassification during training.
- **How hard is to accomplish?** In the general case, it is hard to formulate a linear classifier that minimizes error.
- **So what?** Employ alternative criteria (simple functions of the weights) and iterative optimization methods (e.g. gradient descent).

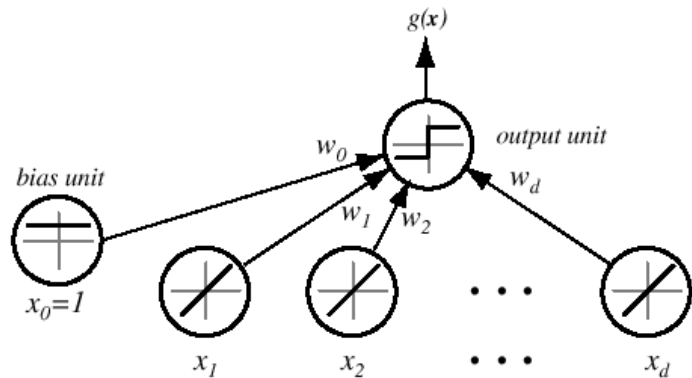


Discriminant Functions and Decision Surfaces

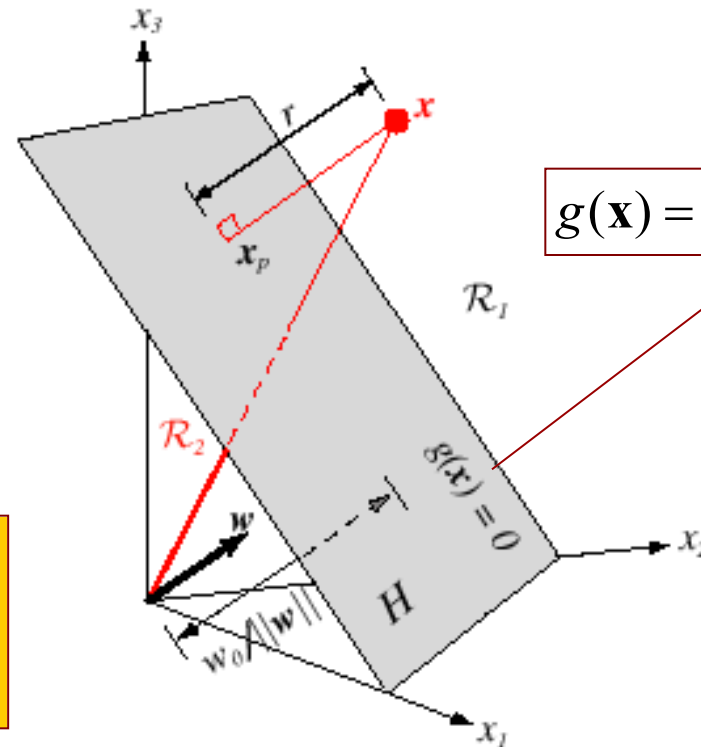
$$g(\mathbf{x}) = \mathbf{w}^T \mathbf{x} + w_0$$

Weight vector

Bias, threshold weight



Weight vector, \mathbf{w} , defines the orientation of the decision hyperplane and the threshold weight, w_0 , defines its relative position with respect to the origin.



$$g(\mathbf{x}) = \mathbf{w}^T \mathbf{x} + w_0 = 0$$

$$\mathbf{x} = \mathbf{x}_p + r \frac{\mathbf{w}}{\|\mathbf{w}\|}, \text{ with } r = \frac{g(\mathbf{x})}{\|\mathbf{w}\|}$$



Discriminant Functions and Decision Surfaces

If \mathbf{x}_1 and \mathbf{x}_2 are both on the decision surface, then

$$\mathbf{w}^t \mathbf{x}_1 + w_0 = \mathbf{w}^t \mathbf{x}_2 + w_0$$

$$\mathbf{w}^t (\mathbf{x}_1 - \mathbf{x}_2) = 0,$$

\mathbf{w} normal to any vector on the hyperplane,
hence normal to the hyperplane

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

$$\Rightarrow r = \frac{g(\mathbf{x})}{\|\mathbf{w}\|}.$$

$$\begin{aligned} \mathbf{w}^t \mathbf{x} + w_0 &= \\ \mathbf{w}^t \left(\mathbf{x}_p + r \frac{\mathbf{w}}{\|\mathbf{w}\|} \right) + w_0 &= \\ \mathbf{w}^t \mathbf{x}_p + r \frac{\mathbf{w}^t \mathbf{w}}{\|\mathbf{w}\|} + w_0 &= r \|\mathbf{w}\| \end{aligned}$$



The multi-class Case

Linear Machine:

Decision boundaries:

Distance of \mathbf{x} from H_{ij} :

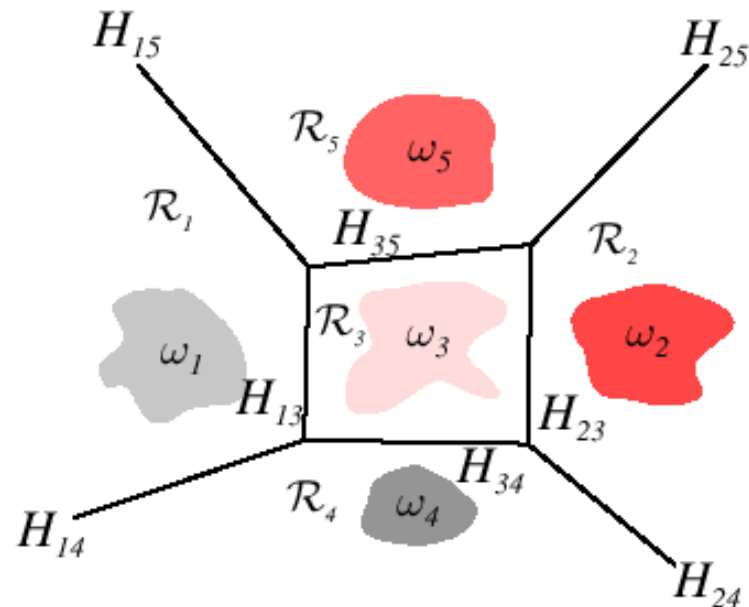
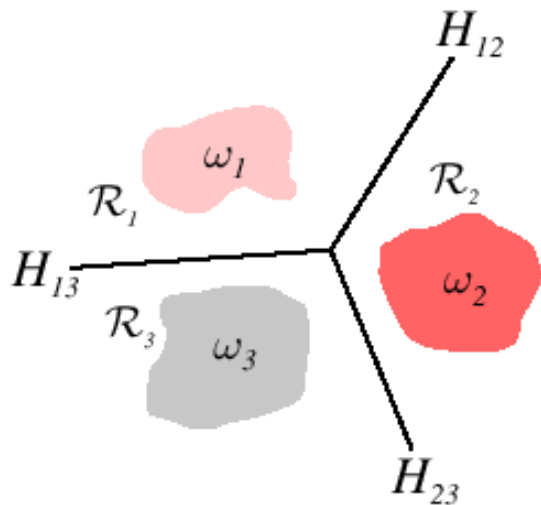
Convex decision regions.

\mathbf{x} in ω_i if $g_i(\mathbf{x}) > g_j(\mathbf{x})$

H_{ij} : $g_i(\mathbf{x}) = g_j(\mathbf{x}) \rightarrow (\mathbf{w}_i - \mathbf{w}_j)^T \mathbf{x} + (w_{i0} - w_{j0}) = 0$
 - part of hyperplane normal to vector $\mathbf{w}_i - \mathbf{w}_j$.

$(g_i(\mathbf{x}) - g_j(\mathbf{x})) / \|\mathbf{w}_i - \mathbf{w}_j\|$

- differences in weight vectors play a role.

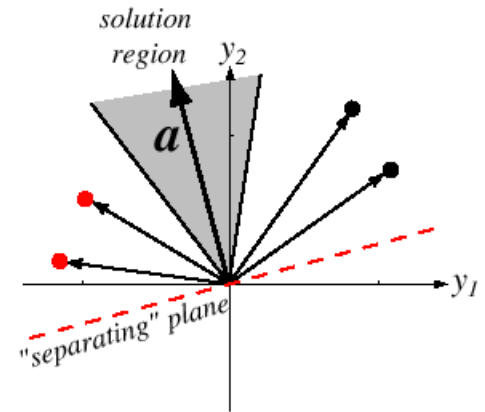
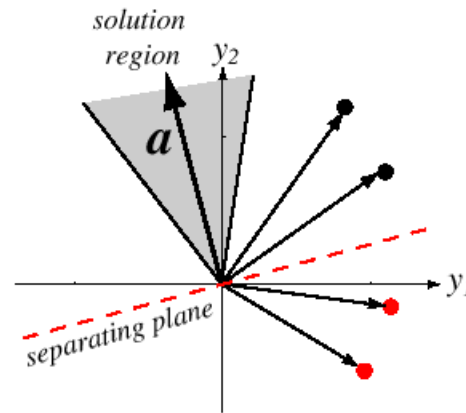




Linearly Separable Classes: Vectors and Solution Regions

Augmented Weight and Feature Vectors:

$$\mathbf{a} = \begin{bmatrix} w_0 \\ w_1 \\ w_2 \\ \vdots \\ w_d \end{bmatrix} = \begin{bmatrix} w_0 \\ \mathbf{w} \end{bmatrix}, \quad \mathbf{y} = \begin{bmatrix} 1 \\ x_1 \\ x_2 \\ \vdots \\ x_d \end{bmatrix} = \begin{bmatrix} x_0 \\ \mathbf{x} \end{bmatrix}$$

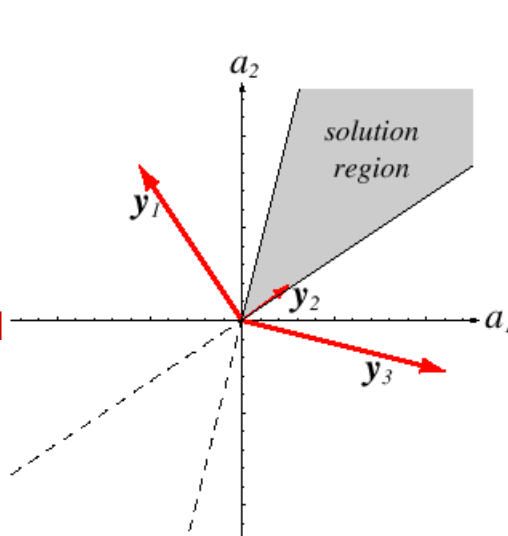


$$H: g(\mathbf{y}) = \mathbf{a}^T \mathbf{y} = 0$$

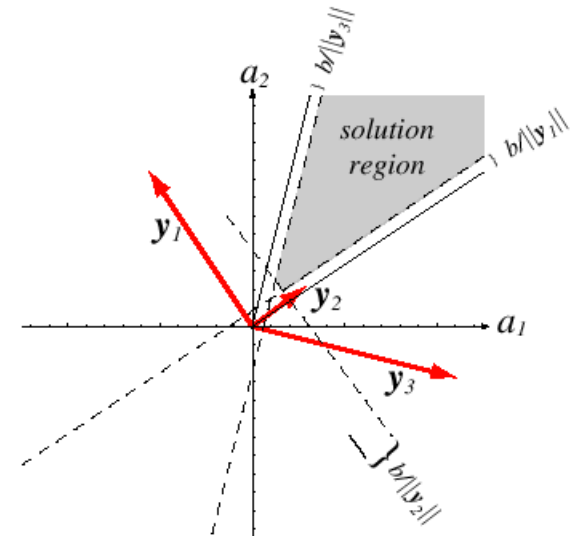
$$\text{Distance of } \mathbf{y} \text{ from } H: \frac{|\mathbf{a}^T \mathbf{y}|}{\|\mathbf{a}\|}$$

Normalization: replace all training samples from class ω_2 with their negatives, and find the separating vectors that satisfy the relation:

$$\mathbf{a}^T \mathbf{y}_i > 0 \quad \forall i$$



$$\mathbf{a}^T \mathbf{y}_i > 0 \quad \forall i$$



$$\mathbf{a}^T \mathbf{y}_i \geq b > 0 \quad \forall i$$



Optimization Procedures

- Task: Find \mathbf{a} that satisfies the set of linear inequalities $\mathbf{a}^t \mathbf{y}_i > 0$ for all $i=1, \dots, n$.
- How an appropriate solution is found?
 - ↳ Define a criterion function, $J(\mathbf{a})$, and minimize it to obtain \mathbf{a} as a solution vector.
 - ↳ Accordingly, we transform the exhaustive search problem to that of minimizing a real-valued function.
- How is $J(\mathbf{a})$ minimized?
 - ↳ Select an initial point \mathbf{a}_1 and compute $J(\mathbf{a}_1)$.
 - ↳ Compute the derivative of $J(\mathbf{a}_1)$: $\nabla J(\mathbf{a}_1)$.
 - ↳ Find the next point \mathbf{a}_2 in the direction of steepest descent, $-\nabla J(\mathbf{a}_1)$, using a step $\eta(k)$, that is, the learning rate or the stepsize.

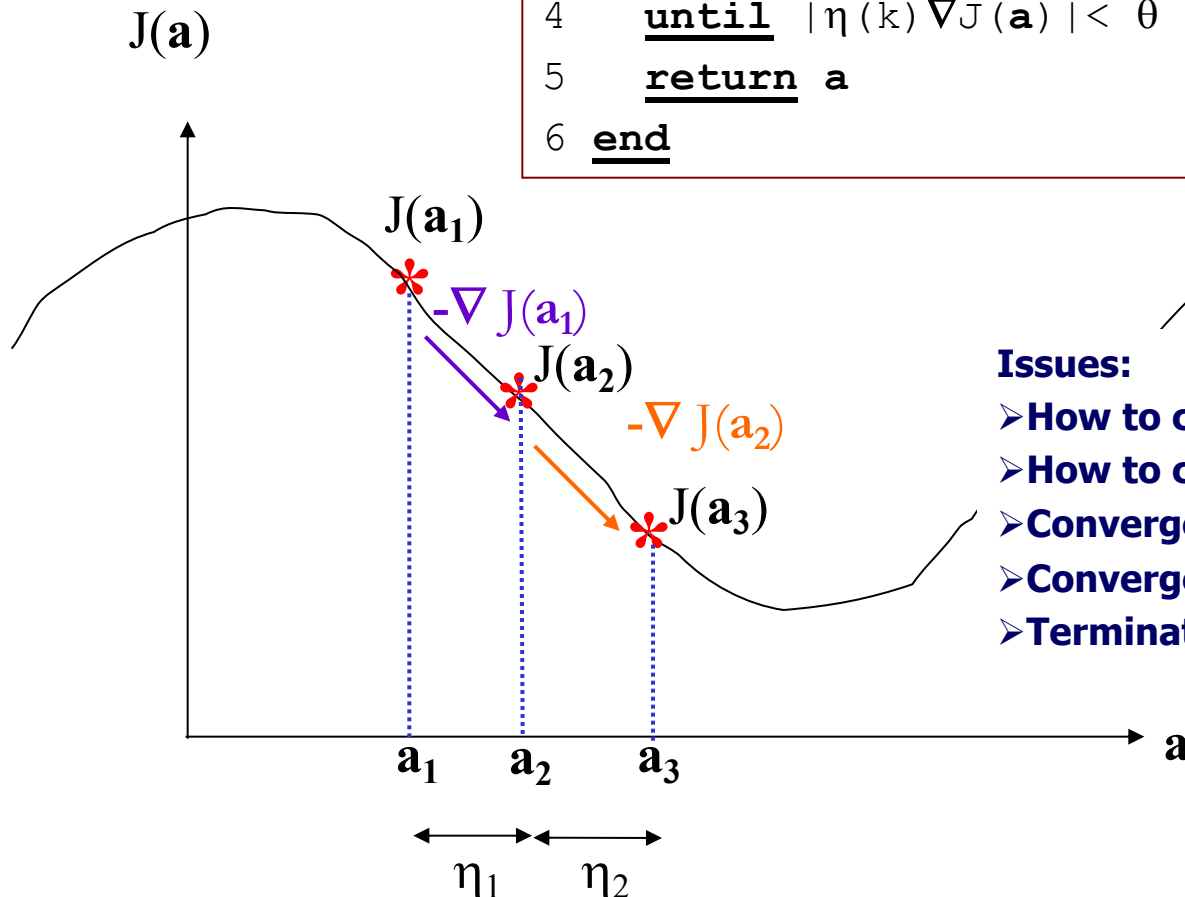


Steepest Descent Algorithm

Algorithm 1. Steepest Descent

```
1 begin initialize  $\mathbf{a}$ , threshold  $\theta$ ,  $\eta(0) > 0$ ,  $k=0$   
2   do  $k \leftarrow k+1$   
3    $\mathbf{a} \leftarrow \mathbf{a} - \eta(k) \nabla J(\mathbf{a})$   
4   until  $|\eta(k) \nabla J(\mathbf{a})| < \theta$   
5   return  $\mathbf{a}$   
6 end
```

$$\mathbf{a}_{k+1} = \mathbf{a}_k - \eta_k \nabla J(\mathbf{a}_k)$$



Issues:

- How to choose criterion function?
- How to choose learning rate $\eta(k)$?
- Convergence in local/global minimum?
- Convergence rate/speed/smoothness?
- Termination criterion?

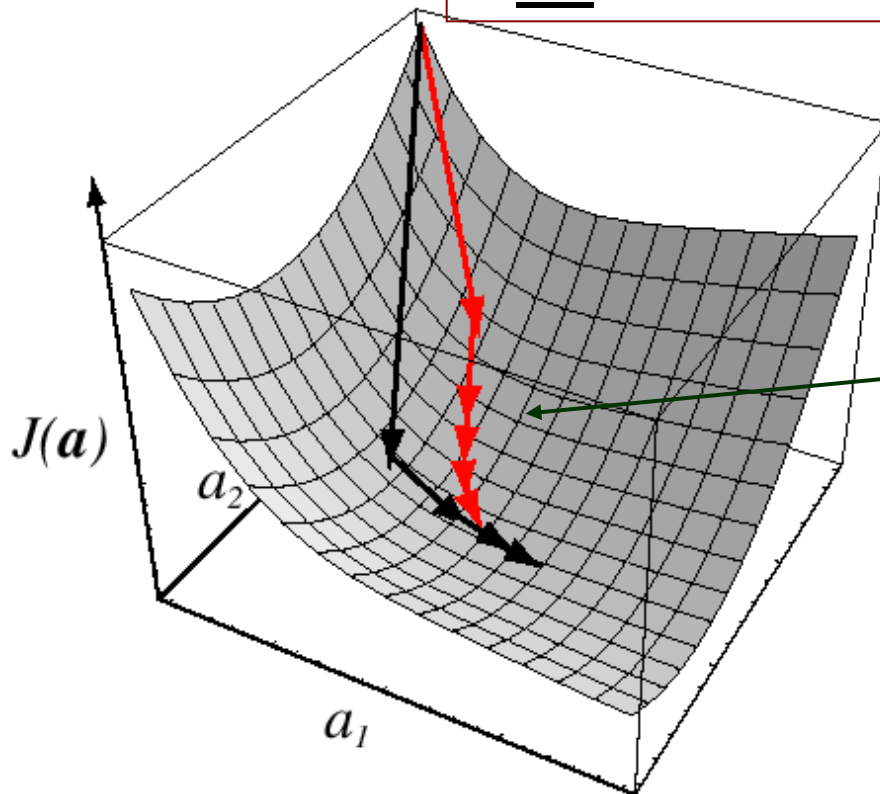


Newton Descent Algorithm

Algorithm 2. Newton Descent

```
1 begin initialize  $\mathbf{a}$ , threshold  $\theta$   
2    $\mathbf{a} \leftarrow \mathbf{a} - \mathbf{H}^{-1} \nabla J(\mathbf{a})$   
3   until  $|\mathbf{H}^{-1} \nabla J(\mathbf{a})| < \theta$   
4   return  $\mathbf{a}$   
5 end
```

$$\mathbf{a}_{k+1} = \mathbf{a}_k - \mathbf{H}_k^{-1} \nabla J(\mathbf{a}_k)$$



$$\mathbf{H}_k = \left[\partial^2 J(\mathbf{a}) / \partial a_i \partial a_j \right]_{\mathbf{a}=\mathbf{a}_k}$$

Red: Steepest Descent

Black: Newton Descent

Newton: greater improvement in each step at the computational cost of the inversion of the Hessian matrix \mathbf{H} .



Perceptron Criterion

➤ Choices of the Criterion Function

- ↪ A first choice: No of erroneously classified training samples.
However: non-continuous function, hence non-differentiable.
- ↪ A better choice: The perceptron criterion function:

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

where $\mathcal{Y}(\mathbf{a})$ is the set of training samples that are erroneously classified by \mathbf{a} .

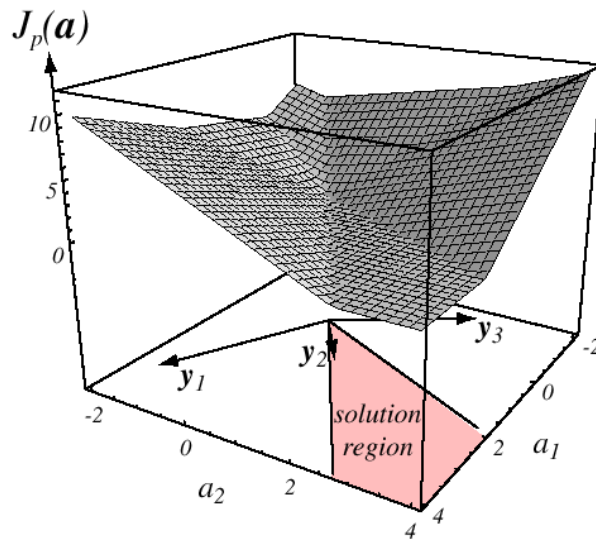
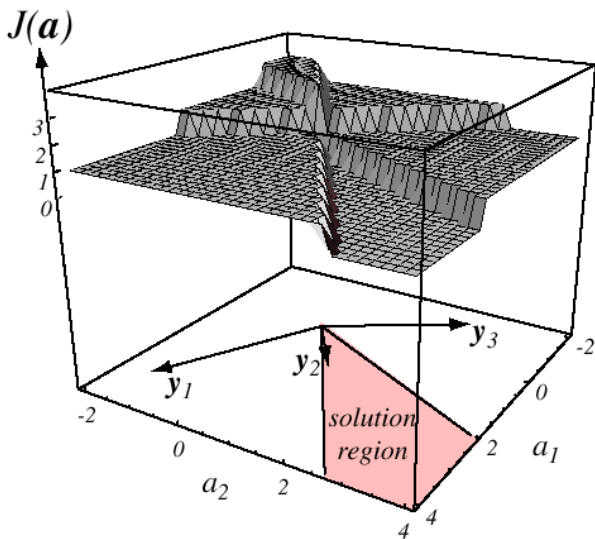
- ↪ If $\mathcal{Y}(\mathbf{a})$ is empty, then $J_p(\mathbf{a})=0$. Given $\mathbf{a}^t \mathbf{y} < 0$ when \mathbf{y} erroneously classified, $J_p(\mathbf{a})$ is never negative and becomes zero when \mathbf{a} is a solution vector.
- ↪ Geometric interpretation: $J_p(\mathbf{a})$ is proportional to the sum of the distances of the erroneously classified samples from the decision boundary.



Four Cost Functions

No of
erroneous
classifications

Bad



Perceptron
Criterion

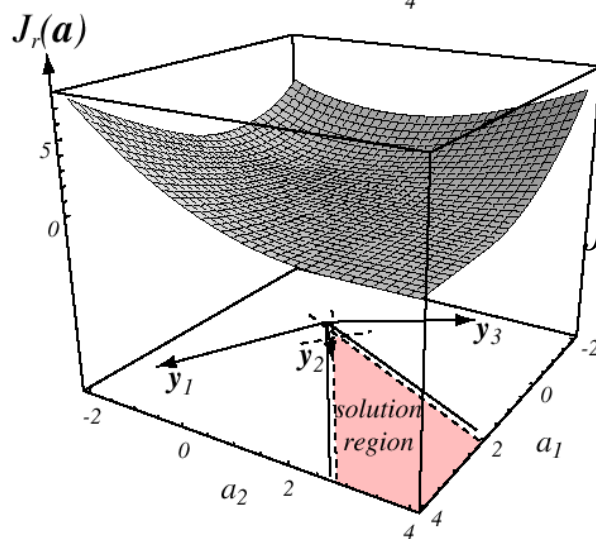
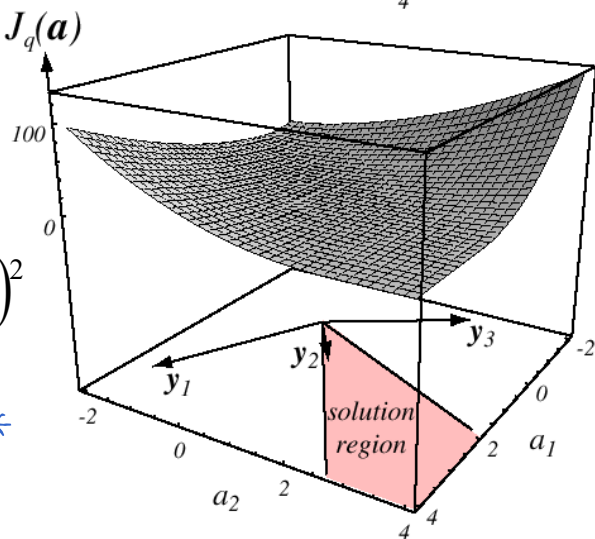
$$J_p(\mathbf{a}) = \sum_{\mathbf{y} \in \mathcal{U}} (-\mathbf{a}^t \mathbf{y})$$

Good!

Total square
error (TSE)

$$J_q(\mathbf{a}) = \sum_{\mathbf{y} \in \mathcal{U}} (\mathbf{a}^t \mathbf{y})^2$$

Better*



TSE with
margin

$$J_r(\mathbf{a}) = \frac{1}{2} \sum_{\mathbf{y} \in \mathcal{U}} \frac{(\mathbf{a}^t \mathbf{y} - b)^2}{\|\mathbf{y}\|}$$

Best*

* Αλλά θα μπορούσε να έχει μεγάλο υπολογιστικό κόστος



Batch Perceptron Algorithm

⇒ After differentiating:

$$\nabla J_p(\mathbf{a}) = \left[\frac{\partial J_p}{\partial a_i} \right] = \sum_{y \in \mathcal{U}} -\mathbf{y}$$

⇒ Recursive formulation:

$$\mathbf{a}(k+1) = \mathbf{a}(k) + \eta(k) \sum_{y \in \mathcal{U}_k} \mathbf{y}$$

where \mathcal{U}_k is the set of erroneously classified samples by \mathbf{a}_k .

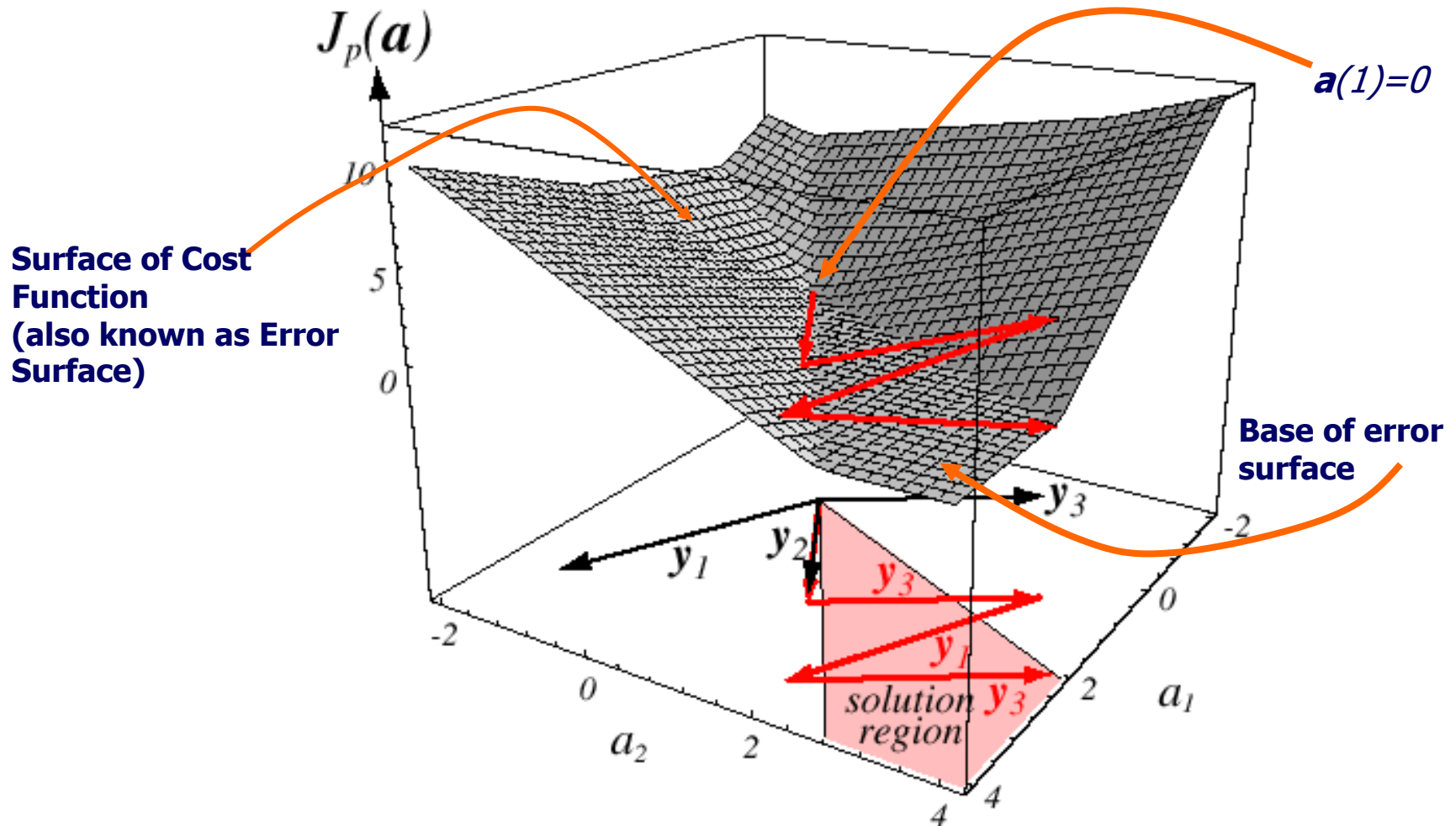
⇒ Next weight vector (w.v.) estimated as the sum of the current w.v. and a multiple of the sum of the erroneously classified samples.

Algorithm 3. Batch Perceptron

```
1 begin initialize  $\mathbf{a}$ , criterion  $\theta$ ,  $\eta(0) > 0$ ,  $k=0$   
2   do  $k \leftarrow k+1$   
3      $\mathbf{a} \leftarrow \mathbf{a} + \eta(k) \sum_{y \in \mathcal{U}_k} \mathbf{y}$   
4   until  $\left| \eta(k) \sum_{y \in \mathcal{U}_k} \mathbf{y} \right| < \theta$   
5   return  $\mathbf{a}$   
6 end
```



Batch Perceptron Steps





Fixed-Increment Single-Sample Perceptron

- Instead of applying the weight vector $\mathbf{a}(k)$ to all samples and then correcting it based on the \mathcal{Y}_k set of the erroneously classified samples, we use the samples one at a time and update or not the weight vector based on the classification result.
- Moreover, if we employ a constant step $\eta(k)$, then we derive the following algorithm:

Algorithm 4. Fixed-Increment Single-Sample Perceptron

```
1 begin initialize  $\mathbf{a}$ ,  $k=0$   
2   do  $k \leftarrow (k+1) \bmod n$   
3     If  $\mathbf{y}^k$  is misclassified by  $\mathbf{a}$ , then  $\mathbf{a} \leftarrow \mathbf{a} + \mathbf{y}^k$   
4   until all patterns properly classified  
5   return  $\mathbf{a}$   
6 end
```

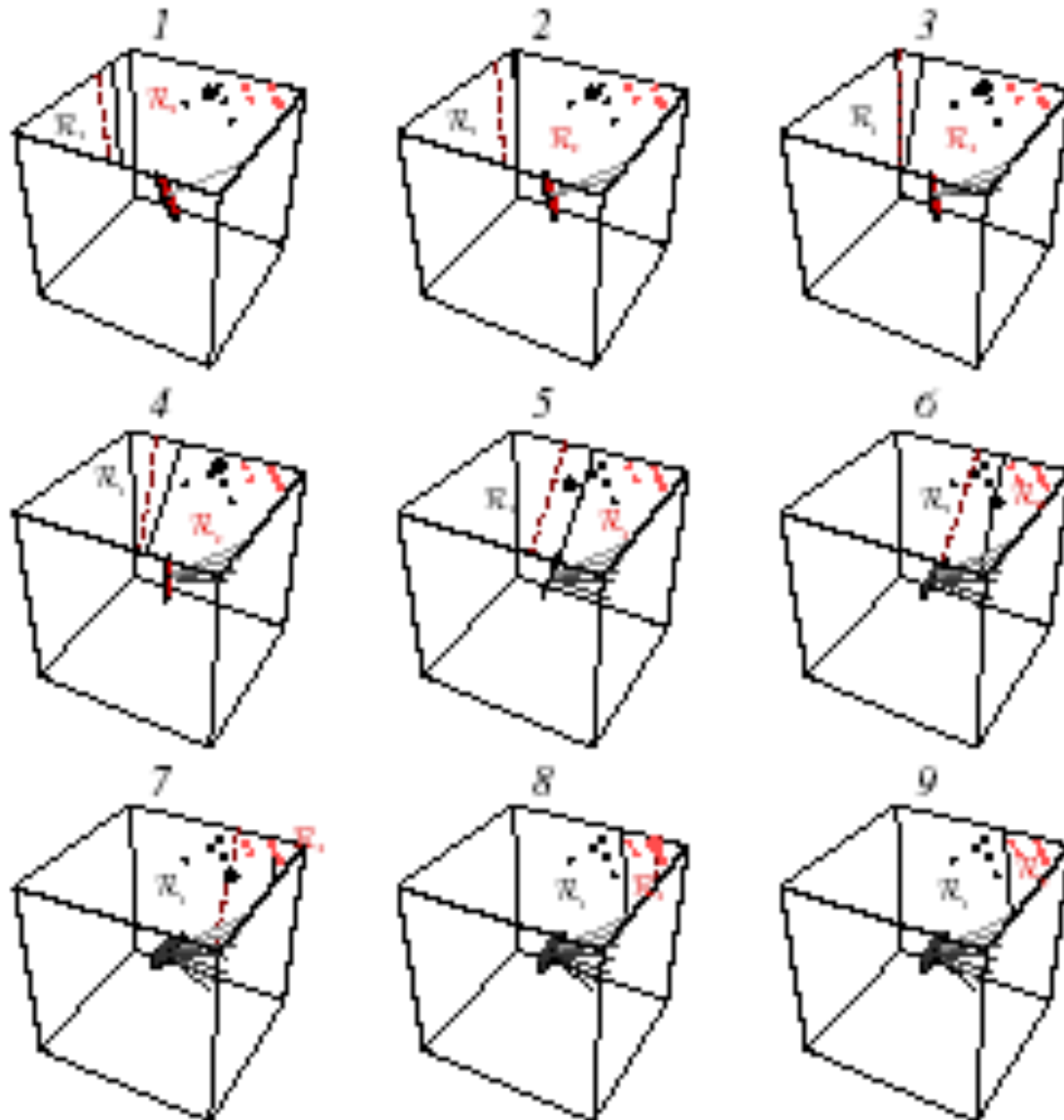
Circular order of samples (green signifies the erroneously classified training samples):

\mathbf{y}_1 \mathbf{y}_2 \mathbf{y}_3 \mathbf{y}_4 \mathbf{y}_1 \mathbf{y}_2 \mathbf{y}_3 \mathbf{y}_4 \mathbf{y}_1 \mathbf{y}_2 \mathbf{y}_3 \mathbf{y}_4

➔ $\mathbf{y}^1 \mathbf{y}^2 \mathbf{y}^3 \mathbf{y}^0 \mathbf{y}^1 \dots = \mathbf{y}_2 \mathbf{y}_1 \mathbf{y}_3 \mathbf{y}_2 \mathbf{y}_3$



Convergence





Variable-Increment Perceptron with Margin

- Use the samples one at a time and update the weight vector $\mathbf{a}(k)$ when its inner product with \mathbf{y}^k is less than a preset threshold b : $\mathbf{a}(k)\mathbf{y}^k < b$.
- Moreover, if we employ a variable step $\eta(k)$, we get the following algorithm:

Algorithm 5. Variable-Increment Perceptron w. Margin

```
1 begin initialize  $\mathbf{a}$ , threshold  $\theta$ , margin  $b$ ,  $\eta(0)$ ,  $k=0$   
2   do  $k \leftarrow (k+1) \bmod n$   
3     if  $\mathbf{a}^t \mathbf{y}^k < b$ , then  $\mathbf{a} \leftarrow \mathbf{a} + \eta(k) \mathbf{y}^k$   
4   until  $\mathbf{a}^t \mathbf{y}^k > b$  for all  $k$   
5   return  $\mathbf{a}$   
6 end
```

Convergence Conditions:

$$\eta(k) \geq 0, \quad \lim_{m \rightarrow \infty} \sum_{k=1}^m \eta(k) = \infty, \quad \lim_{m \rightarrow \infty} \frac{\sum_{k=1}^m \eta^2(k)}{\left(\sum_{k=1}^m \eta(k)\right)^2} = 0$$



Relaxation Procedures

➤ Criterion Function:

$$J_r(\mathbf{a}) = \frac{1}{2} \sum_{\mathbf{y} \in \mathcal{U}} \frac{(\mathbf{a}^t \mathbf{y} - b)^2}{\|\mathbf{y}\|^2}$$

where $\mathcal{U}(\mathbf{a})$ is the set of samples for which $\mathbf{a}^t \mathbf{y} < b$.

↪ If $\mathcal{U}(\mathbf{a})$ is empty, then $J_r(\mathbf{a})=0$. $J_r(\mathbf{a})$ can never become negative and is zero if and only if $\mathbf{a}^t \mathbf{y} > b$ for all the training samples.

↪ Gradient Vector:

$$\nabla J_r(\mathbf{a}) = \left[\frac{\partial J_r}{\partial a_i} \right] = \sum_{\mathbf{y} \in \mathcal{U}} \frac{\mathbf{a}^t \mathbf{y} - b}{\|\mathbf{y}\|^2} \mathbf{y}$$

↪ Update Rule:

$$\mathbf{a}(k+1) = \mathbf{a}(k) + \eta(k) \sum_{\mathbf{y} \in \mathcal{U}_k} \frac{b - \mathbf{a}^t \mathbf{y}}{\|\mathbf{y}\|^2} \mathbf{y}$$



Batch Relaxation with Margin

Algorithm 6. Batch Relaxation with Margin

```
1 begin initialize  $\mathbf{a}$ , margin  $b$ ,  $\eta(0)$ ,  $k=0$ 
2   do  $k \leftarrow (k+1) \bmod n$ 
3      $\mathcal{Y}_k = \{ \}$ 
4      $j=0$ 
5     do  $j \leftarrow j+1$ 
6       if  $\mathbf{a}^t \mathbf{y}^k < b$ , then append  $\mathbf{y}^j$  to  $\mathcal{Y}_k$ 
7     until  $j=n$ 
8      $\mathbf{a} \leftarrow \mathbf{a} + \eta(k) \sum_{\mathbf{y} \in \mathcal{Y}_k} \frac{b - \mathbf{a}^t \mathbf{y}}{\|\mathbf{y}\|^2} \mathbf{y}$ 
9   until  $\mathcal{Y}_k = \{ \}$ 
10  return  $\mathbf{a}$ 
11 end
```



Single-Sample Relaxation with Margin

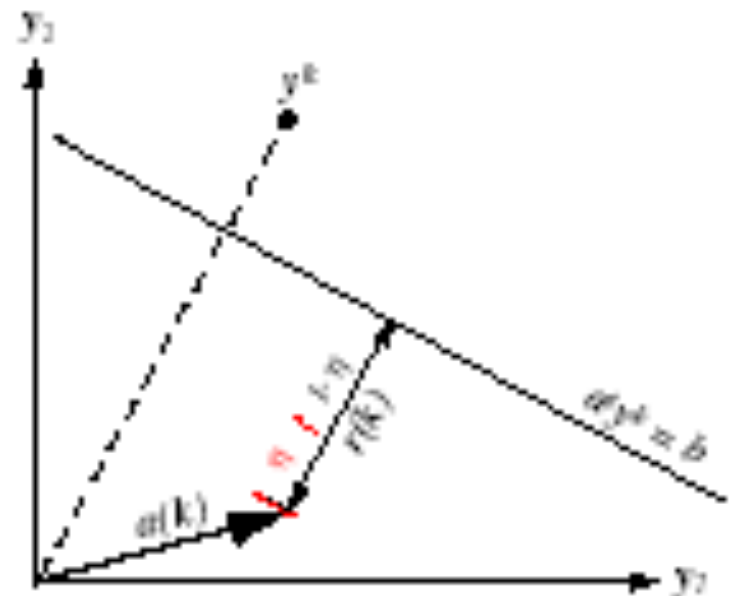
Algorithm 7. Single-Sample Relaxation with Margin

```

1 begin initialize  $\mathbf{a}$ , margin  $b$ ,  $\eta(0)$ ,  $k=0$ 
2   do  $k \leftarrow (k+1) \bmod n$ 
3     if  $\mathbf{a}^t \mathbf{y}^k < b$ , then  $\mathbf{a} \leftarrow \mathbf{a} + \eta(k) \frac{b - \mathbf{a}^t \mathbf{y}^k}{\|\mathbf{y}^k\|^2} \mathbf{y}^k$ 
4   until  $\mathbf{a}^t \mathbf{y}^k > b$  for all  $\mathbf{y}^k$ 
5   return  $\mathbf{a}$ 
6 end
  
```

At each step, weight vector $\mathbf{a}(k)$, is moved towards the hyperplane $\mathbf{a}^t \mathbf{y}^k = b$ by a fraction, $\eta(k)$, of its distance, $r(k)$, from that.

$\eta(k) < 1 \rightarrow$ underrelaxation
 $\eta(k) > 1 \rightarrow$ overrelaxation
 $0 < \eta(k) < 2$ for convergence





Minimum Square-Error – MSE

- **Goal:** Good performance in both cases, linearly separable and non-linearly separable.
- **How:** Criterion that involves all patterns. Moreover,
 - ↳ Before: Find \mathbf{a} such as $\mathbf{a}^t \mathbf{y}_i > 0$ for each pattern \mathbf{y}_i .
 - ↳ Now: Find \mathbf{a} such as $\mathbf{a}^t \mathbf{y}_i = b_i$ for each pattern \mathbf{y}_i . (b_i positive constants).
- **Accordingly:** Instead of solving a problem of linear inequalities, solve one of linear equations.
- **Notations:**

$$\mathbf{Y}_{n \times \hat{d}} = \begin{bmatrix} \mathbf{y}_1^t \\ \mathbf{y}_2^t \\ \vdots \\ \mathbf{y}_n^t \end{bmatrix}$$

Pattern Matrix

$$\mathbf{b}_{n \times 1} = \begin{bmatrix} b_1 \\ b_2 \\ \vdots \\ b_n \end{bmatrix}$$

Vector of positive constants

$$\mathbf{a}_{\hat{d} \times 1} = \begin{bmatrix} a_0 \\ a_1 \\ \vdots \\ a_d \end{bmatrix}$$

Vector of weights

- **Formulation:** Find \mathbf{a} such as: $\mathbf{Y}\mathbf{a}=\mathbf{b}$.



Minimum Square-Error – MSE

➤ Commonly $n > d+1$, number of patterns $>$ dimensions \rightarrow overdetermined system, hence no exact solution.

➤ **Therefore:** Minimization of the square of the error vector length, $\mathbf{e} = \mathbf{Y}\mathbf{a} - \mathbf{b}$:

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

➤ **Formulation:**

$$\nabla J_s(\mathbf{a}) = \mathbf{0} \Rightarrow \mathbf{Y}^t \mathbf{Y} \mathbf{a} = \mathbf{Y}^t \mathbf{b}$$

Normal Equations

➤ If $\mathbf{Y}^t \mathbf{Y}$ can be inverted,

$$\mathbf{a} = \underbrace{(\mathbf{Y}^t \mathbf{Y})^{-1} \mathbf{Y}^t}_{\text{Pseudo-inverse}} \mathbf{b}$$



Widrow-Hoff (Least-mean squares-LMS)

- $J_s(\mathbf{a})$ can be minimized via recursive algorithms that do not require matrix inversions.
- Tangent vector: $\nabla J_s(\mathbf{a}) = 2\mathbf{Y}^t(\mathbf{Y}\mathbf{a} - \mathbf{b})$
- Basic recursive formulae: $\mathbf{a}(k+1) = \mathbf{a}(k) + \eta(k)\mathbf{Y}^t(\mathbf{b} - \mathbf{Y}\mathbf{a}(k))$
- By employing one sample at a step, the Widrow-Hoff (LMS) algorithm is derived:

Αλγόριθμος 8. Widrow-Hoff (LMS)

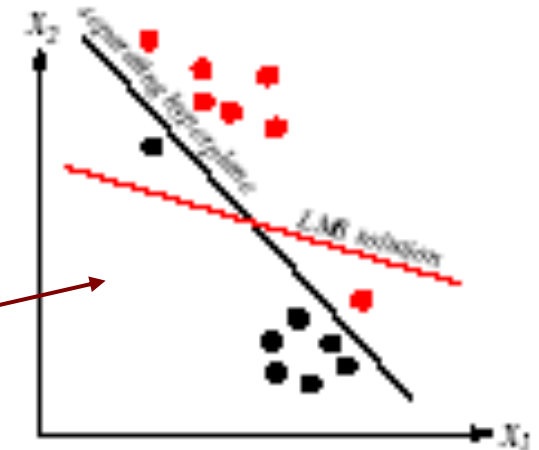
```
1 begin initialize  $\mathbf{a}$ ,  $\mathbf{b}$ , κριτήριο  $\theta$ ,  $\eta()$ ,  $k=0$   
2 do  $k \leftarrow (k+1) \bmod n$   
3  $\mathbf{a} \leftarrow \mathbf{a} + \eta(k)(b_k - \mathbf{a}^t \mathbf{y}^k) \mathbf{y}^k$   
4 until  $|\eta(k)(b_k - \mathbf{a}^t \mathbf{y}^k) \mathbf{y}^k| < \theta$   
5 return  $\mathbf{a}$   
6 end
```



Ho-Kashyap Method

- Perceptron and Relaxation procedures find separating Weight Vectors when the samples are linearly separable, but do not converge for non-separable classes.

- Least-mean squares offers always a solution vector (the one that minimizes $\|\mathbf{Y}\mathbf{a}-\mathbf{b}\|^2$) which is not necessarily separating vector in the separable case.



- The Ho-Kashyap algorithm recursively solves the minimization problem:

$$\min_{\mathbf{a}, \mathbf{b}} J_s(\mathbf{a}, \mathbf{b}) = \|\mathbf{Y}\mathbf{a} - \mathbf{b}\|^2 \quad s.t. \quad \mathbf{b} > \mathbf{0}$$

- Algorithmically it achieves \mathbf{b} not to converge to $\mathbf{0}$, by setting all the positive components of the tangent vector $\nabla_{\mathbf{b}} J_s$ equal to zero.



Ho-Kashyap Algorithm

Algorithm 9. Ho-Kashyap

```
1 begin initialize  $a, b, \eta() < 1, \text{threshold } b_{\min}, k_{\max}$   
2   do  $k \leftarrow (k+1) \bmod n$   
3      $e \leftarrow Ya - b$   
4      $e^+ \leftarrow (e + |e|) / 2$   
5      $b \leftarrow b + 2\eta(k)e^+$   
6      $a \leftarrow (Y^t Y)^{-1} Y^t b$   
7     if  $\text{Abs}(e) < b_{\min}$  then return  $a, b$  and exit  
8   until  $k = k_{\max}$   
9   print "No solution found"  
10 end
```




Multiple Classes Kesler Structure

- Goal: Linear separation of multiple classes:

If $\mathbf{y} \sim \omega_1$, then $\mathbf{a}_1^t \mathbf{y} - \mathbf{a}_j^t \mathbf{y} > 0$ for all $j=2, \dots, c$.

- The system of $c-1$ inequalities can be described as:

The $c\hat{d} - D$ weight vector $\hat{\mathbf{a}}$ classifies correctly all $c-1$ $c\hat{d} - D$ patterns $\boldsymbol{\eta}_{12}, \boldsymbol{\eta}_{13}, \dots, \boldsymbol{\eta}_{1c}$:
 $\hat{\mathbf{a}}^t \boldsymbol{\eta}_{1j} > 0 \quad \forall \quad j = 2, \dots, c$ where:

$$\hat{\mathbf{a}}_{c\hat{d} \times 1} = \begin{bmatrix} \mathbf{a}_1 \\ \mathbf{a}_2 \\ \vdots \\ \mathbf{a}_c \end{bmatrix} \quad \boldsymbol{\eta}_{12} = \begin{bmatrix} \mathbf{y} \\ -\mathbf{y} \\ \mathbf{0} \\ \vdots \\ \mathbf{0} \end{bmatrix} \quad \boldsymbol{\eta}_{13} = \begin{bmatrix} \mathbf{y} \\ \mathbf{0} \\ -\mathbf{y} \\ \vdots \\ \mathbf{0} \end{bmatrix}, \dots, \quad \boldsymbol{\eta}_{1c} = \begin{bmatrix} \mathbf{y} \\ \mathbf{0} \\ \vdots \\ \mathbf{0} \\ -\mathbf{y} \end{bmatrix} \quad \text{Kesler Structure}$$

- Generally: $\hat{\mathbf{a}}^t \boldsymbol{\eta}_{ij} > 0 \quad \forall \quad j \neq i$, with $\boldsymbol{\eta}_{ij} = \begin{bmatrix} \vdots \\ \mathbf{y} \\ \vdots \\ -\mathbf{y} \\ \vdots \end{bmatrix} \begin{matrix} \leftarrow i \\ \leftarrow j \end{matrix}$



Multi-Class Perceptron Classification

- Let $\mathbf{y}_1, \dots, \mathbf{y}_n$ patterns from c classes, linearly separable. Let L_k a linear machine $\mathbf{a}_1(k), \dots, \mathbf{a}_c(k)$. We seek to formulate a series of linear machines L_1, \dots, L_k, \dots that converges to a separation (classification) machine L .
- Let \mathbf{y}^k the k -th sample that is amenable to correction (correct classification). If $\mathbf{y}^k \sim \omega_i$, then at least one $j \neq i$ exists for which $\mathbf{a}_i^t(k) \mathbf{y}^k < \mathbf{a}_j^t(k) \mathbf{y}^k$.
- L_k correction rules specifies (perceptron with constant unit-step) :

$$\mathbf{a}(k+1) = \mathbf{a}(k) + \boldsymbol{\eta}_{ij}^k \quad \text{where} \quad \mathbf{a}^t(k) \boldsymbol{\eta}_{ij}^k \leq 0 \quad \text{with} \quad \mathbf{a}(k) = \begin{bmatrix} \mathbf{a}_1(k) \\ \vdots \\ \mathbf{a}_c(k) \end{bmatrix} \quad \text{and} \quad \boldsymbol{\eta}_{ij}^k = \begin{bmatrix} \vdots \\ \mathbf{y}^k \\ \vdots \\ -\mathbf{y}^k \\ \vdots \end{bmatrix} \begin{matrix} \leftarrow i \\ \\ \leftarrow j \end{matrix}$$

Therefore:

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

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

$$\mathbf{a}_l(k+1) = \mathbf{a}_l(k) \quad l \neq i \quad \text{and} \quad l \neq j$$