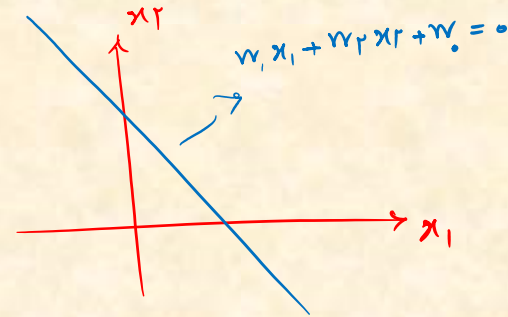# LINEAR CLASSIFIERS

در این فصل : طراحی طبقه بندهای خطی صرفنظر از توزیع‌های توصیف کننده داده‌های آموزشی      برتری : ساده‌تر، کم‌هزینه‌تر

❖ The Problem: Consider a two class task with $\omega_1, \omega_2$

معادله خطی

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

$w_1 x_1 + w_2 x_2 + ... + w_l x_l + w_0 = 0$

$w_1 x_1 + w_2 x_2 + w_0 = 0$

$x_2$

$x_1$

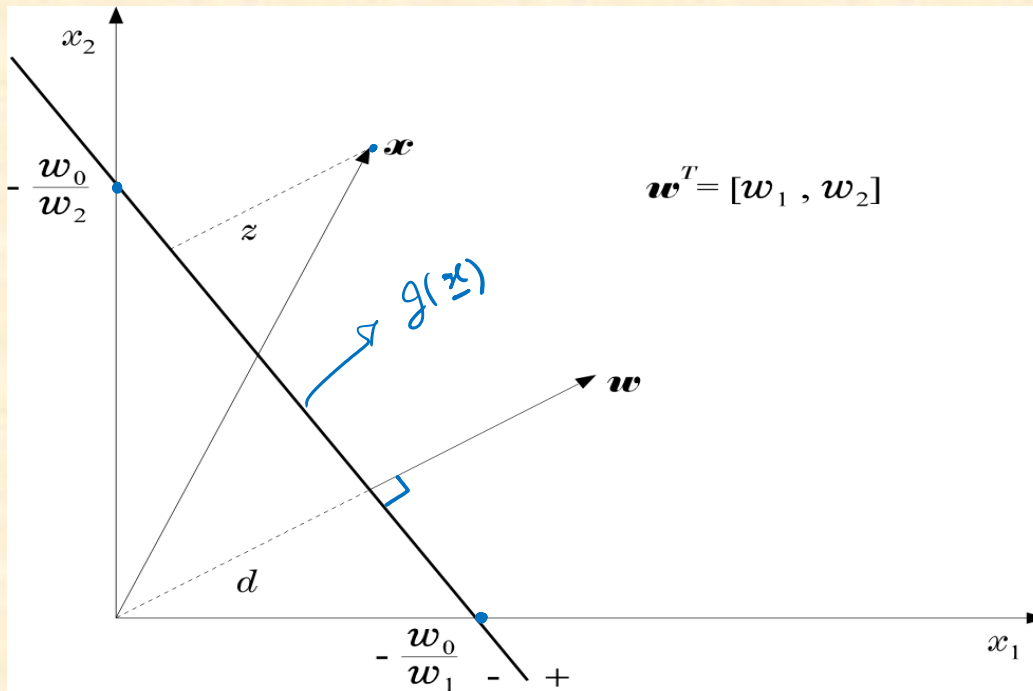➢ Assume $\underline{x}_1, \underline{x}_2$ on the decision hyperplane :

$0 = \underline{w}^T \underline{x}_1 + w_0 = \underline{w}^T \underline{x}_2 + w_0 \Rightarrow$

$\underline{w}^T (\underline{x}_1 - \underline{x}_2) = 0 \quad \forall \underline{x}_1, \underline{x}_2$     بردار وزن‌ها $\underline{w}$ بر ابر صفحه تقسیم عمود است.

1

> Hence:

$$\boxed{\underline{w} \perp \text{ on the hyperplane}}$$

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



$W_1 X_1 + W_2 X_2 + W_0 = 0$

$w^T = [w_1 , w_2]$

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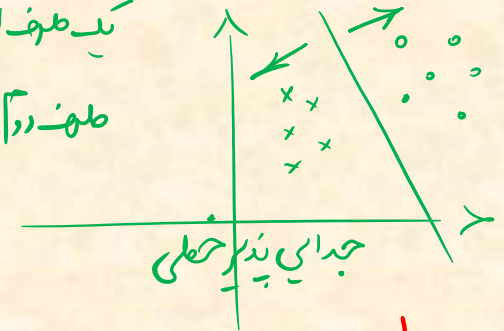
2

❖ The Perceptron Algorithm

الگوریتم Perceptron :

فرض : فطه‌ها جدای پذیر خطی باشند.

جدای پذیر خطی

➤ Assume linearly separable classes, i.e.,

$$\exists \underline{w}^*: \ w^{*T} \underline{x} > 0 \ \forall \underline{x} \in \omega_1$$ یک طرف ابرصفحه

$$\underline{w}^{*T} \underline{x} < 0 \ \forall \underline{x} \in \omega_2$$ طرف دوم ابرصفحه

جدای پذیر خطی

➤ The case $\underline{w}^{*T} \underline{x} + \boxed{w_0^*}$
falls under the above formulation, since

جدای پذیر خطی نیست ✗

• $\underline{w}' \equiv \begin{bmatrix} \underline{w}^* \\ w_0^* \end{bmatrix}, \ \underline{x}' = \begin{bmatrix} \underline{x} \\ 1 \end{bmatrix}$ → $\underline{w}^{*T} \underline{x} + w_0^* = 0$

• $\underline{w}^{*T} \underline{x} + w_0^* = \underline{w}'^T \underline{x}' = 0$

3

هدف : راه حلی برای به دست آوردن ابر صفحهٔ $\underline{w}$ که دو صنف را از هم جدا کنند.

➢ Our goal:  Compute a solution, i.e., a hyperplane $\underline{w}$, so that

$$\underline{w}^T \underline{x} (><)0 \quad \underline{x} \in \quad \begin{array}{c} \omega_1 \\ \\ \omega_2 \end{array}$$

• The steps
   – Define a cost function to be minimized
   – Choose an algorithm to minimize the cost function
   – The minimum corresponds to a solution

➢ The Cost Function

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

تابع هزینه perceptron:

- Where $Y$ is the subset of the vectors wrongly classified by $\underline{w}$. When $Y$=(empty set) a solution is achieved and

۲: زیر مجموعه ای از بردارها که به اشتباه ۰ طبقه بندی شده اند.

وقتی $\phi = Y$ شده به جواب رسیده ایم:

$Y = \phi \rightarrow$

- $J(\underline{w}) = 0$

$\underline{w}^T \underline{x} < 0$

- $\begin{cases} \delta_x = -1 \text{ if } \underline{x} \in Y \text{ and } \underline{x} \in \omega_1 \\ \delta_x = +1 \text{ if } \underline{x} \in Y \text{ and } \underline{x} \in \omega_2 \end{cases}$

$\underline{w}^T \underline{x} > 0$

$\underline{w}^T \underline{x} > 0 \rightarrow x \in \omega_1$

$\underline{w}^T \underline{x} < 0 \rightarrow x \in \omega_2$
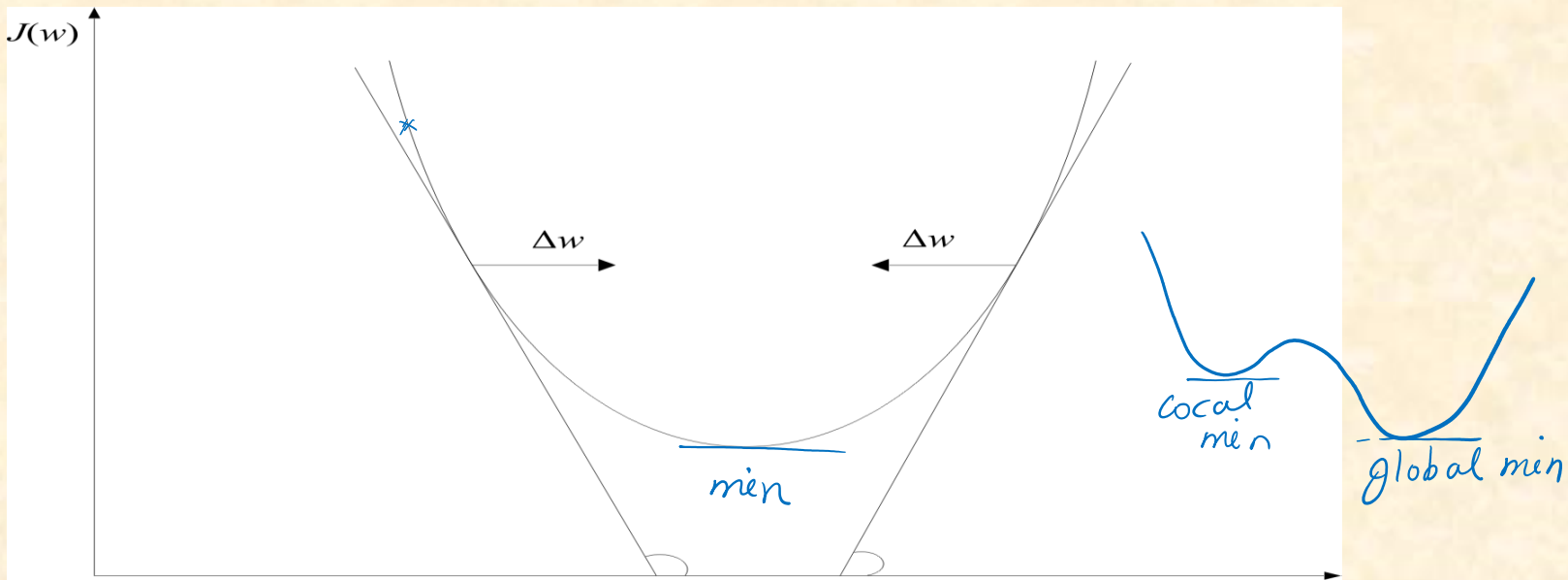
- $J(\underline{w}) \geq 0$

5

- $J(\underline{w})$ is piecewise linear (WHY?)



➢ The Algorithm
  - The philosophy of the gradient descent is adopted.

گرادیان نزولی

$$\underline{w}(\text{new}) = \underline{w}(\text{old}) + \Delta\underline{w}$$

$$\Delta\underline{w} = -\mu \frac{\partial J(\underline{w})}{\partial \underline{w}} \Big|_{\underline{w} = \underline{w}(\text{old})}$$

$$\Delta\underline{w} = -\mu \frac{\partial J(\underline{w})}{\partial\underline{w}} \Big|_{\underline{w} = \underline{w}_{old}}$$

پارامتر تنظیم سرعت همگرایی

- Wherever valid

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

$$\underline{w}_{new} = \underline{w}_{old} + \Delta\underline{w}$$

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

ρ : پارامتر تنظیم
دنباله‌های از اعداد حقیقی مثبت

This is the celebrated Perceptron Algorithm

7

تفسیر هندسی الگوریتم Peraptron :

صفحه جدا کننده w(t)

w(t+1)

$w(t+1) = w(t) + \Delta w$

$\Delta w$ در جهت بردار $\underline{x}$

وزن بهینه جدا کننده $w^*$

**FIGURE 3.2**

Geometric interpretation of the perceptron algorithm. The update of the weight vector is in the direction of $x$ in order to turn the decision hyperplane to include $x$ in the correct class.

8

**FIGURE 3.3**

An example of the perceptron algorithm. After the update of the weight vector, the hyperplane is turned from its initial location (dotted line) to the new one (full line), and all points are correctly classified.

- مقدار دهی اولیه $\underline{w}(0)$ به صورت تصادفی
- انتخاب $\rho_0$
- $t = 0$
- حلقهٔ تکرار
  - $\gamma = \phi$
  - for $i = 1$ to $\underline{N}$
    - if $\delta_{x_i} w(t)^T x_i \geq 0$ then $\gamma = \gamma \cup \{\underline{x_i}\}$
  - end

حلقهٔ for برای تشخیص نمونه های به اشتباه طبقه بندی شده با مدلورژن های $w(t)$

$$\boxed{w(t+1) = w(t) - \rho_t \sum_{x \in \gamma} \delta_x \, \underline{x}}$$

- تنظیم $\rho_t$
- $t = t + 1$
- تا زمانی که $\gamma = \phi$

➤ An example:



$$\underline{w}(t+1) = \underline{w}(t) + \rho_t \underline{x}$$

$$= \underline{w}(t) - \rho_t \delta_x \underline{x} \quad (\delta_x = -1)$$

همگرایی :

➤ The perceptron algorithm **converges** in a **finite** number of iteration steps to a solution if

$$\lim_{t \to \infty} \sum_{k=0}^{t} \rho_k \to \infty, \qquad \lim_{t \to \infty} \sum_{k=0}^{t} \rho_k^2 < +\infty$$

$$\text{e.g.,:} \; \rho_t = \frac{c}{t}$$

min

11

❖ A useful variant of the perceptron algorithm

فرم‌های دیگرالگوریتم پیریرن: ‒ ⌘

⇒ Reward & Punishment تنبیه و پاداش

$$\underline{w}(t+1) = \underline{w}(t) + \rho \underline{x}_{(t)} , \begin{cases} \underline{w}^T(t)\underline{x}_{(t)} \le 0 \\ \\ \underline{x}_{(t)} \in \omega_1 \end{cases}$$

طبقه‌بندی به اشتباه

$$\underline{w}(t+1) = \underline{w}(t) - \rho \underline{x}_{(t)} , \begin{cases} \underline{w}^T(t)\underline{x}_{(t)} \ge 0 \\ \\ \underline{x}_{(t)} \in \omega_2 \end{cases}$$

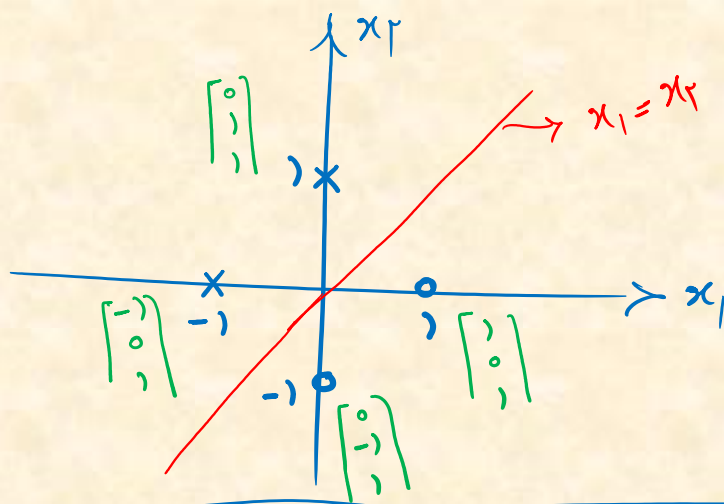$$\underline{w}(t+1) = \underline{w}(t) \quad \text{otherwise} \longrightarrow$$ طبقه‌بندی به درستی

نمونه‌ها را یک به یک به الگوریتم اعمال می‌کنیم ؛ اگر پس از پایان یافتن نمونه‌ها ، الگوریتم همگرا نشد ، دوباره نمونه‌ها را اعمال خواهیم کرد.

➢ It is a ⟦reward and punishment⟧ type of algorithm

مثال 3.2 :

x : ω₁
0 : ω₂

$x_1 = x_2$



طبقه‌بندی خطی با استفاده از الگوریتم پرسپترون به فرم با دانش رشته همگرا کنید.

$$\frac{\omega_1}{\omega_2} \gtrless \underline{w}^T \underline{x} \gtrless 0$$

حل : $\rho = 1$ ، $\underline{w}(0) = [0, 0, 0]$ تصادفی

دور اول:

$\begin{bmatrix} -1 \\ 0 \\ 1 \end{bmatrix}$ : $\underline{w}^T(0)\begin{bmatrix} -1 \\ 0 \\ 1 \end{bmatrix} = 0 \xrightarrow{\times} \underline{w}(1) = \underline{w}(0) + 1 \cdot \begin{bmatrix} -1 \\ 0 \\ 1 \end{bmatrix} = \begin{bmatrix} -1 \\ 0 \\ 1 \end{bmatrix}$

$\begin{bmatrix} 0 \\ 1 \end{bmatrix}$ : $\underline{w}^T(1)\begin{bmatrix} 0 \\ 1 \\ 1 \end{bmatrix} = 1 > 0 \xrightarrow{\checkmark} \underline{w}(2) = \underline{w}(1)$

$\begin{bmatrix} 0 \\ -1 \end{bmatrix}$ : $\omega_2 \leftarrow \quad \underline{w}^T(2)\begin{bmatrix} 0 \\ -1 \\ 1 \end{bmatrix} = 1 > 0 \xrightarrow{\times} \underline{w}(3) = \underline{w}(2) - 1 \cdot \begin{bmatrix} 0 \\ -1 \\ 1 \end{bmatrix} = \begin{bmatrix} -1 \\ 1 \\ 0 \end{bmatrix}$

$\begin{bmatrix} 1 \\ \hat{x}_2 \end{bmatrix}$ : $\omega_2 \leftarrow \quad \underline{w}^T(3)\begin{bmatrix} 1 \\ 0 \\ 1 \end{bmatrix} = -1 < 0 \xrightarrow{\checkmark} \underline{w}(4) = \underline{w}(3)$

$\begin{bmatrix} -1 \\ \cdots \end{bmatrix}$ : $\underline{w}^T(4)\begin{bmatrix} -1 \\ 0 \\ 1 \end{bmatrix} = 1 > 0 \xrightarrow{\checkmark} \underline{w}(5) = \underline{w}(4)$

$\cdots \longrightarrow$ همگرا شده است

$\underline{w}(7) = \underline{w}(6) = \begin{bmatrix} -1 \\ 1 \\ 0 \end{bmatrix}$

$\underline{w}^T \underline{x} = 0 \rightarrow [-1 \ 1 \ 0]\begin{bmatrix} x_1 \\ x_2 \\ 1 \end{bmatrix}$

همگرا $\Rightarrow -x_1 + x_2 = 0$
$x_1 = x_2$

❖ The perceptron

واحد پایه‌ای شبکه‌های عصبی (نورون)



$w^T \underline{x} + w_0 \gtrless 0$

$\omega_1$

$\omega_2$

$w_1 x_1 + w_2 x_2 + \dots + w_0 \gtrless 0$

تابع فعال‌سازی

$f(y) = \begin{cases} -1 & \text{if } y < 0 \\ +1 & \text{if } y > 0 \end{cases}$

از روی مدل شبکه عصبی انسان

$\omega_1$
$\vee$
$\vee$
$\omega_2$

$w_i's$    synapses or synaptic weights

$w_0$      threshold

➢ The network is called perceptron or neuron
➢ It is a learning machine that learns from the training vectors via the perceptron algorithm
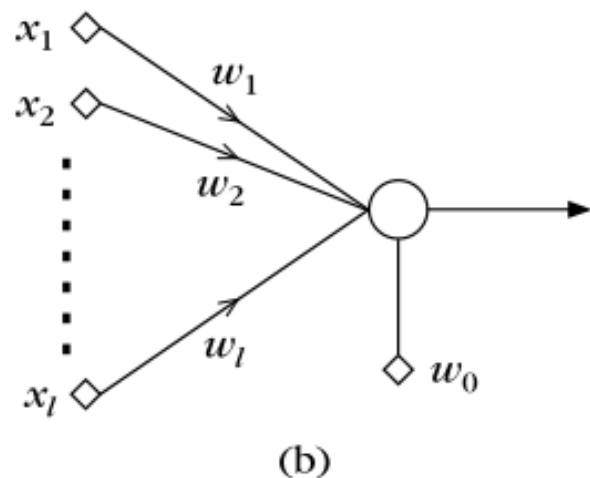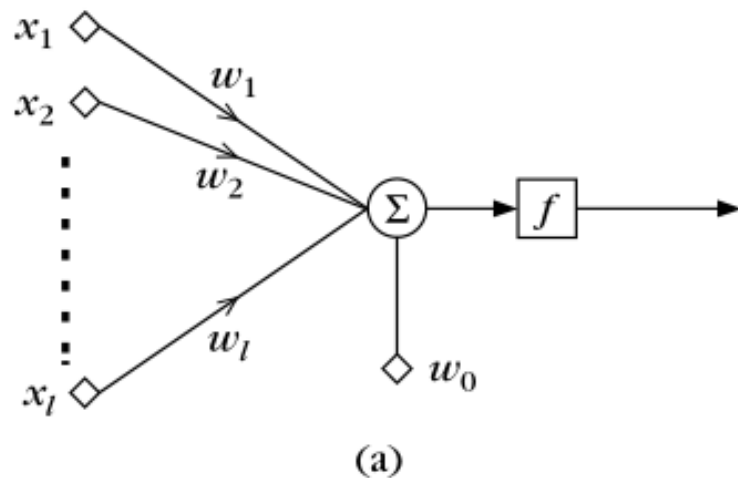
14

**FIGURE 3.5**

The basic perceptron model. (a) A linear combiner is followed by the activation function. (b) The combiner and the activation function are merged together.

➢ Example: At some stage $t$ the perceptron algorithm results in

$w(t)$ → $w_1 = 1, \; w_2 = 1, \; w_0 = -0.5$

$x_1 + x_2 - 0.5 = 0$

The corresponding hyperplane is

$\begin{bmatrix} -0.2 \\ 0.75 \end{bmatrix}$

$\rho = 0.7$

$1.42\,x_1 + 0.51\,x_2 - 0.5 = 0$

$\omega_1$

$\omega_2$

$\begin{bmatrix} 0.4 \\ 0.05 \end{bmatrix}$

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

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

16

الگوریتم Pocket :    (برای حالت دسته‌هایی که جدایی‌پذیر خطی نیستند)

شرط اصلی الگوریتم پرسپترون جدایی‌پذیر بودن کلاس‌هاست. اگر این فرض برقرار نباشد که در عمل معمولاً این چنین است،
الگوریتم پرسپترون همگرا نخواهد شد. فرآهایی دیگر الگوریتم که به راه حل بهینه منجر شوند در این مورد راستا به کار می‌شود.
یکی از آنها الگوریتم Pocket است.

دوگام الگوریتم Pocket :                                      تعداد نمونه‌هایی که به درستی طبقه‌بندی
$W_S = W(0)$                     شده‌اند.

$h_s = 0$  ۱- مقدار دهی اولیه بردار وزن $W(0)$ به صورت تصادفی و ذخیره آن با عنوان $W_S$ ، مقدار دهی شمارنده

۲- در شروع گام t ، با استفاده از قانون پرسپترون $W(t+1)$ را بر می‌کنیم، با استفاده از این $W(t)$ وزن،
بردار ها آموزشی راست می‌کنیم و تعداد بردارهایی که درست طبقه بندی شده‌اند را $h$ می‌نامیم. اگر $h > h_s$
باشد، $W_S$ را با $W(t+1)$ جایگزین می‌کنیم، $h_s$ را با $h$ جایگزین می‌کنیم.
$h_s = h$                                       $W_S = W(t+1)$

17

تاکنون حالت دو کلاسه را بررسی کردیم . تعمیم به حالت $\underline{M}$ کلاسه سرراست است :

تابع جداکننده خطی $w_j$ ; $i = 1, \ldots, M$ و $i \neq j$ برای هر یک از کلاس ها تعریف شده است .

بردار نمونه ها $\underline{x}$ با ابعاد $(\ell+1)$ بعدی به صورت زیر در کلاس $i$ ام طبقه بندی می شود اگر :

$$\underline{w}_i^T \underline{x} > \underline{w}_j^T \underline{x} \quad ; \quad \forall j \neq i$$

این توما به ساختار kesler منجر می شود :

برای هر یک از بردارهای آموزشی متعلق به کلاس $i$ ام که $i = 1, 2, \ldots, M$ می باشد ، $\boxed{(M-1)}$ بردار آموزشی به صورت زیر تعریف می کنیم :

$$\underline{x}_{ij} = [\underline{0}^T, \underline{0}^T, \ldots, \underline{x}^T, \ldots, -\underline{x}^T, \ldots, \underline{0}^T]^T$$
$$(\ell+1)M \times 1$$

همه بلوک ها بردار صفر هستند غیر از بلوک $i$ ام که برابر $\underline{x}^T$ و بلوک $j$ ام که برابر $-\underline{x}^T$ است .

برردار وزن ها :
$$\underline{w} = [\underline{w}_1^T, \ldots, \underline{w}_M^T]^T$$

با این تعریف :
$$\forall j = 1, 2, \ldots, M \atop i \neq j \quad ; \quad \underline{x} \in \omega_i \longrightarrow \underline{w}^T \underline{x}_{ij} > 0$$

$$\frac{w_1}{w_2} \underline{w}^T \underline{x}_{12} \gtrless 0$$

$$\frac{w_1}{w_3} \underline{w}^T \underline{x}_{13} \gtrless 0$$

بنابر این مسأله به طراحی طبقه بند خطی در فضای $M(\ell+1)$ بعدی تبدیل شد .  تعداد نمونه ها : $N(M-1)$

18

مثال : مسأله ۳ کلاسه در فضای دو بُعدی :   بُردارهای آموزشی :

$\omega_1$ : $\boxed{[1,1]^T}$ , $[2,2]^T$ , $[2,1]^T$

$\omega_2$ : $[1,-1]^T$ , $\boxed{[1,-2]^T}$ , $[2,-2]^T$

$\omega_3$ : $[-1,1]^T$ , $[-1,2]^T$ , $\boxed{[-2,1]^T}$

مسأله به جدایی پذیر خطی است.

ساختار kesler :

بُردار وزن‌ها :

$$\begin{cases} \underline{w}_1 = [w_{11}, w_{12}, w_{10}]^T \\ \underline{w}_2 = [w_{21}, w_{22}, w_{20}]^T \\ \underline{w}_3 = [w_{31}, w_{32}, w_{30}]^T \end{cases}$$

$$\rightarrow \underline{w} = [w_1^T, w_2^T, w_3^T]^T$$

با ۱۸ بُردار آموزشی، چون در $\underline{w}$ الگوریتم پرسپترون را اجرا کنیم :

$\underline{w}_1 = [0.13, 3.60, 1.00]^T$

$\underline{w}_2 = [-0.05, -3.16, -0.41]^T$

$\underline{w}_3 = [-3.84, 1.28, 0.69]^T$

$$\frac{(\ell+1)}{3} \underset{3}{M} \times 1 = 9 \times 1$$

$$\begin{cases} \underline{x} = [1,1,1]^T \\ \omega_1 \rightarrow j=1 \end{cases}$$

$$x_{12} = [\underbrace{1,1,1}_{x^T}, \underbrace{-1,-1,-1}_{-x^T}, 0,0,0]^T_{9\times1}$$

$$x_{13} = [\underbrace{1,1,1}_{x^T}, 0,0,0, \underbrace{-1,-1,-1}_{-x^T}]^T_{9\times1}$$

$$\begin{cases} \underline{x} = [1,-2,1]^T \\ \omega_2 \rightarrow j=2 \end{cases}$$

$$x_{21} = [\underbrace{-1,2,-1}_{-x^T}, \underbrace{1,-2,1}_{x^T}, 0,0,0]^T_{9\times1}$$

$$x_{23} = [0,0,0, \underbrace{1,-2,1}_{x^T}, \underbrace{-1,2,-1}_{-x^T}]^T_{9\times1}$$

$$\begin{cases} \underline{x} = [-2,1,1]^T \\ \omega_3 \rightarrow j=3 \end{cases}$$

$$x_{31} = [2,-1,-1, 0,0,0, -2,1,1]^T_{9\times1}$$

$$x_{32} = [0,0,0, 2,-1,-1, -2,1,1]^T_{9\times1}$$

جدائیت طبقه‌بند های خطی در آنهاست بنابراین در بیاری از موارد با وجو اینکه فرض جدای پذیر خطی برقرار نیست . از این طبقه‌بندها استفاده می‌شود . در این موقع طبقه‌بندهای خطی بهترین حل زیر بهینه (sub optimal) منتر می‌شود .

مسئله دو کلاسه :    خروجی    $y = \pm 1$

حل با استفاده از خطاها

بردار ورودی $\underline{x}$ ———— طبقه‌بند $\underline{w}^T \underline{x}$

NSE : Mean Square Error

$J(\underline{w}) = E\left[ |y - \underline{w}^T \underline{x}|^2 \right]$    تابع هزینه

NSE بین خروجی های واقعی و مطلوب
desired    true

خروجی واقعی    که خروجی مطلوب
$y(\underline{x}) = \hat{y} = \pm 1$

minimize $\longrightarrow$ $\hat{\underline{w}} = \underset{\underline{w}}{\arg\min} \, J(\underline{w})$

$3.28 \longrightarrow J(\underline{w}) = P(\omega_1) \int (1 - \underline{x}^T \underline{w})^2 p(\underline{x}|\omega_1) \, d\underline{x} + P(\omega_2) \int (1 + \underline{x}^T \underline{w})^2 p(\underline{x}|\omega_2) d\underline{x}$

$\underbrace{\qquad\qquad\qquad\qquad}_{y = +1}$    $\underbrace{\qquad\qquad\qquad\qquad}_{y = -1}$

$\frac{\partial J(\underline{w})}{\partial \underline{w}} = 0 \longrightarrow 2E\left[ \underline{x}(y - \underline{x}^T \underline{w}) \right] = 0 \rightarrow \hat{\underline{w}} = R_x^{-1} E[\underline{x}y]$ [20]

❖ **Least Squares Methods**

➢ If classes are linearly separable, the perceptron output results in $\pm 1$

➢ If classes are <u>NOT</u> linearly separable, we shall compute the weights $w_1, w_2, ..., w_0$

so that the difference between

- The actual output of the classifier, $\underline{w}^T \underline{x}$, and

- The desired outputs, e.g.

$$+1 \text{ if } \underline{x} \in \omega_1$$
$$-1 \text{ if } \underline{x} \in \omega_2$$

to be SMALL

➤ SMALL, in the mean square error sense, means to choose $\underline{w}$ so that the cost function

- $J(\underline{w}) \equiv E[(y - \underline{w}^T \underline{x})^2]$ is minimum
- $\hat{\underline{w}} = \arg\min_{\underline{w}} J(\underline{w})$
- $y$ the corresponding desired responses

➢ Minimizing

$J(\underline{w})$ w.r. to $\underline{w}$ results in :

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

$$= 2E[\underline{x}(y - \underline{x}^T \underline{w})] \Rightarrow$$

$$E[\underline{x}\underline{x}^T]\underline{w} = E[\underline{x}y] \Rightarrow$$

$$\boxed{\hat{\underline{w}} = R_x^{-1} E[\underline{x}y]}$$

where $R_x$ is the autocorrelation matrix

$$R_x \equiv E[\underline{x}\underline{x}^T] = \begin{bmatrix} E[x_1 x_1] & E[x_1 x_2]... & E[x_1 x_l] \\ ........... & ............... & ........... \\ E[x_l x_1] & E[x_l x_2]... & E[x_l x_l] \end{bmatrix}$$

and $E[\underline{x}y] = \begin{bmatrix} E[x_1 y] \\ ... \\ E[x_l y] \end{bmatrix}$ the crosscorrelation vector
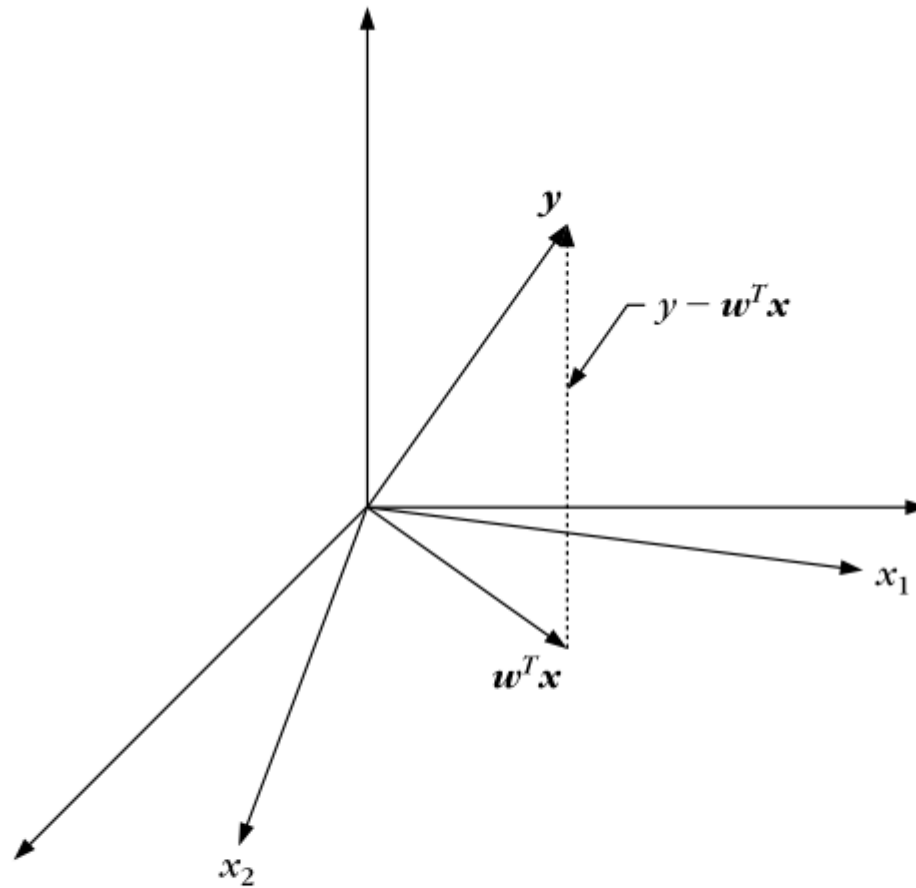
23

**FIGURE 3.6**

Interpretation of the MSE estimate as an orthogonal projection on the input vector elements' subspace.

➤ Multi-class generalization

    • The goal is to compute $\underline{M}$ linear discriminant functions:

$$g_i(\underline{x}) = \underline{w}_i^T \underline{x} \qquad \text{;} \; i = 1, 2, \ldots, M$$

    according to the MSE.

    • Adopt as desired responses $y_i$: خروجی های مطلوب

$$y_i = 1 \quad \text{if} \quad \underline{x} \in \omega_i$$
$$y_i = 0 \quad \text{otherwise}$$

    • Let

$$\underline{y} = \left[ y_1, y_2, \ldots, y_M \right]^T$$

    • And the matrix

ماتریس وزن

$$W = \left[ \underline{w}_1, \underline{w}_2, \ldots, \underline{w}_M \right]$$

تعمیم چند کلاسه

- The goal is to compute $W$:

$$\hat{W} = \arg\min_W E\left[\left\|\underline{y} - W^T \underline{x}\right\|^2\right] = \arg\min_W E\left[\sum_{i=1}^{M}\left(y_i - \underline{w}_i^T \cdot \underline{x}\right)^2\right]$$

حل این M تا مسئله کمینه سازی MSE

- The above is equivalent to a number $M$ of MSE minimization problems. That is:

  Design each $\underline{w}_i$ so that its desired output is 1 for $\underline{x} \in \omega_i$ and 0 for any other class.

➤ Remark: The MSE criterion belongs to a more general class of cost function with the following important property:

- The value of $g_i(\underline{x})$ is an estimate, in the MSE sense, of the a-posteriori probability $P(\omega_i \mid \underline{x})$, provided that the desired responses used during training are $y_i = 1, \underline{x} \in \omega_i$ and 0 otherwise.

26

➤ **Mean square error regression**: Let $y \in \mathfrak{R}^M$, $\underline{x} \in \mathfrak{R}^\ell$ be jointly distributed random vectors with a joint pdf $p(\underline{x}, \underline{y})$

- The goal: Given the value of $\underline{x}$ estimate the value of $\underline{y}$. In the pattern recognition framework, given $\underline{x}$ one wants to estimate the respective label $y = \pm 1$.

- The MSE estimate $\hat{\underline{y}}$ of $\underline{y}$ given $\underline{x}$ is defined as:

$$\hat{\underline{y}} = \arg \min_{\tilde{y}} E\left[ \left\| y - \tilde{y} \right\|^2 \right]$$

- It turns out that:

$$\hat{\underline{y}} = E\left[ \underline{y} \mid \underline{x} \right] \equiv \int_{-\infty}^{+\infty} \underline{y} p(\underline{y} \mid \underline{x}) d\underline{y}$$

The above is known as the regression of $\underline{y}$ given $\underline{x}$ and it is, in general, a non-linear function of $\underline{x}$. If $p(\underline{x}, \underline{y})$ is Gaussian the MSE regressor is linear.

# Sum of Error Squares Estimation

معیاری نزدیک به MSE ، بعنوان مجموع مربعات خطا یا LS    (Least Squares)

$$J(\underline{w}) = \sum_{i=1}^{N} (y_i - \underline{x}_i^T \underline{w})^2 \equiv \sum_{i=1}^{N} e_i^2 \qquad (\text{هزینه})$$

در MSE میانگین خطاها در LS مجموع خطا

در MSE برای $E\{\cdot\}$ به pdf نیاز است ولی در LS این نیاز مطرح نمی‌شود.

$$\frac{\partial J(\underline{w})}{\partial \underline{w}} = 0 \quad \text{Minimizing} \qquad \sum_{i=1}^{N} \underline{x}_i (y_i - \underline{x}_i^T \underline{\hat{w}}) = 0 \longrightarrow \left( \sum_{i=1}^{N} \underline{x}_i \underline{x}_i^T \right) \hat{w} = \sum_{i=1}^{N} \underline{x}_i y_i$$

ماتریس $X = \begin{bmatrix} \underline{x}_1^T \\ \underline{x}_2^T \\ \vdots \\ \underline{x}_N^T \end{bmatrix} = \begin{bmatrix} x_{11} & x_{12} & \cdots & x_{1\ell} \\ x_{21} & x_{22} & \cdots & x_{2\ell} \\ \vdots & & \ddots & \vdots \\ x_{N1} & x_{N2} & & x_{N\ell} \end{bmatrix}_{N \times \ell}$ , $\underline{y} = \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_N \end{bmatrix}_{N \times 1}$

$$(X^T X) \hat{w} = X^T \underline{y} \longrightarrow \boxed{\hat{w} = (X^T X)^{-1} X^T \underline{y}}$$

$X \hat{w} = \underline{y}$

دستگاه معادلات خطی

$X^+$ شبه وارون ماتریس $X$     اگر $X$ مربعی باشد: $X^+ = X^{-1}$

پیچ موارد اینگونه نیست   $N = \ell$

❖ SMALL in the sum of error squares sense means

$$J(\underline{w}) = \sum_{i=1}^{N} (y_i - \underline{w}^T \underline{x}_i)^2$$

$(y_i, \underline{x}_i)$: training pairs   that is, the input $\underline{x}_i$ and its corresponding class label $y_i$ (±1).

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

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

$+1 \rightarrow \omega_1$
$-1 \rightarrow \omega_2$

❖ Pseudoinverse Matrix

➢ Define

$$X = \begin{bmatrix} \underline{x}_1^T \\ \underline{x}_2^T \\ ... \\ \underline{x}_N^T \end{bmatrix} \quad \text{(an } Nxl \text{ matrix)}$$

$$\underline{y} = \begin{bmatrix} y_1 \\ ... \\ y_N \end{bmatrix} \quad \text{corresponding desired responses}$$

➢ $\quad X^T = [\underline{x}_1, \underline{x}_2, ..., \underline{x}_N] \quad \text{(an } lxN \text{ matrix)}$

➢ $\quad X^T X = \sum_{i=1}^{N} \underline{x}_i \underline{x}_i^T$

➢ $\quad X^T \underline{y} = \sum_{i=1}^{N} \underline{x}_i y_i$

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

$$(X^T X)\hat{\underline{w}} = X^T \underline{y} \Rightarrow$$

$$\hat{\underline{w}} = (X^T X)^{-1} X^T \underline{y}$$

$$= X^{\neq} \underline{y}$$

$$\boxed{X^{\neq} \equiv (X^T X)^{-1} X^T}$$ Pseudoinverse of $X$

➢ Assume $N=l \Rightarrow X$ square and invertible. Then

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

$$\boxed{X^{\neq} = X^{-1}}$$

31

➢ Assume $N > l$.   Then, in general, there is no solution to satisfy all equations simultaneously:

$$X\,\underline{w} = \underline{y}:\qquad \begin{array}{l} \underline{x}_1^T\,\underline{w} = y_1 \\[4pt] \underline{x}_2^T\,\underline{w} = y_2 \\[2pt] \cdots \\[2pt] \underline{x}_N^T\,\underline{w} = y_N \end{array}\qquad N \text{ equations} > l \text{ unknowns}$$

➢ The "solution"  $\underline{w} = X^{\neq}\,\underline{y}$  corresponds to the minimum sum of squares solution

> Example:

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

$$\omega_2 : \begin{bmatrix} 0.4 \\ 0.6 \end{bmatrix}, \begin{bmatrix} 0.6 \\ 0.2 \end{bmatrix}, \begin{bmatrix} 0.7 \\ 0.4 \end{bmatrix}, \begin{bmatrix} 0.8 \\ 0.6 \end{bmatrix}, \begin{bmatrix} 0.7 \\ 0.5 \end{bmatrix}$$
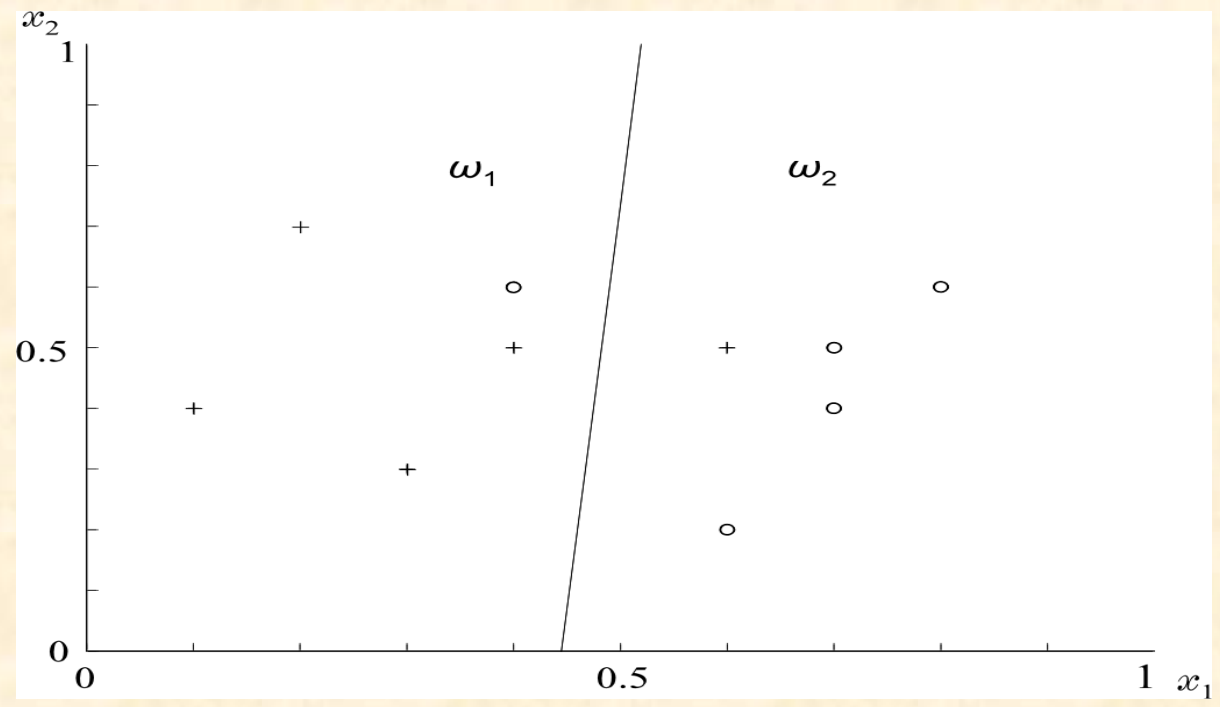
مسأله دسته بندی دو بعدی

$$W_1 x_1 + W_2 x_2 + W_0 = 0$$

حل:

$$\hat{W} = (X^T X)^{-1} X^T \underline{y}$$

طبقه بند خطی بهینه با معیار جمع مربعات خطا (LS) طراحی کنید.

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



33

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

$$\underline{w} = (X^T X)^{-1} X^T \underline{y} = \begin{bmatrix} -3.13 \\ 0.24 \\ 1.34 \end{bmatrix}$$

❖ The Bias – Variance Dilemma

A classifier $g(\underline{x})$ is a learning machine that tries to predict the class label $y$ of $\underline{x}$ . In practice, a finite data set $D$ is used for its training. Let us write $g(\underline{x}; D)$. Observe that:

➢ For some training sets, $D = \{(y_i, \underline{x}_i), i = 1, 2, ..., N\}$, the training may result to good estimates, for some others the result may be worse.

➢ The average performance of the classifier can be tested against the MSE optimal value, in the mean squares sense, that is:

$$E_D\left[\left(g(\underline{x}; D) - E[y \mid \underline{x}]\right)^2\right]$$

where $E_D$ is the mean over all possible data sets $D$.

- The above is written as:

$$E_D\left[\left(g(\underline{x};D)-E[y\,|\,\underline{x}]\right)^2\right]=$$

$$\left(E_D\left[g(\underline{x};D)\right]-E[y\,|\,\underline{x}]\right)^2+E_D\left[\left(g(\underline{x};D)-E_D\left[g(\underline{x};D)\right]\right)^2\right]$$

- In the above, the first term is the contribution of the bias and the second term is the contribution of the variance.

- For a finite $D$, there is a trade-off between the two terms. Increasing bias it reduces variance and vice verse. This is known as the bias-variance dilemma.

- Using a complex model results in low-bias but a high variance, as one changes from one training set to another. Using a simple model results in high bias but low variance.

❖ LOGISTIC DISCRIMINATION

➤ Let an $M$-class task, $\omega_1, \omega_2, ..., \omega_M$. In logistic discrimination, the logarithm of the likelihood ratios are modeled via linear functions, i.e.,

$$\ln\left(\frac{P(\omega_i \mid \underline{x})}{P(\omega_M \mid \underline{x})}\right) = w_{i,0} + \underline{w}_i^T \underline{x}, \ i = 1, \ 2, \ ..., \ M\text{-}1$$

➤ Taking into account that

$$\sum_{i=1}^{M} P(\omega_i \mid \underline{x}) = 1$$

it can be easily shown that the above is equivalent with modeling posterior probabilities as:

$$P(\omega_M \mid \underline{x}) = \frac{1}{1 + \sum\limits_{i=1}^{M-1} \exp\left(w_{i,0} + \underline{w}_i^T \underline{x}\right)}$$

$$P(\omega_i \mid \underline{x}) = \frac{\exp\left(w_{i,0} + \underline{w}_i^T \underline{x}\right)}{1 + \sum\limits_{i=1}^{M-1} \exp\left(w_{i,0} + \underline{w}_i^T \underline{x}\right)}, \iota = 1,2,...M-1$$

➢ For the two-class case it turns out that

$$P(\omega_2 \mid \underline{x}) = \frac{1}{1 + \exp\left(w_0 + \underline{w}^T \underline{x}\right)}$$

$$P(\omega_1 \mid \underline{x}) = \frac{\exp\left(w_0 + \underline{w}^T \underline{x}\right)}{1 + \exp\left(w_0 + \underline{w}^T \underline{x}\right)}$$

➢ The unknown parameters $\underline{w}_i,\, w_{i,0},\; i=1,\, 2,\, ...,\, M\text{-}1$ are usually estimated by maximum likelihood arguments.

➢ Logistic discrimination is a useful tool, since it allows linear modeling and at the same time ensures posterior probabilities to add to one.

Support Vector Machine (SVM)

3.7.1 : کلاس های جدایی پذیر خطی          Seperable Classes

فرض کنید $N$ بردار آموزشی $x_i$ , $i = 1, 2, ..., N$     از کلاس های $\omega_1$ و $\omega_2$ جدایی پذیر خطی

هدف : طراحی ابر صفحه ای که همه ی بردارهای آموزشی را به درستی طبقه بندی کند.

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

همانطور که قبلا گفته شد این ابر صفحه جدا کننده مها یکتا نیست . برای مثال الگوریتم پرسپترون به یکی از این ابر صفحه ها
همگرا می شود.



طبقه بند بهتری است

زیرا فضای بیشتری را در اختیار هر کلاس
قرار می دهد.

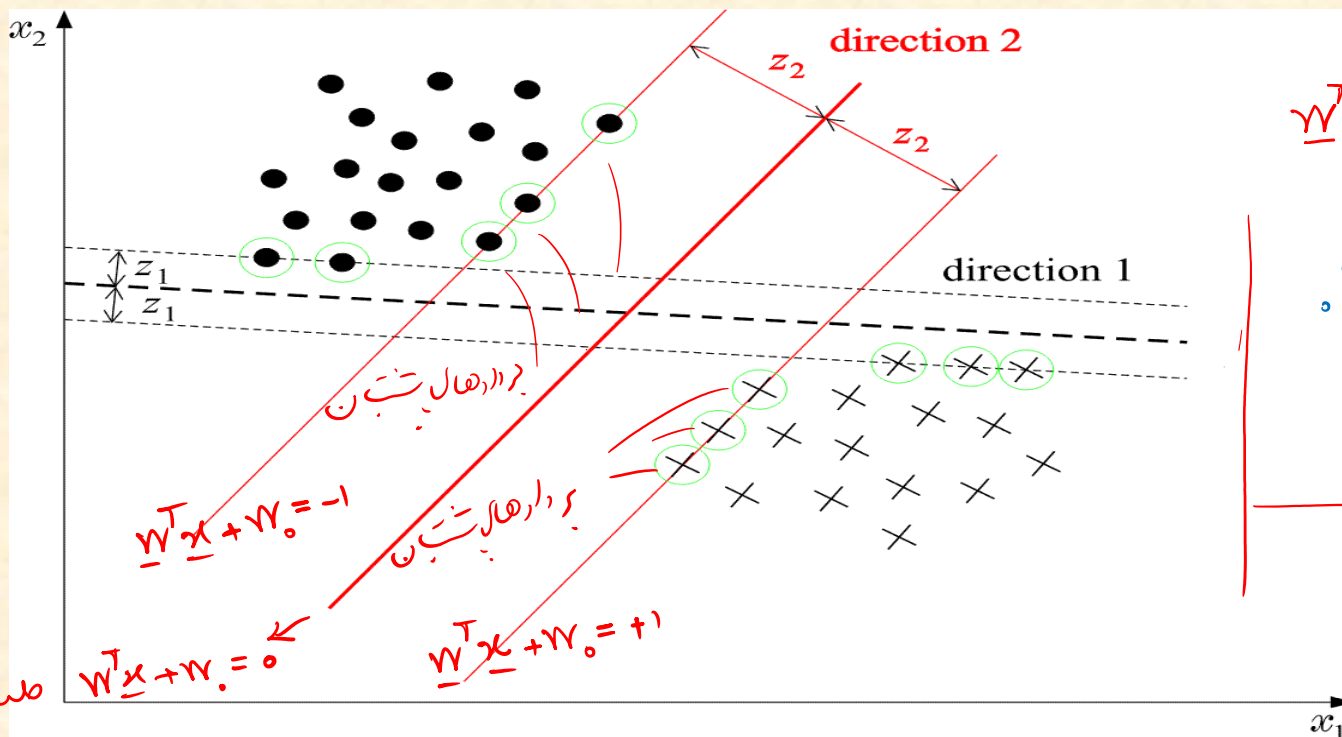generalization performance          کارایی تعمیم یافته

40

❖ Support Vector Machines

➢ The goal: Given two linearly separable classes, design the classifier

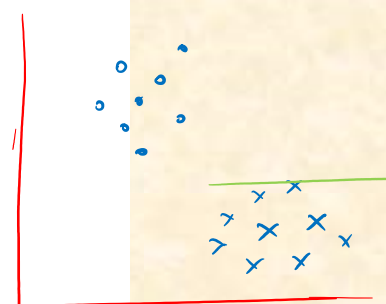$$g(\underline{x}) = \underline{w}^T \underline{x} + w_0 = 0$$

that leaves the maximum margin from both classes



41

➢ 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

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

فاصله نقطه $\hat{x}$ از ابر صفحه $g(\underline{x}) = 0$:

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

➢ Scale, $\underline{w}, \underline{w}_0$, so that at the nearest points from each class the discriminant function is ±1:

$$\left| g(\underline{x}) \right| = 1 \; \left\{ g(\underline{x}) = +1 \, \text{for} \, \omega_1 \, \text{and} \, g(\underline{x}) = -1 \, \text{for} \, \omega_2 \right\}$$

$\underline{w}$ و $w_0$ را طوری مقیاس دهی می‌کنیم که نزدیک‌ترین نقاط هر کلاس مقدار $g(\underline{x}) = \pm 1$ داشته باشند.

➢ Thus the margin is given by

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

➢ Also, the following is valid

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

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

43

➤ SVM (linear) classifier

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

بهينه سازی مقيد

maximize $\frac{2}{\|\underline{w}\|}$

➤ Minimize

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

➤ Subject to

قيد را ببين ⟶ $y_i (\underline{w}^T \underline{x}_i + w_0) \geq 1, \ i = 1, 2, ..., N$

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

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

➤ The above is justified since by minimizing $\|\underline{w}\|$

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

44

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

KKT

بسم (c)

تابع لاگرانژ

ضرایب لاگرانژ

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

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

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

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

- Where $L(\bullet, \bullet, \bullet)$ is the Lagrangian

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

45

➢ 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,

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

    where $N_s \le N_0$ , corresponding to positive Lagrange multipliers

  – From constraint (4) above, i.e.,

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

    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.

– Although the solution is unique, the resulting Lagrange multipliers are not 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
  $$\underline{w} = \sum_{i=1}^{N} \lambda_i y_i \underline{x}_i$$

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

  $$\underline{\lambda} \geq \underline{0}$$

- Combine the above to obtain

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

  - subject to

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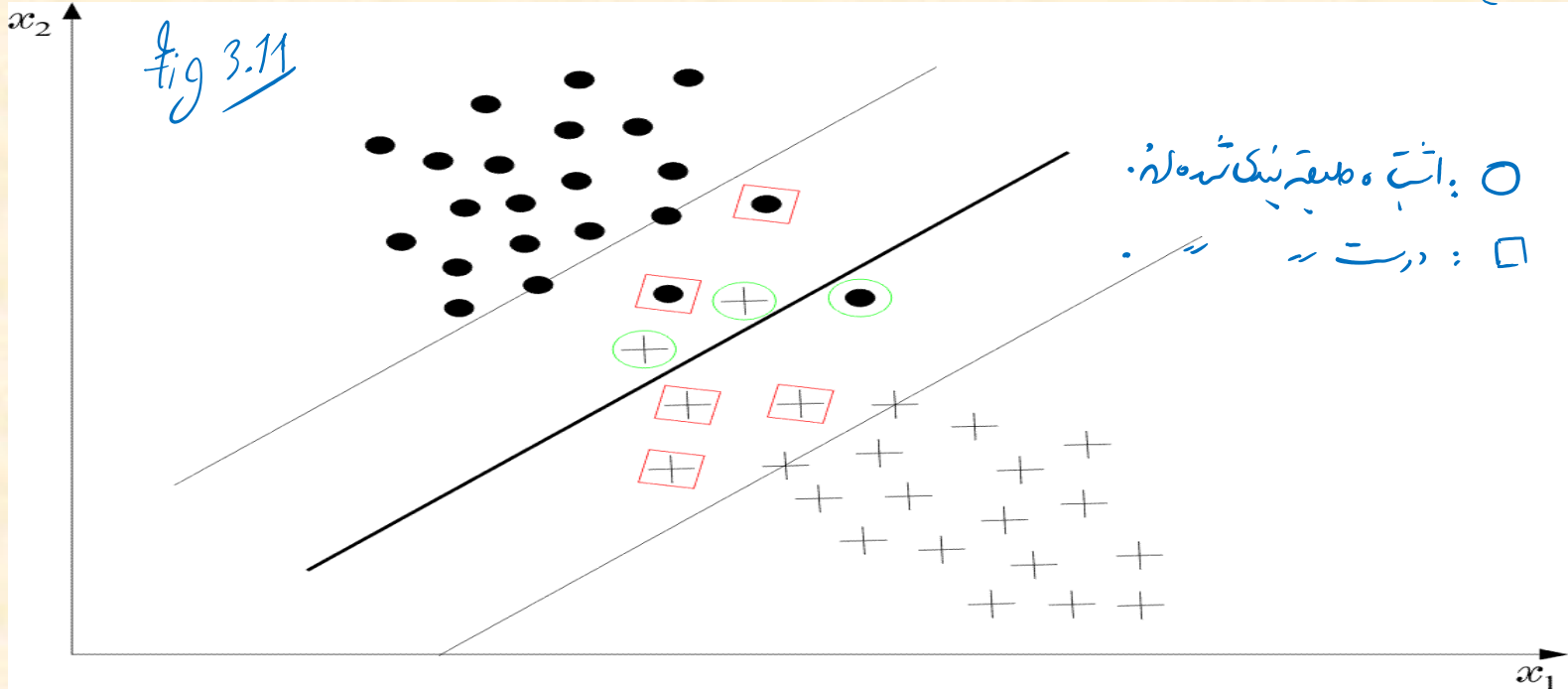
  $$\underline{\lambda} \geq \underline{0}$$

> Remarks:

- Support vectors enter via inner products

> Non-Separable classes

حالت دوم : کلاس های جدایی ناپذیر خطی :

در این حالت از آنجا که هیچ ابرصفحه ای وجب دندارد که مثلۀ کلاس ها را از یکدیگر جدا کند . خمولده نمونه های بیشتر به اشتباه طبقه بندی می شود .



fig 3.11

○ : اشتباه طبقه بندی شده اند .

□ : درست ...

51

داخل حاشیه‌ها ، تعدادی نمونه قرار گرفته‌اند :

سه دسته بردار وجود دارد :

۱- بردارهایی که خارج از باند قرار گرفته‌اند و به درستی طبقه‌بندی شده‌اند .

۲- بردارهایی که داخل باند قرار گرفته‌اند و درست طبقه‌بندی شده‌اند ( $\square$ ) و معادله زیر را برآورده می‌کنند :

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

۳- بردارهایی که داخل باند قرار گرفته‌اند و به اشتباه طبقه‌بندی شده‌اند ( $O$ ) و معادله زیر را برآورده می‌کنند :

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

هر سه دسته را می‌توان با رابطه زیر بیان کرد :

$$y_i \left[ \underline{w}^T \underline{x} + w_0 \right] \geq 1 - \xi_i$$

$$\left\{ \begin{array}{l} \text{برای دسته اول} : \xi_i = 0 \\ \text{دوم} : 1 \leq \xi_i < 0 \\ \text{سوم} : 1 < \xi_i \end{array} \right.$$

به $\xi_i$ ها متغیرهای slack گفته می‌شود .

In this case, there is no hyperplane such that

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

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

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

➤ The training vectors belong to <u>one</u> of <u>three</u> possible categories

1) Vectors outside the band which are correctly classified, i.e.,

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

2) Vectors inside the band, and correctly classified, i.e.,

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

3) Vectors misclassified, i.e.,

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

➤ All three cases above can be represented as

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

1) $\quad \rightarrow \xi_i = 0$

2) $\quad \rightarrow 0 < \xi_i \leq 1$

3) $\quad \rightarrow 1 < \xi_i$

$\xi_i$ are known as slack variables

➢ The goal of the optimization is now two-fold

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

  One way to achieve this goal is via the cost

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

where $C$ is a constant and

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

*Regularization Parameter*

- $I(.)$ is 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 a similar procedure as before we obtain

56

حل مسئله :

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

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

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

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

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

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

➤ The associated dual problem

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

subject to

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

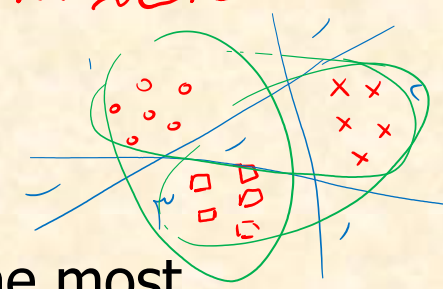$$\sum_{i=1}^{N} \lambda_i y_i = 0$$

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

➢ Training the SVM

A major problem is the high computational cost. To this end, decomposition techniques are used. The rationale behind them consists of the following:

- Start with an arbitrary data subset (working set) that can fit in the memory. Perform optimization, via a general purpose optimizer.

- Resulting support vectors remain in the working set, while others are replaced by new ones (outside the set) that violate severely the KKT conditions.

- Repeat the procedure.

- The above procedure guarantees that the cost function decreases.

- Platt's SMO algorithm chooses a working set of two samples, thus analytic optimization solution can be obtained.

> Multi-class generalization



One against All:

   Although theoretical generalizations exist, the most popular in practice is to look at the problem as $M$ two-class problems (one against all).

One against One:

   Binary Classifiers

— مسئلهٔ $M$ کلاسه :

— یکی در برابر همه :

در این حالت برای هر یک از کلاس ها به دنبال $m$ رای تابع جداکننده میتند $g_i(\underline{x})$ هستیم :

$$g_i(\underline{x}) > g_j(\underline{x}) \quad ; \quad i=1,\dots,M \quad \text{یعنی}$$

$$\forall j \neq i \quad ; \quad \text{if} \quad \underline{x} \in \omega_i$$

در SVM : $g(\underline{x}) = 0$ را طول برداری طراحی کنیم که از همهٔ کلاس های دیگر فاصله برابر داشته باشد ؛

$$\text{assign} \quad \underline{x} \text{ in } \omega_i \quad \text{if} \quad i = \underset{k}{\arg\max} \left[ g_k(\underline{x}) \right]$$

دو مشکل : ۱- نواحی نامشخص : نواحی که چند $g_i(\underline{x})$ دارند.

assymetric training

۲- آموزش نامتوازن

تعداد نمونه هال منفی از مثبت بیشتر شود چون $m$ از تعداد کلاس ها زیاد

— یکی در برابر دیگری ؛ برای هر دو کلاس یک طبقه بند طراحی می کنیم

پرسش : تعداد طبقه بندها ؟

تعداد طبقه بندها : $\dfrac{M(M-1)}{2}$

$\dbinom{M}{2}$

$M$ کلاس

$\omega_1 : [1,1]^T, [1,-1]^T$

$\omega_r : [-1,1]^T, [-1,-1]^T$



$w_1 = 1$

$w_r = w_o = 0$

بر رَسی SVM می خواهیم نشان دهیم

به حَقّ $x_1 = 0$ ابر صفحهٔ

جداکننده است و این نتیجه به ازای

محدوده های مختلف از ضرایب لاگرانژ قابل

دست یابی است.

$g(x) = w_1 x_1 + w_r x_r + w_o = 0$

**FIGURE 3.12**

In this example all four points are support vectors. The margin associated with $g_1(x) = 0$ is smaller compared to the margin defined by the optimal $g(x) = 0$.

61

قيود هر داده : $\quad W_1 + W_2 + W_0 - 1 \geqslant 0$

$$W_1 - W_2 + W_0 - 1 \geqslant 0$$

مقيد است $\quad W_1 - W_2 - W_0 - 1 \geqslant 0$

$$W_1 + W_2 - W_0 - 1 \geqslant 0$$

$$\mathcal{L}(W_1, W_2, W_0, \lambda) = \frac{W_1^2 + W_2^2}{2} - \lambda_1 (W_1 + W_2 + W_0 - 1)$$
$$- \lambda_2 (W_1 - W_2 + W_0 - 1)$$
$$- \lambda_3 (W_1 - W_2 - W_0 - 1)$$
$$- \lambda_4 (W_1 + W_2 - W_0 - 1)$$

$KKT$ :

$$\frac{\partial \mathcal{L}}{\partial W_1} = 0 \longrightarrow W_1 = \lambda_1 + \lambda_2 + \lambda_3 + \lambda_4$$

$$\frac{\partial \mathcal{L}}{\partial W_2} = 0 \longrightarrow W_2 = \lambda_1 - \lambda_2 - \lambda_3 + \lambda_4$$

$$\frac{\partial \mathcal{L}}{\partial W_0} = 0 \longrightarrow \lambda_1 + \lambda_2 - \lambda_3 - \lambda_4 = 0$$

$$W_1 = 1 \qquad \begin{cases} \lambda_1 + \lambda_2 + \lambda_3 + \lambda_4 = 1 \\ \lambda_1 - \lambda_2 - \lambda_3 + \lambda_4 = 0 \\ \lambda_1 + \lambda_2 - \lambda_3 - \lambda_4 = 0 \end{cases}$$
$$W_2 = W_0 = 0 \longrightarrow$$

$$\lambda_1 (W_1 + W_2 + W_0 - 1) = 0$$
$$\lambda_2 (W_1 - W_2 + W_0 - 1) = 0$$
$$\lambda_3 (W_1 - W_2 - W_0 - 1) = 0$$
$$\lambda_4 (W_1 + W_2 - W_0 - 1) = 0$$

$$\lambda_1, \lambda_2, \lambda_3, \lambda_4 \geqslant 0$$

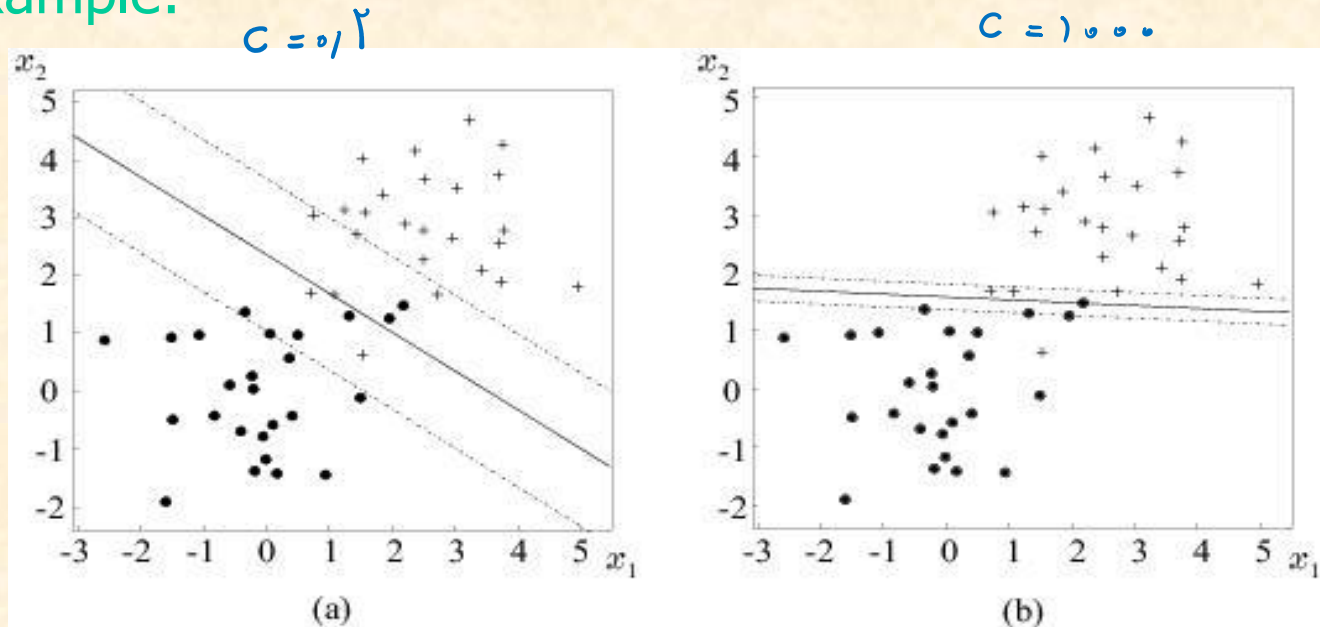حل دستگاه معادلات
جواب نهایی

مثال: تأثیر پارامتر تنظیم C درحالتِ کلاسهای جدایی ناپذیر خطی

➢ Example:

$C = 0/2$          $C = 1000$



(a)          (b)

➢ Observe the effect of different values of $C$ in the case of non-separable classes.

$C$ کمتر ← margin بیشتر

$C$ بیشتر ← تعداد نمونه‌های داخل باند کمتر