

SVM

2020年10月12日 13:33

$$\{x_i, y_i\}, y_i \in \{-1, 1\}$$

Hard Margin (Perfectly Linear Separable).

- Margin: distance from the separator to the closest point

$$\text{Margin length } \frac{1}{\|w\|_2}$$

Goal: $y_i(w^T x_i - b) \geq 1$ (at the same time minimize $\|w\|_2$)

If data not perfectly linearly separable \Rightarrow Allow mistakes

naive answer.

\Rightarrow minimize $\|w\|_2 + \lambda \sum \xi_i$.

minimize $\|w\|_2 + \lambda \cdot \# \text{ mistakes}$

$$\text{s.t. } \forall i, y_i(w^T x_i - b) \geq 1 - \xi_i$$

but indicator function is hard to optimize. (NP-hard)

$$\xi_i \geq 0 \quad (\text{slack variable})$$

$L_0 \xrightarrow{\text{relax}} L_1$

Hinge loss: $\max\{0, 1 - ty\}$

($t = w^T x_i - b$ is the output).

Solve:

Primal - Dual

$$\min_w \|w\|_2 + \lambda \sum \xi_i$$

$$\text{s.t. } \forall i, y_i(w^T x_i) \geq 1 - \xi_i$$

$$\xi_i \geq 0$$

\Rightarrow

$$\max_a \sum_i a_i - \frac{1}{2} \sum_i \sum_j y_i y_j a_i a_j \langle x_i, x_j \rangle$$

$$\text{s.t. } \forall i, 0 \leq a_i \leq \frac{1}{2\lambda}$$

$$\sum_i y_i a_i = 0$$

$$w = \sum_i a_i x_i y_i$$

Kernel Method

Input Space (not separable)

$$x \longrightarrow \phi(x)$$

feature space (usually much higher dim.)
separable.

With SVM:

$$\min_w \|w\|_2 + \lambda \cdot \sum \xi_i$$

$$\text{s.t. } \forall i: y_i (w^T \phi(x_i)) \geq 1 - \xi_i \Rightarrow$$

$$\xi_i \geq 0$$

$$\max_a \sum_i a_i - \frac{1}{2} \sum_i \sum_j y_i y_j a_i a_j \langle \phi(x_i), \phi(x_j) \rangle$$

$$\text{s.t. } \forall i \quad 0 \leq a_i \leq \frac{1}{2n\lambda}$$

$$\sum_i y_i a_i = 0$$

$$\Rightarrow w = \sum_i a_i y_i \phi(x_i)$$

$\phi(x_i)$ hard to compute

but $\langle \phi(x_i), \phi(x_j) \rangle = K(x_i, x_j)$ may be computed easily

↳ kernel trick.

predict:

$$w^T \phi(x) = \sum_i a_i y_i \langle \phi(x_i), \phi(x) \rangle$$

⇒ no need for computing $\phi(x)$

Mercer's Theorem

$$K = \begin{bmatrix} K(x_1, x_1) & K(x_1, x_2) & \dots \\ K(x_2, x_1) & \ddots & \\ \vdots & & \end{bmatrix}$$

if K semi-definite for any $\{x_i\}$, then $\exists \phi$ such that K is a kernel for ϕ .

(We do not even need ϕ in computation!).