




$$d = \frac{b}{\|w\|}$$

margin $M = \frac{2}{\|w\|} = \frac{2}{\sqrt{w^T w}}$

$$\left. \begin{array}{l} w^T x_i + b \geq +1 \quad \text{for } y_i = +1 \\ w^T x_i + b \leq -1 \quad \text{for } y_i = -1 \end{array} \right\} \forall (w^T x_i + b) \geq 1$$

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

Lagrangian

$$L(\theta, \alpha) = f_0(\theta) + \sum_{i=1}^n \alpha_i f_i(\theta)$$

$$\min_{\theta \in \mathbb{R}^d} f_0(\theta) = f_0(\theta^*) \geq f_0(\theta^*) + \sum_{i=1}^n \underbrace{\alpha_i}_{\geq 0} \underbrace{f_i(\theta^*)}_{\leq 0} = L(\theta^*, \alpha) \geq \min_{\theta \in \mathbb{R}^d} L(\theta, \alpha)$$

$\underbrace{\phantom{\sum_{i=1}^n \alpha_i f_i(\theta^*)}_{\leq 0}}_{g(\alpha)}$

Hence, $\forall \alpha \ f_0(\theta^*) \geq g(\alpha)$ $g(\alpha) \rightarrow$ is a lower bound

Lagrangian dual problem

maximized $g(\alpha)$

subject to $\alpha_i \geq 0, i=1, \dots, m$

(always)

weak duality

$$g(\alpha) \leq p^*$$

$$\alpha_i \leq p^*$$

strong duality

under certain conditions

$$g(\alpha) = p^*$$

$$\alpha^* = p^*$$

SVM

1. Calculate the Lagrangian (constrained)

$$L(w, b, \alpha) = \underbrace{\frac{1}{2} w^T w}_{\text{Obj}} - \sum_{i=1}^N \alpha_i \underbrace{[y_i (w^T x_i + b) - 1]}_{f_i(0)}$$

2. Minimize $L(w, b, \alpha)$ w.r.t. w and b $g(\alpha) = \min_{w, b} L(w, b, \alpha)$

$$\nabla_w L(w, b, \alpha) = w - \sum_{i=1}^N \alpha_i y_i x_i = 0 \quad ; \quad w = \sum_{i=1}^N \alpha_i y_i x_i$$

$$\frac{\partial L}{\partial b} = \sum_{i=1}^N \alpha_i y_i = 0$$

The weights are a linear combination of the training samples

Substituting both relations back into $L(w, b, \alpha)$ gives the Lagrange dual function $g(\alpha)$

$$\begin{aligned} L(w, b, \alpha) &= \frac{1}{2} w^T w - \sum_{i=1}^N \alpha_i [y_i (w^T x_i + b) - 1] \\ &= \frac{1}{2} w^T \left(\sum_{i=1}^N \alpha_i y_i x_i \right) - \sum_{i=1}^N \alpha_i y_i w^T x_i - \sum_{i=1}^N \alpha_i y_i b + \sum_{i=1}^N \alpha_i \\ &= -\frac{1}{2} w^T \left(\sum_{i=1}^N \alpha_i y_i x_i \right) - b \sum_{i=1}^N \alpha_i y_i + \sum_{i=1}^N \alpha_i \\ &= -\frac{1}{2} \sum_{j=1}^N \sum_{i=1}^N \alpha_j y_j x_j^T \left(\sum_{i=1}^N \alpha_i y_i x_i \right) + \sum_{i=1}^N \alpha_i , \text{ since } \sum_{i=1}^N \alpha_i y_i = 0 \\ &= g(\alpha) \end{aligned}$$

maximize $g(\alpha)$

subject to $\sum_{i=1}^N \alpha_i y_i = 0$

$\alpha_i \geq 0, \text{ for } i=1, \dots, N$

Rewrite $g(\alpha)$ in vector form

$$g(\alpha) = \frac{1}{2} \alpha^T Q \alpha + \alpha^T 1_n$$

$Q \rightarrow$ symmetric negative (semi) definite matrix

→ To solve this, there is an efficient algorithm Sequential minimal optimization (SMO)

Recover w and b from dual solution α^*

$$w = \sum_{i=1}^N \alpha_i^* y_i x_i$$

$$w^T x_i + b = y_i$$

$$\alpha_i^* f_i(\theta^*) = 0$$

$$b = y_i - w^T x_i$$

$$y_i \begin{cases} -1 \\ 1 \end{cases}$$

$$f_i(\theta^*) = y_i (w^T x_i + b) - 1 \rightarrow y_i (w^T x_i + b) = 1 ; (w^T x_i + b) = \frac{1}{y_i} ;$$

$$w^T x_i + b = y_i ; b = y_i - w^T x_i$$

A training sample x_i can only contribute to the weight vector

($\alpha_i \neq 0$) if it lies on the margin, that is $y_i (w^T x_i + b) = 1$

A training sample x_i with $\alpha_i \neq 0$ is called support vector

Slack variable $\epsilon_i \geq 0$ to control noisy data points

$\ell_i = 0 \Rightarrow$ No violation

$$\begin{cases} w^T x_i + b \geq 1 - \epsilon_i & \text{for } y_i = +1 \\ w^T x_i + b \leq -1 + \epsilon_i & \text{for } y_i = -1 \end{cases} \rightarrow y_i(w^T x_i + b) \geq 1 - \epsilon_i$$

The new cost function is, $J_0(w, b, \epsilon) = \frac{1}{2} w^T w + C \sum_{i=1}^n \epsilon_i$

The factor $C > 0$ determines how heavy a violation is punished
 $C \rightarrow \infty$ recovers hard-margin SVM

$$\min f_0(w, b, \epsilon) = \frac{1}{2} w^T w + C \sum_{i=1}^n \epsilon_i$$

$$\text{subject to } \begin{cases} y_i (w^T x_i + b) - 1 + \epsilon_i \geq 0 \\ \epsilon_i \geq 0 \end{cases}$$

1. Lagrangian

$$L(w, b, \epsilon, \alpha, \mu) = \frac{1}{2} w^T w + C \sum_{i=1}^N \epsilon_i$$

f_0(\alpha)

$$\sum_{i=1}^n d_i [y_i (\mathbf{w}^\top \mathbf{x}_i + b) - 1 + \varepsilon_i] - \sum_{i=1}^n \mu_i \varepsilon_i \approx f_i(\theta)$$

2. Minimize

$$\nabla_w L(w, b, \epsilon, \alpha, \mu) = w - \sum_{i=1}^n \alpha_i y_i x_i \stackrel{!}{=} 0 \quad ; \quad \frac{\partial L}{\partial b} = \sum_{i=1}^n \alpha_i y_i \stackrel{!}{=} 0$$

$$\frac{\partial L}{\partial \varepsilon_i} = C - d_i - \mu_i \stackrel{!}{=} 0 \quad \text{for } i=1, \dots, n$$

From $\alpha_i = C - M_i$, and dual feasibility, $M_i \geq 0$, $\alpha_i \geq 0$ we get $0 \leq \alpha_i \leq C$

$$\text{max } g(\alpha) = \sum_{i=1}^n \alpha_i - \frac{1}{2} \sum_{i=1}^n \sum_{j=1}^n y_i y_j \alpha_i \alpha_j x_i^\top x_j$$

subject to $\begin{cases} \sum_{i=1}^N \alpha_i y_i = 0 \\ 0 \leq \alpha_i \leq C \quad i=1, \dots, N \end{cases}$

Only the constraint $\alpha_i \leq C$ is Neur, it ensures that α_i is bounded and cannot go to infinity

C big \rightarrow small margin (less datapoints violating in margin)

C small \rightarrow big margin (more datapoints violating in margin)
margin is violated

$$\epsilon_i = \begin{cases} 1 - y_i(w^\top x_i + b) & \text{if } y_i(w^\top x_i + b) \leq 1 \\ 0 & \text{else} \end{cases}$$

since we are minimizing over ϵ

$$\max(0, 1 - y_i(w^\top x_i + b))$$

we can rewrite the objective function as an unconstrained optimization problem known as the Hinge loss formulation

$$\min_{w,b} \frac{1}{2} w^\top w + C \sum_{i=1}^n \max\{0, 1 - y_i(w^\top x_i + b)\}$$

Kernel trick \rightarrow to construct non linear classifiers

$$g(\alpha) = \sum_{i=1}^n \alpha_i - \frac{1}{2} \sum_{i=1}^n \sum_{j=1}^n y_i y_j \alpha_i \alpha_j \underbrace{\phi(x_i)^T \phi(x_j)}_{K(x_i, x_j)}$$

Kernel function $K: \mathbb{R}^D \times \mathbb{R}^D \rightarrow \mathbb{R}$

$$K(x_i, x_j) := \phi(x_i)^T \phi(x_j)$$

Kernel is valid if it corresponds to an inner product in some feature space

Mercer's theorem \rightarrow A kernel is valid if it gives rise to a symmetric, positive semi definite kernel matrix K for any input data X

a new point can be classified as

$$h(x) = \text{sign} \left(\sum_{\substack{j \\ x_j \in S}} \alpha_j y_j K(x_j, x) + b \right)$$