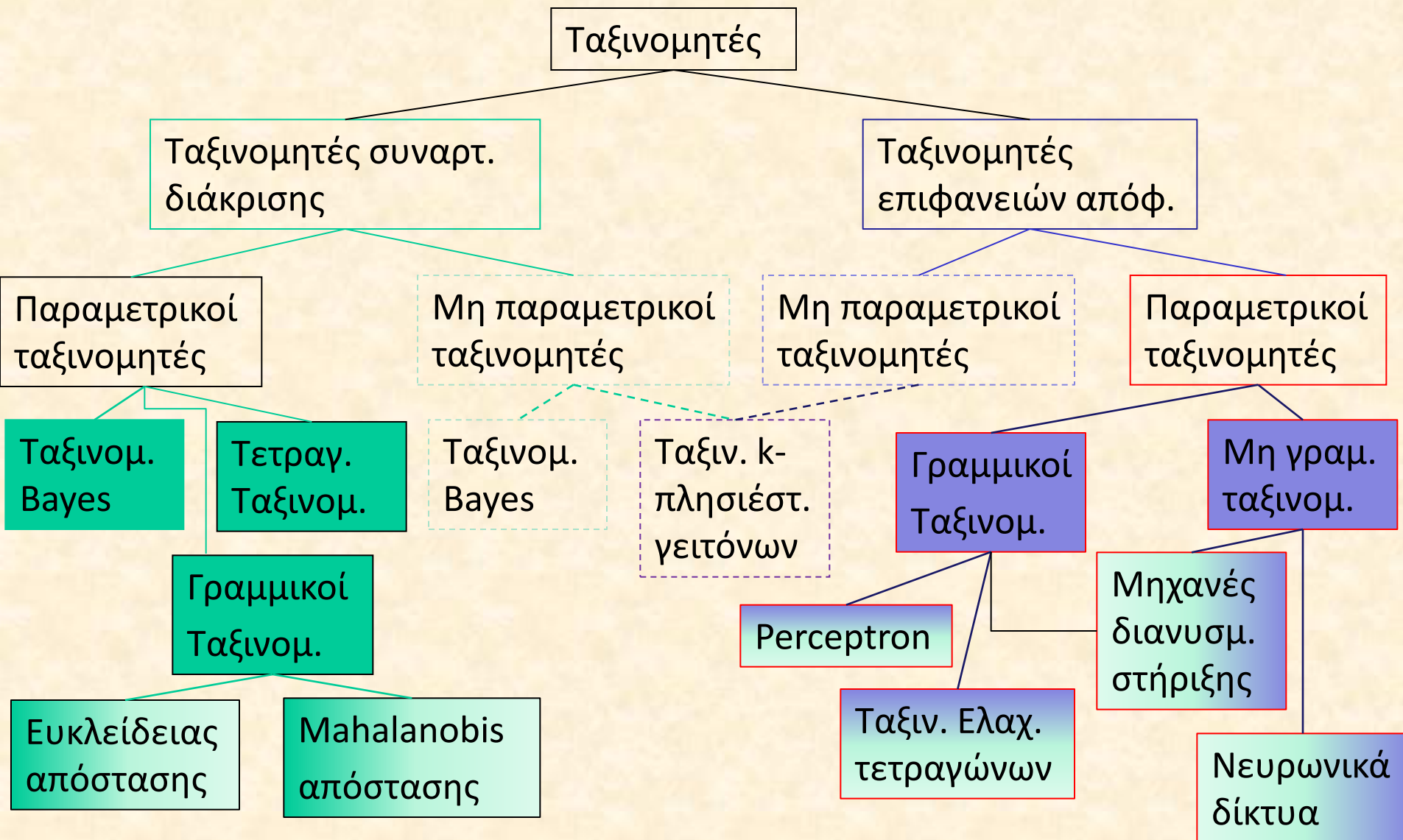


❖ ΑΝΑΓΝΩΡΙΣΗ ΠΡΟΤΥΠΩΝ ❖ (PATTERN RECOGNITION)

Σέργιος Θεοδωρίδης
Κωνσταντίνος Κουτρούμπας

“ΧΑΡΤΟΓΡΑΦΗΣΗ” ΤΟΥ ΧΩΡΟΥ ΤΩΝ ΤΑΞΙΝΟΜΗΤΩΝ



“ΧΑΡΤΟΓΡΑΦΗΣΗ” ΤΟΥ ΧΩΡΟΥ ΤΩΝ ΤΑΞΙΝΟΜΗΤΩΝ

Υπενθ.: X είναι το σύνολο των δεδομένων σημείων όλων των κλάσεων

X_j είναι το υποσύνολο του X που περιέχει τα διανύσματα της κλάσης ω_j ,

$$X = X_1 \cup \dots \cup X_M$$

Ταξινομητές με βάση τις συναρτήσεις διάκρισης

Ταξινομ. Bayes

- $g_j(x) = f(P(\omega_j)p(x|\omega_j))$
- Εκτιμ. $p(x|\omega_j) \approx \hat{p}(x|\omega_j; \mathcal{G}_j)$
- Εκτιμ. \mathcal{G}_j , με βάση το X_j

$$(ML, EM): X_j \rightarrow \hat{\mathcal{G}}_j$$

- $x_i \rightarrow p(x_i|\omega_j) \approx \hat{p}(x_i|\omega_j; \hat{\mathcal{G}}_j)$

Τετραγωνικός
ταξινομητής

- $g_j(x) = (x - \mu_j)^T \Sigma_j^{-1} (x - \mu_j)$
- Υπόθεση: $N(\mu_j, \Sigma_j)$
- Εκτιμ. μ_j, Σ_j , με βάση το X_j

$$(ML): X_j \rightarrow \hat{\mu}_j, \hat{\Sigma}_j$$

- $x_i \rightarrow g_j(x_i) = (x_i - \hat{\mu}_j)^T \hat{\Sigma}_j^{-1} (x_i - \hat{\mu}_j)$

Γραμμικός
Ευκλείδειος
ταξινομητής

- $g_j(x) = (x - \mu_j)^T (x - \mu_j)$
- Υπόθεση: $N(\mu_j, I)$
- Εκτίμ. μ_j , με βάση το X_j

$$(ML): X_j \rightarrow \hat{\mu}_j$$

- $x_i \rightarrow g_j(x_i) = (x_i - \hat{\mu}_j)^T (x_i - \hat{\mu}_j)$

Γραμμικός
Mahalanobis
ταξινομητής

- $g_j(x) = (x - \mu_j)^T \Sigma^{-1} (x - \mu_j)$
- Υπόθεση: $N(\mu_j, \Sigma)$
- Εκτίμ. μ_j, Σ , με βάση το X_j

$$(ML): X_j \rightarrow \hat{\mu}_j, \hat{\Sigma}$$

- $x_i \rightarrow g_j(x_i) = (x_i - \hat{\mu}_j)^T \hat{\Sigma}^{-1} (x_i - \hat{\mu}_j)$

“ΧΑΡΤΟΓΡΑΦΗΣΗ” ΤΟΥ ΧΩΡΟΥ ΤΩΝ ΤΑΞΙΝΟΜΗΤΩΝ

Υπενθ.: X είναι το σύνολο των δεδομένων σημείων όλων των κλάσεων

X_j είναι το υποσύνολο του X που περιέχει τα διανύσματα της κλάσης ω_j ,

$$X = X_1 \cup \dots \cup X_M$$

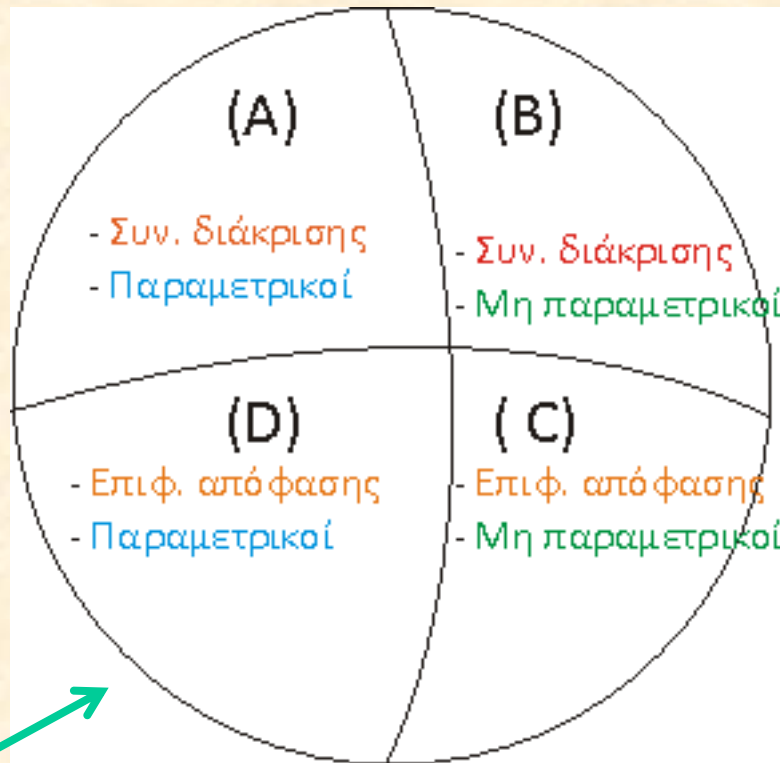
Μη παραμετρικοί ταξινομητές με βάση τις συναρτήσεις διάκρισης

Ταξινομητής
Bayes

– $g_j(x) = f(P(\omega_j)p(x|\omega_j))$
– $x_i \rightarrow p(x_i|\omega_j) \approx \hat{p}(x_i|\omega_j; X_j)$
(παράθυρα Parzen,
εκτίμ. πυκν. βάσει των k -πλησ. γειτ.)

Ταξινομητής k -
πλησιέστερων
γειτόνων

$$- x_i \rightarrow g_j(x_i) = k_i^j$$



ΠΑΡΑΜΕΤΡΙΚΟΙ ΤΑΞΙΝΟΜΗΤΕΣ ΕΠΙΦΑΝΕΙΩΝ ΑΠΟΦΑΣΗΣ

Μερικά προκαταρκτικά:

- Στην πράξη έχουμε στη διάθεσή μας ένα σύνολο δεδομένων (σύν. εκπαίδευσης)

$$X = \{(x_i, d_i), x_i \in R^l, d_i \in \{1, 2, \dots, M\}, i = 1, \dots, N\}$$

όπου

x_i είναι η l -διάστατη αναπαράσταση του i -στού από τις N οντότητες

(διάνυσμα εκπαίδευσης)

d_i είναι η ετικέτα της κλάσης όπου ανήκει το x_i (1 για την ω_1 , 2 για την ω_2, \dots).

- Εστιάζουμε κυρίως στην περίπτωση των δύο κλάσεων.

- Δεν υιοθετούμε καμία υπόθεση σχετικά με τις pdfs που μοντελοποιούν τις διάφορες κλάσεις.

- Αναζητούμε την επιφάνεια (γραμμική ή μη γραμμική) που επιτυγχάνει τον “βέλτιστο” διαχωρισμό των (διανυσμάτων δεδομένων των) κλάσεων.

ΠΑΡΑΜΕΤΡΙΚΟΙ ΤΑΞΙΝΟΜΗΤΕΣ ΕΠΙΦΑΝΕΙΩΝ ΑΠΟΦΑΣΗΣ

Επιφάνεια απόφασης (C): Περιγράφεται από εξίσωση της μορφής $h(\mathbf{x})=0$, ή $h(\mathbf{x};\mathbf{w})=0$, όπου \mathbf{w} είναι το διάνυσμα των παραμέτρων που ορίζουν την επιφάνεια (C).

Παραδείγματα:

1. Αν η (C) είναι **καμπύλη 1^{ου} βαθμού** (υπερεπίπεδο – γραμμικός διαχωρισμός), τότε

$$h(x, w) = w_1 x_1 + \dots + w_l x_l + w_0 = \sum_{k=1}^l w_k x_k + w_0 = \mathbf{w}^T \mathbf{x} + w_0 = 0$$

όπου $\mathbf{w}=[w_1, \dots, w_l]^T$, $\mathbf{x}=[x_1, \dots, x_l]^T$. Οι παράμετροι είναι το διάνυσμα \mathbf{w} και το w_0 .

2. Αν η (C) είναι **καμπύλη 2^{ου} βαθμού** (π.χ. Υπερέλλειψη – μη γραμμικός διαχωρισμός), τότε

$$h(x, w) = \sum_{k=1}^l \sum_{q=k}^l w_{kq} x_k x_q + \sum_{k=1}^l w_k x_k + w_0 = 0$$

όπου $\mathbf{w}=[w_0, w_1, \dots, w_l, w_{11}, \dots, w_{1l}, w_{22}, \dots, w_{2l}, \dots, w_{ll}]^T$

Σημείωση: Σε αρκετούς μη γραμμικούς ταξινομητές (π.χ. στα νευρωνικά δίκτυα) η μορφή της (C) δεν μπορεί να εκφραστεί άμεσα (explicitly).

ΠΑΡΑΜΕΤΡΙΚΟΙ ΤΑΞΙΝΟΜΗΤΕΣ ΕΠΙΦΑΝΕΙΩΝ ΑΠΟΦΑΣΗΣ

Υπόθεση: Αν δεν ορίζεται διαφορετικά, θεωρούμε την περίπτωση των **δύο κλάσεων**, δηλ., $\omega_1 (+1)$ και $\omega_2 (-1)$.

Ορισμός προβλήματος: Δοθέντος ενός συνόλου δεδομένων X προσδιόρισε μία **επιφάνεια** που επιτυγχάνει το **“βέλτιστο” δυνατό διαχωρισμό** των (διανυσμάτων των) δύο κλάσεων.

Στρατηγική αντιμετώπισης του προβλήματος:

1. Υιοθέτησε μια **συγκεκριμένη** (“άμεση” ή “έμμεση”) **παραμετρική μορφή** για την επιφάνεια $h(\mathbf{x}; \mathbf{w})=0$.
2. Όρισε κατάλληλη συνάρτηση (**συνάρτηση κόστους - cost function**) του \mathbf{w} , $J(\mathbf{w})$, η οποία περιλαμβάνει επίσης τα διανύσματα του X , έτσι ώστε **τα βέλτιστά της (ελάχιστα ή μέγιστα) να αντιστοιχούν στις καλύτερες δυνατές επιφάνειες για το υπό μελέτη πρόβλημα**.
3. Βελτιστοποίησε την $J(\mathbf{w})$ ως προς το \mathbf{w} . Η θέση του βέλτιστου αυτής ορίζει την επιφάνεια απόφασης.

Σημαντική παρατήρηση: Το νόημα της φράσης **“βέλτιστη δυνατή καμπύλη”** διαφέρει για **διαφορετικές επιλογές της $J(\mathbf{w})$** .

ΒΑΣΙΚΑ ΣΤΟΙΧΕΙΑ ΘΕΩΡΙΑΣ ΒΕΛΤΙΣΤΟΠΟΙΗΣΗΣ

Έστω $J(w)$ συνεχής συνάρτηση του w .

Πρόβλημα (P1): Προσδιόρισε τη **θέση w^*** όπου η συνάρτηση $J(w)$ λαμβάνει την **ελάχιστη** τιμή της.

Μια απλή μέθοδος για την επίλυση του **(P1)** είναι αυτή της **οξύτερης καθόδου (gradient descent - GD)**.

- Αρχικοποίησε $w=w(0)$

- $t=0$

- Επανάλαβε

$$- w(t+1) = w(t) - \mu \frac{\partial J(w)}{\partial w} \Big|_{w=w(t)}$$

- $t=t+1$

- Έως ότου επιτευχθεί σύγκλιση

ΒΑΣΙΚΑ ΣΤΟΙΧΕΙΑ ΘΕΩΡΙΑΣ ΒΕΛΤΙΣΤΟΠΟΙΗΣΗΣ

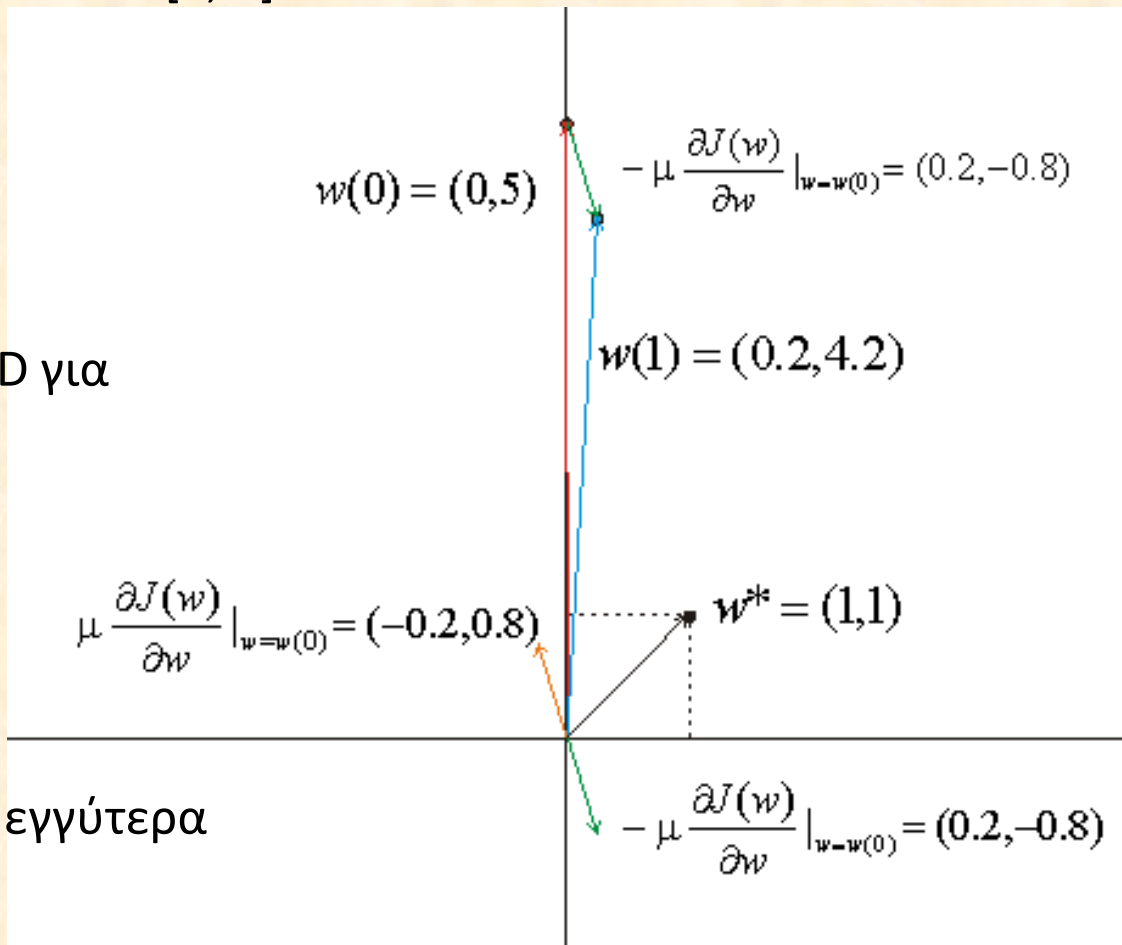
-Ένα παράδειγμα: Έστω $\mathbf{w}=[w_1, w_2]^T$ και $J(\mathbf{w})=(w_1-1)^2+(w_1-1)^2$. Είναι σαφές ότι η $J(\mathbf{w})$ λαμβάνει την ελάχιστη τιμή της στο $\mathbf{w}^*=[1, 1]^T$.

-Είναι
$$\frac{\partial J(\mathbf{w})}{\partial \mathbf{w}} = \begin{bmatrix} 2w_1 - 2 \\ 2w_2 - 2 \end{bmatrix}$$

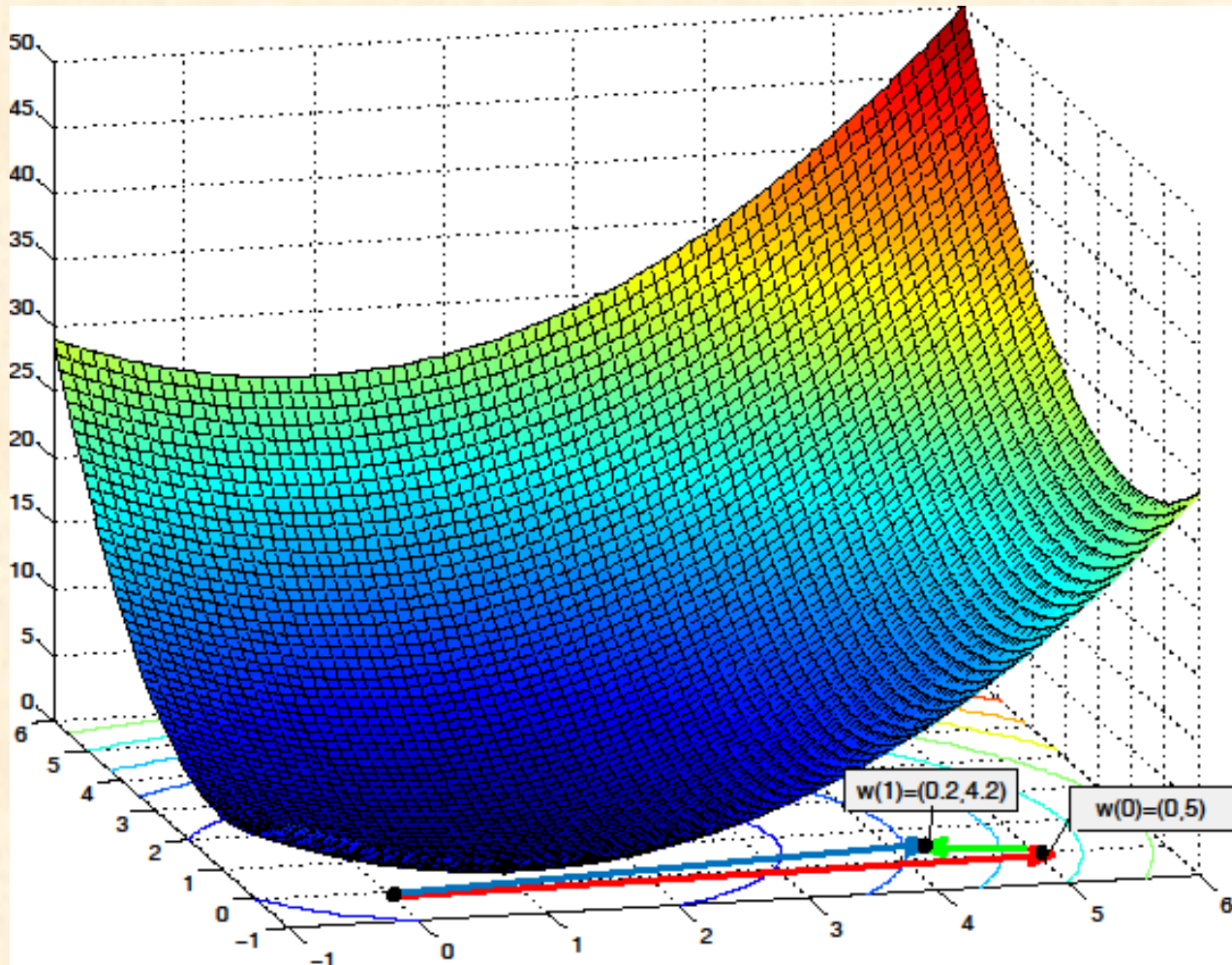
-Εφαρμόζοντας τον αλγόριθμο GD για $\mathbf{w}(0)=[0, 5]^T$, και $\mu=0.1$, έχουμε

$$\mathbf{w}(1) = \begin{bmatrix} 0 \\ 5 \end{bmatrix} - 0.1 \begin{bmatrix} -2 \\ 8 \end{bmatrix} = \begin{bmatrix} 0.2 \\ 4.2 \end{bmatrix}$$

-Βλέπουμε ότι το $\mathbf{w}(1)$ βρίσκεται εγγύτερα στο \mathbf{w}^* .



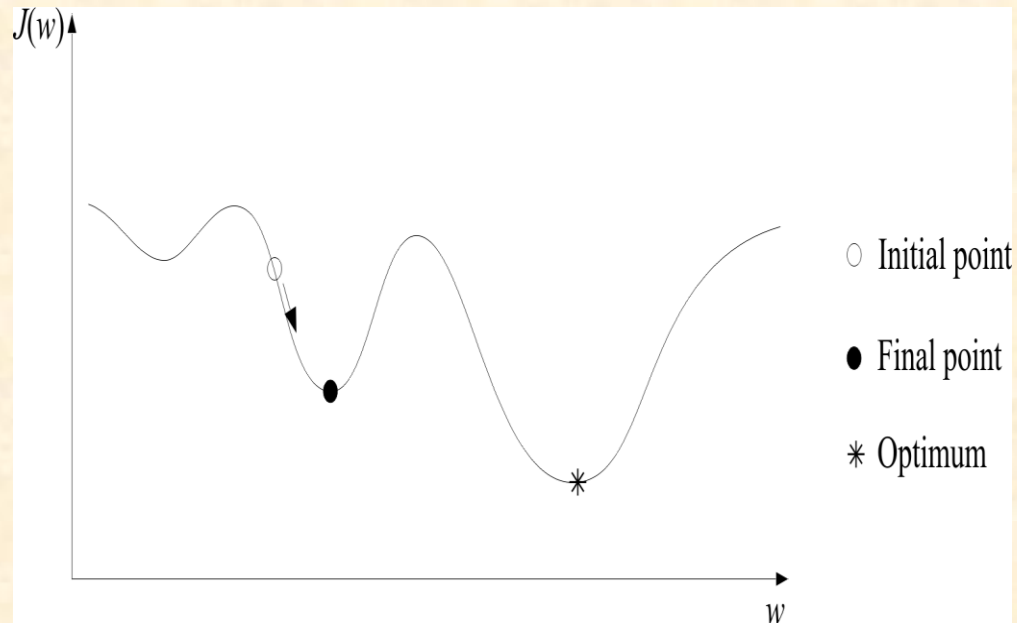
ΒΑΣΙΚΑ ΣΤΟΙΧΕΙΑ ΘΕΩΡΙΑΣ ΒΕΛΤΙΣΤΟΠΟΙΗΣΗΣ



ΒΑΣΙΚΑ ΣΤΟΙΧΕΙΑ ΘΕΩΡΙΑΣ ΒΕΛΤΙΣΤΟΠΟΙΗΣΗΣ

Σχόλια για τον αλγόριθμο GD:

- Η τιμή του μ πρέπει να είναι **όχι πολύ μεγάλη**, ώστε να αποφεύγονται ταλαντώσεις γύρω από το ελάχιστο της συνάρτησης και **όχι πολύ μικρή**, ώστε να αποφεύγονται άσκοπες καθυστερήσεις στη σύγκλιση.
- Αν η $J(\mathbf{w})$ έχει **περισσότερα από ένα ελάχιστα**, ο αλγόριθμος GD θα συγκλίνει (γενικά) σε εκείνο που βρίσκεται εγγύτερα στο $\mathbf{w}(0)$.
- Αν ο αλγόριθμος παγιδευτεί σε ένα **τοπικό ελάχιστο που αντιστοιχεί σε μια “κακή” λύση**, ο μόνος τρόπος για να το **αποφύγουμε** είναι να **επανεκκινήσουμε** τον αλγόριθμο από μια διαφορετική αρχική θέση.
- Μπορεί να αποδειχθεί ότι, κάτω από κάποιες συνθήκες, ο αλγόριθμος **συγκλίνει ασυμπτωτικά** σε ένα **τοπικό ελάχιστο** της $J(\mathbf{w})$.



ΒΑΣΙΚΑ ΣΤΟΙΧΕΙΑ ΘΕΩΡΙΑΣ ΒΕΛΤΙΣΤΟΠΟΙΗΣΗΣ

Έστω $J(\mathbf{w})$ μια συνεχής συνάρτηση του \mathbf{w} .

Πρόβλημα (P2): Προσδιόρισε τη **θέση \mathbf{w}^*** όπου η συνάρτηση $J(\mathbf{w})$ λαμβάνει την **ελάχιστη** τιμή της, υπό την προϋπόθεση ότι το \mathbf{w} ικανοποιεί μερικούς **περιορισμούς ισότητας (equality constraints)**.

Για **γραμμικούς περιορισμούς ισότητας**, το πρόβλημα διατυπώνεται ως εξής

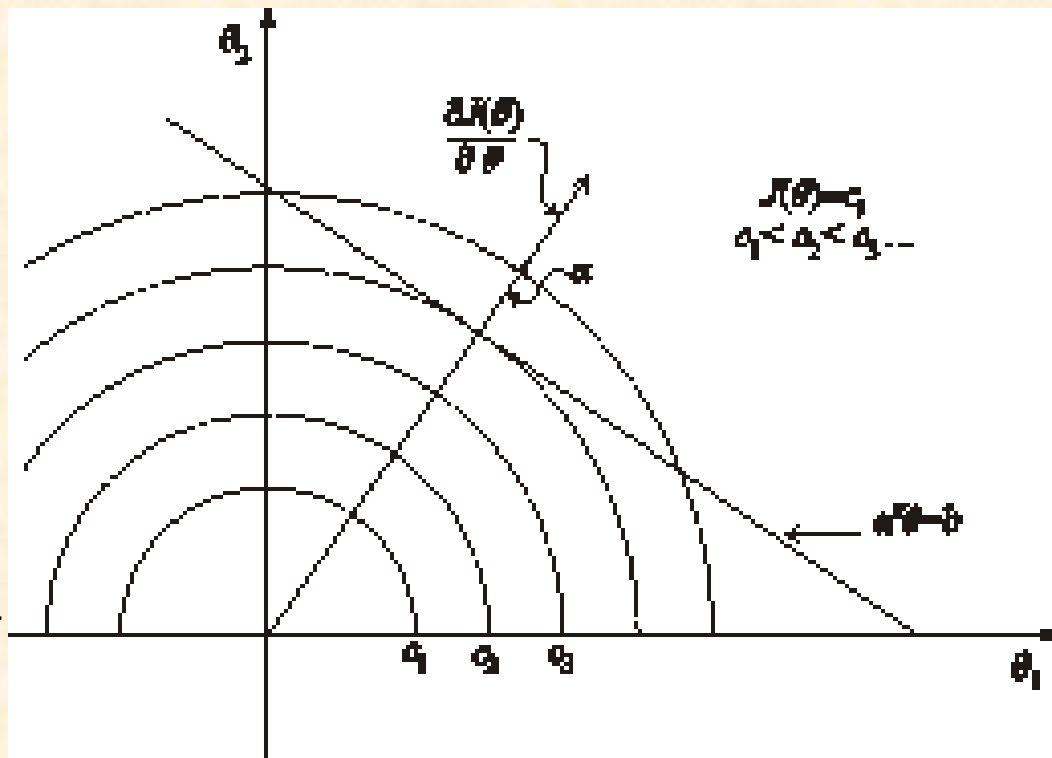
- Ελαχιστοποίησε την $J(\mathbf{w})$
- Υπό τους περιορισμούς $A\mathbf{w}=\mathbf{b}$, όπου A ένας $m \times l$ πίνακας και \mathbf{b} ένα m -διάστατο διάνυσμα.

Λύση: Πολ/στές Lagrange

Ελαχιστοποίησε την

$$-L(\mathbf{w})=J(\mathbf{w})+\boldsymbol{\lambda}^T(A\mathbf{w}-\mathbf{b})$$

- $\boldsymbol{\lambda}$ είναι ένα m -διάστατο διάνυσμα που εκτιμάται μέσω των περιορισμών $A\mathbf{w}=\mathbf{b}$



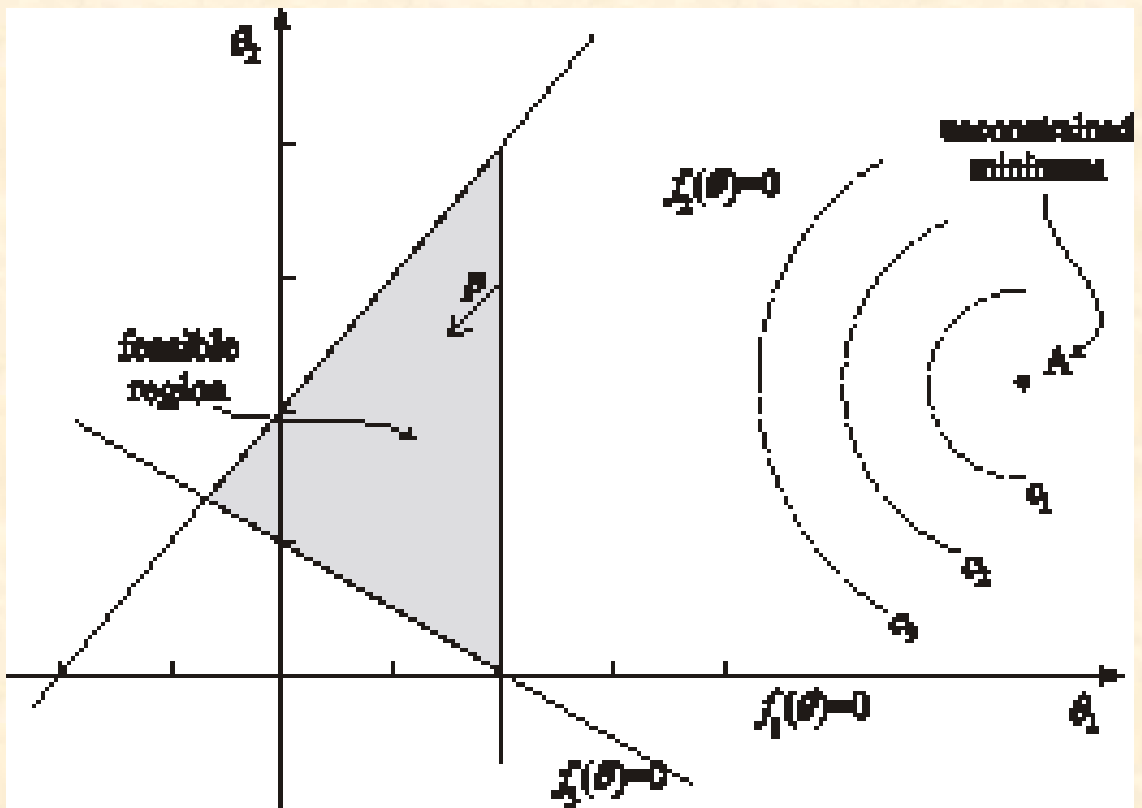
ΒΑΣΙΚΑ ΣΤΟΙΧΕΙΑ ΘΕΩΡΙΑΣ ΒΕΛΤΙΣΤΟΠΟΙΗΣΗΣ

Έστω $J(\mathbf{w})$ μια συνεχής συνάρτηση του \mathbf{w} .

Πρόβλημα (P3): Προσδιόρισε τη **θέση \mathbf{w}^*** όπου η συνάρτηση $J(\mathbf{w})$ λαμβάνει την **ελάχιστη** τιμή της, υπό τις την προϋπόθεση ότι το \mathbf{w} ικανοποιεί μερικούς **περιορισμούς ανισότητας**.

Για **γραμμικούς περιορισμούς ανισότητας**, το πρόβλημα διατυπώνεται ως εξής

- Ελαχιστοποίησε την $J(\mathbf{w})$
- Υπό τους περιορισμούς $A\mathbf{w} \geq \mathbf{b}$, όπου A ένας $m \times l$ πίνακας και \mathbf{b} ένα m -διάστατο διάνυσμα.



ΓΡΑΜΜΙΚΟΙ ΠΑΡΑΜΕΤΡΙΚΟΙ ΤΑΞΙΝΟΜΗΤΕΣ ΕΠΙΦΑΝΕΙΩΝ ΑΠΟΦΑΣΗΣ

- Στην περίπτωση αυτή η επιφάνεια είναι ένα **υπερεπίπεδο** (H) της μορφής

$$(H) : h(x, w) = w_1 x_1 + \dots + w_l x_l + w_0 = \sum_{k=1}^l w_k x_k + w_0 = w^T x + w_0 = 0$$

όπου $w = [w_1, \dots, w_l]^T$, $x = [x_1, \dots, x_l]^T$

- Το **(H)** ορίζεται πλήρως από τα **w** και **w₀**.

- Μερικά στοιχεία από τη Γεωμετρία:

- Το διάνυσμα **w** είναι **κάθετο** στο **(H)**.

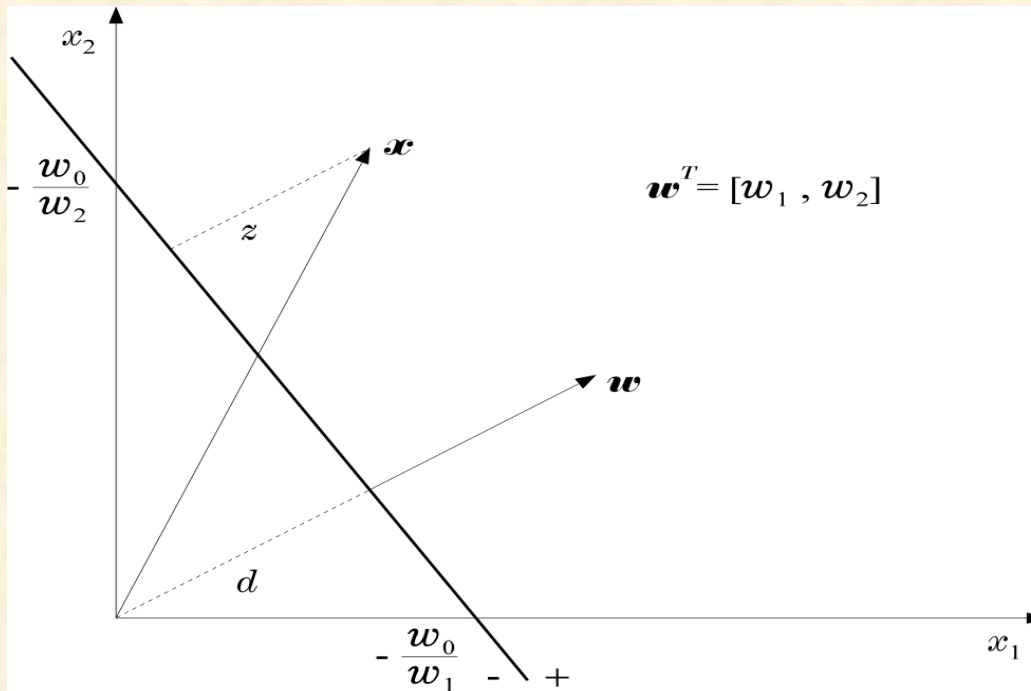
- Το **w** ορίζει τον **προσανατολισμό** (orientation) του **(H)**.

- Όλα τα **υπερεπίπεδα** που είναι **παράλληλα** με το **(H)** έχουν τον **ίδιο προσανατολισμό** με αυτό.

- Το **w** είναι **κάθετο** σε **όλα** τα **υπερεπίπεδα** που είναι **παράλληλα** με το **(H)**.

- Το **w₀** είναι η παράμετρος που **ταυτοποιεί μοναδικά** ένα συγκεκριμένο υπερεπίπεδο, ανάμεσα από υπερεπίπεδα ίδιου προσανατολισμού.

ΓΡΑΜΜΙΚΟΙ ΠΑΡΑΜΕΤΡΙΚΟΙ ΤΑΞΙΝΟΜΗΤΕΣ ΕΠΙΦΑΝΕΙΩΝ ΑΠΟΦΑΣΗΣ



$$d = \frac{|w_0|}{\sqrt{w_1^2 + w_2^2}}, \quad z = \frac{|g(\underline{x})|}{\sqrt{w_1^2 + w_2^2}}$$

ΓΡΑΜΜΙΚΟΙ ΠΑΡΑΜΕΤΡΙΚΟΙ ΤΑΞΙΝΟΜΗΤΕΣ ΕΠΙΦΑΝΕΙΩΝ ΑΠΟΦΑΣΗΣ

ΜΙΑ ΣΥΜΒΑΣΗ: Προκειμένου να καταστήσουμε πιο συμπαγή το συμβολισμό

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

ορίζουμε

$$\mathbf{w}^* = [\mathbf{w}^T, w_0]^T \text{ and } \mathbf{x}^* = [\mathbf{x}^T, 1]^T.$$

Έτσι έχουμε

$$(H) : \mathbf{w}^{*T} \mathbf{x}^* = 0$$

ΣΗΜΕΙΩΣΗ: Στη συνέχεια, εκτός αν ορίζεται διαφορετικά, θα γράφουμε \mathbf{w} και \mathbf{x} , εννοώντας τα \mathbf{w}^* και \mathbf{x}^* , αντίστοιχα.

ΓΡΑΜΜΙΚΟΙ ΠΑΡΑΜΕΤΡΙΚΟΙ ΤΑΞΙΝΟΜΗΤΕΣ ΕΠΙΦΑΝΕΙΩΝ ΑΠΟΦΑΣΗΣ

Ο αλγόριθμος perceptron

Υπόθεση: (Τα διανύσματα των) δύο κλάσεων ω_1 και ω_2 είναι γραμμικώς διαχωρίσιμα.

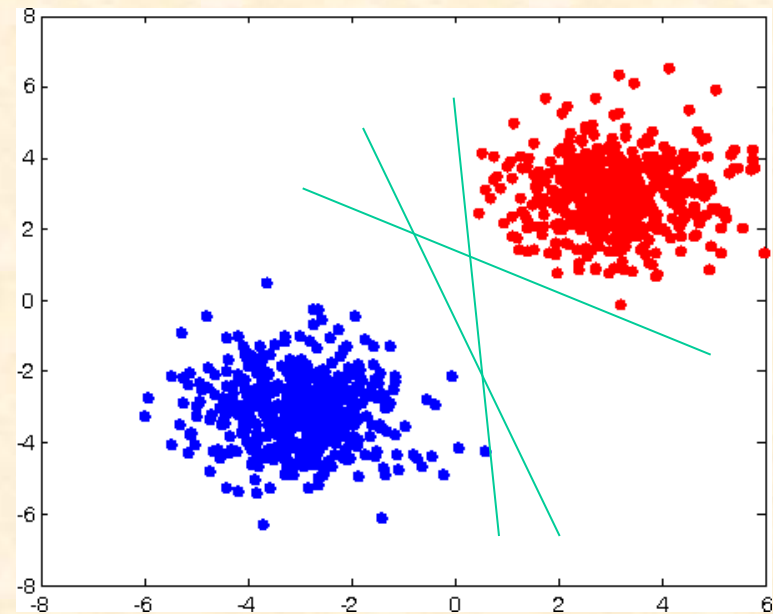
Έτσι, υπάρχει ένα υπερεπίπεδο (H^*): $\mathbf{w}^{*T}\mathbf{x}=0$, ώστε

$$\mathbf{w}^{*T} \mathbf{x} > 0, \quad \forall \mathbf{x} \in \omega_1$$

$$\mathbf{w}^{*T} \mathbf{x} < 0, \quad \forall \mathbf{x} \in \omega_2$$

Το κριτήριο στο perceptron: Προσδιόρισε ένα υπερεπίπεδο το οποίο διαχωρίζει τέλεια τα διανύσματα από τις δύο κλάσεις.

ΣΗΜΕΙΩΣΗ: Αν οι κλάσεις είναι γραμμικώς διαχωρίσιμες, υπάρχουν άπειρα υπερεπίπεδα που ικανοποιούν το κριτήριο perceptron.



ΓΡΑΜΜΙΚΟΙ ΠΑΡΑΜΕΤΡΙΚΟΙ ΤΑΞΙΝΟΜΗΤΕΣ ΕΠΙΦΑΝΕΙΩΝ ΑΠΟΦΑΣΗΣ

Ο αλγόριθμος perceptron

Η συνάρτηση κόστους:

$$J(\underline{w}) = \sum_{x \in Y(\underline{w})} (\delta_x \underline{w}^T x)$$

όπου

- $Y(\underline{w})$ το σύνολο των **λανθασμένα ταξινομημένων** διανυσμάτων από το \underline{w} και

$$\begin{aligned} \delta_x &= -1 \text{ if } x \in \omega_1 \\ \delta_x &= +1 \text{ if } x \in \omega_2 \end{aligned}$$

Σημειώσεις:

- Είναι $J(\underline{w}) \geq 0$.

Πράγματι, αν $x \in Y(\underline{w})$ και $x \in \omega_1$ τότε $\underline{w}^T x < 0$ και $\delta_x = -1$. Επομένως $\delta_x \underline{w}^T x > 0$.

Επίσης, αν $x \in Y(\underline{w})$ και $x \in \omega_2$ τότε $\underline{w}^T x > 0$ και $\delta_x = +1$. Επομένως $\delta_x \underline{w}^T x > 0$.

- Είναι $J(\underline{w}) = 0$ μόνον όταν $Y = \emptyset$, πράγμα που σημαίνει ότι προσδιορίστηκε **μία λύση που ικανοποιεί το κριτήριο perceptron**.

ΓΡΑΜΜΙΚΟΙ ΠΑΡΑΜΕΤΡΙΚΟΙ ΤΑΞΙΝΟΜΗΤΕΣ ΕΠΙΦΑΝΕΙΩΝ ΑΠΟΦΑΣΗΣ

Ο αλγόριθμος perceptron

Η συνάρτηση κόστους:

Η $J(w)$ είναι τμηματικά γραμμική (piecewise linear) (γιατί;)



Η **ελαχιστοποίηση** της $J(w)$ επιτυγχάνεται μέσω διαδικασίας που ακολουθεί τη φιλοσοφία της **GD**.

Όπου είναι έγκυρο, ορίζουμε
$$\frac{\partial J(w)}{\partial w} = \frac{\partial}{\partial w} \left(\sum_{x \in Y(w)} \delta_x w^T x \right) = \sum_{x \in Y(w)} \delta_x x$$

ΓΡΑΜΜΙΚΟΙ ΠΑΡΑΜΕΤΡΙΚΟΙ ΤΑΞΙΝΟΜΗΤΕΣ ΕΠΙΦΑΝΕΙΩΝ ΑΠΟΦΑΣΗΣ

Ο αλγόριθμος perceptron

- Αρχικοποίησε $w=w(0)$

- $t=0$

- Επανάλαβε

- Προσδιόρισε τα **λανθασμένα ταξινομημένα διανύσματα** του X από το $w(t)$ και συσσωρεύσέ τα στο $Y(w_t)$.

$$w(t+1) = w(t) - \rho_t \sum_{x \in Y(w_t)} \delta_x x$$

- $t=t+1$

- Έως ότου επιτευχθεί σύγκλιση

Σημειώσεις:

• Πρόκειται για αλγόριθμο επεξεργασίας **κατά συρροή** (batch algorithm), όπου η ενημέρωση των παραμέτρων σε κάθε επανάληψη πραγματοποιείται μετά την επεξεργασία όλων των διανυσμάτων του X .

• Ο αλγόριθμος perceptron **συγκλίνει μόνον όταν** οι κλάσεις είναι **γραμμικώς διαχωρίσιμες**. Διαφορετικά, αποτυγχάνει να συγκλίνει.

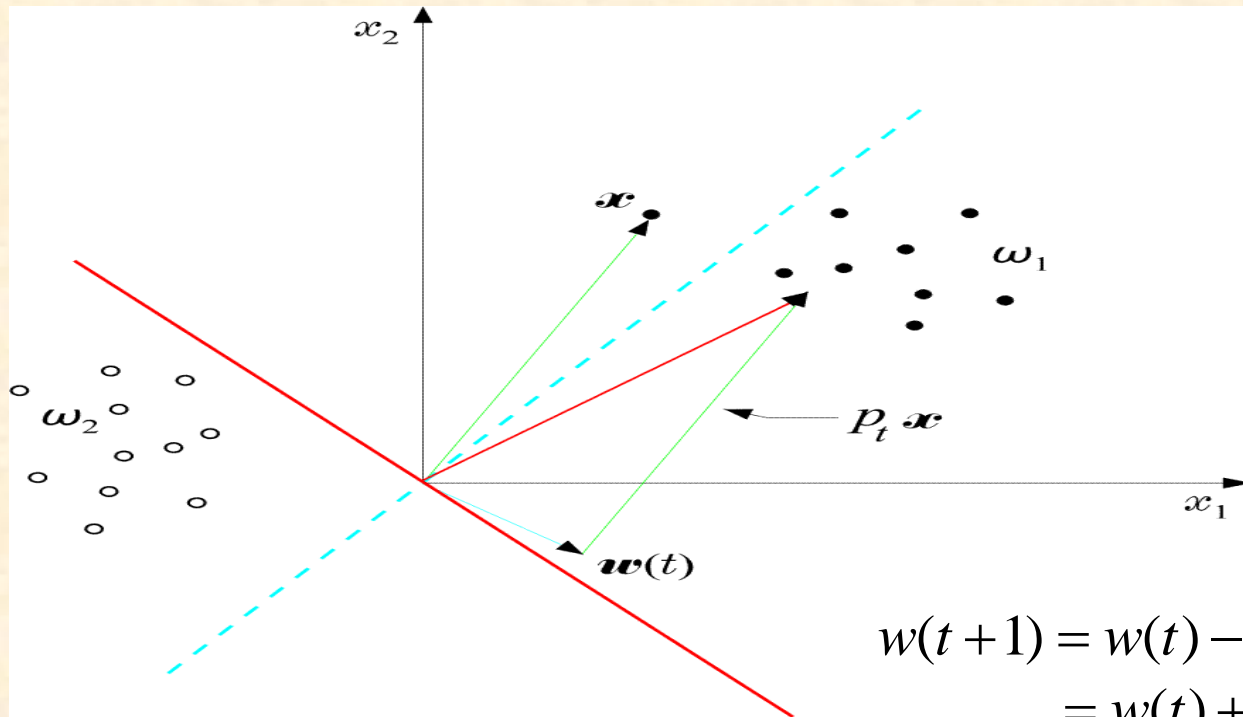
• Αν οι κλάσεις είναι **γραμμικώς διαχωρίσιμες**, ο αλγόριθμος **συγκλίνει** σε **πεπερασμένο αριθμό βημάτων**

$$\lim_{t \rightarrow \infty} \sum_{k=0}^t \rho_k \rightarrow \infty, \quad \lim_{t \rightarrow \infty} \sum_{k=0}^t \rho_k^2 < +\infty \text{ π.χ. } \rho_t = c/t$$

ΓΡΑΜΜΙΚΟΙ ΠΑΡΑΜΕΤΡΙΚΟΙ ΤΑΞΙΝΟΜΗΤΕΣ ΕΠΙΦΑΝΕΙΩΝ ΑΠΟΦΑΣΗΣ

Ο αλγόριθμος perceptron

Παράδειγμα 1:



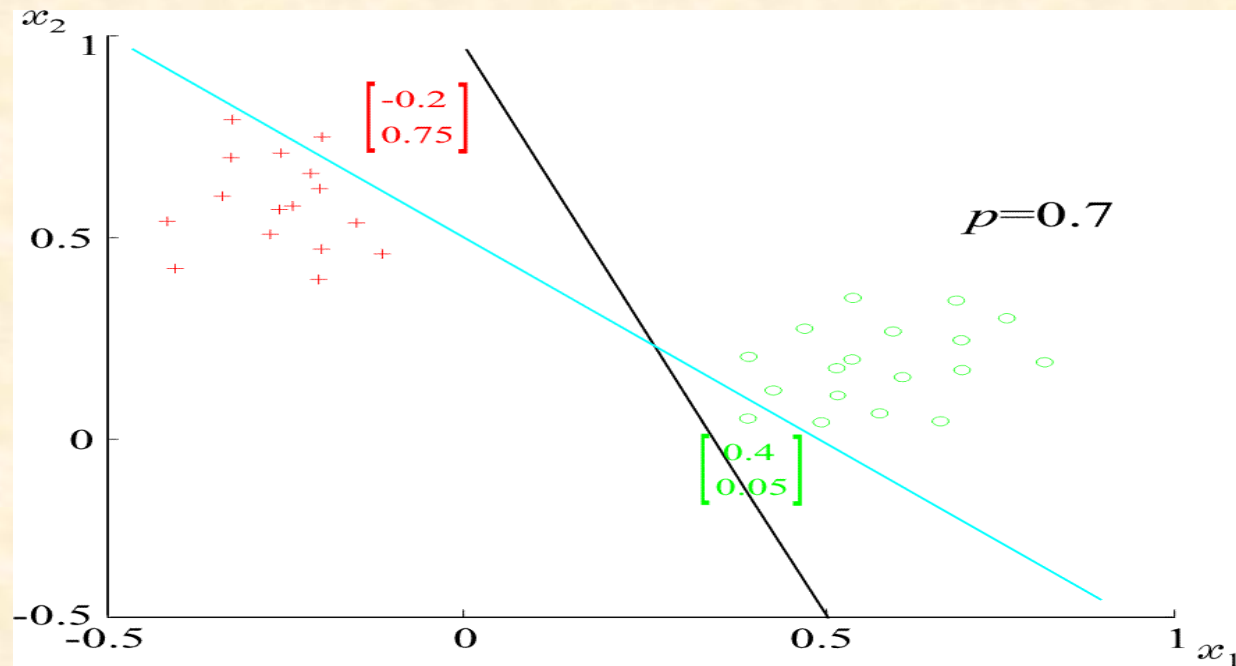
ΓΡΑΜΜΙΚΟΙ ΠΑΡΑΜΕΤΡΙΚΟΙ ΤΑΞΙΝΟΜΗΤΕΣ ΕΠΙΦΑΝΕΙΩΝ ΑΠΟΦΑΣΗΣ

Ο αλγόριθμος perceptron

Παράδειγμα 2: Σε κάποια επανάληψη t ο αλγόριθμος perceptron δίνει

$w_1 = 1$, $w_2 = 1$, $w_0 = -0.5$ που αντιστοιχεί στην ευθεία $x_1 + x_2 - 0.5 = 0$

Το υπερέπιπεδο της επόμενης επανάληψης είναι



$$\underline{w}(t+1) = \begin{bmatrix} 1 \\ 1 \\ -0.5 \end{bmatrix} - 0.7(-1) \begin{bmatrix} 0.4 \\ 0.05 \\ 1 \end{bmatrix} - 0.7(+1) \begin{bmatrix} -0.2 \\ 0.75 \\ 1 \end{bmatrix} = \begin{bmatrix} 1.42 \\ 0.51 \\ -0.5 \end{bmatrix}$$

ΓΡΑΜΜΙΚΟΙ ΠΑΡΑΜΕΤΡΙΚΟΙ ΤΑΞΙΝΟΜΗΤΕΣ ΕΠΙΦΑΝΕΙΩΝ ΑΠΟΦΑΣΗΣ

Παραλλαγές του αλγόριθμου perceptron:

-Ο αλγόριθμος τσέπης (rocket algorithm).

- Κρατά την καλύτερη λύση που βρέθηκε κατά την εκτέλεση του αλγόριθμου perceptron για συγκεκριμένο αριθμό επαναλήψεων.
- Κατάλληλος για προβλήματα μη γραμμικώς διαχωρίσιμων κλάσεων.

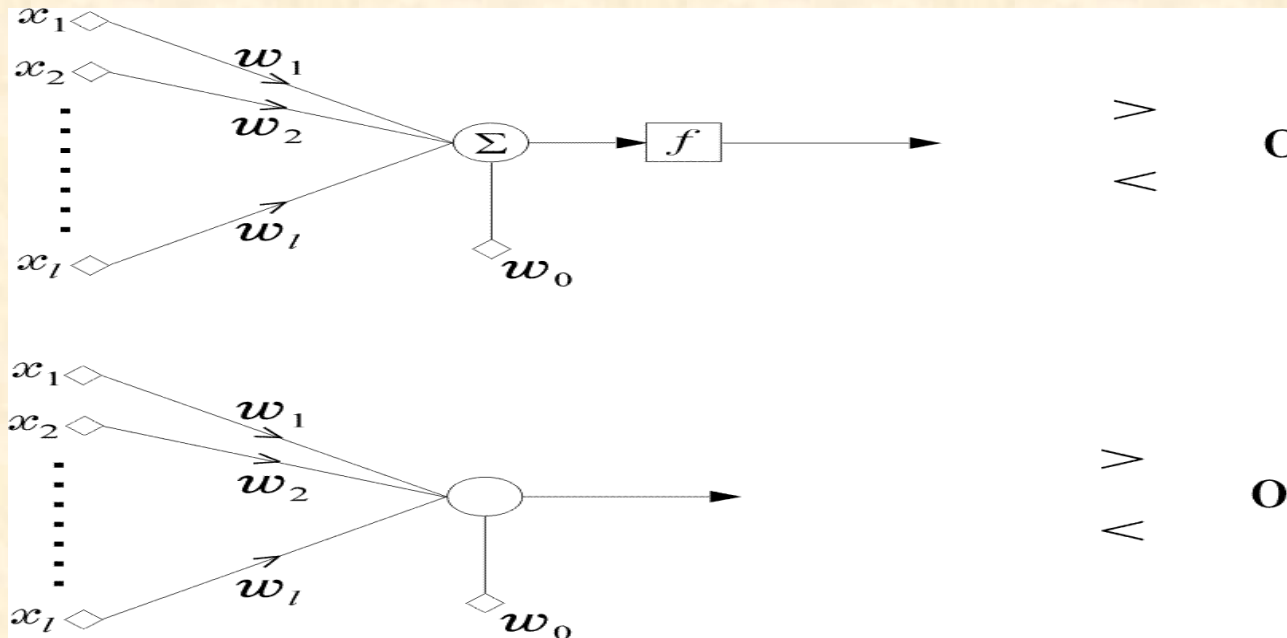
-Ο on-line αλγόριθμος perceptron:

- Η ενημέρωση των παραμέτρων λαμβάνει χώρα μετά την μεμονωμένη επεξεργασία κάθε διανύσματος του X .

$$\begin{aligned} w(t+1) &= w(t) + \rho x_{(t)}, & w^T(t)x_{(t)} &\leq 0 \\ & & x_{(t)} &\in \omega_1 \\ w(t+1) &= w(t) - \rho x_{(t)}, & w^T(t)x_{(t)} &\geq 0 \\ & & x_{(t)} &\in \omega_2 \\ w(t+1) &= w(t) & \text{otherwise} & \end{aligned}$$

ΓΡΑΜΜΙΚΟΙ ΠΑΡΑΜΕΤΡΙΚΟΙ ΤΑΞΙΝΟΜΗΤΕΣ ΕΠΙΦΑΝΕΙΩΝ ΑΠΟΦΑΣΗΣ

Η δομή perceptron



w_i 's synapses or synaptic weights
 w_0 threshold

- Η δομή αυτή καλείται **perceptron** ή **νευρώνας (neuron)**
- Είναι μια **μηχανή** που μπορεί να **“μάθει”** (δηλ. να προσδιορίσει τις τιμές των εμπλεκόμενων παραμέτρων) από τα **διανύσματα εκπαίδευσης** μέσω του **αλγόριθμου perceptron**.

ΓΡΑΜΜΙΚΟΙ ΠΑΡΑΜΕΤΡΙΚΟΙ ΤΑΞΙΝΟΜΗΤΕΣ ΕΠΙΦΑΝΕΙΩΝ ΑΠΟΦΑΣΗΣ

Μέθοδοι ελαχίστων τετραγώνων (Least square methods)

- Αν οι κλάσεις είναι γραμμικώς διαχωρίσιμες, το perceptron θα δώσει σαν έξοδο ± 1
- Αν οι κλάσεις ΔΕΝ είναι γραμμικώς διαχωρίσιμες, θα υπολογίσουμε τα βάρη (παραμέτρους) w_1, w_2, \dots, w_0

έτσι ώστε η **διαφορά** ανάμεσα

- Στην **πραγματική απόκριση** του **ταξινομητή**, $w^T x$, και
- Την αντίστοιχη **επιθυμητή απόκριση**, δηλ.
$$\begin{cases} +1 & \text{if } x \in \omega_1 \\ -1 & \text{if } x \in \omega_2 \end{cases}$$

να είναι όσο το δυνατόν **ΜΙΚΡΟΤΕΡΗ** για όλα τα διανύσματα του X .

ΓΡΑΜΜΙΚΟΙ ΠΑΡΑΜΕΤΡΙΚΟΙ ΤΑΞΙΝΟΜΗΤΕΣ ΕΠΙΦΑΝΕΙΩΝ ΑΠΟΦΑΣΗΣ

Μέθοδοι ελαχίστων τετραγώνων – μικρά ως προς το MSE

– **ΜΙΚΡΟΤΕΡΑ**, ως προς το κριτήριο του μέσου τετραγωνικού σφάλματος (**mean square error - MSE**), σημαίνει επιλογή του \mathcal{W} ώστε η συνάρτηση κόστους

- $J(w) \equiv E[(d - w^T x)^2]$ να ελαχιστοποιηθεί
 $\hat{w} = \arg \min_w J(w)$

όπου d οι αντίστοιχες επιθυμητές αποκρίσεις

Ελαχιστοποιώντας την $J(w)$ ως προς το w έχουμε:

$$\frac{\partial J(w)}{\partial w} = \frac{\partial}{\partial w} E[(d - w^T x)^2] = 0$$
$$= 2E[x(d - x^T w)] \Rightarrow$$

R_x : Πίν. αυτοσυσχέτισης - **autocorrelation matrix**

$E[xd]$: Διάν. Ετεροσυσχέτ.-**cross-correlation vector**

$$E[xx^T]w = E[xd] \Rightarrow \hat{w} = R_x^{-1}E[xd]$$

$$E[xd] = \begin{bmatrix} E[x_1 d] \\ \dots \\ E[x_l d] \end{bmatrix}$$

$$R_x \equiv E[xx^T] = \begin{bmatrix} E[x_1 x_1] & E[x_1 x_2] \dots & E[x_1 x_l] \\ \dots & \dots & \dots \\ E[x_l x_1] & E[x_l x_2] \dots & E[x_l x_l] \end{bmatrix}$$

ΓΡΑΜΜΙΚΟΙ ΠΑΡΑΜΕΤΡΙΚΟΙ ΤΑΞΙΝΟΜΗΤΕΣ ΕΠΙΦΑΝΕΙΩΝ ΑΠΟΦΑΣΗΣ

Μέθοδοι ελαχίστων τετραγώνων – μικρά ως προς το MSE: Περ. από 2 κλάσεις

Στην περίπτωση που το υπό μελέτη πρόβλημα περιλαμβάνει M κλάσεις, λύνουμε M προβλήματα σαν το παραπάνω. Συγκεκριμένα:

-Για την j -στή κλάση

- Θέσε $d_i=1$, αν $\mathbf{x}_i \in \omega_j$ και $d_i=0$, διαφορετικά.

- Λύσε το πρόβλημα δύο κλάσεων που προκύπτει και έστω \mathbf{w}^j το αντίστοιχο διάνυσμα παραμέτρων.

Μετά τον προσδιορισμό των \mathbf{w}^j 's:

-Για δεδομένο \mathbf{x}_i :

- Υπολόγισε τις ποσότητες $g^j(\mathbf{x}_i) = \mathbf{w}^j T \mathbf{x}_i$, $j=1, \dots, M$.

- Καταχώρησε το \mathbf{x}_i στην κλάση ω_q για την οποία $g^q(\mathbf{x}_i) = \max_{j=1, \dots, M} g^j(\mathbf{x}_i)$

Σημείωση: Το κριτήριο MSE ανήκει σε μια γενικότερη κλάση συναρτήσεων κόστους με την ακόλουθη σημαντική ιδιότητα:

Η τιμή $g^j(\mathbf{x}_i)$ είναι μια εκτίμηση, ως προς το κριτήριο MSE, της εκ των υστέρων πιθανότητας, $P(\omega_j | \mathbf{x})$, υπό την προϋπόθεση ότι οι επιθυμητές αποκρίσεις που χρησιμοποιούνται κατά την εκπαίδευση είναι $d_i=1$, αν $\mathbf{x}_i \in \omega_j$ και $d_i=0$, διαφορετικά.

ΓΡΑΜΜΙΚΟΙ ΠΑΡΑΜΕΤΡΙΚΟΙ ΤΑΞΙΝΟΜΗΤΕΣ ΕΠΙΦΑΝΕΙΩΝ ΑΠΟΦΑΣΗΣ

Μέθοδοι ελαχίστων τετραγώνων – μικρά ως προς το SSE

- **ΜΙΚΡΟΤΕΡΟ**, ως προς το κριτήριο του **αθροίσματος των τετραγώνων των σφαλμάτων**, σημαίνει επιλογή του \mathcal{W} που **ελαχιστοποιεί** τη συνάρτηση κόστους

$$\bullet J(\mathbf{w}) = \sum_{i=1}^N (d_i - \mathbf{w}^T \mathbf{x}_i)^2$$

όπου $d_i = +1$, αν $\mathbf{x}_i \in \omega_1$ και $d_i = -1$, αν $\mathbf{x}_i \in \omega_2$.

Minimizing $J(\mathbf{w})$ with respect to \mathbf{w} results in:

$$\frac{\partial J(\mathbf{w})}{\partial \mathbf{w}} = \frac{\partial}{\partial \mathbf{w}} \sum_{i=1}^N (d_i - \mathbf{w}^T \mathbf{x}_i)^2 = 0 \Rightarrow$$

$$\left(\sum_{i=1}^N \mathbf{x}_i \mathbf{x}_i^T \right) \mathbf{w} = \sum_{i=1}^N \mathbf{x}_i d_i$$

ΓΡΑΜΜΙΚΟΙ ΠΑΡΑΜΕΤΡΙΚΟΙ ΤΑΞΙΝΟΜΗΤΕΣ ΕΠΙΦΑΝΕΙΩΝ ΑΠΟΦΑΣΗΣ

Μέθοδοι ελαχίστων τετραγώνων – μικρά ως προς το SSE

Ένας εναλλακτικός τρόπος διατύπωσης:

Ορίζουμε τα ακόλουθα

$$X = \begin{bmatrix} x_1^T \\ x_2^T \\ \dots \\ x_N^T \end{bmatrix} \text{ } N \times l \text{ πίνακας} \quad \underline{d} = \begin{bmatrix} d_1 \\ \dots \\ d_N \end{bmatrix} \text{ Αντίστοιχες επιθυμητές αποκρίσεις}$$

$$X^T = [x_1, x_2, \dots, x_N] \text{ } l \times N \text{ πίνακας}$$

$$X^T X = \sum_{i=1}^N x_i x_i^T$$

$$X^T d = \sum_{i=1}^N x_i d_i$$

ΓΡΑΜΜΙΚΟΙ ΠΑΡΑΜΕΤΡΙΚΟΙ ΤΑΞΙΝΟΜΗΤΕΣ ΕΠΙΦΑΝΕΙΩΝ ΑΠΟΦΑΣΗΣ

Μέθοδοι ελαχίστων τετραγώνων – μικρά ως προς το SSE

Ένας εναλλακτικός τρόπος διατύπωσης:

Τότε

$$\left(\sum_{i=1}^N x_i^T x_i \right) \hat{w} = \left(\sum_{i=1}^N x_i d_i \right) \Leftrightarrow (X^T X) \hat{w} = X^T d \Rightarrow \hat{w} = (X^T X)^{-1} X^T d = X^\# d$$

$$X^\# \equiv (X^T X)^{-1} X^T \quad \text{Ψευδοαντίστροφος (pseudoinverse) του } X$$

Έστω $N=I$. Τότε ο X είναι τετραγωνικός και (γενικά) αντιστρέψιμος. Τότε έχουμε

$$(X^T X)^{-1} X^T = X^{-1} X^{-T} X^T = X^{-1} \Rightarrow X^\# = X^{-1}$$

ΓΡΑΜΜΙΚΟΙ ΠΑΡΑΜΕΤΡΙΚΟΙ ΤΑΞΙΝΟΜΗΤΕΣ ΕΠΙΦΑΝΕΙΩΝ ΑΠΟΦΑΣΗΣ

Μέθοδοι ελαχίστων τετραγώνων – μικρά ως προς το SSE

Ένας εναλλακτικός τρόπος διατύπωσης:

Έστω $N > l$. Τότε, γενικά, **δεν υπάρχει λύση** που να ικανοποιεί ταυτόχρονα όλους τους περιορισμούς

$$Xw = d : \begin{array}{l} x_1^T w = d_1 \\ x_2^T w = d_2 \\ \dots \\ x_N^T w = d_N \end{array} \quad N \text{ equations} > l \text{ unknowns}$$

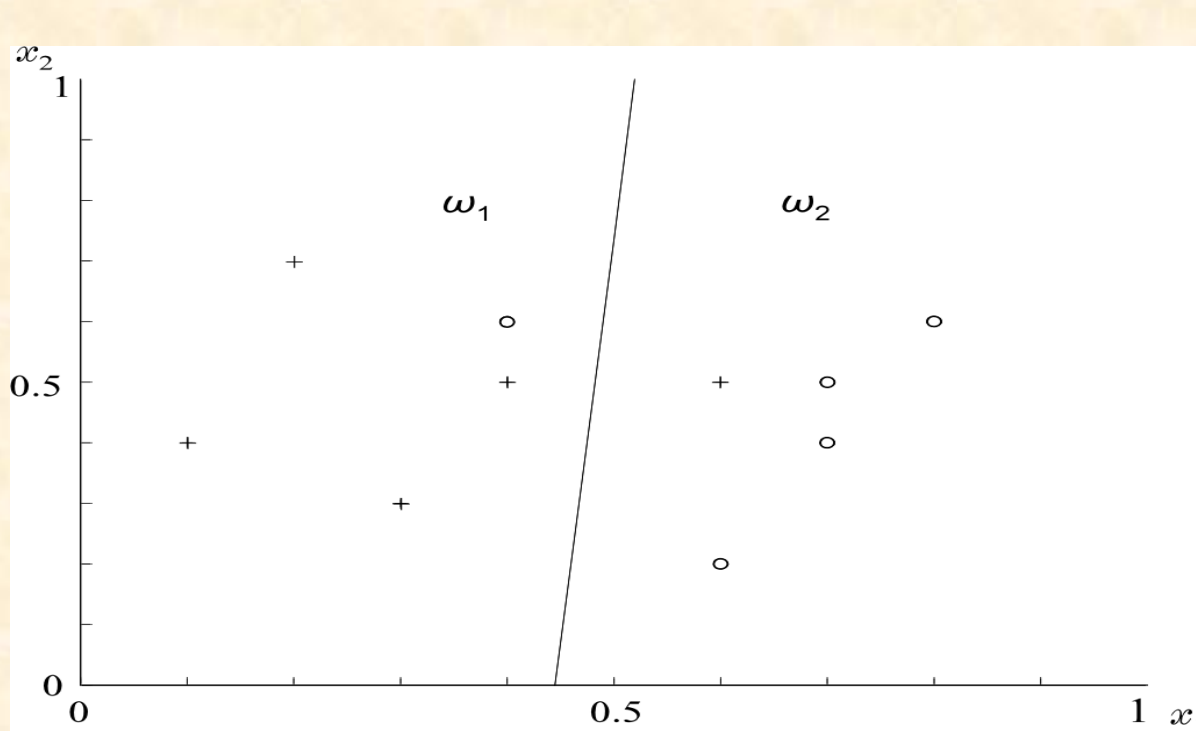
Σ' αυτή την περίπτωση, η **“λύση”** $w = X^\# d$ ελαχιστοποιεί το **άθροισμα των τετραγώνων των σφαλμάτων**.

ΓΡΑΜΜΙΚΟΙ ΠΑΡΑΜΕΤΡΙΚΟΙ ΤΑΞΙΝΟΜΗΤΕΣ ΕΠΙΦΑΝΕΙΩΝ ΑΠΟΦΑΣΗΣ

Μέθοδοι ελαχίστων τετραγώνων – μικρά ως προς το SSE

Ένα παράδειγμα:

$$\omega_1 : \begin{bmatrix} 0.4 \\ 0.5 \end{bmatrix}, \begin{bmatrix} 0.6 \\ 0.5 \end{bmatrix}, \begin{bmatrix} 0.1 \\ 0.4 \end{bmatrix}, \begin{bmatrix} 0.2 \\ 0.7 \end{bmatrix}, \begin{bmatrix} 0.3 \\ 0.3 \end{bmatrix} \quad \omega_2 : \begin{bmatrix} 0.4 \\ 0.6 \end{bmatrix}, \begin{bmatrix} 0.6 \\ 0.2 \end{bmatrix}, \begin{bmatrix} 0.7 \\ 0.4 \end{bmatrix}, \begin{bmatrix} 0.8 \\ 0.6 \end{bmatrix}, \begin{bmatrix} 0.7 \\ 0.5 \end{bmatrix}$$



$$X = \begin{bmatrix} 0.4 & 0.5 & 1 \\ 0.6 & 0.5 & 1 \\ 0.1 & 0.4 & 1 \\ 0.2 & 0.7 & 1 \\ 0.3 & 0.3 & 1 \\ 0.4 & 0.6 & 1 \\ 0.6 & 0.2 & 1 \\ 0.7 & 0.4 & 1 \\ 0.8 & 0.6 & 1 \\ 0.7 & 0.5 & 1 \end{bmatrix} \quad \underline{d} = \begin{bmatrix} 1 \\ 1 \\ 1 \\ 1 \\ 1 \\ -1 \\ -1 \\ -1 \\ -1 \\ -1 \end{bmatrix}$$

$$X^T X = \begin{bmatrix} 2.8 & 2.24 & 4.8 \\ 2.24 & 2.41 & 4.7 \\ 4.8 & 4.7 & 10 \end{bmatrix}, \quad X^T \underline{d} = \begin{bmatrix} -1.6 \\ 0.1 \\ 0.0 \end{bmatrix} \quad \underline{w} = (X^T X)^{-1} X^T \underline{d} = \begin{bmatrix} -3.13 \\ 0.24 \\ 1.34 \end{bmatrix}$$

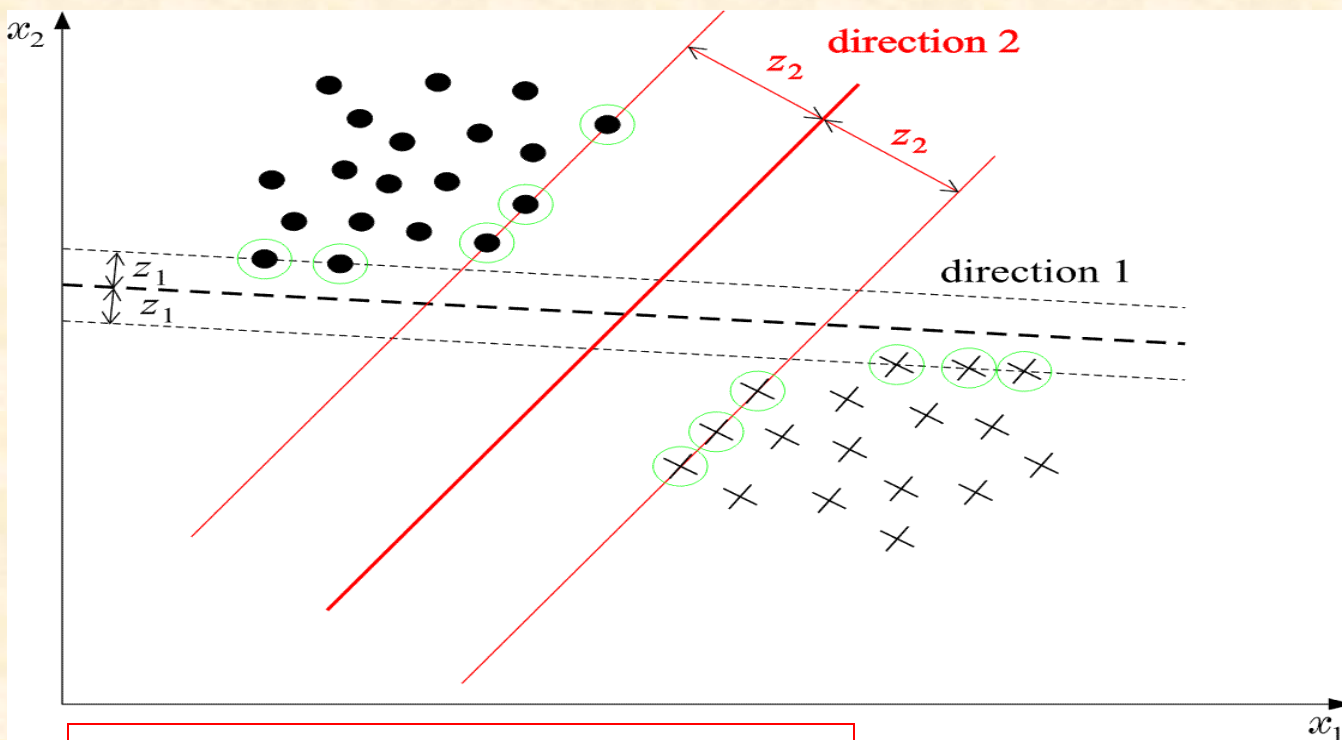
ΓΡΑΜΜΙΚΟΙ ΠΑΡΑΜΕΤΡΙΚΟΙ ΤΑΞΙΝΟΜΗΤΕΣ ΕΠΙΦΑΝΕΙΩΝ ΑΠΟΦΑΣΗΣ

Support vector machines – η γραμμική περίπτωση – γραμ. διαχ. κλάσεις

Ο στόχος: Δοθέντων δύο **γραμμικώς διαχωρίσιμων** κλάσεων, προσδιόρισε τον ταξινομητή

$$g(x) = w^T x + w_0 = 0$$

που αφήνει το **μέγιστο περιθώριο** από τις δύο κλάσεις.



Support vector machines:
Μηχανές διανυσματικής στήριξης

ΓΡΑΜΜΙΚΟΙ ΠΑΡΑΜΕΤΡΙΚΟΙ ΤΑΞΙΝΟΜΗΤΕΣ ΕΠΙΦΑΝΕΙΩΝ ΑΠΟΦΑΣΗΣ

Support vector machines – η γραμμική περίπτωση – γραμ. διαχ. κλάσεις

- Περιθώριο (Margin): Κάθε υπερεπίπεδο χαρακτηρίζεται από
 - Τον **προσανατολισμό** του στο χώρο, δηλ. το διάνυσμα \mathbf{w}
 - Τη **θέση** του χώρο, i.e., w_0
- Για **ΚΑΘΕ** προσανατολισμό, \mathbf{w} , επέλεξε το υπερεπίπεδο που **αφήνει την ΙΔΙΑ απόσταση** από τα **εγγύτερα** σημεία από κάθε κλάση. Το περιθώριο είναι το διπλάσιο αυτής της απόστασης.

ΓΡΑΜΜΙΚΟΙ ΠΑΡΑΜΕΤΡΙΚΟΙ ΤΑΞΙΝΟΜΗΤΕΣ ΕΠΙΦΑΝΕΙΩΝ ΑΠΟΦΑΣΗΣ

Support vector machines – η γραμμική περίπτωση – γραμ. διαχ. κλάσεις

– Η απόσταση ενός σημείου \hat{x} από ένα υπερεπίπεδο είναι

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

– Διαβάθμισε τα w , w_0 , ώστε για τα **εγγύτερα σημεία** από κάθε κλάση να είναι:

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

– Τότε το **περιθώριο** είναι

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

– Επίσης, ισχύουν και οι ακόλουθοι περιορισμοί

$$\begin{aligned} w^T x + w_0 &\geq 1 \quad \forall x \in \omega_1 \\ w^T x + w_0 &\leq -1 \quad \forall x \in \omega_2 \end{aligned}$$

(δηλ. όλα τα στοιχεία της κλάσης +1 (-1) βρίσκονται στη θετική (αρνητική) πλευρά του υπερεπιπέδου)

ΓΡΑΜΜΙΚΟΙ ΠΑΡΑΜΕΤΡΙΚΟΙ ΤΑΞΙΝΟΜΗΤΕΣ ΕΠΙΦΑΝΕΙΩΝ ΑΠΟΦΑΣΗΣ

Support vector machines – η γραμμική περίπτωση – γραμ. διαχ. κλάσεις

– SVM (γραμμικός) ταξινομητής

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

– (Πρόβλημα SVM) Ελαχιστοποίησε την

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

– Υπό τους περιορισμούς

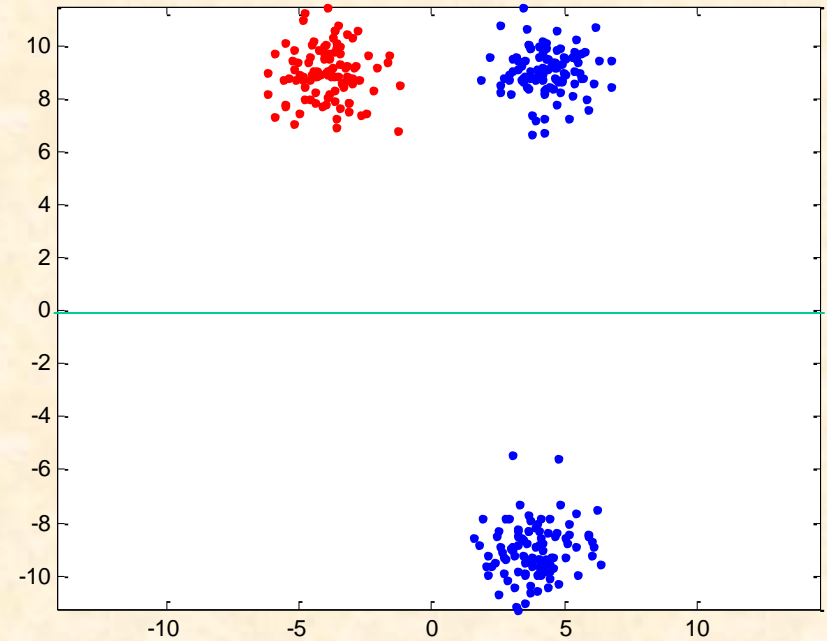
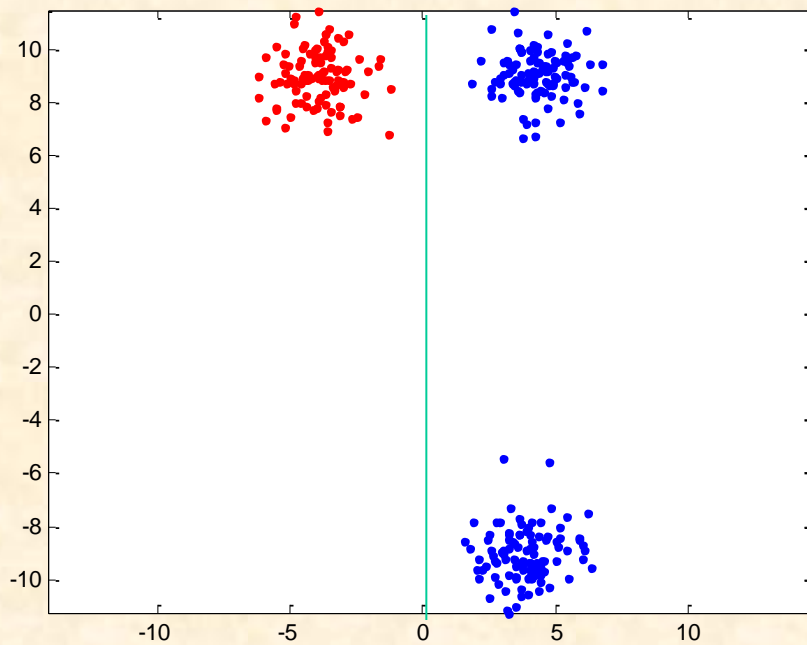
$$d_i(\underline{w}^T \underline{x}_i + w_0) \geq 1, \quad i = 1, 2, \dots, N \quad \begin{array}{l} d_i = 1, \text{ for } \underline{x}_i \in \omega_1, \\ d_i = -1, \text{ for } \underline{x}_i \in \omega_2 \end{array}$$

– Τα παραπάνω δικαιολογούνται από το γεγονός ότι ελαχιστοποιώντας το $\|\underline{w}\|^2$

μεγιστοποιείται το περιθώριο $\frac{2}{\|\underline{w}\|}$.

ΓΡΑΜΜΙΚΟΙ ΠΑΡΑΜΕΤΡΙΚΟΙ ΤΑΞΙΝΟΜΗΤΕΣ ΕΠΙΦΑΝΕΙΩΝ ΑΠΟΦΑΣΗΣ

Support vector machines – η γραμμική περίπτωση – γραμ. διαχ. κλάσεις



Οι επιπλέον περιορισμοί ισότητας αποθαρρύνουν λύσεις σαν αυτή του 2^{ου} σεναρίου πιο πάνω.

ΒΕΛΤΙΣΤΟΠΟΙΗΣΗ ΣΥΝΑΡΤΗΣΗΣ ΥΠΟ ΓΡΑΜ. ΑΝΙΣΟΤΙΚΟΥΣ ΠΕΡΙΟΡΙΣΜΟΥΣ

Έστω το πρόβλημα

$$\text{Min } J(\vartheta)$$

Υπό τους περιορισμούς $f_i(\vartheta) \geq 0, i=1, \dots, m.$

Ορίζουμε τη **συνάρτηση Lagrange**

$$L(\vartheta, \lambda) = J(\vartheta) - \sum \lambda_i f_i(\vartheta)$$

Συνθήκες KKT (Karush-Kuhn-Tacker): Για τη θέση του **ελαχίστου** ισχύουν τα ακόλουθα

$$\frac{\partial}{\partial \vartheta} L(\vartheta, \lambda) = 0$$

$$\lambda_i \geq 0, i = 1, 2, \dots, m$$

$$\lambda_i f_i(\vartheta) = 0, i = 1, 2, \dots, m$$

ΒΕΛΤΙΣΤΟΠΟΙΗΣΗ ΣΥΝΑΡΤΗΣΗΣ ΥΠΟ ΓΡΑΜ. ΑΝΙΣΟΤΙΚΟΥΣ ΠΕΡΙΟΡΙΣΜΟΥΣ

Έστω

$$L^*(\vartheta) = \max_{\lambda} L(\vartheta, \lambda) = \max_{\lambda} (J(\vartheta) - \sum_i \lambda_i f_i(\vartheta))$$

Αφού $\lambda_i \geq 0$, $f_i(\vartheta) \geq 0$ έχουμε

$$L^*(\vartheta) = J(\vartheta)$$

Άρα

$$\min_{\vartheta} J(\vartheta) = \min_{\vartheta} L^*(\vartheta) = \min_{\vartheta} \max_{\lambda \geq 0} L(\vartheta, \lambda)$$

Θεώρημα:

- Έστω (α) η $J(\vartheta)$ είναι **κυρτή** και (β) οι $f_i(\vartheta)$ είναι **γραμμικές**.
- Έστω ϑ^* μία θέση ελαχίστου για το πρόβλημα ελαχιστοποίησης και λ^* το αντίστοιχο διάνυσμα των πολ/στών Lagrange.
- Τότε το (ϑ^*, λ^*) είναι ένα σαγματικό σημείο (saddle point) της συνάρτησης Lagrange, για το οποίο ισχύει

$$L(\vartheta^*, \lambda^*) = \min_{\vartheta} \max_{\lambda \geq 0} L(\vartheta, \lambda) = \max_{\lambda \geq 0} \min_{\vartheta} L(\vartheta, \lambda)$$

Συνεπώς, σύμφωνα με το θεώρημα, για τον προσδιορισμό του ελαχίστου **δεν παίζει ρόλο η σειρά με την οποία βελτιστοποιούμε την $L(\vartheta, \lambda)$ ως προς τα ϑ και λ .**

Στο SVM πρόβλημα θα βελτιστοποιήσουμε την $L(\cdot)$ πρώτα ως προς ϑ και μετά ως προς λ .

ΓΡΑΜΜΙΚΟΙ ΠΑΡΑΜΕΤΡΙΚΟΙ ΤΑΞΙΝΟΜΗΤΕΣ ΕΠΙΦΑΝΕΙΩΝ ΑΠΟΦΑΣΗΣ

Support vector machines – η γραμμική περίπτωση – γραμ. διαχ. κλάσεις

– Το πρόβλημα SVM είναι μια διαδικασία (κυρτής) τετραγωνικής βελτιστοποίησης (quadratic optimization task), με γραμμικούς περιορισμούς. Οι συνθήκες Karush-Kuhn-Tucker ορίζουν ότι για τη θέση του ελαχίστου ισχύουν τα ακόλουθα:

- (1) $\frac{\partial}{\partial w} L(w, w_0, \lambda) = 0$
- (2) $\frac{\partial}{\partial w_0} L(w, w_0, \lambda) = 0$
- (3) $\lambda_i \geq 0, i = 1, 2, \dots, N$
- (4) $\lambda_i [d_i (w^T x_i + w_0) - 1] = 0, i = 1, 2, \dots, N$
- Όπου $L(\bullet, \bullet, \bullet)$ είναι η συνάρτηση Lagrange

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

ΓΡΑΜΜΙΚΟΙ ΠΑΡΑΜΕΤΡΙΚΟΙ ΤΑΞΙΝΟΜΗΤΕΣ ΕΠΙΦΑΝΕΙΩΝ ΑΠΟΦΑΣΗΣ

Support vector machines – η γραμμική περίπτωση – γραμ. διαχ. κλάσεις

– **Η λύση:** από τα παραπάνω προκύπτει ότι είναι

$$w = \sum_{i=1}^N \lambda_i d_i x_i \quad \sum_{i=1}^N \lambda_i d_i = 0$$

Αντικαθιστώντας τα παραπάνω στην $L(\cdot)$, τα λ_i 's εκτιμώνται ως οι λύσεις του δυϊκού SVM προβλήματος:

(Δυϊκό SVM)

Μεγιστοποίησε την $\left(\sum_{i=1}^N \lambda_i - \frac{1}{2} \sum_{ij} \lambda_i \lambda_j d_i d_j x_i^T x_j \right)$ ως προς το $\lambda = [\lambda_1, \dots, \lambda_N]^T$

Έτσι ώστε $\sum_{i=1}^N \lambda_i d_i = 0$
 $\underline{\lambda} \geq \underline{0}$

ΓΡΑΜΜΙΚΟΙ ΠΑΡΑΜΕΤΡΙΚΟΙ ΤΑΞΙΝΟΜΗΤΕΣ ΕΠΙΦΑΝΕΙΩΝ ΑΠΟΦΑΣΗΣ

Support vector machines – η γραμμική περίπτωση – γραμ. διαχ. κλάσεις

Παρατηρήσεις:

- Οι πολ/στές Lagrange λ_i είναι είτε **θετικοί** είτε **μηδέν**. Έτσι

$$w = \sum_{i=1}^{N_s} \lambda_i d_i x_i$$

όπου N_s είναι ο αριθμός των διανυσμάτων με θετικά λ_i 's.

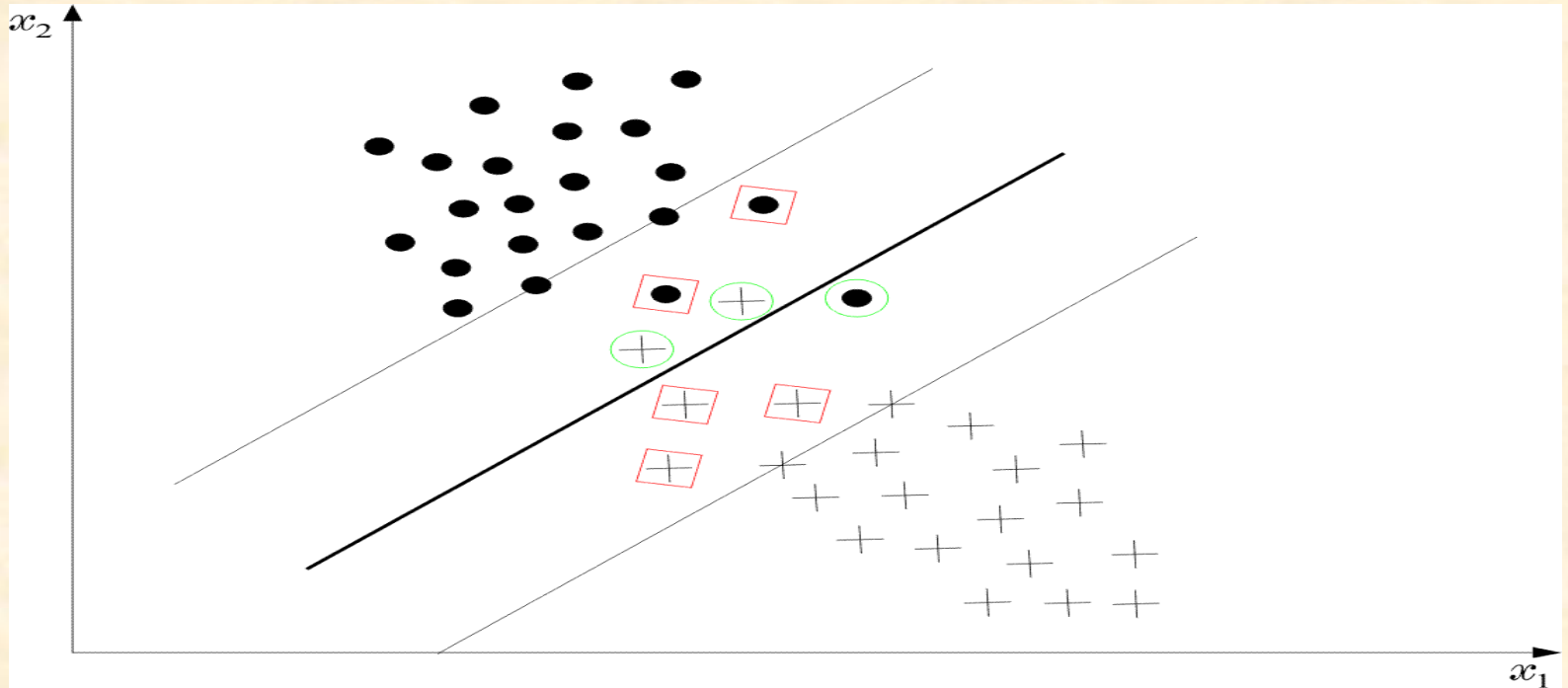
- Θετικά λ_i έχουν τα διανύσματα που ικανοποιούν τη συνθήκη $w^T x_i + w_0 = \pm 1$ λόγω των περιορισμών $\lambda_i [d_i (w^T x_i + w_0) - 1] = 0, \quad i = 1, 2, \dots, N$

Τα διανύσματα αυτά ονομάζονται **ΔΙΑΝΥΣΜΑΤΑ ΣΤΗΡΙΞΗΣ (SUPPORT VECTORS)** και είναι τα εγγύτερα προς το υπερεπίπεδο του ταξινομητή διανύσματα από κάθε κλάση. Αυτά είναι που καθορίζουν το w .

- Μετά τον προσδιορισμό του w , το w_0 προσδιορίζεται από τις συνθήκες (4).
- Το βέλτιστο υπερεπίπεδο ως προς το κριτήριο της μεγιστοποίησης του περιθωρίου είναι **ΜΟΝΑΔΙΚΟ**.
- Παρότι η λύση είναι μοναδική, οι πολ. Lagrange **δεν είναι απαραίτητα** μοναδικοί.

ΓΡΑΜΜΙΚΟΙ ΠΑΡΑΜΕΤΡΙΚΟΙ ΤΑΞΙΝΟΜΗΤΕΣ ΕΠΙΦΑΝΕΙΩΝ ΑΠΟΦΑΣΗΣ

Support vector machines – η γραμμική περίπτωση – μη γραμ. διαχ. κλάσεις



ΓΡΑΜΜΙΚΟΙ ΠΑΡΑΜΕΤΡΙΚΟΙ ΤΑΞΙΝΟΜΗΤΕΣ ΕΠΙΦΑΝΕΙΩΝ ΑΠΟΦΑΣΗΣ

Support vector machines – η γραμμική περίπτωση – μη γραμ. διαχ. κλάσεις

- Στην περίπτωση αυτή, δεν υπάρχει υπερεπίπεδο έτσι ώστε

$$w^T x + w_0 (><)1, \forall x$$

- Υπενθυμίζεται ότι το **περιθώριο** ορίζεται ως το **διπλάσιο της απόστασης** μεταξύ των ακόλουθων δύο υπερεπιπέδων

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

– Για τα **διανύσματα εκπαίδευσης** έχουμε **ένα** από τα ακόλουθα **τρία** δυνατά σενάρια

1) Διανύσματα **εκτός** της ζώνης των δύο επιπέδων που είναι **ορθά** ταξινομημένα

$$d_i(w^T x + w_0) > 1$$

2) Διανύσματα **εντός** της ζώνης των δύο επιπέδων που είναι **ορθά** ταξινομημένα

$$0 \leq d_i(w^T x + w_0) < 1$$

3) Διανύσματα τα οποία είναι **μη ορθά ταξινομημένα**

$$d_i(w^T x + w_0) < 0$$

ΓΡΑΜΜΙΚΟΙ ΠΑΡΑΜΕΤΡΙΚΟΙ ΤΑΞΙΝΟΜΗΤΕΣ ΕΠΙΦΑΝΕΙΩΝ ΑΠΟΦΑΣΗΣ

Support vector machines – η γραμμική περίπτωση – μη γραμ. διαχ. κλάσεις

- Οι παραπάνω περιπτώσεις μπορούν να παρασταθούν με συμπαγή τρόπο ως εξής

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

Όπου για

1^ο σενάριο $\rightarrow \xi_i = 0$

2^ο σενάριο $\rightarrow 0 < \xi_i \leq 1$

3^ο σενάριο $\rightarrow 1 < \xi_i$

Οι μεταβλητές ξ_i είναι γνωστές ως **μεταβλητές χαλαρότητας (slack variables)**

ΓΡΑΜΜΙΚΟΙ ΠΑΡΑΜΕΤΡΙΚΟΙ ΤΑΞΙΝΟΜΗΤΕΣ ΕΠΙΦΑΝΕΙΩΝ ΑΠΟΦΑΣΗΣ

Support vector machines – η γραμμική περίπτωση – μη γραμ. διαχ. κλάσεις

–Ο στόχος της βελτιστοποίησης είναι τώρα διττός

- Μεγιστοποίηση του **περιθωρίου**

- Ελαχιστοποίηση του **αριθμού των διανυσμάτων** με $\xi_i > 0$,

- Ένας τρόπος να επιτύχουμε τον παραπάνω στόχο είναι μέσω της ακόλουθης συνάρτησης κόστους

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

όπου C είναι μια **σταθερά** και

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

- Η συνάρτηση $I(\cdot)$ **ΔΕΝ** είναι **διαφορίσιμη**. Στην πράξη, μπορούμε να

χρησιμοποιήσουμε μια **προσέγγιση**

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

ΓΡΑΜΜΙΚΟΙ ΠΑΡΑΜΕΤΡΙΚΟΙ ΤΑΞΙΝΟΜΗΤΕΣ ΕΠΙΦΑΝΕΙΩΝ ΑΠΟΦΑΣΗΣ

Support vector machines – η γραμμική περίπτωση – μη γραμ. διαχ. κλάσεις

- Η **συνάρτηση Lagrange** για την περίπτωση αυτή γίνεται

$$L(w, w_0, \xi, \lambda, \mu) \equiv \frac{1}{2} w^T w + C \sum_{i=1}^N \xi_i - \sum_{i=1}^N \mu_i \xi_i - \sum_{i=1}^N \lambda_i [d_i (w^T x_i + w_0) - 1 + \xi_i]$$

(τα λ_i και μ_i είναι πολ/στές Lagrange)

- Οι αντίστοιχες **ΚΚΤ συνθήκες** είναι:

$$(1) \quad w = \sum_{i=1}^N \lambda_i d_i x_i$$

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

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

$$(4) \quad \lambda_i [y_i (w^T x_i + w_0) - 1 + \xi_i] = 0, \quad i = 1, 2, \dots, N$$

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

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

➤ Αντικαθιστώντας τα παραπάνω στην $L(\cdot)$, οι πολ/στές Lagrange προσδιορίζονται ως οι λύσεις του ακόλουθου δυϊκού SVM προβλήματος

Μεγιστοποίηση $\lambda_{\geq 0} \left(\sum_{i=1}^N \lambda_i - \frac{1}{2} \sum_{ij} \lambda_i \lambda_j d_i d_j x_i^T x_j \right)$ ως προς το $\lambda = [\lambda_1, \dots, \lambda_N]^T$

υπό τις προϋποθέσεις

$$0 \leq \lambda_i \leq C, \quad i = 1, 2, \dots, N$$
$$\sum_{i=1}^N \lambda_i d_i = 0$$

➤ Σχόλιο:

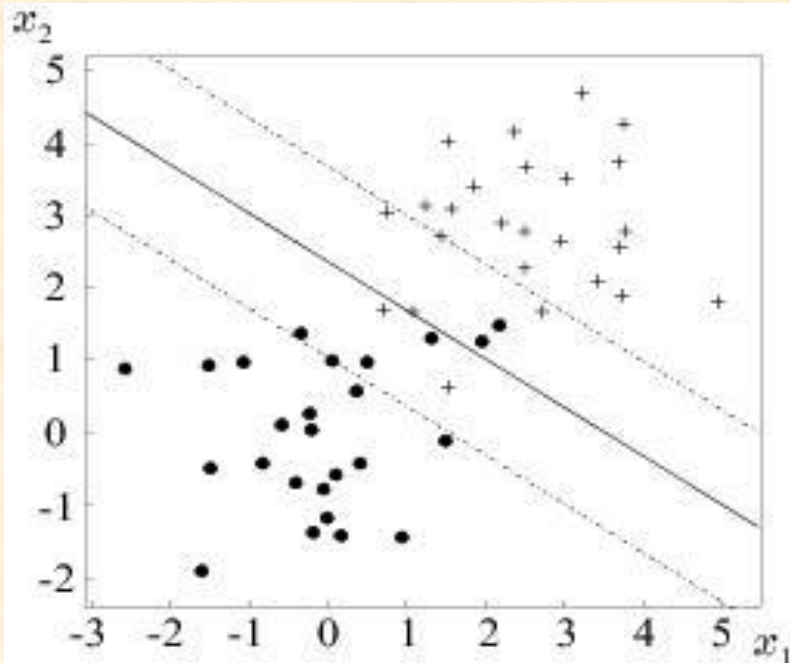
Η μόνη διαφορά με την περίπτωση των διαχωρίσιμων κλάσεων είναι η ύπαρξη του C στους περιορισμούς.

ΓΡΑΜΜΙΚΟΙ ΠΑΡΑΜΕΤΡΙΚΟΙ ΤΑΞΙΝΟΜΗΤΕΣ ΕΠΙΦΑΝΕΙΩΝ ΑΠΟΦΑΣΗΣ

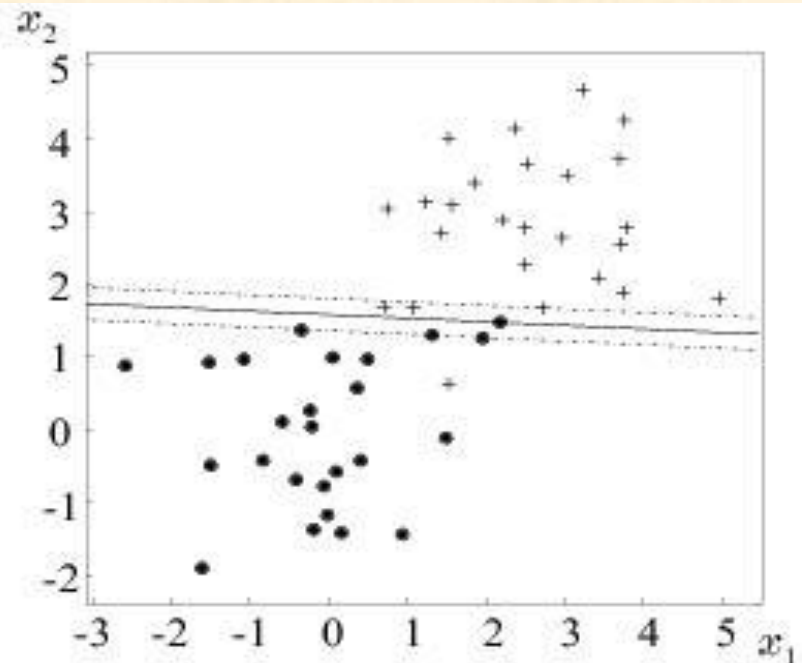
Support vector machines – η γραμμική περίπτωση – μη γραμ. διαχ. κλάσεις

Η λύση του SVM προβλήματος ακολουθεί βήματα ανάλογα με αυτά της προηγούμενης περίπτωσης.

Ωστόσο, στην παρούσα περίπτωση, η παράμετρος C επηρεάζει την επιλογή της τελικής λύσης.



(a)

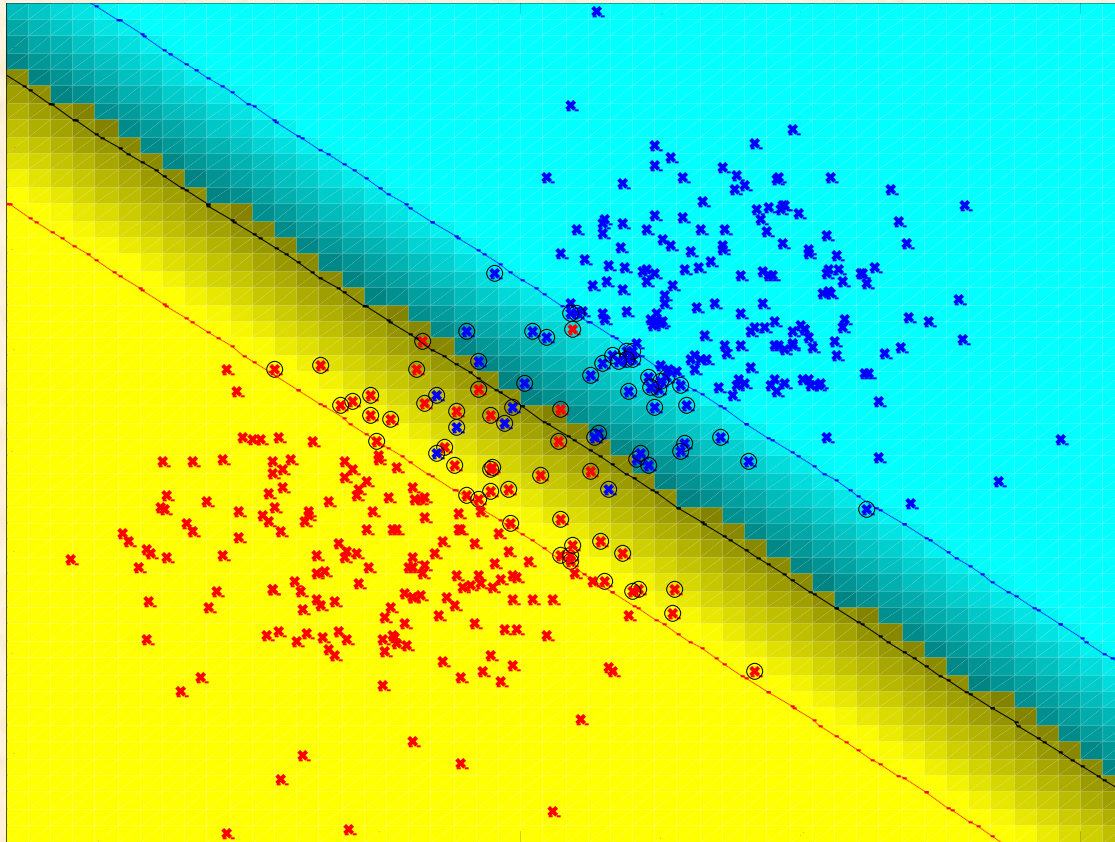


(b)

Στο παραπάνω παράδειγμα η C έχει “μικρότερη τιμή” για το σχήμα (a) και “μεγαλύτερη τιμή” για το σχήμα (b).

ΓΡΑΜΜΙΚΟΙ ΠΑΡΑΜΕΤΡΙΚΟΙ ΤΑΞΙΝΟΜΗΤΕΣ ΕΠΙΦΑΝΕΙΩΝ ΑΠΟΦΑΣΗΣ

Support vector machines – η γραμμική περίπτωση – μη γραμ. διαχ. κλάσεις



C = 0.1

Pe_tr = 0.0225

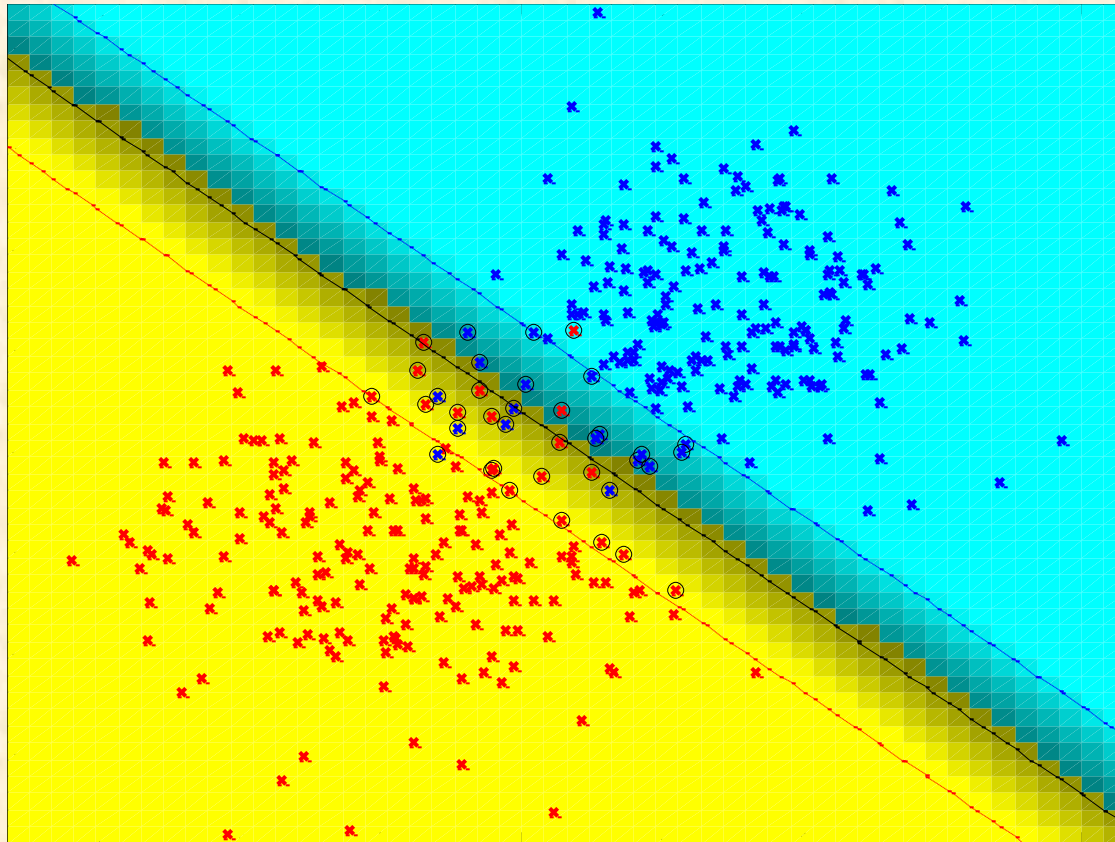
Pe_te = 0.0325

sup_vec = 82

marg = 0.9412

ΓΡΑΜΜΙΚΟΙ ΠΑΡΑΜΕΤΡΙΚΟΙ ΤΑΞΙΝΟΜΗΤΕΣ ΕΠΙΦΑΝΕΙΩΝ ΑΠΟΦΑΣΗΣ

Support vector machines – η γραμμική περίπτωση – μη γραμ. διαχ. κλάσεις



C = 1

Pe_tr = 0.0225

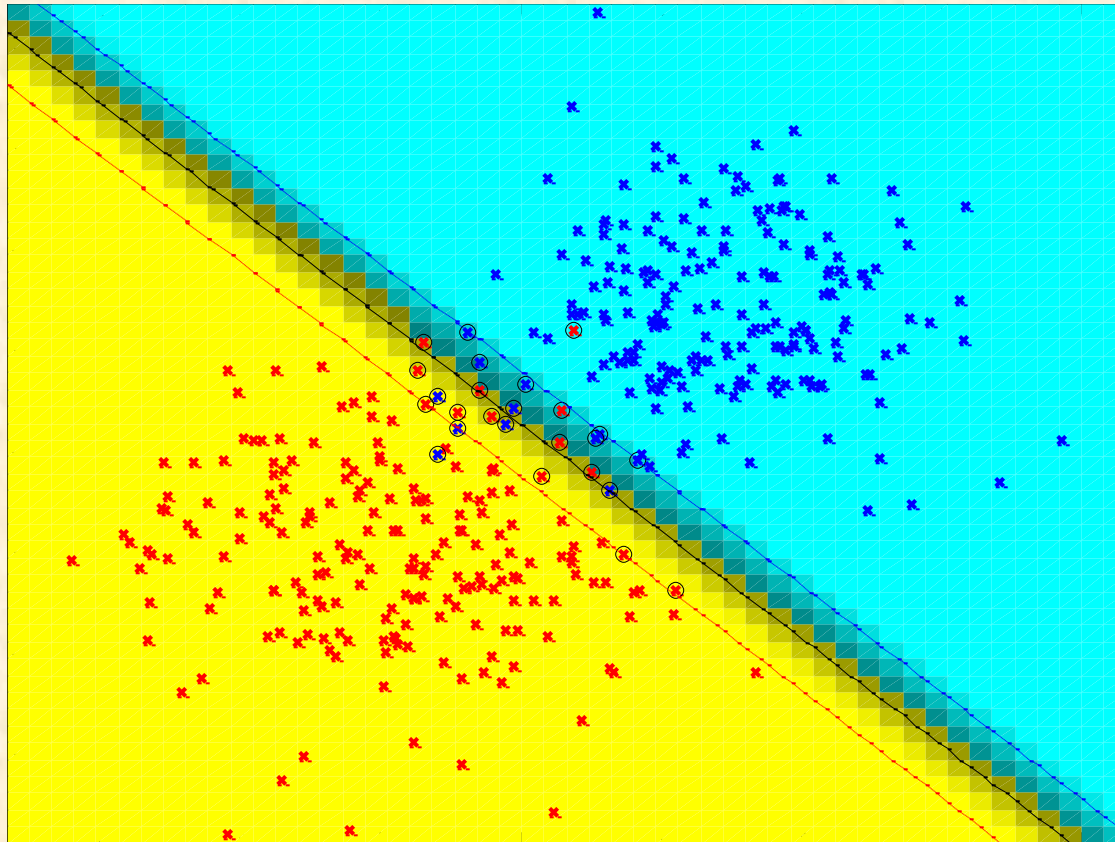
Pe_te = 0.0325

sup_vec = 37

marg = 0.6317

ΓΡΑΜΜΙΚΟΙ ΠΑΡΑΜΕΤΡΙΚΟΙ ΤΑΞΙΝΟΜΗΤΕΣ ΕΠΙΦΑΝΕΙΩΝ ΑΠΟΦΑΣΗΣ

Support vector machines – η γραμμική περίπτωση – μη γραμ. διαχ. κλάσεις



C =20

Pe_tr =0.0225

Pe_te =0.0350

sup_vec =25

marg = 0.3573