

## P28 Hard-Margin SVM - Model Definition

SVM three features: Margin; Duality; kernel

① hard-margin SVM ② soft-margin SVM ③ kernel SVM

$$f(w) = \text{sign}(w^T x + b) \quad \text{Let distance } (w, b, x_i) \text{ be the distance of } x_i \text{ to hyperplane } w^T x + b$$

$$\text{distance } (w, b, x_i) = \frac{1}{\|w\|} |w^T x_i + b|$$

Maximum interval classifier

$$\begin{cases} \hat{w}, \hat{b} = \underset{w, b}{\operatorname{argmax}} \text{margin}(w, b) \\ \text{margin}(w, b) = \min_{x_i \in N} \text{distance}(w, b, x_i) \end{cases} \rightarrow \text{The distance from the nearest point to the hyperplane among all points}$$

$$\text{S.t. } y_i(w^T x_i + b) > 0 \quad i=1 \dots N$$

$$\Rightarrow \begin{cases} \hat{w}, \hat{b} = \underset{w, b}{\operatorname{argmax}} \frac{1}{\|w\|} \min_{x_i \in N} y_i(w^T x_i + b) \\ \text{S.t. } y_i(w^T x_i + b) > 0 \Rightarrow \exists r > 0, \text{s.t. } \min_{x_i \in N} y_i(w^T x_i + b) = r \end{cases}$$

$\because 2w^T x + 2b$  and  $w^T x + b$  are the same plane

$\therefore r$  can take any value. After  $r$  is fixed,  $w, b$  are also determined.

如果不做限定，则 $w$ 可以任意缩放，无限种取值。所以此处只取能使 $r=1$ 的 $w$ 的值

$$\text{设 } r=1 \Rightarrow \begin{cases} \hat{w}, \hat{b} = \underset{w, b}{\operatorname{argmax}} \frac{1}{\|w\|} \\ \text{S.t. } \min_{x_i \in N} y_i(w^T x_i + b) = 1 \end{cases} \Rightarrow \begin{cases} \hat{w}, \hat{b} = \underset{w, b}{\operatorname{argmin}} \frac{1}{2} w^T w \\ \text{S.t. } y_i(w^T x_i + b) \geq 1, i=1 \dots N \end{cases}$$

## P29 Model Solution, Duality Problem

$$L(w, b, \lambda) = \frac{1}{2} w^T w + \sum_{i=1}^N \lambda_i [1 - y_i(w^T x_i + b)] \Rightarrow \begin{cases} \hat{w}, \hat{b} = \underset{w, b}{\operatorname{argmin}} \sum_{i=1}^N \lambda_i \\ \text{S.t. } \lambda_i \geq 0 \end{cases}$$

# Lagrange Multiplier Method hides the  $y_i(w^T x_i + b) \geq 1$  condition in the filtering process of finding minimum

$$\Rightarrow \begin{cases} \text{If } 1 - y_i(w^T x_i + b) > 0, \max_{\lambda} L(w, b, \lambda) = \frac{1}{2} w^T w + \infty = \infty \\ \text{If } 1 - y_i(w^T x_i + b) \leq 0, \max_{\lambda} L(w, b, \lambda) = \frac{1}{2} w^T w \\ \therefore \min_{w, b} \max_{\lambda} L(w, b, \lambda) = \min_{w, b} (\infty, \frac{1}{2} w^T w) = \min_{w, b} \frac{1}{2} w^T w \end{cases}$$

Duality Problem: Weak Duality:  $\min \max f \geq \max \min f$ Strong Duality:  $\min \max f = \max \min f$ 

$$\therefore \begin{cases} \hat{\lambda} = \max_{\lambda} \min_{w, b} L(w, b, \lambda) \\ \text{S.t. } \lambda_i \geq 0, i=1 \dots N \end{cases}$$

$$\frac{\partial L}{\partial b} = -\sum_{i=1}^N \lambda_i y_i = 0 \quad \text{At } \lambda \text{ in } L(w, b, \lambda) \Rightarrow \text{Find it with support vectors in the next section.}$$

$$\frac{\partial L}{\partial w} = w - \sum_{i=1}^N \lambda_i y_i x_i = 0 \quad \hat{w} = \sum_{i=1}^N \lambda_i y_i x_i \quad \text{At } \lambda \text{ in } L(w, b, \lambda)$$

$$\begin{aligned} \therefore L(w, b, \lambda) &= \frac{1}{2} \left( \sum_{i=1}^N \lambda_i y_i x_i \right)^T \left( \sum_{j=1}^N \lambda_j y_j x_j \right) + \sum_{i=1}^N \lambda_i - \sum_{i=1}^N \lambda_i y_i \left( \sum_{j=1}^N \lambda_j y_j x_j \right)^T x_i \\ &= \frac{1}{2} \sum_{i=1}^N \sum_{j=1}^N \lambda_i y_i y_j x_i^T x_j + \sum_{i=1}^N \lambda_i - \sum_{i=1}^N \sum_{j=1}^N \lambda_i \lambda_j y_i y_j x_i^T x_j \\ &= \sum_{i=1}^N \lambda_i - \frac{1}{2} \sum_{i=1}^N \sum_{j=1}^N \lambda_i \lambda_j y_i y_j x_i^T x_j \end{aligned}$$

$$\therefore \begin{cases} \hat{\lambda} = \underset{\lambda}{\operatorname{argmin}} \left( \sum_{i=1}^N \sum_{j=1}^N \lambda_i \lambda_j y_i y_j x_i^T x_j - \sum_{i=1}^N \lambda_i \right) \\ \text{S.t. } \lambda_i \geq 0, \sum_{i=1}^N \lambda_i y_i = 0 \quad i=1 \dots N \end{cases}$$

## P30 Model Solution, KKT Condition (Karush-Kuhn-Tucker condition)

$$\text{KKT condition: } \begin{cases} \frac{\partial L}{\partial b} = 0 = \frac{\partial L}{\partial w} = 0 \\ \lambda_i [1 - y_i(w^T x_i + b)] = 0 \end{cases}$$

?  $\lambda_i [1 - y_i(w^T x_i + b)] = 0 \rightarrow \text{Complementary Slackness (互补松弛)}$

$\lambda_i \geq 0, 1 - y_i(w^T x_i + b) \leq 0 \quad \text{Support Vectors: } \lambda_i \neq 0 ; \text{ Non-support Vectors: } \lambda_i = 0$

Find  $\hat{b}$ :  $\exists (x_k, y_k), \text{s.t. } 1 - y_k(w^T x_k + b) = 0 \quad (\text{i.e. support vectors})$

$$\therefore y_k(w^T x_k + b) = 1 \Rightarrow y_k^2 (w^T x_k + b) = y_k \quad \because y_k = \pm 1 \quad \therefore w^T x_k + b = y_k$$

$$\therefore \hat{w} = \sum_{i=1}^N \lambda_i y_i x_i \quad \therefore \hat{b} = y_k - \sum_{i=1}^N \lambda_i y_i x_i^T x_k$$

## P31 Soft-margin SVM

Allow classifier to make mistakes:  $\hat{w}, \hat{b} = \underset{w, b}{\operatorname{argmin}} \frac{1}{2} w^T w + \text{loss}$ ① Count: loss =  $\sum_{i=1}^N \{y_i(w^T x_i + b) < 1\}$  But the function isn't continuous② Distance: If  $y_i(w^T x_i + b) \geq 1$ , loss<sub>i</sub> = 0; If  $y_i(w^T x_i + b) < 1$ , loss<sub>i</sub> =  $1 - y_i(w^T x_i + b)$ 

$$\therefore \text{loss}_i = \max \{0, 1 - y_i(w^T x_i + b)\}$$

$$\text{Simplify to } \Rightarrow g_i = 1 - y_i(w^T x_i + b), g_i \geq 0 \quad \begin{cases} \hat{w}, \hat{b} = \underset{w, b}{\operatorname{argmin}} \frac{1}{2} w^T w + C \sum_{i=1}^N g_i \\ \text{S.t. } y_i(w^T x_i + b) \geq 1 - g_i, g_i \geq 0 \quad i=1 \dots N \end{cases}$$

## P32 Weak Duality

Generalized constrained optimization problem:  $\min_x f(x)$ 

$$\text{S.t. } m_i(x) < 0, n_j(x) = 0 \quad i=1 \dots M, j=1 \dots N$$

$$L(x, \lambda, \eta) = f(x) + \sum_{i=1}^M \lambda_i m_i(x) + \sum_{j=1}^N \eta_j n_j(x) \quad \Rightarrow \begin{cases} \min_{x, \lambda, \eta} L(x, \lambda, \eta) \\ \text{S.t. } \lambda_i \geq 0 \quad i=1 \dots M \end{cases} \approx \text{Unconstrained form of the original problem}$$

# If  $m_i(x) > 0$ , coefficient is negative

Duality Problem  $\leq$  Original Problem

$\#$  If  $m_i(x) > 0$ , coefficient is negative

Duality Problem  $\leq$  Original Problem

$$\max_{\lambda, \eta} \min_x L(x, \lambda, \eta) \leq \min_x \max_{\lambda, \eta} L(x, \lambda, \eta)$$

$$\text{证: } \min_x L(x, \lambda, \eta) \leq L(x, \lambda, \eta) \leq \max_{\lambda, \eta} L(x, \lambda, \eta)$$

$x$  is constant now, i.e.  $A(\lambda, \eta) \leq B(x)$

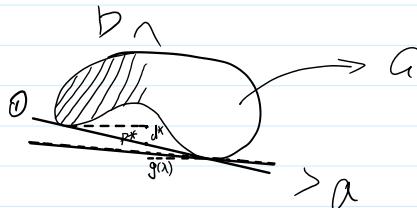
$$\therefore \max A(\lambda, \eta) \leq \min B(x)$$

### P33 Duality Problem - Geometric Interpretation

Simplified optimization problem:  $\begin{cases} \min f(x) \\ \text{s.t. } m(x) \leq 0 \end{cases} \Rightarrow L(x, \lambda) = f(x) + \lambda m(x), \lambda \geq 0$

$$P^* = \min f(x), m(x) \leq 0 \quad (\text{Original problem optimum solution}) \quad d^* = \max \min_x L(x, \lambda) \quad (\text{Duality problem optimum solution})$$

$$G = \{ (m(x), f(x)) \mid x \in D \} \quad \text{设 } m(x) = a, f(x) = b$$



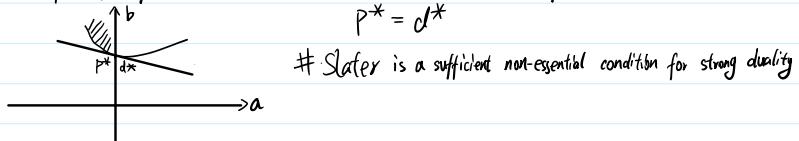
$$P^* = \inf \{ b \mid (a, b) \in G, a \leq 0 \} \quad (\inf \text{ is infimum (下确界)})$$

$$d^* = \max \min_x L(x, \lambda) \quad \text{Set } g(\lambda) = \min_x (b + \lambda a) \quad d^* = \max g(\lambda)$$

$$\underbrace{g(\lambda) = \inf \{ b + \lambda a \mid (a, b) \in G \}}_{\text{i.e. the minimum of the intersection of the line } b + \lambda a = ? \text{ and } G}$$

$d^*$  is the point that maximum  $g(\lambda)$  when  $\lambda$  (slope) changes, i.e. line O

$\therefore d^* \leq P^*$  When the part of the function G that intersects the b-axis is a convex function + Slater condition



### P34 Slater Condition

$$\begin{cases} \min f(x) \\ \text{s.t. } m_i(x) \leq 0, i=1 \dots M \\ n_j(x) = 0, j=1 \dots N \end{cases} \quad D = \{ \text{dom } f \cap \text{dom } \bigcap_{i=1}^M m_i \cap \bigcap_{j=1}^N n_j \}$$

Domain of definition

Slater condition:

$\exists \bar{x} \in \text{relint } D \quad$  The set is the inner part after removing the boundary  $\Rightarrow$

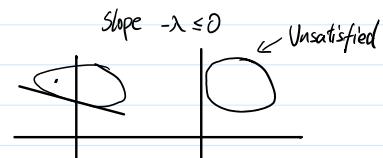
S.t.  $\forall i=1 \dots M, m_i(\bar{x}) < 0 \Rightarrow$  The set G must have points of  $m_i(x) < 0$

① Most convex optimization problems satisfy Slater ② Affine function no need to proof  $m_i(\bar{x}) < 0$

# Affine function: The polynomial function with the highest order is 1 e.g.  $f(x) = Ax + b$

Affine function with a constant term of zero is Linear function e.g.  $f(x) = Ax$

From ②, if  $f(x)$  is a convex function and  $m_i, n_j$  are affine functions, then it satisfies Slater.  $\therefore$  Satisfy strong duality



### P35 KKT Condition

$$\begin{cases} \min f(x) \\ \text{s.t. } m_i(x) \leq 0, i=1 \dots N \\ n_j(x) = 0, j=1 \dots M \end{cases} \quad \begin{aligned} L(x, \lambda, \eta) &= f(x) + \sum_{i=1}^N \lambda_i m_i(x) + \sum_{j=1}^M \eta_j n_j(x) \\ g(\lambda, \eta) &= \min_x L(x, \lambda, \eta) \end{aligned}$$

Convex + Slater  $\Rightarrow$  Strong Duality  $\Leftrightarrow$  KKT

$$\text{KKT} \left\{ \begin{array}{l} \text{Passable condition: } \begin{cases} m_i(x^*) \leq 0 \\ n_j(x^*) = 0 \\ \lambda^* \geq 0 \end{cases} \\ \text{complementary slackness: } \lambda^* m_i(x^*) = 0 \end{array} \right.$$

optimal solution  
 $\downarrow$   
Set  $d^* \rightarrow \lambda^*, \eta^*$   
 $\downarrow$   
 $p^* \rightarrow x^*$

gradient is 0:  $\frac{\partial L(x, \lambda^*, \eta^*)}{\partial x} \Big|_{x=x^*} = 0$   $\Rightarrow$  It's also used in SVM!

Proof KKT by Strong Duality Support vectors with  $\lambda=0$ , and non-support vectors with  $\lambda \neq 0$

$$d^* \stackrel{?}{=} \max_{\lambda, \eta} g(\lambda, \eta) \stackrel{?}{=} g(\lambda^*, \eta^*) = \min_x L(x, \lambda^*, \eta^*) \stackrel{?}{\leq} L(x^*, \lambda^*, \eta^*) \geq$$

$$= f(x^*) + \sum_{i=1}^N \lambda^* m_i(x^*) + \sum_{j=1}^M \eta^* n_j(x^*) \quad \because \lambda^* m_i(x^*) \leq 0 \quad n_j(x^*) = 0$$

Original function  $\stackrel{?}{\geq} f(x^*) \stackrel{?}{\leq} d^*$

$$d^* \triangleq \min_{\lambda, \eta} g(\lambda, \eta) \triangleq g(x^*, \eta^*) = \max_{x, \lambda, \eta} L(x, \lambda^*, \eta^*) \triangleq L(x^*, \lambda^*, \eta^*) \geq$$
$$= f(x^*) + \sum_i \lambda^* m_i(x) + \sum_j \eta^* n_j(x) \quad \because \lambda^* m_i(x) \leq 0 \quad n_j(x) = 0$$

$\therefore$  Original function  $\stackrel{\text{③}}{\leq} f(x^*) \stackrel{\text{由①}}{=} p^*$

$\therefore$  Strong Duality  $\Rightarrow d^* = p^* \quad \therefore$  equations at ①③ hold

$\therefore$  ②  $\max_x L(x, \lambda^*, \eta^*) = L(x^*, \lambda^*, \eta^*) \quad \because L(x, \lambda^*, \eta^*)$  takes the minimum value at  $x^*$ , it can be proved the gradient is 0

$\therefore$  ③  $f(x^*) + \sum_i \lambda^* m_i(x) = f(x^*) \quad \therefore$  Complementary slackness can be proved