# Support Vector Machines

Neural Network Lab.

之前討論只要分界線合理讓所有資料正確即可，但我
們不可能已經有所有資料

$$\vec{w}_0^T\vec{x} + b = +1$$



$$\vec{w}_0^T\vec{x} + b = 0$$

$$\vec{w}_0^T\vec{x} + b = -1$$

類神經網路暨醫學影像處理實驗室

# Support Vector Machines

The optimal hyperplane $\quad \vec{w}_0^T \vec{x} + b = 0$

the distance from $\vec{x}$ to hyperplane is

$$g(x) = \vec{w}_0^T \vec{x} + b$$

Let $x_p$ be the normal projection of x onto the optimal hyperplane

r: be the algebra distance
(r>0 on positive side, r<0 on negative side)



Then $\quad \mathbf{x} = \mathbf{x}_p + \boldsymbol{r} \dfrac{\boldsymbol{w}_0}{\|\boldsymbol{w}_0\|}$

$$g(\vec{x}) = \vec{w}_0^T \vec{x} + b = \boldsymbol{r}\|\boldsymbol{w}_0\|$$

$$=> \boldsymbol{r} = \dfrac{g(\vec{x})}{\|\boldsymbol{w}_0\|}$$

類神經網路暨醫學影像處理實驗室
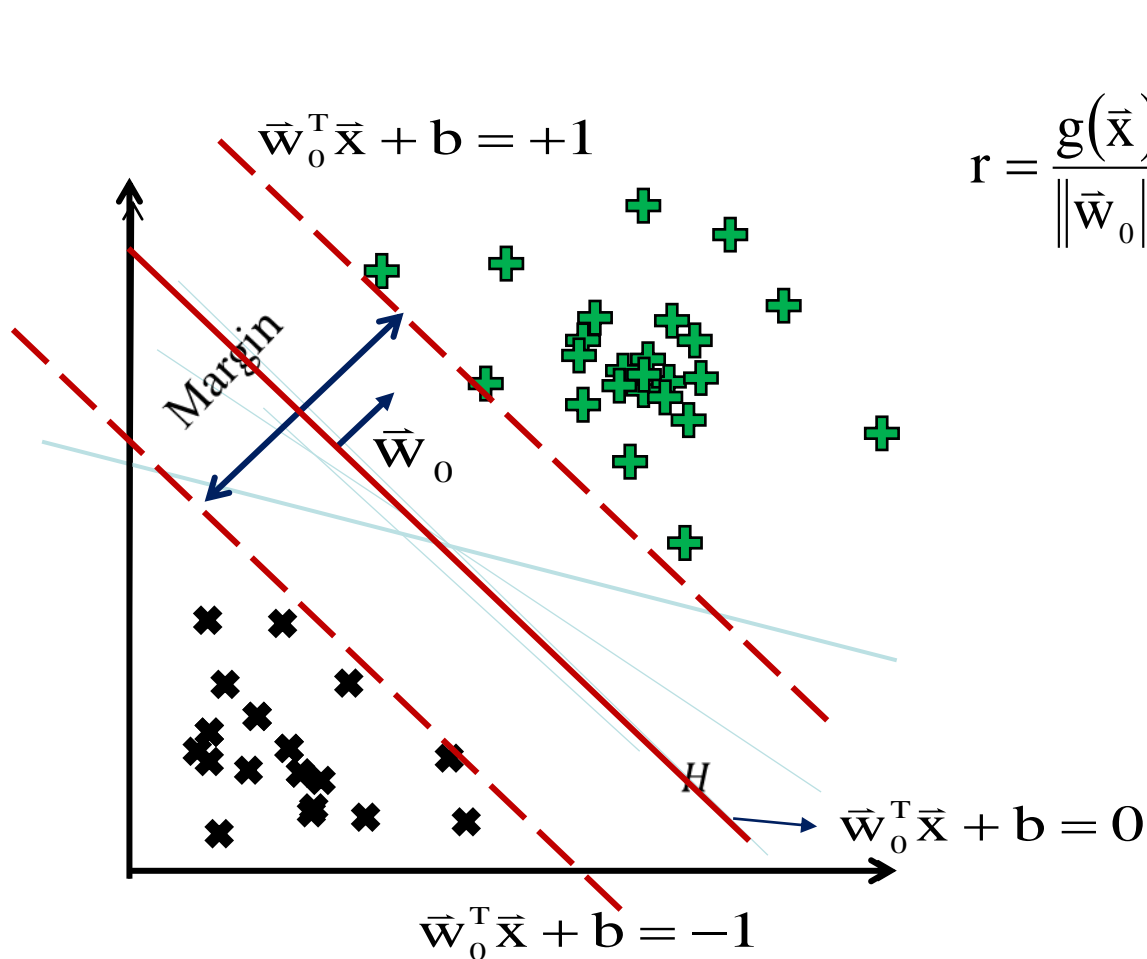
The algebraic distance from support vector to optimal hyperlane is

$$r = \frac{g(\vec{x})}{\|\vec{w}_0\|} = \begin{cases} \dfrac{1}{\|\vec{w}_0\|} & \text{if } d = +1 \\ \dfrac{-1}{\|\vec{w}_0\|} & \text{if } d = -1 \end{cases}$$

Margin of separation between two classes:

$$\rho = 2r = \frac{2}{\|\vec{w}_0\|}$$

➔ Maximum separation implies minimizes $\|\vec{w}_0\|$

Neural Network Lab.

$$r = \frac{g(\vec{x})}{\|\vec{w}_0\|} = \begin{cases} \dfrac{1}{\|\vec{w}_0\|} & \text{if } d = +1 \\ \dfrac{-1}{\|\vec{w}_0\|} & \text{if } d = -1 \end{cases}$$

$$\vec{w}_0^T \vec{x} + b = +1$$

Margin

$$\vec{w}_0$$

$$H$$

$$\vec{w}_0^T \vec{x} + b = 0$$

$$\vec{w}_0^T \vec{x} + b = -1$$

類 神 經 網 路 暨 醫 學 影 像 處 理 實 驗 室

Neural Network Lab.

# Lagrangian function

$$\mathbf{w}_o^T \mathbf{x}_i + b_o \geq 1 \qquad \text{for } d_i = +1 \tag{6.6}$$

$$\mathbf{w}_o^T \mathbf{x}_i + b_o \leq -1 \qquad \text{for } d_i = -1$$

➔ $$d_i(\mathbf{w}^T \mathbf{x}_i + b) \geq 1 \qquad \text{for } i = 1, 2, \ldots, N \tag{6.10}$$

The particular data points $(x_i, d_i)$ for which the first or second line if Eq. (6.6) is satisfied with the quality sign are called *support vectors*.

*Given the training sample* $\{(\mathbf{x}_i, d_i)\}_{i=1}^N$, *find the optimum values of the weight vector* **w** *and bias b such that they satisfy the constraints*

Kuhn-Tucker conditions: $\underline{d_i(\mathbf{w}^T \mathbf{x}_i + b) \geq 1 \qquad \text{for } i = 1, 2, \ldots, N}$

*and the weight vector* **w** *minimizes the cost function:*

$$\Phi(\mathbf{w}) = \frac{1}{2} \mathbf{w}^T \mathbf{w}$$

類神經網路暨醫學影像處理實驗室

# Lagrangian function

➔ Minimize

$$J(\mathbf{w}, b, \alpha) = \frac{1}{2} \mathbf{w}^T \mathbf{w} - \sum_{i=1}^{N} \alpha_i [d_i(\mathbf{w}^T \mathbf{x}_i + b) - 1] \quad (6.11)$$

$\alpha_j$ are called *Lagrange multipliers*.

代入

Taking directive:

Condition 1: $\dfrac{\partial J(\mathbf{w}, b, \alpha)}{\partial \mathbf{w}} = \mathbf{0}$ ⟹ $w = \sum_{i=1}^{N} \alpha_i d_i \vec{x}_i$

Condition 2: $\dfrac{\partial J(\mathbf{w}, b, \alpha)}{\partial b} = 0$ ⟹ $\sum_{i=1}^{N} \alpha_i d_i = \mathbf{0}$

Plug w into J(w,b,$\alpha$), we will have

$$J(w, b, \alpha) = \frac{1}{2} w^T w - \sum_{i=1}^{N} \alpha_i [d_i(w^T x_i + b) - 1]$$

$$= \sum_{i=1}^{N} \alpha_i - \frac{1}{2} \sum_{i=1}^{N} \sum_{j=1}^{M} d_i d_j \alpha_i \alpha_j x_i^T x_j$$

$w^{\mathrm{T}} w$

$b \sum \alpha_i d_i = 0$

= J($\alpha$)

Neural Network Lab.

max $\quad Q(\alpha) = \sum_{i=1}^{N} \alpha_{ij} - \dfrac{1}{2} \sum_{i=1}^{N} \sum_{j=1}^{N} \alpha_i \alpha_j d_i d_j \mathbf{x}_i^T \mathbf{x}_j$ $\qquad$ (6.16)

➔ min $\quad \dfrac{1}{2} \boldsymbol{\alpha}^T \begin{bmatrix} d_1 d_1 x_1^T x_1 & d_1 d_2 x_1^T x_2 & .... & d_1 d_N x_1^T x_N \\ d_2 d_1 x_2^T x_1 & d_2 d_2 x_2^T x_2 & .... & d_2 d_N x_2^T x_N \\ .... & .... & .... & .... \\ d_N d_1 x_N^T x_1 & d_N d_2 x_N^T x_2 & .... & d_N d_N x_N^T x_N \end{bmatrix} \boldsymbol{\alpha} + -1^T \boldsymbol{\alpha}$

Such that $\quad d^T \boldsymbol{\alpha} = 0 \qquad 0 \le \boldsymbol{\alpha} \le \infty$

This can be rewritten as

$$\min \dfrac{1}{2} \boldsymbol{\alpha}^T Q \boldsymbol{\alpha} - 1^T \boldsymbol{\alpha} \quad \text{subject to} \quad d^T \boldsymbol{\alpha} = 0$$

此為標準二次程式 **convex function** 之**optimization**

類神經網路暨醫學影像處理實驗室

*Neural Network Lab.*

Let the optimization solutions be $\alpha_1$ , $\alpha_2$ , $\alpha_3$ ,…, $\alpha_N$

At saddle point, for each Lagrange multiplier, Kuhn-Tucker Conditions of Optimization theory

$$\boldsymbol{\alpha}_i(d_i(\vec{w}_0^T\vec{x}_i + b\ ) - 1) = 0$$

Then we will have

$$\boldsymbol{\alpha}_i = 0 \quad \text{or} \quad d_i(\vec{w}_0^T\vec{x}_i + b\ ) - 1 = 0$$

➔1:  All interior points have $\boldsymbol{\alpha}_i = 0$

2:  if $\alpha_I > 0$, then $x_i$ is a support vector

Then $w = \sum_{i=1}^{N}\boldsymbol{\alpha}_i\boldsymbol{d}_i\vec{x}_i$ is determined only by the support vectors

Once w is obtained, $\boldsymbol{w} = \sum_{i=1}^{N} \boldsymbol{\alpha}_i \boldsymbol{d}_i \vec{\boldsymbol{x}}_i$ by the support vectors

we can plug the w into $\quad \mathrm{d}_i(\vec{\mathrm{w}}_0^{\mathrm{T}} \vec{\mathrm{x}}_i + \mathrm{b}) - 1 = 0$
to obtain the b

From $\quad J(\boldsymbol{\alpha}) = \sum_{i=1}^{N} \boldsymbol{\alpha}_i - \frac{1}{2} \sum_{i=1}^{N} \sum_{j=1}^{M} \boldsymbol{d}_i \boldsymbol{d}_j \boldsymbol{\alpha}_i \boldsymbol{\alpha}_j \boldsymbol{x}_i^T \boldsymbol{x}_j$

## What if the data are not linearly separable?

The minimization is on $x_i^T x_i$ if we can transform to minimize on $z_i^T z_i$, it would still work

# The Dual Problem

*Given the training sample $\{(\mathbf{x}_i, d_i)\}_{i=1}^N$, find the Lagrange multipliers $\{\alpha_i\}_{i=1}^N$ that maximize the objective function*

$$Q(\alpha) = \sum_{i=1}^{N} \alpha_i - \frac{1}{2} \sum_{i=1}^{N} \sum_{j=1}^{N} \alpha_i \alpha_j d_i d_j \mathbf{x}_i^T \mathbf{x}_j \qquad \text{只含}\alpha_i$$

*subject to the constraints*

(1) $\quad \displaystyle\sum_{i=1}^{N} \alpha_i d_i = 0$

(2) $\quad \alpha_i \geq 0 \qquad$ for $i = 1, 2, \dots, N$

由 dual problem 決定 optimal $\alpha_i$ (叫做 $\alpha_{o,i}$)

代入 (6.17) 得到 optimal weight $\boldsymbol{w_0}$ $\qquad \mathbf{w}_o = \displaystyle\sum_{i=1}^{N} \alpha_{o,i} d_i \mathbf{x}_i$

代入 (6.18) 得到 optimal weight $\boldsymbol{b_0}$

$$b_o = 1 - \mathbf{w}_o^T \mathbf{x}^{(s)} \qquad \text{for } d^{(s)} = 1 \qquad\qquad (6.18)$$

類神經網路暨傢處理實驗室

*Neural Network Lab.*

# Optimal hyper-plane for non-separable patterns

之前討論for linear separable，現在為non-separable pattern (允許部分pattern落入partition margin內)：

$$d_i(\mathbf{w}^T\mathbf{x}_i + b) \geq 1 - \xi_i, \qquad i = 1, 2, ..., N \tag{6.22}$$
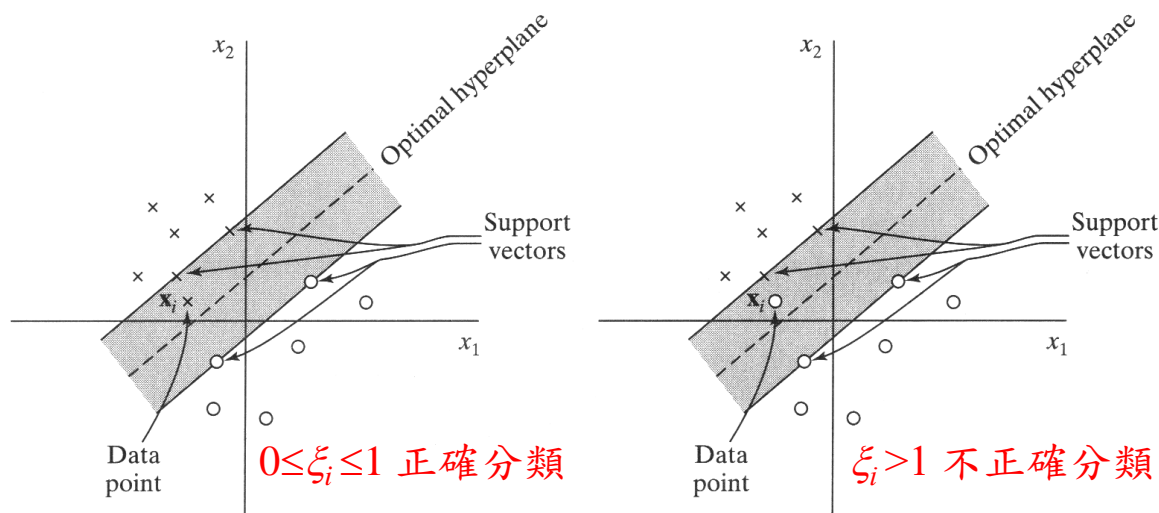
The $\xi_i$ are called slack variables



0≤$\xi_i$≤1 正確分類　　　　　$\xi_i$>1 不正確分類

**FIGURE 6.3** (a) Data point $\mathbf{x}_i$ (belonging to class $\mathscr{C}_1$) falls inside the region of separation, but on the right side of the decision surface. (b) Data point $\mathbf{x}_i$ (belonging to class $\mathscr{C}_2$) falls on the wrong side of the decision surface.

類神經網路暨醫學影像處理實驗室

The misclassification error, averaged on the training set, is minimized

$$\Phi(\xi) = \sum_{i=1}^{N} I(\xi_i - 1)$$

$$I(\xi) = \begin{cases} 0 & \text{if } \xi \leq 0 \\ 1 & \text{if } \xi > 0 \end{cases}$$

<span style="color:red">Correct but maybe inside the margin</span>

<span style="color:red">Incorrect</span>

<span style="color:red">為了計算方便更改為</span> $\Phi(\xi) = \sum_{i=1}^{N} \xi_i$

<span style="color:red">Then</span> $\Phi(\mathbf{w}, \xi) = \frac{1}{2}\mathbf{w}^T\mathbf{w} + C\sum_{i=1}^{N}\xi_i$ (6.23)

The first term in Eq. (6.23) is related to minimizing the VC dimension of the support vector machine.
The second term is an upper bound on the number of the test errors.
The parameter C is user determined (1)*experimentally* (2)*analytically*.

類神經網路暨醫學影像處理實驗室

*Then for the soft classification:*

Given the training sample $\{(\mathbf{x}_i, d_i)\}_{i=1}^{N}$, find the optimum values of the weight vector $\mathbf{w}$ and bias b such that they satisfy the constraint

$$d_i(\mathbf{w}^T\mathbf{x}_i + b) \geq 1 - \xi_i \qquad \text{for } i = 1, 2, \ldots, N$$

$$\xi_i \geq 0 \quad \text{for all } i$$

and such that the weight vector $\mathbf{w}$ and the slack variables $\xi_i$ minimize the cost functional

$$\Phi(\mathbf{w}, \xi) = \frac{1}{2}\mathbf{w}^T\mathbf{w} + C\sum_{i=1}^{N}\xi_i$$

where C is a user-specified positive parameter.

Neural Network Lab.

# Duality

*Given the training sample* $\{(\mathbf{x}_i, d_i)\}_{i=1}^{N}$, *find the Lagrange multipliers* $\{\alpha_i\}_{i=1}^{N}$ *that maximize the objective function*

$$Q(\alpha) = \sum_{i=1}^{N} \alpha_i - \frac{1}{2} \sum_{i=1}^{N} \sum_{j=1}^{N} \alpha_i \alpha_j d_i d_j \mathbf{x}_i^T \mathbf{x}_j$$

*subject to the constraints*

(1) $\displaystyle\sum_{i=1}^{N} \alpha_i d_i = 0$

(2) $0 \leq \alpha_i \leq C \qquad$ for $i = 1, 2, \ldots, N$

*where C is a user-specified positive parameter.*

Neural Network Lab.
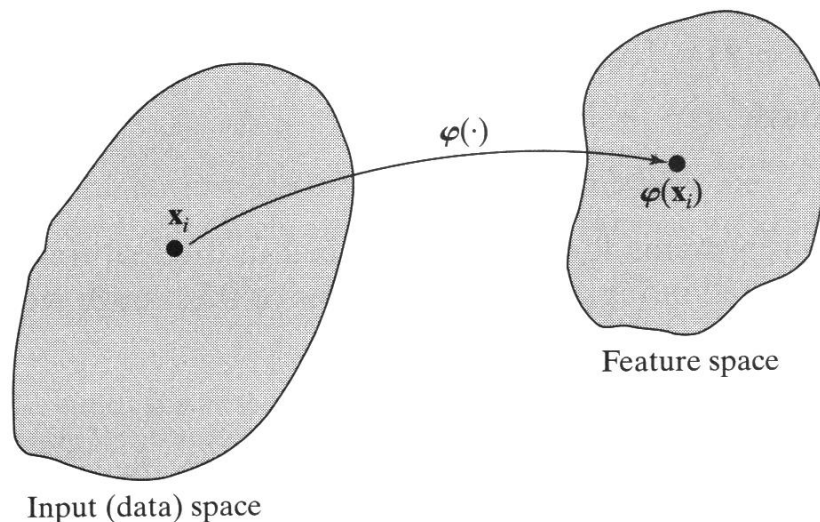
Then the optimum solution for weight vector is

$$w_0 = \sum \alpha_{0i} d_i x_i$$

and the Kuhn - Tucker conditions are :

$$\alpha_i \left[ d_i \left( \vec{w}^T x_i + b \right) - 1 + \xi_i \right] = 0 \qquad i = 1, 2, ..., N$$

$$u_i \, \xi_i = 0 \qquad i = 1, 2, ..., N$$



FIGURE 6.4 Nonlinear map $\varphi(\cdot)$ from the input space to the feature space.

We may define a hyperplane acting as the decision surface as follows

$$\sum_{j=1}^{m_1} w_j \varphi_j(\mathbf{x}) + b = 0 \qquad\qquad (6.29)$$

$$\mathbf{w}^T \boldsymbol{\varphi}(\mathbf{x}) = 0 \qquad\qquad (6.33) \quad \text{Hyperplane}$$

由6.12得知

$$\mathbf{w} = \sum_{i=1}^{N} \alpha_i d_i \boldsymbol{\varphi}(\mathbf{x}_i) \qquad\qquad (6.34)$$

6.34代入6.33

$$\sum_{i=1}^{N} \alpha_i d_i \boldsymbol{\varphi}^T(\mathbf{x}_i) \boldsymbol{\varphi}(\mathbf{x}) = 0 \qquad\qquad (6.35)$$

The term $\boldsymbol{\varphi}^T(\boldsymbol{x_i})\boldsymbol{\varphi}(\boldsymbol{x})$ represents the inner product of two vectors induced in the feature space by the input vector $\boldsymbol{x}$ and the input pattern $\boldsymbol{x_i}$

$$K(\mathbf{x}, \mathbf{x}_i) = \boldsymbol{\varphi}^T(\mathbf{x})\boldsymbol{\varphi}(\mathbf{x}_i)$$

$$= \sum^{m_1} \varphi_j(\mathbf{x})\varphi_j(\mathbf{x}_i) \qquad \text{for } i = 1, 2, \ldots, N \qquad (6.36)$$

$$K(\mathbf{x}, \mathbf{x}_i) = K(\mathbf{x}_i, \mathbf{x}) \qquad \text{for all } i \qquad\qquad (6.37)$$

$$\sum_{i=1}^{N} \alpha_i d_i K(\mathbf{x}, \mathbf{x}_i) = 0 \qquad\qquad (6.38)$$

K is a symmetric function

類神經網路暨醫學影像處理實驗室

Neural Network Lab.

**TABLE 6.1**　Summary of Inner-Product Kernels

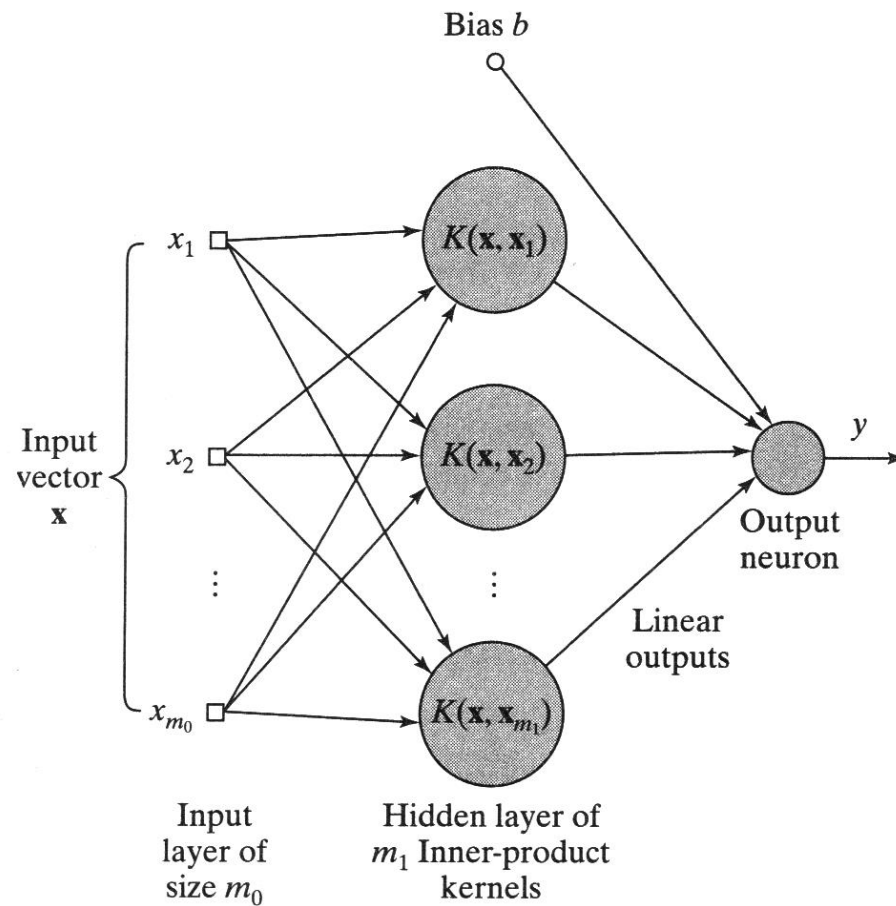| Type of support vector machine | Inner product kernel $K(\mathbf{x}, \mathbf{x}_i),\ i = 1, 2, \ldots, N$ | Comments |
|---|---|---|
| Polynomial learning machine | $(\mathbf{x}^T\mathbf{x}_i + 1)^p$ | Power $p$ is specified *a priori* by the user |
| Radial-basis function network | $\exp\left(-\dfrac{1}{2\sigma^2}\|\mathbf{x} - \mathbf{x}_i\|^2\right)$ | The width $\sigma^2$, common to all the kernels, is specified *a priori* by the user |
| Two-layer perceptron | $\tanh(\beta_0\mathbf{x}^T\mathbf{x}_i + \beta_1)$ | Mercer's theorem is satisfied only for some values of $\beta_0$ and $\beta_1$ |

**FIGURE 6.5** Architecture of support vector machine.