

## P28 Hard-Margin SVM - Model Definition

SVM three features: Margin; Duality; Kernel

① Hard-margin SVM ② Soft-margin SVM ③ Kernel SVM

 $f(w) = \text{sign}(w^T x + b)$  Let distance  $(w, b, x_i)$  be the distance of  $x_i$  to hyperplane  $w^T x + b$ distance  $(w, b, x_i) = \frac{1}{\|w\|} |w^T x_i + b|$ Maximum interval classifier  $\begin{cases} \hat{w}, \hat{b} = \arg \max_{w, b} \text{margin}(w, b) \\ \text{margin}(w, b) = \min_{i=1 \dots N} \text{distance}(w, b, x_i) \end{cases}$  $\Rightarrow$  The distance from the nearest point to the hyperplane among all points  
s.t.  $y_i(w^T x_i + b) > 0 \quad i=1 \dots N$ 

$$\Rightarrow \begin{cases} \hat{w}, \hat{b} = \arg \max_{w, b} \frac{1}{\|w\|} \min_{i=1 \dots N} y_i(w^T x_i + b) \\ \text{s.t. } y_i(w^T x_i + b) > 0 \Rightarrow \exists r > 0, \text{s.t. } \min_{i=1 \dots N} y_i(w^T x_i + b) = r \end{cases}$$

 $\therefore 2w^T x + 2b$  and  $w^T x + b$  are the same plane $\therefore r$  can take any value. After  $r$  is fixed,  $w, b$  are also determined.

$$\text{设 } r=1 \Rightarrow \begin{cases} \hat{w}, \hat{b} = \arg \max_{w, b} \frac{1}{\|w\|} \\ \text{s.t. } \min_{i=1 \dots N} y_i(w^T x_i + b) = 1 \end{cases} \Rightarrow \begin{cases} \hat{w}, \hat{b} = \arg \min_{w, b} \frac{1}{2} w^T w \\ \text{s.t. } y_i(w^T x_i + b) \geq 1, i=1 \dots N \end{cases}$$

## P29 Model Solution, Duality Problem

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

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

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

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

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

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

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

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

$$\therefore \begin{cases} \hat{\lambda} = \arg \min_{\lambda} \left( \sum_{i=1}^N \lambda_i - \frac{1}{2} \sum_{i=1}^N \sum_{j=1}^N \lambda_i \lambda_j y_i y_j x_i^T x_j \right) \\ \text{s.t. } \lambda_i \geq 0, \sum_{i=1}^N \lambda_i y_i = 0 \quad i=1 \dots N \end{cases}$$

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

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

Complementary Slackness (互补松弛)

 $\lambda_i \geq 0, 1 - y_i(w^T x_i + b) \leq 0$  Support Vectors:  $\lambda \neq 0$ ; Non-support Vectors:  $\lambda = 0$ Find  $\hat{b}$ :  $\exists (x_k, y_k)$ , s.t.  $1 - y_k(w^T x_k + b) = 0$  (i.e. support vectors)

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

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

## P31 Soft-margin SVM

Allow classifier to make mistakes:  $\hat{w}, \hat{b} = \arg \min_{w, b} \frac{1}{2} w^T w + \text{loss}$ ① Count:  $\text{loss} = \sum_{i=1}^N \mathbb{I}\{y_i(w^T x_i + b) < 1\}$  But the function is not continuous② Distance: If  $y_i(w^T x_i + b) \geq 1$ ,  $\text{loss}_i = 0$ ; If  $y_i(w^T x_i + b) < 1$ ,  $\text{loss}_i = 1 - y_i(w^T x_i + b)$ 

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

Simplify to  $\Rightarrow \xi_i = 1 - y_i(w^T x_i + b), \xi_i \geq 0$ 

$$\begin{cases} \hat{w}, \hat{b} = \arg \min_{w, b} \frac{1}{2} w^T w + C \sum_{i=1}^N \xi_i \\ \text{s.t. } y_i(w^T x_i + b) \geq 1 - \xi_i, \xi_i \geq 0 \quad i=1 \dots N \end{cases}$$

## P32 Weak Duality

Generalized constrained optimization problem:  $\begin{cases} \min_x f(x) \\ \text{s.t. } m_i(x) \leq 0, n_j(x) = 0 \quad i=1 \dots M, j=1 \dots N \end{cases}$ 

$$\mathcal{L}(x, \lambda, \eta) = f(x) + \sum_{i=1}^M \lambda_i m_i(x) + \sum_{j=1}^N \eta_j n_j(x)$$

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

$$\Rightarrow \begin{cases} \min_x \max_{\lambda, \eta} \mathcal{L}(x, \lambda, \eta) \\ \text{s.t. } \lambda_i \geq 0 \quad i=1 \dots M \end{cases} \approx \text{Unconstrained form of the original problem}$$

In fact there are still constraints

Duality Problem  $\leq$  Original Problem

$\alpha_1, \alpha_2, \dots, \alpha_M = 1, \dots, M, \alpha_i \geq 0$   
 # If  $m_i(x) > 0$ , coefficient is negative  $\left\{ \begin{array}{l} x \text{ s.t. } \lambda_i \geq 0 \quad i=1 \dots M \\ \text{In fact there are still constraints} \end{array} \right.$

Duality Problem  $\leq$  Original Problem

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

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

$x$  is constant now, i.e.  $A(\lambda, \eta) \leq B(x)$   
 $\therefore \max_{\lambda, \eta} A(\lambda, \eta) \leq \min_x B(x)$

### P33 Duality Problem - Geometric Interpretation

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

$p^* = \min f(x), m(x) \leq 0$  (Original problem optimum solution)  $d^* = \max_{\lambda} \min_x \mathcal{L}(x, \lambda)$  (Duality problem optimum solution)  
 $G = \{ (m(x), f(x)) \mid x \in D \}$  设  $m(x) = a, f(x) = b$



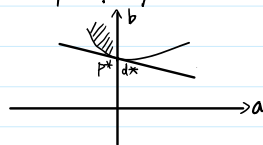
$p^* = \inf \{ b \mid (a, b) \in G, a \leq 0 \}$  (inf is infimum (下确界))

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

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

$d^*$  is the point that maximum  $g(a)$  when  $\lambda$  (slope) changes, i.e. line ①

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



$p^* = d^*$

# Slater is a sufficient non-essential condition for strong duality

### 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:

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

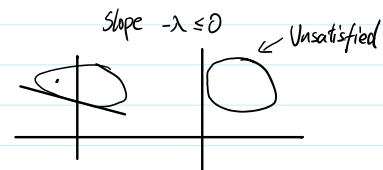
s.t.  $\forall i=1 \dots M, m_i(\tilde{x}) < 0 \Rightarrow$  The set  $G$  must have points of  $m_i(x) < 0$

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

# 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

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

$$g(\lambda, \eta) = \min_x \mathcal{L}(x, \lambda, \eta)$$

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

KKT

$$\begin{cases} \text{Passable condition: } \begin{cases} m_i(x^*) \leq 0 \\ n_j(x^*) = 0 \\ \lambda^* \geq 0 \end{cases} \\ \text{complementary slackness: } \lambda_i m_i = 0 \\ \text{gradient is 0: } \frac{\partial \mathcal{L}(x, \lambda^*, \eta^*)}{\partial x} \Big|_{x=x^*} = 0 \end{cases}$$

optimal solution optimal value  
 Set  $d^* \Rightarrow \lambda^*, \eta^*$   
 $p^* \Rightarrow x^*$

Proof KKT by Strong Duality

$$d^* \triangleq \max_{\lambda, \eta} g(\lambda, \eta) \triangleq g(\lambda^*, \eta^*) = \min_x \mathcal{L}(x, \lambda^*, \eta^*) \triangleq \mathcal{L}(x^*, \lambda^*, \eta^*) \triangleq$$

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

$$\therefore \text{Original function } \mathcal{L}(x^*) \triangleq \mathcal{L}(x^*, \lambda^*, \eta^*) \triangleq d^*$$

It's also used in SVM: Support vectors with  $\lambda = 0$ , and non-support vectors with  $\lambda \neq 0$

$$d^* \triangleq \min_{\lambda, \eta} y(\lambda, \eta) \triangleq y(\lambda^*, \eta^*) = \min_x L(x, \lambda^*, \eta^*) \triangleq L(x^*, \lambda^*, \eta^*) \geq$$

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

$$\therefore \text{Original function} \stackrel{\text{②}}{\leq} f(x^*) \stackrel{\text{①}}{=} p^*$$

$\therefore$  Strong Duality  $\Rightarrow d^* = p^*$   $\therefore$  equations at ② hold

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

$\therefore$  ③  $f(x^*) + \sum_{i=1}^M \lambda_i^* m_i(x) = f(x^*)$   $\therefore$  Complementary slackness can be proved