



# Αναγνώριση Προτύπων

## Γραμμικές Συναρτήσεις Διάκρισης (Linear Discriminant Functions)

*Χριστόδουλος Χαμζάς*

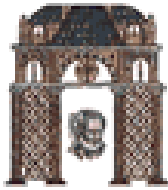
*Τα περιεχόμενα των παρουσιάσεων προέρχονται από τις παρουσιάσεις του αντίστοιχου διδασκτέου μαθήματος του καθ. Παναγιώτη Τσακαλίδη, Τμ. Επιστήμης Υπολογιστών, Παν. Κρήτης και του καθ. Σέργιου Θεοδωρίδη, Τμήμα Πληροφορικής, Πανεπιστήμιο Αθηνών. Βασίζεται στα βιβλία: "Pattern Classification", R.O. Duda, P.E. Hart, D.G. Stork, Wiley, 2<sup>nd</sup> Ed., 2001 και S. Theodoridis, K. Koutroumbas, Pattern Recognition, 3<sup>rd</sup> Edition, Academic Press, 2006*



# Γραμμικές Συναρτήσεις Διάκρισης

## ΠΟΡΕΙΑ ΤΟΥ ΜΑΘΗΜΑΤΟΣ

1. Γνωρίζαμε τις κατανομές και βρίσκαμε τις **συναρτήσεις διάκρισης** (Εκτιμητές Bayesian)
2. Υποθέτουμε ότι γνωρίζουμε την μορφή της κατανομής πιθανότητας και από τα δεδομένα εκπαίδευσης βρίσκουμε τις τιμές των παραμέτρων τους (Parametric Estimation)
3. Από τα δεδομένα εκπαίδευσης βρίσκουμε (εκτιμούμε) τις κατανομές πιθανότητας (Παράθυρα Parzen και KNN)
4. Υποθέτουμε ότι γνωρίζουμε την μορφή των **συναρτήσεων διάκρισης** και θέλουμε να βρούμε τις τιμές των παραμέτρων τους από τα δεδομένα εκπαίδευσης. Δεν χρειαζόμαστε τις υποκείμενες κατανομές (non parametric estimation). **Κίνδυνος** no-robust estimators



# Γραμμικές Συναρτήσεις Διάκρισης

- **Στόχος:** Η σχεδίαση γραμμικών ως προς το διάνυσμα χαρακτηριστικών  $x$  συναρτήσεων διάκρισης που ορίζουν υπερεπίπεδα ως επιφάνειες απόφασης.
- **Γιατί;** Απλή μορφή, εύκολη υλοποίηση, βέλτιστες για Γκαουσιανές σ.π.π.
- **Πώς:** Διατυπώνοντας το πρόβλημα εύρεσης των παραμέτρων (βαρών) ως πρόβλημα βελτιστοποίησης μιας συνάρτησης κριτηρίου (κόστους).
- **Τι είναι η συνάρτηση κριτηρίου;** Μια βαθμωτή συνάρτηση των βαρών που θα πρέπει να ελαχιστοποιηθεί, π.χ. η πιθανότητα λάθος ταξινόμησης κατά την εκπαίδευση.
- **Είναι δύσκολο να επιτευχθεί;** Ναι, γενικώς είναι δύσκολη η σχεδίαση ενός γραμμικού ταξινομητή που να ελαχιστοποιεί το ρίσκο.
- **Επομένως;** Χρησιμοποιούμε εναλλακτικά κριτήρια (απλές συναρτήσεις των βαρών) και επαναληπτικές μεθόδους βελτιστοποίησης (καθόδου κατά την κλίση του κριτηρίου – gradient descent).



# Γιατί γπαμμικοί ταξινομητέρ;

- Πξόβιεκα 2 θαηεγνξηώλ: Αλ ν αξηζκόο Ν ηωλ πξννύπωλ είλαη κηθξόηεξο από ηελ αξηζκό Ν ηωλ ζπληεξεζώλ θάζε πξννύπνπ, ηότε ππάξξεη πάληερα ππεξεπίπεδν πνπ ηα δηαξααίδεη
- Επνκέλω νη γξακκηθελί ηαμηακεηέο είλαη ηξίζηηκη
  - ζε πξνβιήκαηα πνύ κεγάινο δηαζηαμηθήηηαο
  - Σε πξνβιήκαηα κέηξηαο δηαζηαμηθήηηαο, όπνπ ππάξξεη έλαο ζξεηηθάκηθόο αξηζκόο πξννύπωλ εθπαίδεπζεο
- Επηπιένλ ν αξηζκόο ηωλ ζπληεξεζώλ ελόο πξννύπνπ κπνρεί λα απκεζει απζαίξεηα κε ηελ πξνζήθε λέωλ ζπληεξεζώλ πνπ είλαη κε γξακκηθέο ζπλαξηήζεηο ηωλ αξηηώλ ζπληεξεζώλ (π.ρ. πνιπώλκκα)

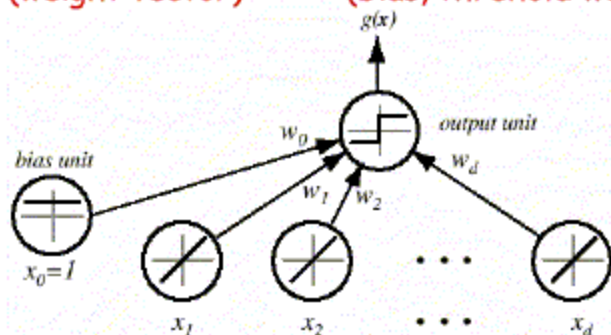


# Συναρτήσεις και Επιφάνειες Διάκρισης

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

Διάνυσμα βαρών  
(weight vector)

Βάρος κατωφλίου  
(bias, threshold weight)

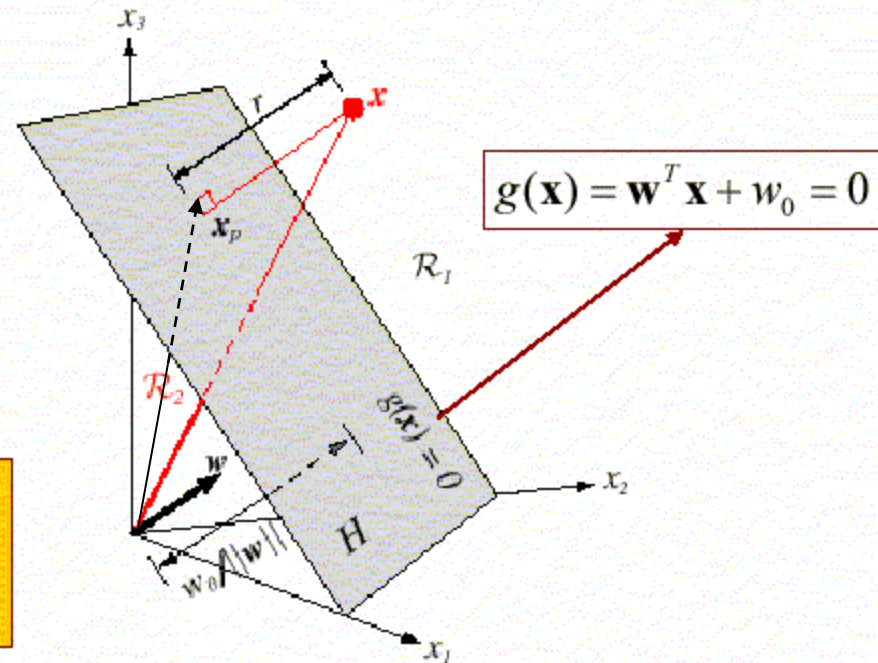


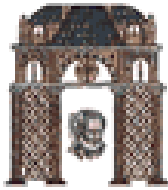
Το διάνυσμα βαρών,  $\mathbf{w}$ , καθορίζει τον προσανατολισμό του υπερεπιπέδου απόφασης και το βάρος κατωφλίου,  $w_0$ , καθορίζει τη σχετική θέση του ως προς την αρχή των αξόνων.

Έστω  $x_p$  η προβολή του  $\mathbf{x}$  στο υπερεπίπεδο. Επειδή  $x_p$  είναι πάνω στο υπερεπίπεδο.

Άρα  $g(x_p) = 0 = \mathbf{w}^T x_p + w_0$  και

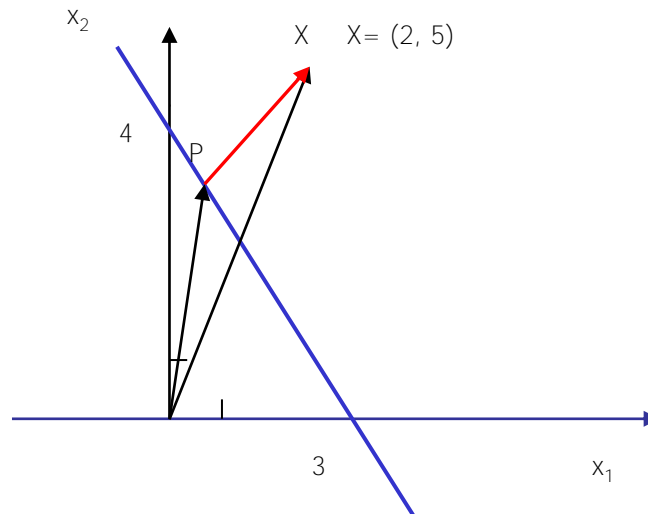
$$\mathbf{x} = \mathbf{x}_p + r \frac{\mathbf{w}}{\|\mathbf{w}\|}, \text{ όπου } r = \frac{g(\mathbf{x})}{\|\mathbf{w}\|}$$





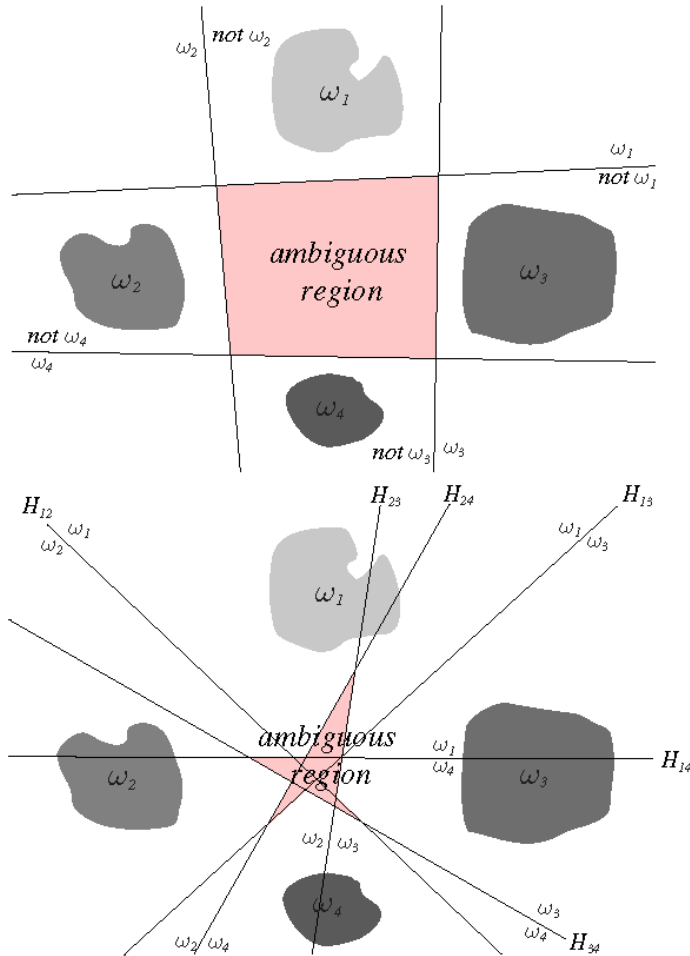
# Παράδειγμα

- $g(x_1, x_2) = x_1 + 0.75x_2 - 3 = [1 \ 0.75][x_1 \ x_2]^T + (-3)$
- $w = (1, .75), w_0 = -3, ||w|| = 5/4$
- $X = (2, 5)$ . Η απόσταση  $r$  του  $X$  από την ευθεία  $g(X) = 0$  είναι
- $r = (PX) = g(2, 5) / ||w|| = (11/4) / (5/4) = 11/5$





# Η Περίπτωση Πολλών Κλάσεων



Περίπτωση διαχωρισμού σε  $\omega_1$  / όχι  $\omega_1$

Μπορεί να υπάρχει περιοχή που δεν καθορίζεται

Περίπτωση διαχωρισμού σε  $\omega_i$  /  $\omega_j$

Μπορεί πάλι να υπάρχει περιοχή που δεν καθορίζεται

Θα ακολουθήσουμε την φιλοσοφία των επιφανειών διαχωρισμού  $g_i(x) = w^T x_i + w_{i0}$  και  $x \in \omega_i$  εάν  $g_i(x) > g_j(x)$  για κάθε  $j \neq i$



# Η Περίπτωση Πολλών Κλάσεων

Γραμμική Μηχανή (linear machine):

$\mathbf{x}$  in  $\omega_i$  αν  $g_i(\mathbf{x}) > g_j(\mathbf{x})$

Σύνορα Απόφασης:

$H_{ij}: g_i(\mathbf{x}) = g_j(\mathbf{x}) \rightarrow (\mathbf{w}_i - \mathbf{w}_j)^T \mathbf{x} + (w_{i0} - w_{j0}) = 0$

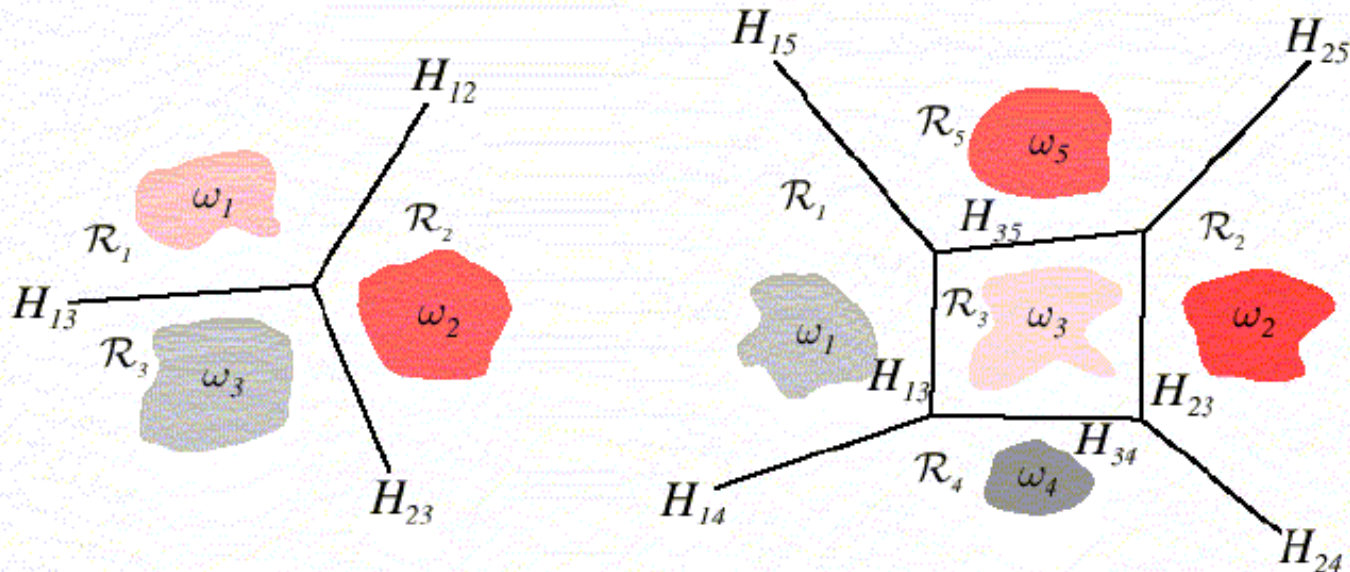
- τμήμα υπερεπιπέδου κάθετο στο διάνυσμα  $\mathbf{w}_i - \mathbf{w}_j$

Απόσταση του  $\mathbf{x}$  από το  $H_{ij}$ :

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

- σημαντικές είναι οι διαφορές των διανυσμάτων βαρών.

Κυρτές περιοχές απόφασης.







# Γραμμικά Διαχωρίσιμες Κλάσεις: Διανύσματα και Περιοχές Λύσης

Επαυξημένα διανύσματα  
χαρακτηριστικών και βαρών  
(Augmented Vectors):

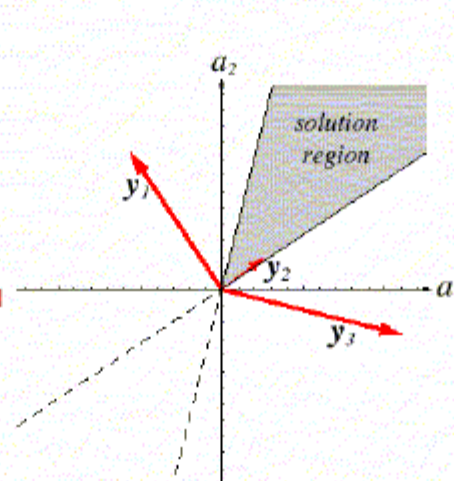
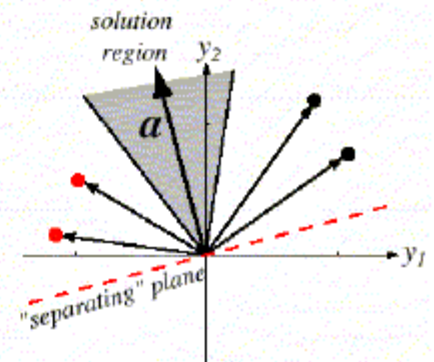
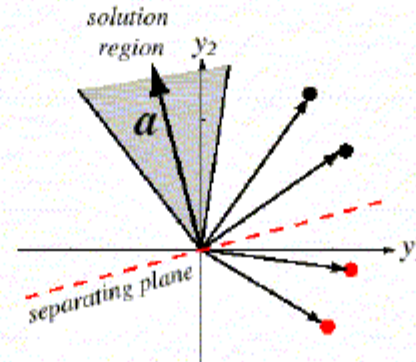
$$\mathbf{a} = \begin{bmatrix} w_0 \\ w_1 \\ w_2 \\ \vdots \\ w_d \end{bmatrix} = \begin{bmatrix} w_0 \\ \mathbf{w} \end{bmatrix}, \quad \mathbf{y} = \begin{bmatrix} 1 \\ x_1 \\ x_2 \\ \vdots \\ x_d \end{bmatrix} = \begin{bmatrix} x_0 \\ \mathbf{x} \end{bmatrix}$$

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

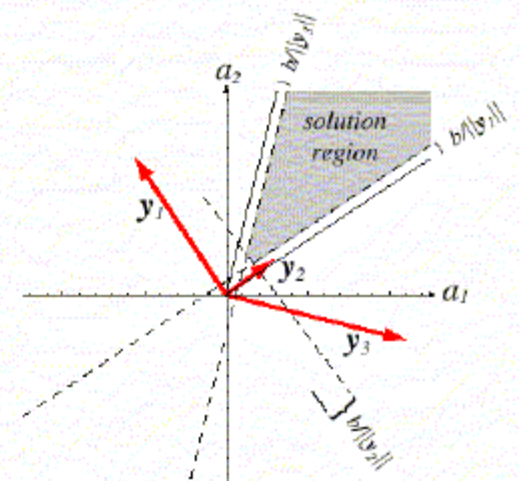
$$\text{απόσταση } \mathbf{y} \text{ από } H: \frac{|\mathbf{a}^T \mathbf{y}|}{\|\mathbf{a}\|}$$

Κανονικοποίηση: αντικαθιστώντας όλα τα δείγματα εκπαίδευσης της κλάσης  $\omega_2$  με τα αρνητικά τους, ζητούμε τα διαχωριστικά διανύσματα (separating vectors) που πρέπει να ικανοποιούν την σχέση:

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



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

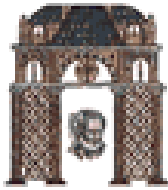


$$\mathbf{a}^T \mathbf{y}_i \geq b > 0 \quad \forall i$$



# Διαδικασίες Βελτιστοποίησης

- Πρόβλημα: Εύρεση του  $\mathbf{a}$  που ικανοποιεί το σύνολο των γραμμικών ανισοτήτων  $\mathbf{a}'\mathbf{y}_i > 0$  για κάθε  $i=1, \dots, n$ .
- Πώς βρίσκουμε την κατάλληλη λύση;
  - ↳ Ορίζουμε μια συνάρτηση κριτηρίου,  $J(\mathbf{a})$ , και την ελαχιστοποιούμε ώστε το  $\mathbf{a}$  να είναι ένα διάνυσμα λύσης.
  - ↳ Με αυτό τον τρόπο μετασχηματίζουμε το πρόβλημα της εξαντλητικής αναζήτησης σε πρόβλημα ελαχιστοποίησης μιας βαθμωτής συνάρτησης.
- Πώς ελαχιστοποιούμε την  $J(\mathbf{a})$ ;
  - ↳ Επιλέγουμε κάποιο αρχικό σημείο  $\mathbf{a}_1$  και υπολογίζουμε την τιμή  $J(\mathbf{a}_1)$ .
  - ↳ Υπολογίζουμε την κλίση στο  $J(\mathbf{a}_1)$ :  $\nabla J(\mathbf{a}_1)$ .
  - ↳ Παιρνουμε το επόμενο σημείο  $\mathbf{a}_2$  κινούμενοι στην κατεύθυνση αρνητικής κλίσης (steepest descent),  $-\nabla J(\mathbf{a}_1)$ , κατά μια ποσότητα  $\eta(k)$ , τον λεγόμενο ρυθμό μάθησης (learning rate) ή το βήμα (stepsize).



# Αλγόριθμος της πιο Απότομης Καθόδου (Steepest Descent)

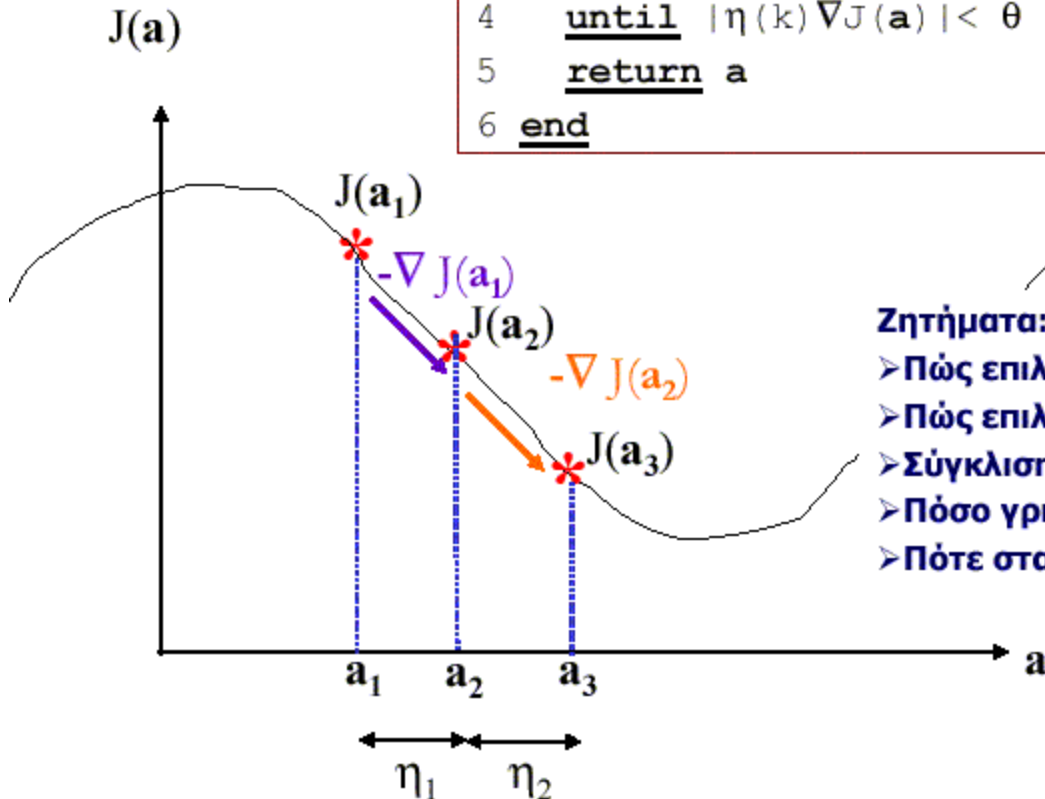
## Αλγόριθμος 1. Πιο Απότομη Κάθοδος (Steepest Descent)

```

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

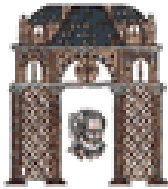
```

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



### Ζητήματα:

- Πώς επιλέγουμε την συνάρτηση κριτηρίου;
- Πώς επιλέγουμε τον ρυθμό μάθησης  $\eta(k)$ ;
- Σύγκλιση σε τοπικό/ολικό ελάχιστο;
- Πόσο γρήγορα συγκλίνουμε, πόσο ομαλά;
- Πότε σταματάμε;



# Αλγόριθμος Καθόδου Newton (Newton Descent)

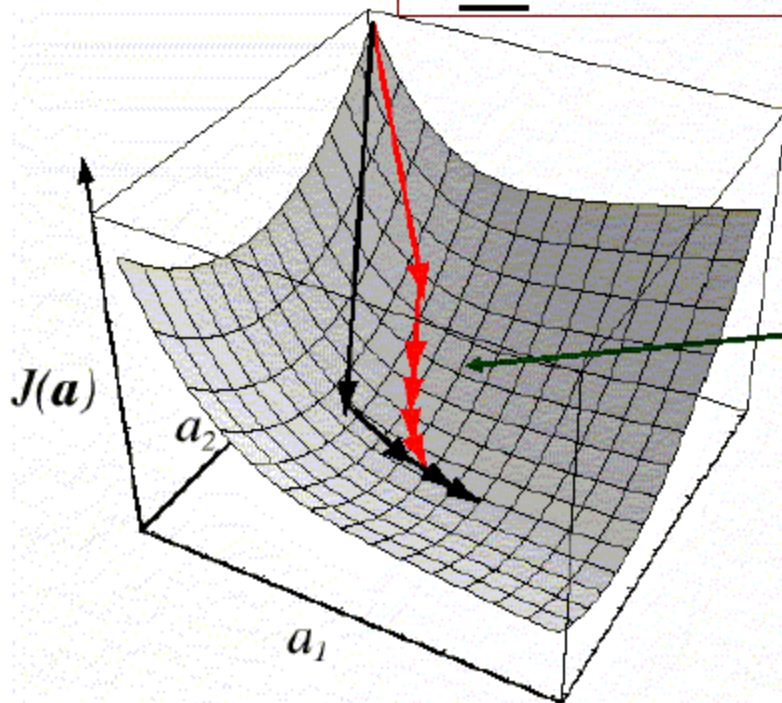
## Αλγόριθμος 2. Κάθοδος Newton (Newton Descent)

```

1 begin initialize 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 a
5 end

```

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



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

**Κόκκινο: Steepest Descent**

**Μαύρο: Newton Descent**

**Newton: μεγαλύτερη βελτίωση σε κάθε βήμα πληρώνοντας το υπολογιστικό κόστος της αντιστροφής του Hessian πίνακα  $H$ .**



# Το κριτήριο Perceptron

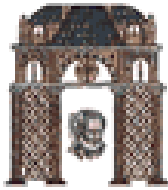
## ➤ Ποια μπορεί να είναι η συνάρτηση κριτηρίου;

- ↳ Μια πρώτη επιλογή: Πλήθος των λάθος ταξινομημένων δειγμάτων εκπαίδευσης.  
Αλλά: Αυτή η συνάρτηση είναι ασυνεχής οπότε δεν είναι διαφορίσιμη.
- ↳ Μια καλύτερη επιλογή: Η συνάρτηση κριτηρίου perceptron:

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

όπου  $\mathcal{Y}(\mathbf{a})$  είναι το σύνολο των δειγμάτων που δεν έχουν ταξινομηθεί σωστά από το  $\mathbf{a}$ .

- ↳ Αν το  $\mathcal{Y}(\mathbf{a})$  είναι κενό, τότε  $J_p(\mathbf{a})=0$ . Αφού  $\mathbf{a}'\mathbf{y}<0$  όταν το  $\mathbf{y}$  δεν είναι σωστά ταξινομημένο, η  $J_p(\mathbf{a})$  δεν είναι ποτέ αρνητική και μηδενίζεται όταν το  $\mathbf{a}$  είναι διάνυσμα λύσης.
- ↳ Γεωμετρικά, η  $J_p(\mathbf{a})$  είναι ανάλογη του αθροίσματος των αποστάσεων των λάθος ταξινομημένων δειγμάτων από το σύνορο απόφασης.



# Ο Αλγόριθμος Batch Perceptron

↪ Το διάνυσμα κλίσεων είναι:  $\nabla J_p(\mathbf{a}) = \left[ \frac{\partial J_p}{\partial a_i} \right] = \sum_{y \in Y} -\mathbf{y}$

↪ Αναδρομική σχέση:  $\mathbf{a}(k+1) = \mathbf{a}(k) + \eta(k) \sum_{y \in Y_k} \mathbf{y}$

όπου  $Y_k$  είναι το σύνολο των δειγμάτων που έχουν ταξινομηθεί λάθος από το  $\mathbf{a}_k$ .

↪ Το επόμενο διάνυσμα βάρους (δ.β.) προκύπτει ως το άθροισμα του τρέχοντος δ.β. και ενός πολλαπλασίου του αθροίσματος των λάθος ταξινομημένων δειγμάτων.

## Αλγόριθμος 3. Batch Perceptron

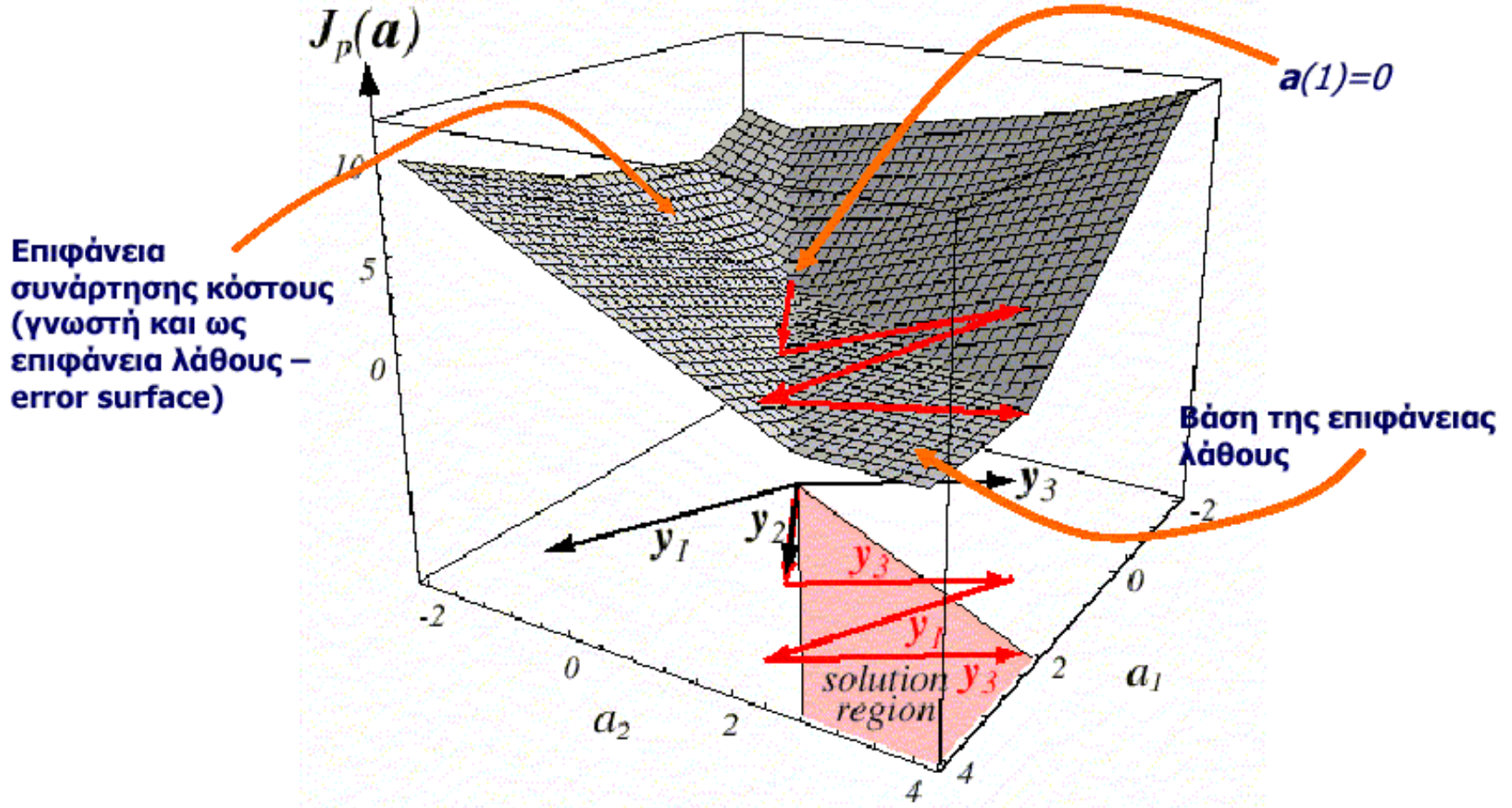
```

1 begin initialize  $\mathbf{a}$ , κριτήριο  $\theta$ ,  $\eta(0) > 0$ ,  $k=0$ 
2   do  $k \leftarrow k+1$ 
3      $\mathbf{a} \leftarrow \mathbf{a} + \eta(k) \sum_{y \in Y_k} \mathbf{y}$ 
4   until  $|\eta(k) \sum_{y \in Y_k} \mathbf{y}| < \theta$ 
5   return  $\mathbf{a}$ 
6 end

```



# Βήματα του Batch Perceptron

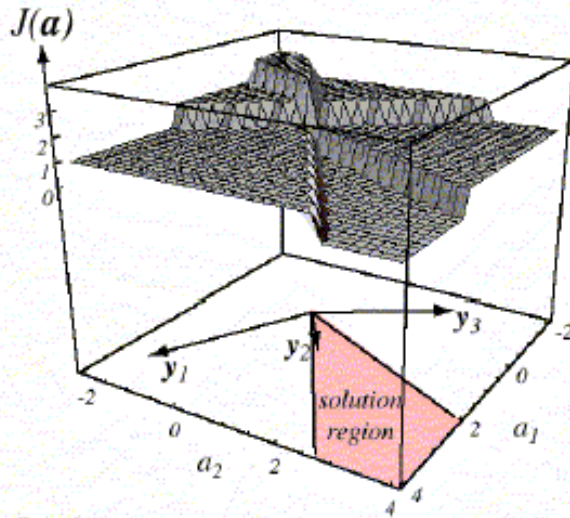




# Τέσσερις Συναρτήσεις Κόστους

Πλήθος  
εσφαλμένων  
ταξινομήσεων

Bad



Κριτήριο  
Perceptron

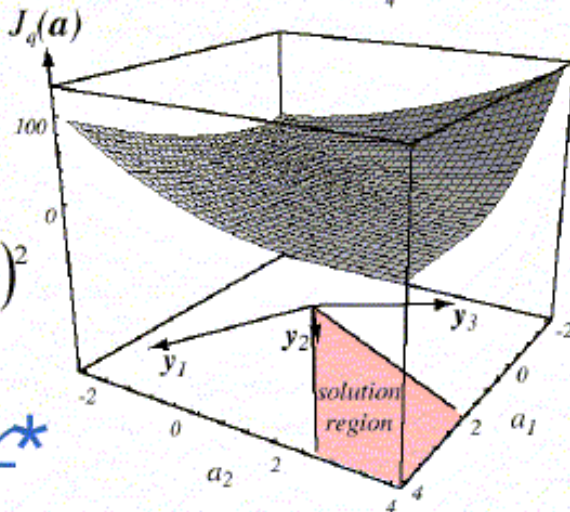
$$J_p(\mathbf{a}) = \sum_{\mathbf{y} \in Y} (-\mathbf{a}'\mathbf{y})$$

Good!

Συνολικό  
τετραγωνικό  
λάθος -  
Total square  
error (TSE)

$$J_q(\mathbf{a}) = \sum_{\mathbf{y} \in Y} (\mathbf{a}'\mathbf{y})^2$$

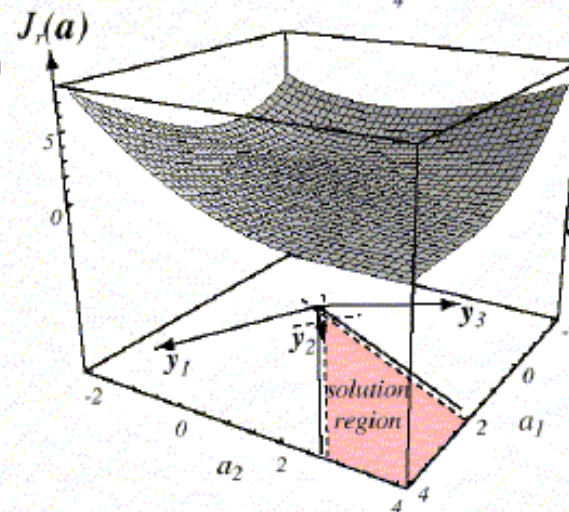
Better\*



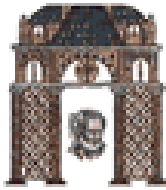
TSE with  
margin

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

Best\*







# Fixed Increment Single-Sample Perceptron

- Αντί να δοκιμάζουμε το διάνυσμα βαρών  $\mathbf{a}(k)$  σε όλα τα δείγματα και να το διορθώνουμε βάσει του συνόλου  $Y_k$  των λάθος ταξινομημένων δειγμάτων, χρησιμοποιούμε τα δείγματα ένα κάθε φορά και ανάλογα με την ταξινόμηση του ανανεώνουμε ή όχι το διάνυσμα βαρών.
- Αν επιπλέον, χρησιμοποιήσουμε ένα σταθερό βήμα  $\eta(k)$ , τότε προκύπτει ο αλγόριθμος:

## Αλγόριθμος 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

```

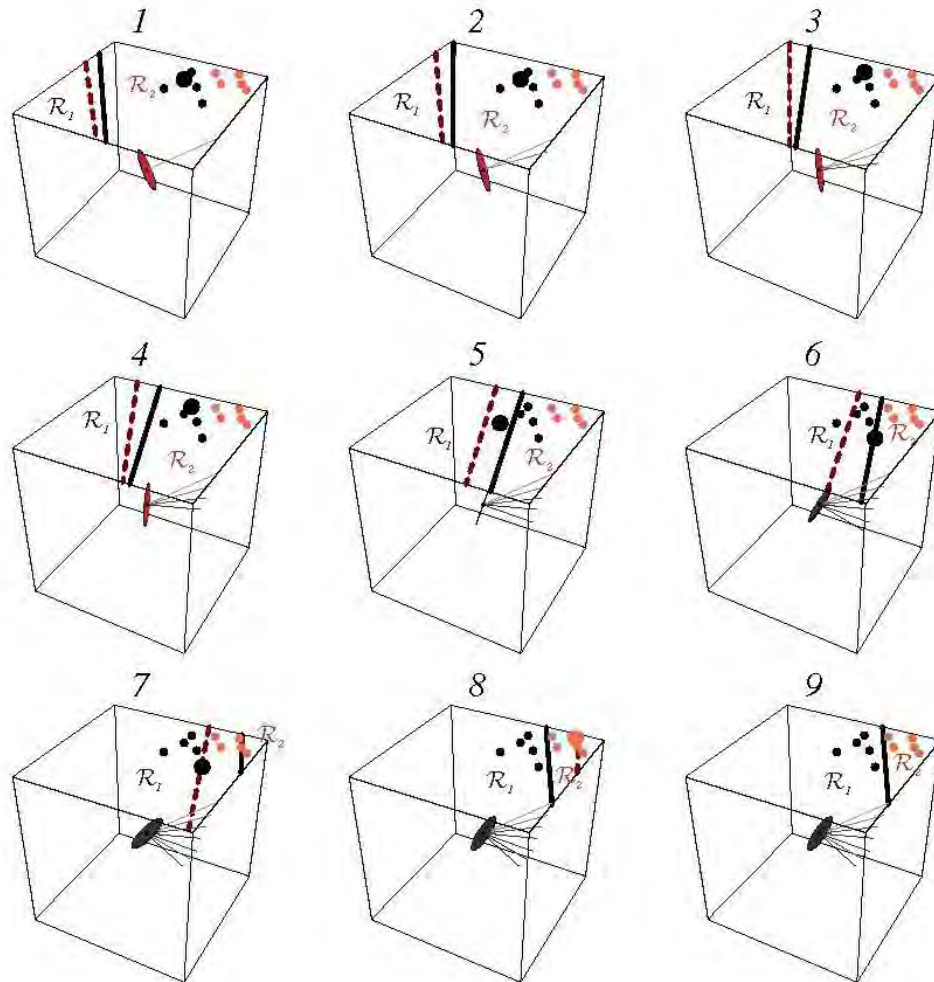
Κυκλική Σειρά Δεδομένων (με πράσινο υποδηλώνονται τα λάθος ταξινομημένα δείγματα):

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

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



# Σύγκλιση

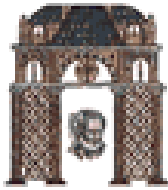


**Πεπερασμένος αριθμός βημάτων,  $k_0$ , για σύγκλιση στη λύση  $a_0$ :**

$$k_0 = \frac{\max_i \|\mathbf{y}_i\|^2 \|\mathbf{a}_0\|^2}{\min_i [\mathbf{y}_i^t \mathbf{a}_0]^2}$$

**Ο παρονομαστής δηλώνει ότι η δυσκολία στη σύγκλιση καθορίζεται από τα δείγματα εκπαίδευσης  $\mathbf{y}_i$ , τα οποία είναι σχεδόν κάθετα στο διάνυσμα λύσης  $\mathbf{a}_0$ : Γραμμικά διαχωρίσιμα προβλήματα είναι δύσκολα επιλύσιμα όταν τα πρότυπα είναι σχεδόν ομοεπίπεδα.**

**ΔΥΣΤΥΧΩΣ ΤΟ ΟΡΙΟ ΕΚΦΡΑΖΕΤΑΙ ΣΥΝΑΡΤΗΣΕΙ ΤΗΣ ΛΥΣΗΣ  $\mathbf{a}_0$**



# Variable-Increment Perceptron with Margin

- Χρησιμοποιούμε τα δείγματα ένα κάθε φορά και διορθώνουμε το διάνυσμα βαρών  $\mathbf{a}(k)$  όταν το εσωτερικό γινόμενο του με το δείγμα  $\mathbf{y}^k$  είναι μικρότερο από κάποιο προκαθορισμένο θετικό όριο  $b$ :  $\mathbf{a}(k)\mathbf{y}^k < b$ .
- Αν επιπλέον, χρησιμοποιήσουμε ένα μεταβαλλόμενο βήμα  $\eta(k)$ , τότε προκύπτει ο αλγόριθμος:

## Αλγόριθμος 5. Variable-Increment Perceptron w. Margin

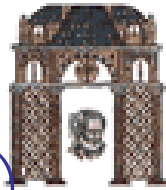
```

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

```

Συνθήκες Σύγκλισης:

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



# Μέθοδοι Χαλάρωσης (Relaxation Procedures)

➤ Συνάρτηση κριτηρίου:

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

όπου  $Y(\mathbf{a})$  είναι το σύνολο των δειγμάτων για τα οποία  $\mathbf{a}'\mathbf{y} < b$ .

↪ Αν το  $Y(\mathbf{a})$  είναι κενό, τότε  $J_r(\mathbf{a})=0$ . Η  $J_r(\mathbf{a})$  δεν είναι ποτέ αρνητική και μηδενίζεται αν και μόνο αν  $\mathbf{a}'\mathbf{y} > b$  για όλα τα δείγματα εκπαίδευσης.

↪ Το διάνυσμα κλίσεων είναι:

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

↪ Αναδρομική σχέση:

$$\mathbf{a}(k+1) = \mathbf{a}(k) + \eta(k) \sum_{\mathbf{y} \in Y_k} \frac{b - \mathbf{a}'\mathbf{y}}{\|\mathbf{y}\|^2} \mathbf{y}$$



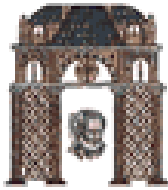
# Batch Relaxation with Margin

## Αλγόριθμος 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  $\mathbf{a}^t \mathbf{y}^j < b$ , 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

## Αλγόριθμος 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

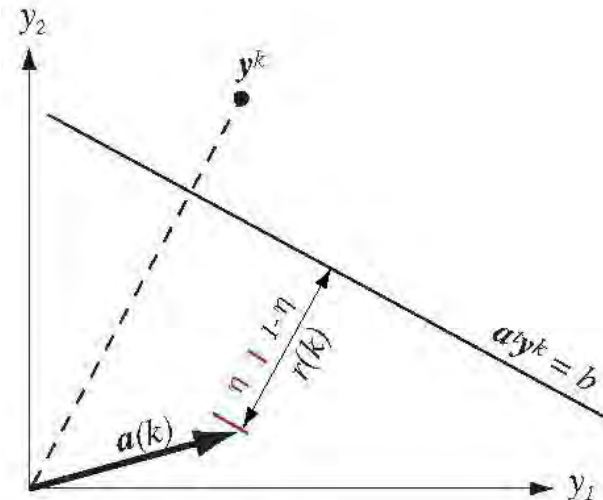
```

Σε κάθε βήμα, το διάνυσμα βαρών  $\mathbf{a}(k)$ , μετατοπίζεται προς το υπερεπίπεδο  $\mathbf{a}^t \mathbf{y}^k = b$  κατά ένα ποσοστό,  $\eta(k)$ , της απόστασής του,  $r(k)$ , από αυτό.

$\eta(k) < 1 \rightarrow$  underrelaxation

$\eta(k) > 1 \rightarrow$  overrelaxation

$0 < \eta(k) < 2$  για σύγκλιση





# Μη Διαχωρίσιμες Κατηγορίες

- Οι προηγούμενοι αλγόριθμοι βασίζονται στην υπόθεση ότι οι κατηγορίες είναι γραμμικά διαχωρίσιμες.
- Αλλά ακόμα και εάν το δείγμα εκπαίδευσης είναι γραμμικά διαχωρίσιμο αυτό δεν εγγυάται καλή συμπεριφορά στην πραγματικότητα

ΠΩΣ ΘΑ ΣΥΜΠΕΡΙΦΕΡΟΝΤΑΙ ΣΕ ΜΗ ΓΡΑΜΜΙΚΑ ΔΙΑΧΩΡΙΣΙΜΕΣ ΚΑΤΗΓΟΡΙΕΣ

Ας βάλλουμε όλα τα δεδομένα με ένα margin  $a^T y_i - b_i > 0$

και ας προσπαθήσουμε να ελαχιστοποιήσουμε μια συνάρτηση κριτηρίου που βασίζεται στα ελάχιστα τετράγωνα ->

Μέθοδο Ελαχίστων Τετραγώνων



# Μέθοδοι Ελαχίστων Τετραγώνων

## (Minimum Square Error – MSE)

- **Στόχος:** Καλή απόδοση τόσο στις γραμμικά διαχωρίσιμες όσο και στις μη γραμμικά διαχωρίσιμες περιπτώσεις.
- **Πώς:** Κριτήριο που περιλαμβάνει όλα τα πρότυπα. Επίσης,
  - ↳ Πριν: Βρες  $\mathbf{a}$  έτσι ώστε  $\mathbf{a}'\mathbf{y}_i > 0$  για κάθε πρότυπο  $\mathbf{y}_i$ .
  - ↳ Τώρα: Βρες  $\mathbf{a}$  έτσι ώστε  $\mathbf{a}'\mathbf{y}_i = b_i$  για κάθε πρότυπο  $\mathbf{y}_i$ . ( $b_i$  θετικές σταθερές).
- **Δηλαδή;** Μετατρέπουμε το πρόβλημα επίλυσης ενός συνόλου γραμμικών ανισοτήτων σε πρόβλημα επίλυσης γραμμικών εξισώσεων.
- **Συμβολισμοί:**

$\mathbf{Y}_{n \times \hat{d}} = \begin{bmatrix} \mathbf{y}_1^t \\ \mathbf{y}_2^t \\ \vdots \\ \mathbf{y}_n^t \end{bmatrix}$ <p>Πίνακας Προτύπων</p>	$\mathbf{b}_{n \times 1} = \begin{bmatrix} b_1 \\ b_2 \\ \vdots \\ b_n \end{bmatrix}$ <p>Διάνυσμα θετικών παραμέτρων</p>	$\mathbf{a}_{\hat{d} \times 1} = \begin{bmatrix} a_0 \\ a_1 \\ \vdots \\ a_d \end{bmatrix}$ <p>Διάνυσμα βαρών</p>
--	--	---

- **Πρόβλημα:** Εύρεση  $\mathbf{a}$  έτσι ώστε:  $\mathbf{Y}\mathbf{a}=\mathbf{b}$ .





# Μέθοδοι Ελαχίστων Τετραγώνων (Minimum Square Error – MSE)

- Συνήθως  $n > d+1$ , περισσότερα πρότυπα από διαστάσεις  $\rightarrow$  σύστημα υπερπροσδιορισμένο, δεν έχει ακριβή λύση.
- **Επομένως;** Ελαχιστοποίηση του τετραγώνου του μήκους του διανύσματος σφάλματος,  $\mathbf{e} = \mathbf{Y}\mathbf{a} - \mathbf{b}$ :

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

- **Κλασσικό Πρόβλημα:**

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

**Κανονικές Εξισώσεις**  
**Normal Equations**

- Αν ο  $\mathbf{Y}'\mathbf{Y}$  είναι ομαλός,

$$\mathbf{a} = \underbrace{(\mathbf{Y}'\mathbf{Y})^{-1}}_{\text{ψευδοαντίστροφος}} \mathbf{Y}' \mathbf{b}$$



# Widrow-Hoff (Least Mean Squares - LMS)

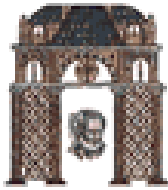
- Η  $J_s(\mathbf{a})$  μπορεί να ελαχιστοποιηθεί μέσω αναδρομικών αλγορίθμων που δεν απαιτούν την αντιστροφή πινάκων.
- Διάνυσμα κλίσεων:  $\nabla J_s(\mathbf{a}) = 2\mathbf{Y}'(\mathbf{Y}\mathbf{a} - \mathbf{b})$
- Βασική αναδρομική σχέση:  $\mathbf{a}(k+1) = \mathbf{a}(k) + \eta(k)\mathbf{Y}'(\mathbf{b} - \mathbf{Y}\mathbf{a}(k))$
- Χρησιμοποιώντας ένα δείγμα σε κάθε βήμα, προκύπτει ο αλγόριθμος Widrow-Hoff (LMS):

## 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      $\mathbf{a} \leftarrow \mathbf{a} + \eta(k)(b_k - \mathbf{a}'\mathbf{y}^k)\mathbf{y}^k$ 
4   until  $|\eta(k)(b_k - \mathbf{a}'\mathbf{y}^k)\mathbf{y}^k| < \theta$ 
5   return  $\mathbf{a}$ 
6 end

```



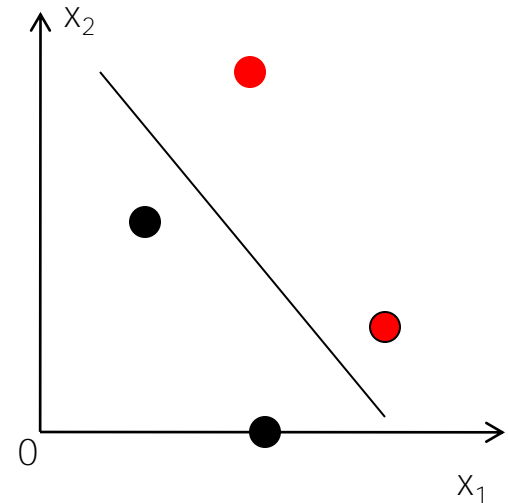
# Παράδειγμα

- Υποθέσουμε ότι έχουμε τα παρακάτω σημεία 2 διαστάσεων που ανήκουν σε 2 κατηγορίες:
- $\omega_1$ : (1, 2) και (2, 0), και  $\omega_2$ : (3, 1) και (2, 3), τα πρώτα είναι με μαύρο και τα δεύτερα με κόκκινο χρώμα στο επόμενο σχήμα.
- Ο πίνακας  $\mathbf{Y}$  είναι επομένως
- $\mathbf{Y} = (1 \ 1 \ 2, 1 \ 2 \ 0, -1 \ -3 \ -1, -1 \ -2 \ -3)$
- και μετά από μερικές απλές πράξεις έχουμε ότι ο ψευδοαντίστροφος πίνακας είναι:

$$\mathbf{Y} = \begin{pmatrix} 1 & 1 & 2 \\ 1 & 2 & 0 \\ -1 & -3 & -1 \\ -1 & -2 & -3 \end{pmatrix} \quad \mathbf{a} = (\mathbf{Y}^T \mathbf{Y})^{-1} \mathbf{Y}^T \mathbf{b} = \begin{pmatrix} 5/4 & 13/12 & 3/4 & 7/12 \\ -1/2 & -1/6 & -1/2 & -1/6 \\ 0 & -1/3 & 0 & -1/3 \end{pmatrix} \mathbf{b}$$

Και διαλέγοντας αυθαίρετα  $\mathbf{b}^T = (1, 1, 1, 1)$ , έχουμε  $\mathbf{a}^T = (11/3, -4/3, -2/3)$   
 άρα

$$g(\mathbf{x}) = (11/3) - (4/3)x_1 - (2/3)x_2$$

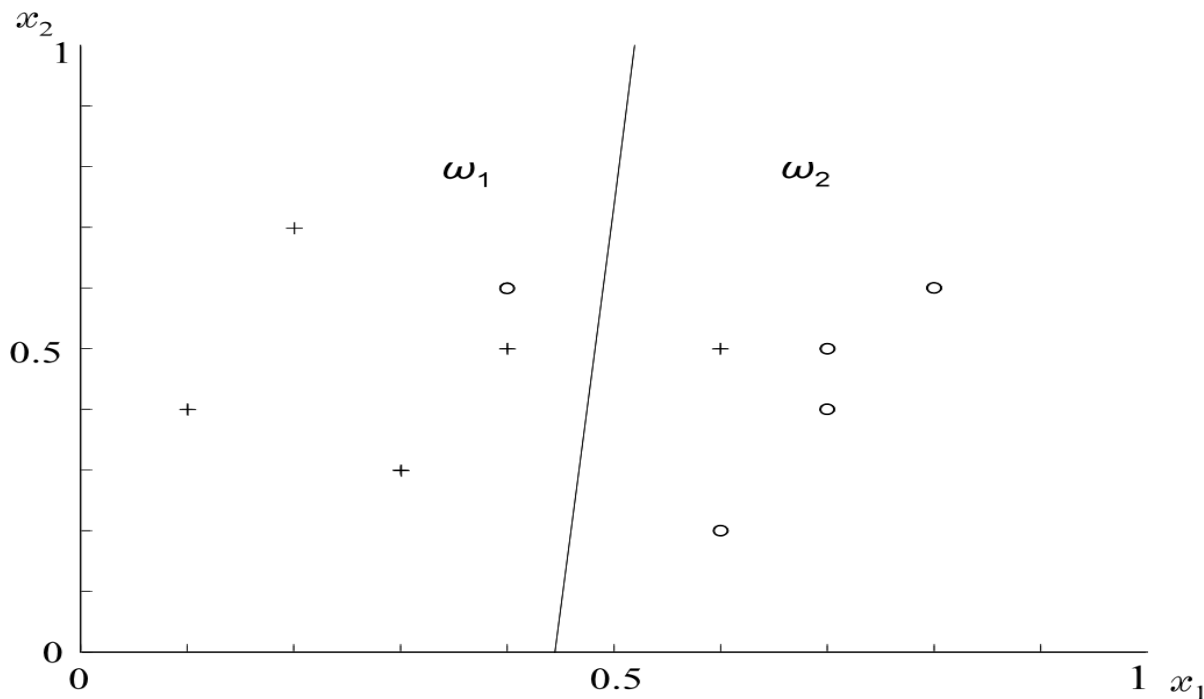




# Παράδειγμα 2

$$\omega_1 : \begin{bmatrix} 0.4 \\ 0.5 \end{bmatrix}, \begin{bmatrix} 0.6 \\ 0.5 \end{bmatrix}, \begin{bmatrix} 0.1 \\ 0.4 \end{bmatrix}, \begin{bmatrix} 0.2 \\ 0.7 \end{bmatrix}, \begin{bmatrix} 0.3 \\ 0.3 \end{bmatrix}$$

$$\omega_2 : \begin{bmatrix} 0.4 \\ 0.6 \end{bmatrix}, \begin{bmatrix} 0.6 \\ 0.2 \end{bmatrix}, \begin{bmatrix} 0.7 \\ 0.4 \end{bmatrix}, \begin{bmatrix} 0.8 \\ 0.6 \end{bmatrix}, \begin{bmatrix} 0.7 \\ 0.5 \end{bmatrix}$$



$$Y = \begin{bmatrix} 1 & 0.4 & 0.5 \\ 1 & 0.6 & 0.5 \\ 1 & 0.1 & 0.4 \\ 1 & 0.2 & 0.7 \\ 1 & 0.3 & 0.3 \\ -1 & -0.4 & -0.6 \\ -1 & -0.6 & -0.2 \\ -1 & -0.7 & -0.4 \\ -1 & -0.8 & -0.6 \\ -1 & -0.7 & -0.5 \end{bmatrix}^T$$



$$\underline{b} = [1 \ 1 \ 1 \ 1 \ 1 \ 1 \ 1 \ 1 \ 1 \ 1]$$

$$Y^T Y = \begin{bmatrix} 10 & 4.8 & 4.7 \\ 4.8 & 2.8 & 2.24 \\ 4.7 & 2.24 & 2.41 \end{bmatrix}, Y^T \underline{b} = \begin{bmatrix} 0.0 \\ -1.6 \\ 0.1 \end{bmatrix}$$

$$\underline{a} = (Y^T Y)^{-1} Y^T \underline{b} = \begin{bmatrix} 1.43 \\ -3.21 \\ 0.24 \end{bmatrix}$$

$$g(\mathbf{x}) = ((Y^T Y)^{-1} Y^T \underline{b}) \begin{bmatrix} 1 \\ x_1 \\ x_2 \end{bmatrix} = [1.43 \quad -3.21 \quad 0.24] \begin{bmatrix} 1 \\ x_1 \\ x_2 \end{bmatrix} = 1.434 - 3.21x_1 + 0.24x_2$$

**Σημείωση:** Στο MSE μπορούμε αντί να αλλάξουμε το πρόσημο στα  $Y$  για την κλάση  $\omega_2$  να αλλάξουμε το πρόσημο στο  $\underline{b}$  για τα δεδομένα που ανήκουν στη κλάση  $\omega_2$ , δηλαδή

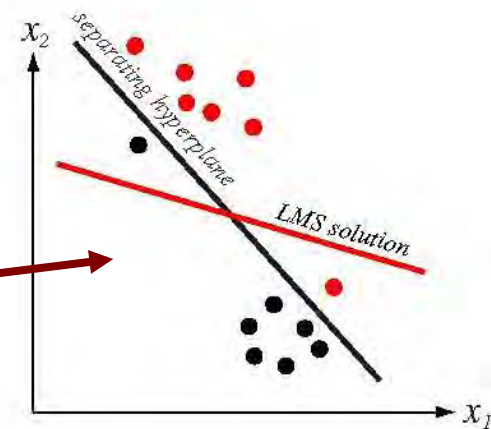
$$Y = \begin{bmatrix} 1 & 0.4 & 0.5 \\ 1 & 0.6 & 0.5 \\ 1 & 0.1 & 0.4 \\ 1 & 0.2 & 0.7 \\ 1 & 0.3 & 0.3 \\ 1 & 0.4 & 0.6 \\ 1 & 0.6 & 0.2 \\ 1 & 0.7 & 0.4 \\ 1 & 0.8 & 0.6 \\ 1 & 0.7 & 0.5 \end{bmatrix}^T \quad \underline{b}^T = \begin{bmatrix} 1 \\ 1 \\ 1 \\ 1 \\ 1 \\ -1 \\ -1 \\ -1 \\ -1 \\ -1 \end{bmatrix}$$



# Μέθοδος Ho-Kashyap

- Το perceptron και οι τεχνικές χαλάρωσης βρίσκουν διαχωριστικά διανύσματα βαρών αν τα δείγματα είναι γραμμικά διαχωρίσιμα, αλλά δεν συγκλίνουν για μη διαχωρίσιμες κλάσεις.

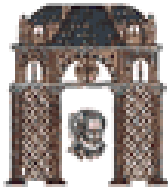
- Οι τεχνικές ελαχίστων τετραγώνων δίνουν πάντα ένα διάνυσμα λύσης (αυτό που ελαχιστοποιεί το  $\|\mathbf{Y}\mathbf{a}-\mathbf{b}\|^2$ ) το οποίο όμως δεν είναι απαραίτητα διαχωριστικό στην διαχωρίσιμη περίπτωση.



- Ο αλγόριθμος Ho-Kashyap λύνει αναδρομικά το εξής πρόβλημα ελαχιστοποίησης:

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

- Είναι ένας αλγόριθμος που φροντίζει ώστε το  $\mathbf{b}$  να μην συγκλίνει στο  $\mathbf{0}$ , θέτοντας όλες τις θετικές συνιστώσες του διανύσματος κλίσης  $\nabla_{\mathbf{b}} J_s$  ίσες με το μηδέν.



# Αλγόριθμος Ho-Kashyap

## Αλγόριθμος 9. Ho-Kashyap

```

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

```

**Βασικά βρίσκουμε πρώτα το gradient descent ως προς b και μετά εφαρμόζουμε ελάχιστα τετράγωνα**



- Έως ηώρα είδακε ην πξόβι εκ α κόλν γηα2 θαηεγνξίεο. Τη κπνύκε λα θάλνπκε εάλ έρνπκε πνιύο θαηεγνξίεο πξνηύπωλ?
- Θα δύκε ην πξόβιεκα π πνζέ ην ληαοόηηηο θιάζεηο είλαη γξακκηθά απωξίζηηεο. Βαζηδύκε ληηηηο γξακκηθέο ζπλαξηήζεηο δηάθεηεο, ζέινπκε λα βξύκε έλα βξύκε έλα ζύλνι ν  $g_i(x)$  γξακκηθώλ ζπλαξηήζεωλ δηάθεηεο έηη ώζηε εάλ  $x \in \omega_i$ , ήθηε  $g_i(x) > g_j(x)$  γηα θάζε  $j \neq i$ .
- Απνηπώλνπκε ηώρα ην πξόβιεκα ζ ηλ ρώξν ηωλ  $y$ , (νη ζπλαξηήζεηο δηάθεηεο είλαη πεξίπεδα πνπ πεξινύλ από ηελ αξρή ηωλ αμώλ).





## Ταξινόμηση Πολλαπλών Κλάσεων με MSE

➤ Έστω  $\mathbf{y}_1, \dots, \mathbf{y}_n$ , πρότυπα από  $c > 2$  κλάσεις. Θέλουμε να βρούμε ένα ταξινομητή που να αποτελείται από γραμμικές συναρτήσεις διάκρισης  $g_i(\mathbf{x}) = w_i \mathbf{x} + w_{i0}$  (μία για κάθε κλάση), ούτως ώστε εάν  $\mathbf{x}$  ανήκει στην  $i$  κλάση τότε  $g_i(\mathbf{x}) > g_j(\mathbf{x}) \quad \forall j \neq i$ . Συνεπώς στον επαυξημένο χώρο ψάχνουμε για τα βάρη  $\mathbf{a}_i$  ώστε εάν  $\mathbf{y}_k$  ανήκει στην  $i$  κλάση τότε

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

➤ Για την επίλυση του προβλήματος συνεπώς υπολογίζουμε το βέλτιστο  $\mathbf{a}_i$  ώστε να διαχωρίζει την κλάση  $i$  από όλες τις υπόλοιπες κλάσεις  $j \neq i$ .

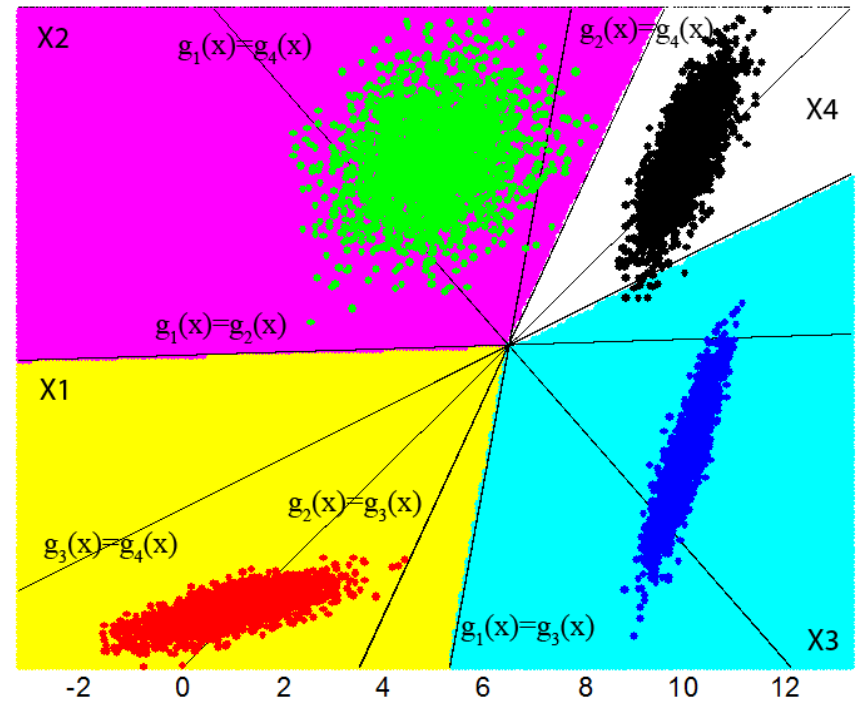
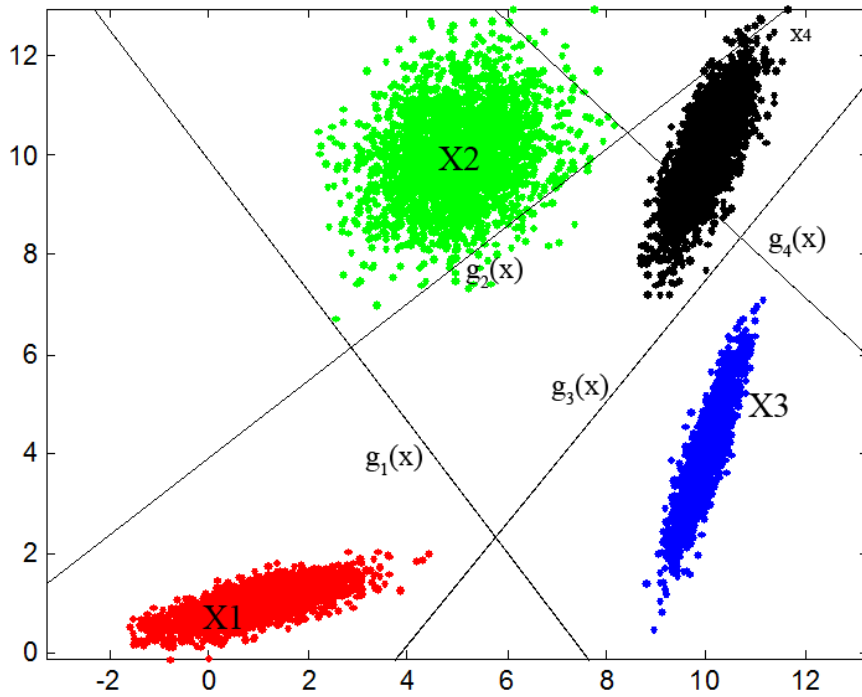
$$\mathbf{a}_i = (\mathbf{Y}^T \mathbf{Y})^{-1} \mathbf{b}_i$$

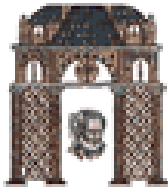
$\mathbf{Y} = [\mathbf{y}_1, \dots, \mathbf{y}_n]$  και  $\mathbf{b}_i = [\mathbf{b}_{jk}]$   $k=1, \dots, n$  και  $b_{jk} = 1$  εάν  $\mathbf{y}_j$  ανήκει στην κλάση  $\omega_i$  και  $b_{jk} = -1$  εάν  $\mathbf{y}_j$  **δεν** ανήκει στην κλάση  $\omega_i$

➤ Σημειώνεται ότι δεν χρειάζεται να αλλάξουμε το πρόσημο στα  $\mathbf{y}_j$  αλλά μπορούμε να το αλλάξουμε στο  $\mathbf{b}_i$



# MSE ταξινόμηση





# Πολλαπλές Κλάσεις

## Δομή του Kesler

- Στόχος: Γραμμικός διαχωρισμός πολλαπλών κλάσεων:

Αν  $y \sim \omega_1$ , τότε  $\mathbf{a}_j' \mathbf{y} - \mathbf{a}_j' \mathbf{y} > 0$  για κάθε  $j=2, \dots, c$ .

- Αυτό το σύστημα των  $c-1$  ανισοτήτων μπορεί να περιγραφεί ως εξής:  
Το  $c\hat{d} - D$  διάνυσμα βαρών  $\hat{\mathbf{a}}$  ταξινομεί ορθά όλα τα  $c-1$   $c\hat{d} - D$  πρότυπα  $\boldsymbol{\eta}_{12}, \boldsymbol{\eta}_{13}, \dots, \boldsymbol{\eta}_{1c}$ :  $\hat{\mathbf{a}}' \boldsymbol{\eta}_{1j} > 0 \quad \forall \quad j=2, \dots, c$  όπου:

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

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



# Ταξινόμηση Πολλαπλών Κλάσεων Perceptron

- Έστω  $y_1, \dots, y_n$  πρότυπα από  $c$  κλάσεις, γραμμικά διαχωρίσιμα. Έστω  $L_k$  μία γραμμική μηχανή  $\mathbf{a}_1(k), \dots, \mathbf{a}_c(k)$ . Θέλουμε να κατασκευάσουμε μία ακολουθία γρ. μηχ.  $L_1, \dots, L_k, \dots$  που να συγκλίνει σε μία διαχωριστική μηχανή  $L$ .
- Έστω  $\mathbf{y}^k$  το  $k$ -στό δείγμα που ζητά διόρθωση (σωστή ταξινόμηση). Αν  $\mathbf{y}^k \sim \omega_i$ , σημαίνει ότι υπάρχει τουλάχιστον ένα  $j \neq i$  για το οποίο  $\mathbf{a}_i^t(k) \mathbf{y}^k < \mathbf{a}_j^t(k) \mathbf{y}^k$ .
- Ο κανόνας διόρθωσης του  $L_k$  (perceptron με σταθερό μοναδιαίο βήμα) λέει:

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

Δηλαδή:

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

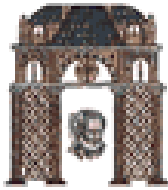
$$\mathbf{a}_j(k+1) = \mathbf{a}_j(k) - \mathbf{y}^k$$

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



---

## ■ ΤΕΛΟΣ ΠΑΡΟΥΣΙΑΣΗΣ

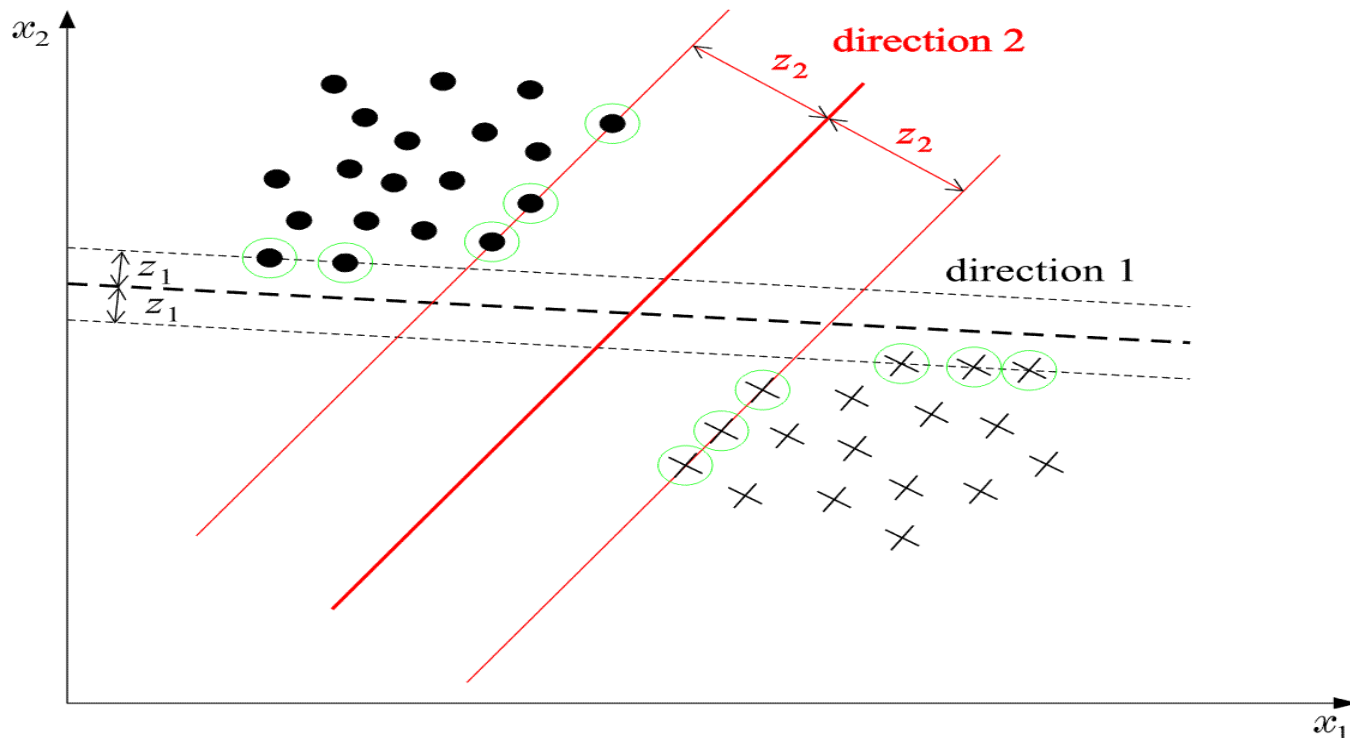


# SUPPORT VECTOR MACHINES

- Στόχος: Εάν έχουμε, δύο κατηγορίες που είναι γραμμικά διαχωρίσιμες, να βρούμε την συνάρτηση διάκρισης

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

- Η οποία αφήνει την μέγιστη απόσταση, περιθώριο, (**maximum margin**) και από τις δύο κατηγορίες from both classes





- Margin: Each hyperplane is characterized by
  - Its direction in space, i.e.,  $\underline{w}$
  - Its position in space, i.e.,  $w_0$
  - For **EACH** direction,  $\underline{w}$  choose the hyperplane that **leaves the SAME distance** from the **nearest** points from each class. The margin is twice this distance.



- The distance of a point  $\hat{x}$  from a hyperplane is given by

$$z_{\hat{x}} = \frac{g(\hat{x})}{\|\underline{w}\|}$$

- Scale,  $\underline{w}$ ,  $w_0$ , so that at the nearest points from each class the discriminant function is  $\pm 1$ :

$$|g(\underline{x})| = 1 \quad \{g(\underline{x}) = +1 \text{ for } \omega_1 \text{ and } g(\underline{x}) = -1 \text{ for } \omega_2\}$$

- Thus the **margin** is given by

$$\frac{1}{\|\underline{w}\|} + \frac{1}{\|\underline{w}\|} = \frac{2}{\|\underline{w}\|}$$

- Also, the following is valid

$$\underline{w}^T \underline{x} + w_0 \geq 1 \quad \forall \underline{x} \in \omega_1$$

$$\underline{w}^T \underline{x} + w_0 \leq -1 \quad \forall \underline{x} \in \omega_2$$





- SVM (linear) classifier

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

- Minimize

$$J(\underline{w}) = \frac{1}{2} \|\underline{w}\|^2$$

- Subject to

$$y_i (\underline{w}^T \underline{x}_i + w_0) \geq 1, \quad i = 1, 2, \dots, N$$

$$y_i = 1, \text{ for } \underline{x}_i \in \omega_1,$$

$$y_i = -1, \text{ for } \underline{x}_i \in \omega_2$$

- This is because minimizing  $\|\underline{w}\|$

the margin  $\frac{2}{\|\underline{w}\|}$  is maximised



- The above is a **quadratic optimization task** subject to a set of linear inequality constraints. The **Karush-Kuhn-Tucker** conditions, state that the **minimizer** satisfies:

- (1) 
$$\frac{\partial}{\partial \underline{w}} L(\underline{w}, w_0, \underline{\lambda}) = 0$$

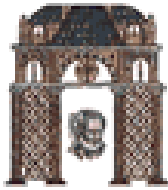
- (2) 
$$\frac{\partial}{\partial w_0} L(\underline{w}, w_0, \underline{\lambda}) = 0$$

- (3) 
$$\lambda_i \geq 0, i = 1, 2, \dots, N$$

- (4) 
$$\lambda_i [y_i (\underline{w}^T \underline{x}_i + w_0) - 1] = 0, i = 1, 2, \dots, N$$

- Where  $L(.,.,.)$  is the **Lagrangian**

$$L(\underline{w}, w_0, \underline{\lambda}) \equiv \frac{1}{2} \underline{w}^T \underline{w} - \sum_{i=1}^N \lambda_i [y_i (\underline{w}^T \underline{x}_i + w_0)]$$



- The solution: from the above, it turns out that

- $$\underline{w} = \sum_{i=1}^N \lambda_i y_i \underline{x}_i$$

- $$\sum_{i=1}^N \lambda_i y_i = 0$$



## Remarks:

- The Lagrange multipliers can be either zero or positive. Thus,

$$\underline{w} = \sum_{i=1}^{N_s} \lambda_i y_i \underline{x}_i$$

where  $N_s \leq N_\theta$  corresponding to positive Lagrange multipliers

- From constraint (4) above, i.e.,

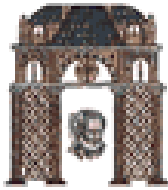
$$\lambda_i [y_i (\underline{w}^T \underline{x}_i + w_0) - 1] = 0, \quad i = 1, 2, \dots, N$$

- the vectors contributing to  $\underline{w}$  satisfy

$$\underline{w}^T \underline{x}_i + w_0 = \pm 1$$



- These vectors are known as **SUPPORT VECTORS** and are the **closest vectors**, from each class, to the classifier.
- Once  $\underline{w}$  is computed,  $w_0$  is determined from conditions (4).
- The optimal hyperplane classifier of a support vector machine is **UNIQUE**.



# Dual Problem Formulation

- The SVM formulation is a convex programming problem, with
  - Convex cost function
  - Convex region of feasible solutions
- Thus, its solution can be achieved by its dual problem, i.e.,

- maximize  $\underset{\underline{\lambda}}{L(\underline{w}, w_0, \underline{\lambda})}$

- subject to

$$\underline{w} = \sum_{i=1}^N \lambda_i y_i \underline{x}_i \quad \sum_{i=1}^N \lambda_i y_i = 0 \quad \underline{\lambda} \geq \underline{0}$$



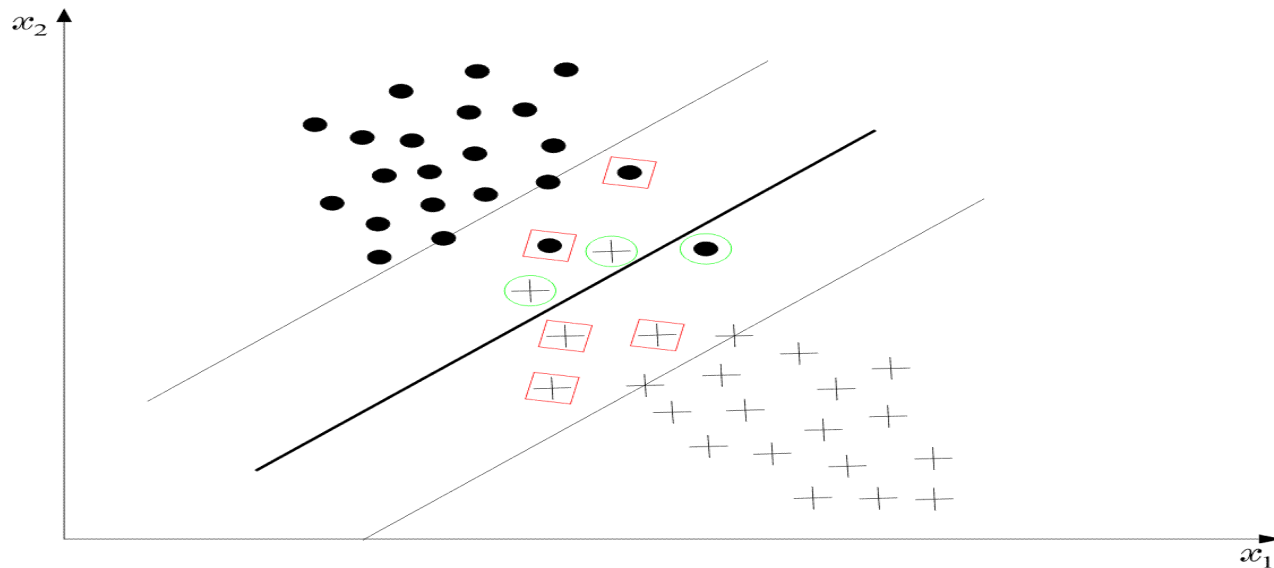
- Combine the above to obtain

- maximize  $\underline{\lambda}$   $\left( \sum_{i=1}^N \lambda_i - \frac{1}{2} \sum_{ij} \lambda_i \lambda_j y_i y_j \underline{x}_i^T \underline{x}_j \right)$

- subject to  $\sum_{i=1}^N \lambda_i y_i = 0$   
 $\underline{\lambda} \geq \underline{0}$



- Remarks:
  - Support vectors enter via **inner products**
  - Although the solution,  $\underline{w}$ , is unique, the Lagrange multipliers ARE NOT.
  
- Non-Separable classes







- In this case there is no hyperplane, such that

$$\underline{w}^T \underline{x} + w_0 (> <) 1, \quad \forall \underline{x}$$

- Recall that the margin is defined as twice the distance between the following two hyperplanes

$$\underline{w}^T \underline{x} + w_0 = 1$$

and

$$\underline{w}^T \underline{x} + w_0 = -1$$



- The training vectors belong to one of three possible categories

- A. Vectors **outside** the band which are **correctly** classified, i.e.,

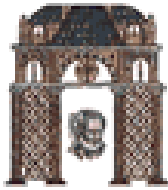
$$y_i(\underline{w}^T \underline{x} + w_0) > 1$$

- B. Vectors **inside** the band, and **correctly** classified, i.e.,

$$0 \leq y_i(\underline{w}^T \underline{x} + w_0) < 1$$

- C. Vectors **misclassified**, i.e.

$$y_i(\underline{w}^T \underline{x} + w_0) < 0$$



- All three cases above can be represented as

$$y_i(\underline{w}^T \underline{x} + w_0) \geq 1 - \xi_i$$

- A.  $\rightarrow \xi_i = 0$
- B.  $\rightarrow 0 < \xi_i \leq 1$
- C.  $\rightarrow 1 < \xi_i$

$\xi_i$  are known as **slack variables**



The goal of the optimization is now two-fold

- Maximize margin
- Minimize the number of patterns with  $\xi_i > 0$ ,

That is

$$J(\underline{w}, w_0, \underline{\xi}) = \frac{1}{2} \|\underline{w}\|^2 + C \sum_{i=1}^N I(\xi_i)$$

where C a constant and

$$I(\xi_i) = \begin{cases} 1 & \xi_i > 0 \\ 0 & \xi_i = 0 \end{cases}$$

- I(.) not differentiable. In practice, we use an approximation

- $$J(\underline{w}, w_0, \underline{\xi}) = \frac{1}{2} \|\underline{w}\|^2 + C \sum_{i=1}^N \xi_i$$

- Following similar procedure as before we obtain



# ■ KKT conditions

$$(1) \underline{w} = \sum_{i=1}^N \lambda_i y_i \underline{x}_i$$

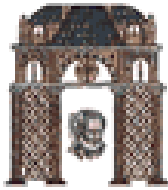
$$(2) \sum_{i=1}^N \lambda_i y_i = 0$$

$$(3) C - \mu_i - \lambda_i = 0, i = 1, 2, \dots, N$$

$$(4) \lambda_i [y_i (\underline{w}^T \underline{x}_i + w_0) - 1 + \xi_i] = 0, i = 1, 2, \dots, N$$

$$(5) \mu_i \xi_i = 0, i = 1, 2, \dots, N$$

$$(6) \mu_i, \lambda_i \geq 0, i = 1, 2, \dots, N$$



- The associated dual problem

Maximize 
$$\underline{\lambda} \left( \sum_{i=1}^N \lambda_i - \frac{1}{2} \sum_{i,j} \lambda_i \lambda_j y_i y_j \underline{x}_i^T \underline{x}_j \right)$$

subject to

$$0 \leq \lambda_i \leq C, \quad i = 1, 2, \dots, N$$

$$\sum_{i=1}^N \lambda_i y_i = 0$$

- Remarks:

The only difference with the separable class case is the existence of  $C$  in the constraints