# Theoretical principles of Support Vector Machines

Aim: We try to understand the kernel trick from a mathematical perspective.

> book: Support Vector Machines for Pattern Classification (Chapter 2); Abe., S. (2010)

## 1. Set the scene

Let us start with understanding duality. Later on we want to bring the constrained problem of finding the optimal separating hyperplane into its dual form.

Suppose we want to minimize

$$f(x,y) = x^2 + 2y$$

subject to an equality constraint

$$3x + 2y + 1 = 0.$$

①

We define a new function called
the Lagrangian :    *Lagrangian multiplier*

$$g(x,y,\lambda) = f(x,y) - \lambda \cdot [3x + 2y + 1]$$

$$= [x^2 + 2y] - \lambda \cdot [3x + 2y + 1]$$

Let us compute the partial derivatives:

$$\frac{d}{dx} g(x,y,\lambda) = 2x - 3\lambda \qquad \overset{!}{=} 0$$

$$\frac{d}{dy} g(x,y,\lambda) = 2 - 2\lambda \qquad \overset{!}{=} 0$$

$$\frac{d}{d\lambda} g(x,y,\lambda) = -[3x + 2y + 1] \overset{!}{=} 0$$
$$\Rightarrow -3x = 1 + 2y$$

$$\Rightarrow 2x - 3\lambda = 2 - 2\lambda$$

$$\Rightarrow \lambda = 2[x - 1]$$

$$\Rightarrow \lambda =$$

②

$$x = + \frac{1 + 2y}{3}$$

$$\Rightarrow \lambda = 2\left[-\frac{1}{3} - \frac{2}{3}y - 1\right]$$

$$\lambda = -\frac{8}{3} - \frac{4}{3}y$$

$$\Rightarrow 2 - 2\left[-\frac{8}{3} - \frac{4}{3}y\right] = 0$$

$$\Rightarrow \frac{6}{3} + \frac{16}{3} + \frac{8}{3}y = 0$$

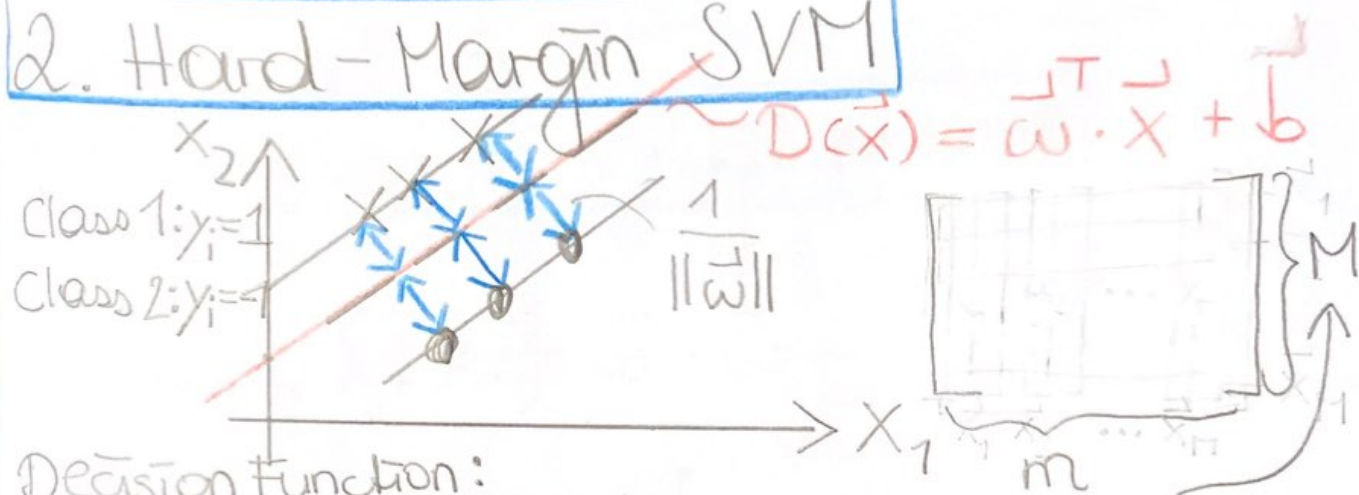$$\Rightarrow -\frac{22}{8} = y \quad \Rightarrow \quad \boxed{y = -\frac{11}{4}}$$

$$\Rightarrow x = -\frac{1}{3} - \frac{2}{3}\left[-\frac{11}{4}\right]$$

$$= -\frac{1}{3} + \frac{22}{12} = -\frac{4}{12} + \frac{22}{12} = \frac{\overset{3}{\cancel{18}}}{\underset{2}{\cancel{12}}} =$$

$$\Rightarrow \boxed{x = \frac{3}{2}} \qquad \Rightarrow \boxed{\lambda = 2\left[\frac{3}{2} - \frac{3}{3}\right] =}$$

$$= 2\left[\frac{9}{6} - \frac{6}{6}\right] = \boxed{1} \quad ③$$

This method can be generalized to inequality constraints (Karush-Kuhn-Tucker (KKT) conditions)

## 2. Hard-Margin SVM

$D(\vec{x}) = \vec{\omega}^T \cdot \vec{x} + \vec{b}$



Class 1: $y_i = 1$
Class 2: $y_i = -1$

$\frac{1}{\|\vec{\omega}\|}$

Decision function:

$$D(\vec{x}) = \vec{\omega}^T \cdot \vec{x} + b \quad ; \underline{\underline{M}} \text{ training data}$$

If linearly separable:     bias term

$$\vec{\omega}^T \cdot \vec{x}_{\cdot i} + b \begin{cases} \geq 1 & \text{for } y_i = 1 \\ < -1 & \text{for } y_i = -1 \end{cases}$$

$$\Leftrightarrow \boxed{y_i \left[ \vec{\omega}^T \cdot \vec{x}_i + b \right] \geq 1}, \quad i = 1 \ldots M$$

training sample $\vec{x}$, distance to hyperplane

$$\frac{|D(\vec{x})|}{\|\vec{\omega}\|} \Rightarrow \text{optimal separating hyperplane, which minimizes } \|\omega\|$$

④

minimize $Q(\vec{\omega}, b) = \frac{1}{2} \|\vec{\omega}\|^2$  

optimal separating hyperplane

subject to $y_i [\vec{\omega}^T \cdot \vec{x}_i + b] \geq 1$

has a unique solution!

## 2.1 Convert into dual problem

$$Q(\vec{\omega}, b, \vec{\lambda}) = \frac{1}{2} \vec{\omega}^T \cdot \vec{\omega} - \sum_{i=1}^{M} \lambda_i \{ y_i (\vec{\omega}^T \cdot \vec{x}_i + b) - \text{...} \quad (1)$$

$\vec{\lambda} = (\lambda_1, \ldots, \lambda_M)^T$ non-negative Lagrangian multipliers

Karush-Kuhn-Tucker (KKT) conditions:

$$\frac{\partial Q}{\partial \vec{\omega}} = 0 \quad (2)$$

$$\frac{\partial Q}{\partial b} = 0 \quad (3)$$

$$\lambda_i \{ y_i (\vec{\omega}^T \cdot \vec{x}_i + b) - 1 \} = 0 \quad (4)$$

$$\lambda_i \geq 0 \quad (5)$$

⑤

$$\lambda_i \neq 0 \Rightarrow y_i(\vec{\omega}^T \cdot \vec{x_i} + b) = 1$$

$$\Rightarrow \vec{x_i} \text{ with } \lambda_i \neq 0 \text{ the support vector}$$

from (1) and (3):

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

from (1) and (2): $\left( \frac{1}{2} \vec{\omega} + \frac{1}{2} \vec{\omega} - \sum \lambda_i y_i x_i \right)$

$$\vec{\omega} = \sum_{i=1}^{M} \lambda_i y_i \vec{x_i}$$

Dual problem:

$$Q(\lambda) = \sum_{i=1}^{M} \lambda_i - \frac{1}{2} \sum_{i,j=1}^{M} \lambda_i y_i \lambda_j y_j \vec{x_i}^T \vec{x_j}$$

$$\text{with } \sum_{i=1}^{M} \lambda_i y_i = 0 \quad ; \quad \lambda_i \geqslant 0, \ i=1..M$$
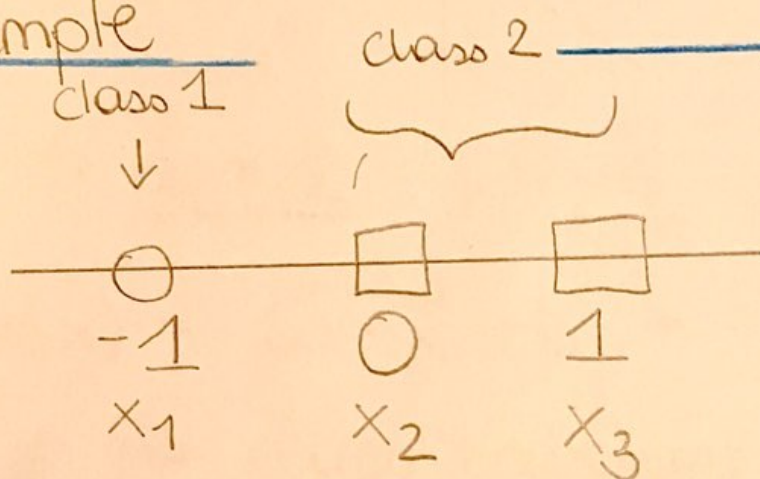
Decision function: $D(\vec{x}) = \sum_{i \in S} \lambda_i y_i \vec{x_i}^T \vec{x} + b$

⑥

Then

$$D(\vec{x}) = \sum_{i \in S} \lambda_i y_i \vec{x}_i^T \cdot \vec{x} + b$$

$$b = \sum_{i \in S} (y_i - \vec{w}^T \cdot \vec{x}_i) \frac{1}{|S|}$$

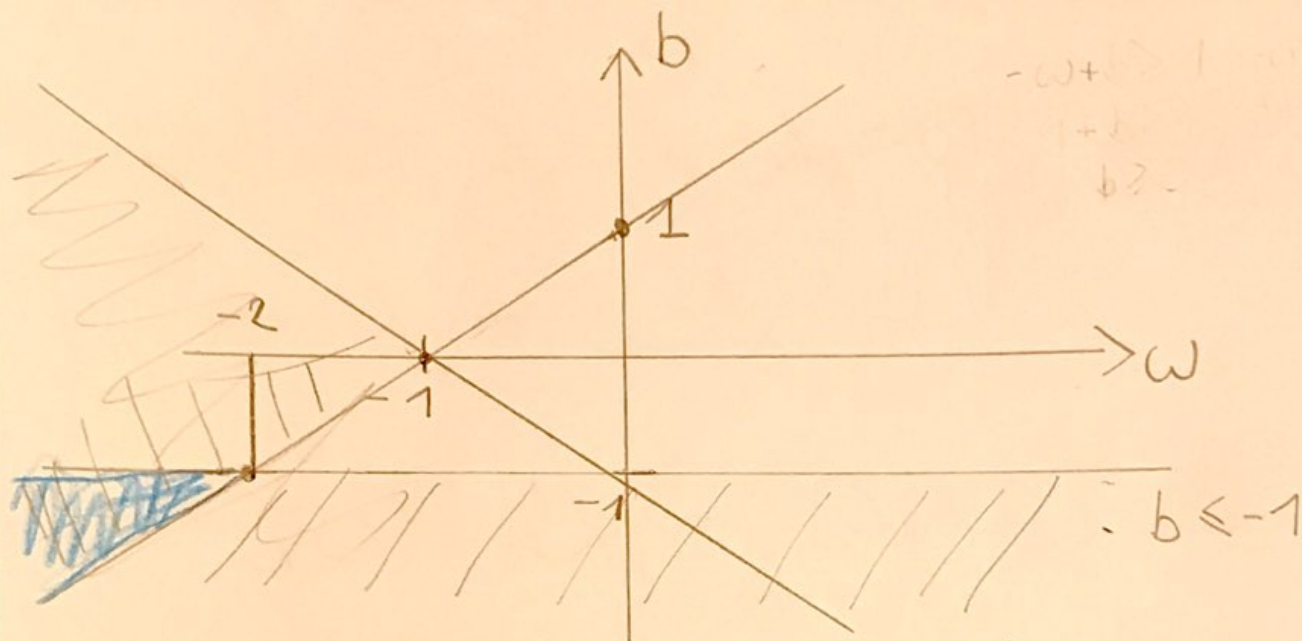## 2.2. Understand it better via an example

class 1

class 2



$$-1 \qquad 0 \qquad 1$$
$$x_1 \qquad x_2 \qquad x_3$$

$$y_i [\vec{w}^T \cdot \vec{x}_i + b] \geqslant 1 \quad , \quad i = 1 \dots M$$

$$\Rightarrow \qquad w \cdot (-1) + b \geqslant -1 \quad \}(y_i = 1)$$
$$(-1) \cdot b \geqslant 1 \quad \}(y_i = -1)$$
$$(-1) [w + b] \geqslant 1 \quad \}$$

Our solution has to minimize $\|\vec{\omega}\|^2$



$\Rightarrow \quad \underline{\omega = -2} \quad \underline{b = -1}$

$$D(x) = \vec{\omega}^T \cdot \vec{x} + b = -2 \cdot x - 1$$

$\Rightarrow$ class boundary at $x = \frac{1}{2}$

$x = 0$ and $x = -1$ are
support vector

Dual problem:

$$Q(\vec{\lambda}) = \lambda_1 + \lambda_2 + \lambda_3 - \frac{1}{2}\left[\lambda_1 \cdot 1 \cdot \lambda_1 \cdot 1 (-1) \cdot (-1)\right.$$
$$\left. + \lambda_1 \cdot 1 \cdot \lambda_2 \cdot (-1) \cdot (-1)(0) + \ldots\right]$$

⑧

$$= \lambda_1 + \lambda_2 + \lambda_3 - \frac{1}{2}\left[\lambda_1^2 + 2\lambda_1\lambda_3 + \lambda_3^2\right]$$

subject to $\quad \lambda_1 - \lambda_2 - \lambda_3 = 0 \; , \quad \alpha_i \geqslant 0$

$$\Rightarrow \quad \lambda_2 = \lambda_1 - \lambda_3$$

$$\Rightarrow Q(\lambda) = 2\lambda_1 - \frac{1}{2}\left[\lambda_1 + \lambda_3\right]^2 ; \; \alpha_i \geqslant 0$$

$$\Rightarrow \quad \lambda_3 = 0 \quad Q(\lambda) = 2\lambda_1 - \frac{1}{2}\lambda_1^2$$

$$= -\frac{1}{2}\left[\lambda_1 - 2\right]^2 + 2$$

$$\Rightarrow \lambda_1 = 2$$

$$\Rightarrow \underbrace{\lambda_1 = 2, \; \lambda_2 = 2, \; \lambda_3 = 0}$$

support vectors

$$\vec{w} = \sum \lambda_i y_i \vec{x}_i = \lambda_1 y_1 x_1 + \lambda_2 y_2 x_2 + \lambda_3 y_3 x_3$$

$$= (2)\cdot(1)\cdot(-1) = -2$$

$$b = -1$$

If $x_3 = 1$ would be in Class 1 then this problem is not linearly separable in this space :
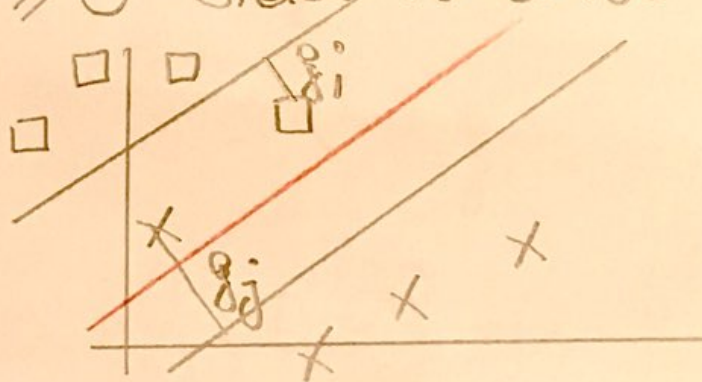
$$(-1)[\omega + b] \geq 1$$

$$\downarrow$$

$$(+1)[\omega + b] \geq 1$$

## 2.3 Soft-Margin Support Vector Machines

To allow inseparable solutions (when hard-margin SVM is unsolvable):

$$y_i [\vec{\omega}^T \cdot \vec{x}_i + b] \geq 1 - \xi_i , \quad i = 1...M$$

$\xi_i \geq 0$ slack variables

With $\xi_i$ feasible solutions always exist.

minimize $Q(\vec{\omega}, b, \xi) = \frac{1}{2}\|\vec{\omega}\|^2 + \frac{C}{p}\sum\limits_{i=1}^{M}\xi_i^p$

subject to $y_i[\vec{\omega}^T \cdot \vec{x}_i + b] \geqslant 1 - \xi_i$

$p = 1$    L1 norm
$p = 2$    L2 norm

$C$ is the margin parameter (trade-off between maximization of margin and minimization of classification error)

$Q(\vec{\omega}, b, \vec{\xi}, \vec{\alpha}, \vec{\beta}) = \frac{1}{2}\|\vec{\omega}\|^2 + C\sum\xi_i -$
(L1 norm)
$\qquad\qquad - \sum\alpha_i(y_i[\vec{\omega}^T\cdot\vec{x}_i + b] - 1 - \xi_i)$

$\qquad\qquad\qquad - \sum\beta_i\xi_i$

$\frac{\partial Q}{\partial \omega} = 0, \frac{\partial Q}{\partial b} = 0, \frac{\partial Q}{\partial \xi} = 0, \begin{array}{l}\alpha_i(y_i[\ldots]) = 0 \\ \beta_i\xi_i = 0\end{array}$   ⑪

$$\Rightarrow \quad \vec{w} = \sum_{i=1}^{M} \alpha_i y_i \vec{x}_i$$

$$\sum_{i=1}^{M} \alpha_i y_i = 0, \qquad \alpha_i + \beta_i = C$$
$$i = 1 \ldots M$$

Dual Problem:

$$\text{maximize} \quad Q(\vec{\alpha}) = \sum \alpha_i - \frac{1}{2} \sum \alpha_i \alpha_j y_i y_j \vec{x}_i^{T} \vec{x}_j$$

$$\text{subject to} \quad \sum_{i=1}^{M} y_i \alpha_i = 0 \qquad C \geq \alpha_i \geq 0$$
$$\text{(unbounded)} \qquad \underbrace{\qquad\qquad}_{\substack{\text{only difference} \\ \text{to hard-margin} \\ \text{SVM}}}$$

$$0 < \alpha_i < C \quad \forall \text{ support vectors}$$

$$\alpha_i = C \quad \text{bounded support vector}$$

$$D(\vec{x}) = \sum_{i \in S} \alpha_i y_i \vec{x}_i^{T} \vec{x} + b \, ;$$

$$i \in S \Leftarrow \text{support vector} \quad b = \frac{1}{|U|} \sum_{i \in U} \left( y_i - \vec{w}^{T} \cdot \vec{x}_i \right)$$

(12)

## 3. Kernel Tricks

$$[X_1 \ldots X_M] \longmapsto [\phi_1(\vec{x}), \ldots, \phi_\ell(\vec{x})]$$

$$D(\vec{x}) = \vec{\omega}^T \cdot \vec{\phi}(\vec{x}) + b$$

Note: $\vec{\omega}$ is now $[\omega_1 \ldots, \omega_\ell]$ $\ell$-dim!

Hilbert-Schmidt theory
$\downarrow$

Let us assume (kernel trick)

$$K(\vec{x}, \vec{x}') = \phi^T(\vec{x}) \, \phi(\vec{x}')$$

$\underbrace{\hphantom{K(\vec{x}, \vec{x}')}}$

kernel (no need to treat high-dim.
feature space explicitly)

Dual problem:

$$\max \quad Q(\vec{\lambda}) = \sum_{i=1}^{M} \lambda_i - \frac{1}{2} \sum_{i=1}^{M} \lambda_i \lambda_j y_i y_j K(\vec{x}_i, \vec{x}_j)$$

$$\text{subject to} \quad \sum_{i=1}^{M} y_i \lambda_i = 0 \qquad 0 \le \lambda_i \le C$$

KKT - conditions:

$$\lambda_i \left( y_i \left[ \sum_j y_j \lambda_j K(\vec{x}_i, \vec{x}_j) + b \right] - 1 - \xi_i \right) = 0$$

$$(C - \lambda_i) \xi_i = 0$$

$$(\alpha_i \geqslant 0), \quad (\xi_i \geqslant 0)$$

Decision function:

$$D(\vec{x}) = \sum_i \lambda_i y_i K(\vec{x}_i, \vec{x}) + b$$

$$b = y_j - \sum_i \alpha_i y_i K(\vec{x}_i, \vec{x}_j)$$

Linear kernel: $K(\vec{x}, \vec{x}') = \vec{x}^T \cdot \vec{x}'$

Polynomial kernel: $K(\vec{x}, \vec{x}') = [\vec{x}^T \cdot \vec{x}' + 1]^d$

$$d = 2, \quad m = 2$$

$$\vec{x} = \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} \qquad \vec{x}' = \begin{bmatrix} x_1' \\ x_2' \end{bmatrix}$$

$$K(\vec{x}, \vec{x}') = \left[ \begin{bmatrix} x_1 & x_2 \end{bmatrix} \begin{bmatrix} x_1' \\ x_2' \end{bmatrix} + 1 \right]^2 =$$
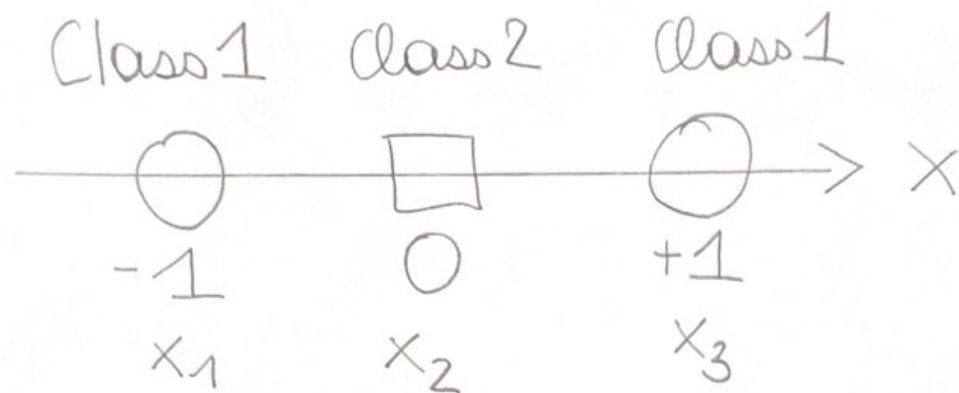
$$= \left[ x_1 x_1' + x_2 x_2' + 1 \right]^2 =$$

$$= (x_1 x_1')^2 + 2 x_1 x_1' x_2 x_2' + 1 + (x_2 x_2')^2$$
$$+ 2 x_1 x_1' + 2 x_2 x_2' = \phi^T(\vec{x}) \phi(\vec{x}')$$

$$= \underbrace{\begin{bmatrix} 1, & \sqrt{2}\, x_1, & \sqrt{2}\, x_2, & \sqrt{2}\, x_1 x_2, & x_1^2, & x_2^2 \end{bmatrix}}_{\phi^T(\vec{x})} \underbrace{\begin{bmatrix} 1 \\ \sqrt{2}\, x_1' \\ \sqrt{2}\, x_2' \\ \sqrt{2}\, x_1' x_2' \\ x_1'^2 \\ x_2'^2 \end{bmatrix}}_{\phi(\vec{x}')}$$

(15)

# 3.1 Our Example

Class 1     Class 2     Class 1

$$\begin{array}{ccc} \bigcirc & \square & \bigcirc \\ -1 & 0 & +1 \\ x_1 & x_2 & x_3 \end{array} \longrightarrow X$$

$$Q(\vec{\alpha}) = \alpha_1 + \alpha_2 + \alpha_3 -$$

$$- \frac{1}{2} \sum_{i=1}^{3} \sum_{j=1}^{3} \alpha_i \alpha_j \, y_i \, y_j \left[ \vec{x}_i^T \cdot \vec{x}_j + 1 \right]^2$$

$i=1, j=1 \quad \alpha_1 \alpha_1 (1)(1) \left[ (-1)^2 + 1 \right]^2 = \alpha_1^2 \cdot 4$

$i=1, j=2 \quad \alpha_1 \alpha_2 (1)(-1) \left[ 0 + 1 \right]^2 = -\alpha_1 \alpha_2$

$i=1, j=3 \quad \alpha_1 \alpha_3 (1)(1) \left[ (-1)(1) + 1 \right]^2 = 0$

$i=2, j=1 \quad \alpha_2 \alpha_1 (-1)(1) \left[ 0 + 1 \right]^2 = -\alpha_1 \alpha_2$

$i=2, j=2 \quad \alpha_2^2$

$i=2, j=3 \quad \alpha_2 \alpha_3 (-1)(1) \left[ 0 + 1 \right]^2 = +\alpha_2 \alpha_3$

$i=3, j=1 \quad 0$

$i=3, j=2 \quad -\alpha_2 \alpha_3$

$i=3, j=3 \quad \alpha_3 \alpha_3 (1)(1) \left[ 1 + 1 \right]^2 = 4\alpha_3^2$

(16)

$$Q(\vec{\alpha}) = \alpha_1 + \alpha_2 + \alpha_3$$
$$-\frac{1}{2}\left[\alpha_1^2 \cdot 4 + \alpha_3^2 \cdot 4 + \alpha_2^2 - 2\alpha_1\alpha_2 - 2\alpha_2\alpha_3\right]$$

$$(1)$$

subject to

$$\sum_i y_i \alpha_j = \alpha_1 - \alpha_2 + \alpha_3 \stackrel{!}{=} 0 \quad (2)$$

$$c \geqslant \alpha_i \geqslant 0$$

From (2).    $\alpha_2 = \alpha_1 + \alpha_3$

into (1):

$$Q(\vec{\alpha}) = 2\alpha_1 + 2\alpha_3 - \left[2\alpha_1^2 + 2\alpha_3^2 - \frac{1}{2}(\alpha_1 + \alpha_3)^2\right]$$

$$\frac{\partial Q}{\partial \alpha_1} = 2 - 4\alpha_1 + \frac{1}{2}2(\alpha_1 + \alpha_3)$$
$$\stackrel{!}{=} 0 \quad (1)$$
$$= 2 - 3\alpha_1 + \alpha_3$$

$$\frac{\partial Q}{\partial \alpha_3} = 2 - 4\alpha_3 + \alpha_1 + \alpha_3 \quad (2)$$
$$= 2 - 3\alpha_3 + \alpha_1 = 0 \quad \textcircled{17}$$

$$\Rightarrow {}^{(1)} \quad 3\alpha_1 = 2 + \alpha_3$$

$$\alpha_1 = \frac{2}{3} + \frac{\alpha_3}{3}$$

$$(2) \quad 2 - 3\alpha_3 + \frac{2}{3} + \frac{\alpha_3}{3} = 0$$

$$\frac{8}{3} - \frac{8}{3}\alpha_3 = 0$$

$$\Rightarrow \underline{\alpha_3 = 1} \quad \Rightarrow \underline{\alpha_1 = 1}$$

$$\underline{\alpha_2 = \alpha_1 + \alpha_3 = 2}$$

$$C \geqslant 2, \quad x = -1, 0, 1 \text{ are support vectoo}$$

$$b = y_j - \sum_{i \in S} \alpha_i y_i \underbrace{k(\vec{x}_i, \vec{x}_j)}_{[\vec{x}_i^T x_j + 1]^2}$$

$$= (1) - \alpha_1 (1) [1]^2 - \alpha_2 (-1) + \cdots$$

$$= -1$$

⑱

$$D(\vec{x}) = \sum_i \alpha_i y_i \left[ \vec{x_i}^T \cdot \vec{x} + 1 \right]^2 + b$$

$$= (1)(1)\left[ -\vec{x} + 1 \right]^2 + 1$$
$$+ (2)(-1)\left[ 1 \right]^2 +$$
$$+ (1)(1)\left[ \vec{x} + 1 \right]^2 - 1$$

$$= \left[ x - 1 \right]^2 + \left[ x + 1 \right]^2 - 3$$

$$= 2x^2 - 1$$

boundaries $\quad x = \pm \dfrac{1}{\sqrt{2}}$



class 1 | Class 2 | Class 1

$D(x)$

$-1 \qquad 0 \qquad 1 \qquad X$

$-1$