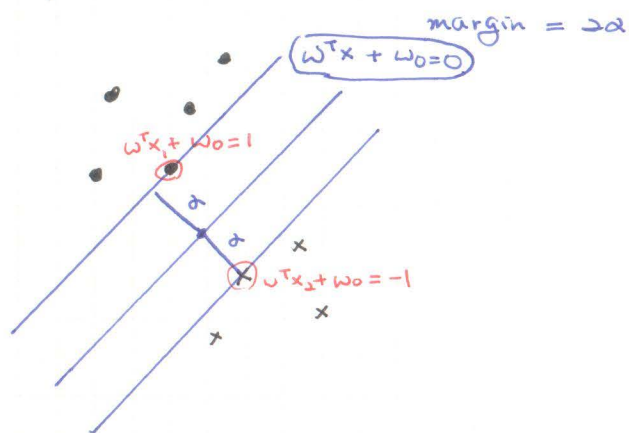


Support Vector Machines

1) Separable cases



• We want to find the equation $w^T x + w_0 = 0$ such that the margin 2α is maximized.

• Normalize $w^T x + w_0$ such that

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

• Problem

$$\max 2\alpha = \max \frac{2}{\|w\|}$$

subject to

$$w^T x + w_0 \geq 1 \quad x \in \omega_1$$

$$w^T x + w_0 \leq -1 \quad x \in \omega_2$$

• Alternatively

①

$$\min \frac{1}{2} \|w\|^2$$

$$\text{subject to } y_i (w^T x_i + w_0) \geq 1 \text{ for all } i.$$

$$\text{where } y_i = \begin{cases} 1 & \text{if } x_i \in \omega_1 \\ -1 & \text{if } x_i \in \omega_2 \end{cases}$$

- * It is a convex programming problem.
- * We can use KKT conditions and the Wolfe dual representation technique.
- * Recall

- KKT Conditions

$$L(\theta, \lambda) = J(\theta) - \sum_i \lambda_i f_i(\theta)$$

$$\textcircled{1} \quad \frac{\partial L}{\partial \theta} = 0$$

$$\textcircled{2} \quad \lambda \geq 0$$

$$\textcircled{3} \quad \lambda_i f_i(\theta) = 0$$

$$- \min_{\theta} \max_{\lambda \geq 0} L(\theta, \lambda)$$

$$= \max_{\lambda \geq 0} L(\theta, \lambda) \text{ subject to } \frac{\partial L}{\partial \theta} = 0.$$

$$f(w, w_0, \lambda) = \frac{1}{2} w^T w - \sum \lambda_i [y_i (w^T x_i + w_0) - 1]$$

(a) KKT Conditions

$$\textcircled{1} \frac{\partial f}{\partial w} = w - \sum \lambda_i y_i x = 0$$

$$\frac{\partial f}{\partial w_0} = - \sum \lambda_i y_i = 0$$

$$\boxed{\begin{aligned} w &= \sum \lambda_i y_i x \\ \sum \lambda_i y_i &= 0 \end{aligned}}$$

$$\textcircled{2} \lambda_i \geq 0$$

$$\textcircled{3} \lambda_i [y_i (w^T x_i + w_0) - 1] = 0$$

* For $\lambda_i \neq 0$ (i.e. $\lambda_i > 0$)

x_i is called a support vector.

$$w^T x_i + w_0 = \pm 1$$

Support vectors are the training vectors that are closest to the linear classifier.

(b) Wolfe Dual Representation Form

$$\max_{\lambda \geq 0} f(w, w_0, \lambda)$$

$$\text{subject to } \frac{\partial f}{\partial w} = 0 \quad \text{i.e. } w = \sum \lambda_i y_i x$$

$$\left\{ \begin{aligned} \frac{\partial f}{\partial w_0} = 0 \quad \text{i.e. } \sum \lambda_i y_i = 0 \end{aligned} \right.$$

or

$$\max f(w, w_0, \lambda)$$

$$\text{subject to } \begin{cases} \lambda_i \geq 0 \\ w = \sum \lambda_i y_i x \\ \sum \lambda_i y_i = 0 \end{cases}$$

Note that

$$f(w, w_0, \lambda)$$

$$= \frac{1}{2} (\sum_i \lambda_i y_i x_i)^T \sum_j \lambda_j y_j x_j$$

$$- \sum \lambda_i \left\{ y_i \left[(\sum_j \lambda_j y_j x_j)^T x_i + w_0 \right] - 1 \right\}$$

$$= - \frac{1}{2} \sum_{i,j} \lambda_i \lambda_j y_i y_j x_i^T x_j - \sum \lambda_i y_i w_0 + \sum \lambda_i$$

$$= \sum_i \lambda_i - \frac{1}{2} \sum_{i,j} \lambda_i \lambda_j y_i y_j x_i^T x_j$$

Therefore, we solve

$$\max_{\lambda} \left(\sum \lambda_i - \frac{1}{2} \sum \lambda_i \lambda_j y_i y_j x_i^T x_j \right)$$

$$\text{subject to } \lambda_i \geq 0$$

① λ 결정

$$\left\{ \begin{aligned} \sum \lambda_i y_i &= 0 \end{aligned} \right.$$

It does not depend on the dimensionality of x !

After solving this

$$\textcircled{2} w = \sum_{i=1}^n \lambda_i y_i x_i$$

For $\lambda_i \neq 0$

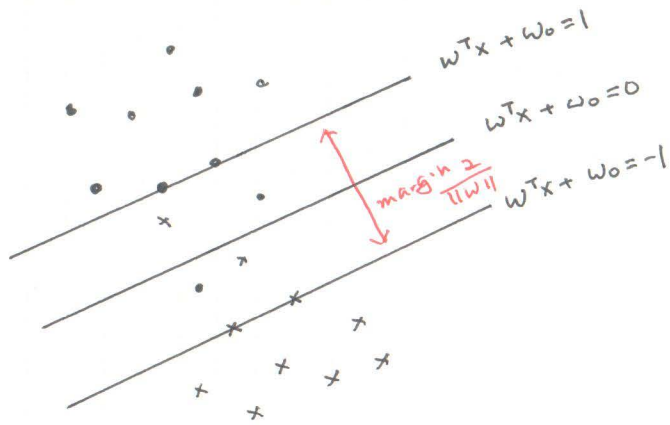
$$y_i (w^T x_i + w_0) = 1.$$

②

③

②

2) Nonseparable cases



* Slack variables $\xi_i > 0$

$$y_i [w^T x + w_0] \geq 1 - \xi_i$$

* New formulation

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

$$\text{subject to } y_i [w^T x + w_0] \geq 1 - \xi_i \quad \text{for each } i$$

$$\xi_i \geq 0 \quad "$$

- It is a convex programming problem

(3)

Lagrangian

$$f(w, w_0, \xi, \lambda, u)$$

$$= \frac{1}{2} \|w\|^2 + C \sum \xi_i - \sum u_i \xi_i$$

$$- \sum \lambda_i [y_i (w^T x + w_0) - 1 + \xi_i]$$

KKT Conditions

$$(1) \quad \frac{\partial f}{\partial w} = w - \sum \lambda_i y_i x = 0 \quad \therefore w = \sum \lambda_i y_i x_i$$

$$\frac{\partial f}{\partial w_0} = - \sum \lambda_i y_i = 0 \quad \therefore \sum \lambda_i y_i = 0$$

$$\frac{\partial f}{\partial \xi_i} = C - u_i - \lambda_i = 0$$

$$(2) \quad \lambda_i \geq 0, \quad u_i \geq 0$$

$$(3) \quad u_i \xi_i = 0, \quad \lambda_i [y_i (w^T x + w_0) - 1 + \xi_i] = 0$$

Wolfe Dual Representation

$$\min \max_{\substack{u \geq 0 \\ \lambda \geq 0}} f = \max_{\substack{u \geq 0 \\ \lambda \geq 0}} \min f$$

$$= \max_{\substack{u \geq 0 \\ \lambda \geq 0}} f \quad \text{subject to} \quad w = \sum \lambda_i y_i x_i$$

$$\sum \lambda_i y_i = 0$$

$$C - u_i - \lambda_i = 0$$

This is simplified to

$$f = \frac{1}{2} \sum_i \sum_j \lambda_i \lambda_j y_i y_j x_i^T x_j + C \sum \xi_i - \sum u_i \xi_i - \sum \lambda_i \xi_i$$

$$- \sum \lambda_i [y_i (\omega^T x_i + \omega_0) - 1]$$

$$= - \frac{1}{2} \sum_i \sum_j \lambda_i \lambda_j y_i y_j x_i^T x_j - \sum \lambda_i (y_i \omega_0 - 1)$$

$$= \sum_i \lambda_i - \frac{1}{2} \sum_i \sum_j \lambda_i \lambda_j y_i y_j x_i^T x_j$$

subject to $\lambda_i \geq 0$

$$u_i = C - \lambda_i \geq 0$$

\therefore $0 \leq \lambda_i \leq C$ \leftarrow The only difference.

$$\sum \lambda_i y_i = 0$$

• If x_i is inside the margin or wrongly classified

$$\xi_i > 0 \Rightarrow u_i = 0 \Rightarrow \lambda_i = C.$$