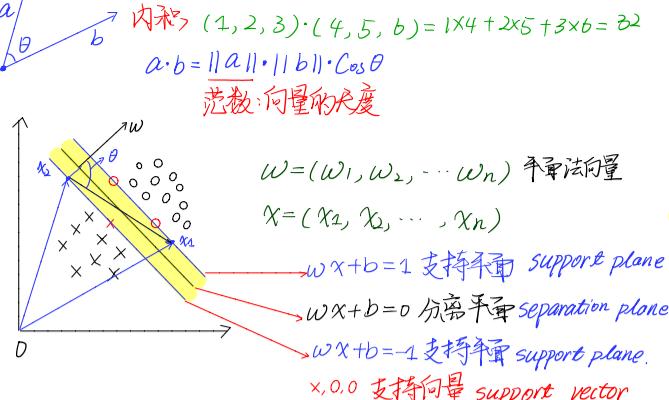


## (1) Linear Separable Condition SVM



Compute the distance between 2 support plane

$$\begin{aligned} & \text{(1)} \\ & \text{(2)} \Rightarrow w(x_1 - x_2) = 2 \Rightarrow \|w\| \cdot \|x_1 - x_2\| \cos \theta = 2 \Rightarrow d = \frac{2}{\|w\|} \end{aligned}$$

间隔要求最大 max

Decision function  $\left\{ \begin{array}{l} w^T x_i + b \geq 1, y_i = 1 \\ \text{上支持平面以上} \\ w^T x_i + b \leq -1, y_i = -1 \\ \text{下支持平面以下} \end{array} \right\}$  Constraint

### (4) Convex Optimization Problem

$$\min_{w} \frac{\|w\|^2}{2}, \text{ subject to } y_i(w^T x_i + b) \geq 1, i=1, \dots, N$$

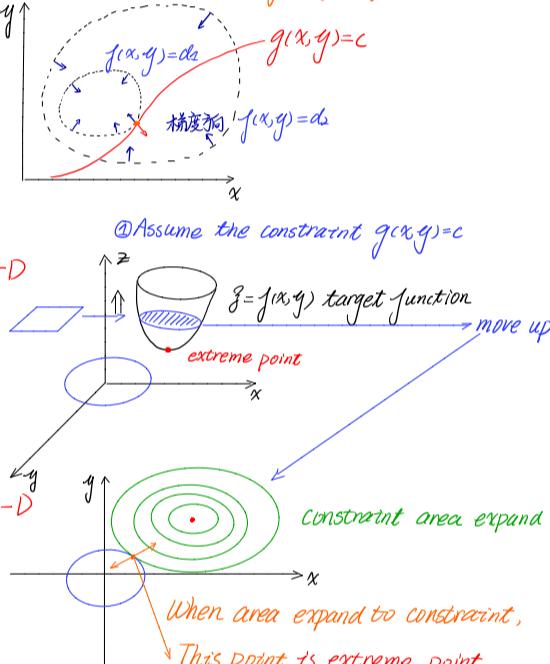
#### (1) Lagrange Multiplier Method

$$L(w, b, \alpha) = \frac{1}{2} \|w\|^2 - \sum_{i=1}^n \alpha_i [y_i(w^T x_i + b) - 1]$$

$$\frac{\partial L}{\partial w} = 0 \Rightarrow w = \sum_{i=1}^n \alpha_i y_i x_i$$

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

#### Principle of Lagrange



$$\min f(x, y), \text{ s.t. } g(x, y) = c$$

$$L(x, y) = f(x, y) + \lambda(g(x, y) - c)$$

Lagrange multiplier

$$\nabla L(x, y) = \nabla f + \lambda \nabla g$$

与约束首次接触的点即为极值点

$$\Rightarrow \nabla f = -\lambda \nabla g$$

### (2) KKT 最优化: Lagrange multiplier method 的推广

True for Lagrange multiplier method,

$$\min_{w} \frac{\|w\|^2}{2}, \text{ subject to } y_i(w^T x_i + b) \geq 1, i=1, \dots, N$$

So, introduce KKT 解带有双等式约束的凸优化问题

$$\begin{aligned} \min f(x) \quad & \text{s.t.} \quad h_j(x) = 0, j=0, 1, \dots, p \\ & \text{受制于} \quad g_k(x) \leq 0, k=0, \dots, q \\ & x \in X \subset R^n \end{aligned}$$

Can convert it into

$$\begin{aligned} & \text{1. } h_j(x) = 0, j=0, \dots, p \quad \text{2. } \nabla f(x) + \sum_{j=1}^p \lambda_j \nabla h_j(x) + \sum_{k=1}^q \mu_k \nabla g_k(x) = 0 \\ & g_k(x) \leq 0, k=0, \dots, q \quad \lambda_j \geq 0, \mu_k \geq 0, \mu_k g_k(x) = 0 \end{aligned}$$

### (3) Take result of gradient from Lagrangian method into original function $\Rightarrow$ Dual problem

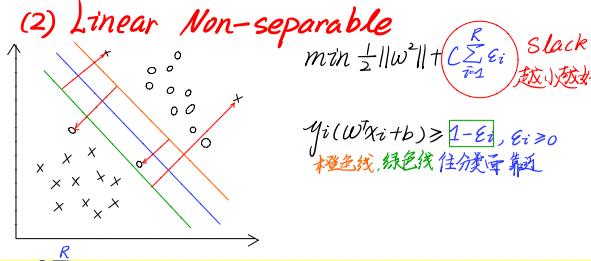
$$\text{Finally } L(w, b, \alpha) = \frac{1}{2} \sum_{i=1}^n \alpha_i - \frac{1}{2} \sum_{i=1}^n \sum_{j=1}^p \alpha_i \alpha_j y_i y_j (x_i^T x_j)$$

$$\text{s.t. } \alpha_i \geq 0, i=1, \dots, n$$

$$\sum_{i=1}^n \alpha_i y_i = 0$$

Assume  $\alpha = \{\alpha_i\}, i=1, \dots, n$  is solution of,

$$w = \sum_{i=1}^n \alpha_i y_i x_i; b = y_j - \sum_{i=1}^n \alpha_i y_i (x_i \cdot x_j)$$



(1)  $C \sum_{i=1}^n \epsilon_i$  is a penalty item based on linear non-separable condition.

1. When  $x_i$  in right side,  $\epsilon_i = 0, R = \text{size}(X), C$  is assigned by user, 表示对错的高加多少的惩罚。

2. When  $C$  is large, 错分少; When  $C$  is small, 错分多, 多有 overfitting 模型又太正确

### (2) Dual problem in linear non-separable condition

$$\max_{\alpha} \sum_{i=1}^n \alpha_i - \frac{1}{2} \sum_{i=1}^n \sum_{j=1}^n \alpha_i \alpha_j y_i y_j x_i^T x_j \quad \text{s.t. } C \geq \alpha_i \geq 0, i=1, \dots, n$$

$$\sum_{i=1}^n \alpha_i y_i = 0$$

### (3) SMO Algorithm (Sequential Minimal optimization) 角 lagrange dual

1. Update 2 multipliers each time, recursion  $\Rightarrow$  get result.

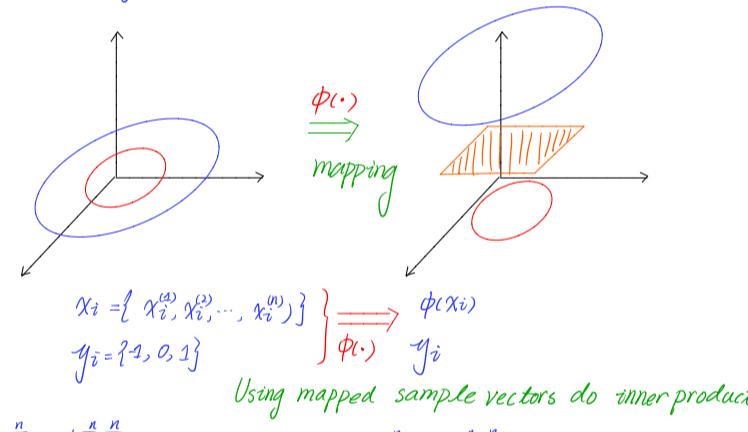
Repeat till convergence?

1. Select some pair  $\alpha_i, \alpha_j$  to update next! Using a heuristic that tries to pick the two that will allow us to make the biggest progress towards the global maximum

2. Reoptimize  $W, b$  with respect to  $\alpha_i, \alpha_j$ , while holding all the other  $\alpha_i$ 's fixed }

### (4) Non-linear Condition — Dimension transform

Convert linear non-separable data into linear separable data in a high dimension space.



If original space is  $n$ -dimensional, mapping to  $n^2$ -dimensional space. Time complexity from  $O(n)$   $\rightarrow$   $O(n^2)$

### (5) Kernel Function

① Example: Assume  $x, z$  are  $N$ -dimensional variable.

Using Kernel function  $K(x, z) = (x^T z)^2$

$$\text{Expand it: } K(x, z) = (x^T z)^2 = \left( \sum_{i=1}^n x_i z_i \right)^2 = \sum_{i=1}^n \sum_{j=1}^n x_i x_j z_i z_j$$

$$= \sum_{i=1}^n \sum_{j=1}^n (x_i z_j) (z_i z_j) = \phi(x)^T \phi(z)$$

So, we can compute the inner product of  $x, z$ , then  $(x^T z)^2$ , decreasing the time complexity

② If for any mapping function, we can find a kernel function,

making it possible that we could compute inner product in the mapped high dimensional space, into the low level dimensional space. we know  $\phi$ , and we can find the  $K$ ,

In fact, we can find a kernel function first, for example the problem in 3-D space is non-linear. Need to map the problem into a high-D space, making it a linear problem. we can replace the inner product of high level dimension with the value computed by the kernel function.

不用找  $\phi$ , 只需要在低维空间找出一个合理的核函数, 并在数学上证明: 一定存在  $\phi$  可以把低维空间中非线性映射成为高维空间的线性问题。

低维空间计算核函数的值可代替高维空间的内积。

### ③ Necessary and sufficient condition of kernel function

#### Mercer Theorem

If function  $K$  is a mapping from  $R^n \times R^n \rightarrow R$  (2 n-dimensional vector mapped to real number). If  $K$  is a mercer kernel function then if for training sample  $\{x_i\}_{i=1, \dots, N}$ , its corresponding kernel function matrix is symmetric semi-positive.

④ Kernel function matrix  $\begin{bmatrix} \cdots & \cdots \\ \cdots & \cdots \\ k_{ij} & = K(x_i, x_j) \end{bmatrix}$

⑤ Positive define matrix: Assume  $M$  is a  $n$ -th order matrix, for any non-zero vector  $z$ ,  $z^T M z > 0$ , calling  $M$  a positive define matrix

⑥ Semi-positive define matrix: Assume  $M$  is a  $n$ -th order matrix, for any non-zero vector  $z$ ,  $z^T M z \geq 0$ , calling  $M$  a positive define matrix

#### ⑦ Commonly used kernel function

$$\begin{aligned} & 1. K(x_i, x_j) = (x_i \cdot x_j)^k \\ & \text{RBF2. } K(x_i, x_j) = \exp(-\frac{\|x_i - x_j\|^2}{2\sigma^2}) \\ & 3. K(x_i, x_j) = \tanh(k x_i \cdot x_j - \delta) \end{aligned} \quad \left. \begin{array}{l} \text{逼近 Approach} \\ \text{逼近 Approach} \end{array} \right\}$$

#### (b) Using kernel function in lagrange function

$$\begin{aligned} & \text{① maximize } L(\alpha) = \sum_{i=1}^n \alpha_i - \frac{1}{2} \sum_{i=1}^n \sum_{j=1}^n \alpha_i \alpha_j y_i y_j \phi(x_i)^T \phi(x_j) \\ & \text{s.t. } \sum_{i=1}^n \alpha_i y_i = 0 \end{aligned}$$

低维空间中计算: 此时未知数为  $\alpha$ , 可用 SMO 解出。

#### ② Using kernel function to do classification

$$w^T x + b = \left( \sum_{i=1}^n \alpha_i y_i \phi(x_i) \right)^T x + b = \sum_{i=1}^n \alpha_i y_i K(x_i, x) + b$$

#### (7) Algorithm

Input: Linear separable training data set

$$D = \{(x_i, y_i)\}, i=1, \dots, N; x_i \in \mathcal{X} = R^n; y_i \in \{-1, 1\}, i=1, \dots, N$$

Output: Maximum margin separation hyperplane

#### ① Linear separable SVM — margin maximum

$$\begin{aligned} \text{step(1): } & \min_{w, b} \frac{1}{2} \|w\|^2; \text{s.t. } y_i(w^T x_i + b) - 1 \geq 0, i=1, \dots, N \\ & \text{compute optimal solution } w^*, b^* \end{aligned}$$

$$\text{step(2) Hyperplane } w^* x + b^* = 0$$

$$\text{Decision function } f(x) = \text{sign}(w^* x + b^*)$$

#### ② Linear separable SVM

$$\begin{aligned} \text{Step(1)} & \min_{\alpha} \frac{1}{2} \sum_{i=1}^n \sum_{j=1}^n \alpha_i \alpha_j y_i y_j (x_i^T x_j) - \sum_{i=1}^n \alpha_i \\ & \text{s.t. } \sum_{i=1}^n \alpha_i y_i = 0; \alpha_i \geq 0, i=1, \dots, N \end{aligned}$$

compute the optimal solution  $\alpha^* = (\alpha_i^*)^T, i=1, \dots, N$

$$\text{Step(2)} \quad \text{Compute } w^* = \sum_{i=1}^n \alpha_i^* y_i x_i$$

choose any  $\alpha_i^* > 0$  in  $\alpha^*$ , and compute

$$b^* = y_j - \sum_{i=1}^n \alpha_i^* y_i (x_i \cdot x_j)$$

$$\text{Step(3) Hyperplane } w^* x + b = 0$$

$$\text{Decision function } f(x) = \text{sign}(w^* x + b)$$

#### ③ Linear SVM

step(1) Choose penalty parameter  $C$ , construct convex quadratic programming problem

$$\min_{\alpha} \frac{1}{2} \sum_{i=1}^n \sum_{j=1}^n \alpha_i \alpha_j y_i y_j (x_i^T x_j) - \sum_{i=1}^n \alpha_i; \text{s.t. } \sum_{i=1}^n \alpha_i y_i = 0$$

$$\text{Solution } \alpha^* = (\alpha_i^*)^T, i=1, \dots, N$$

$$\text{Step(2) Compute } w^* = \sum_{i=1}^n \alpha_i^* y_i x_i, \text{ choose any } \alpha_i^* \in (0, C), \text{ compute}$$

$$b = y_j - \sum_{i=1}^n \alpha_i^* y_i (x_i \cdot x_j)$$

$$\text{Step(3) Hyperplane } w^* x + b = 0$$

$$\text{Decision function } f(x) = \text{sign}(w^* x + b)$$

#### ④ Non-linear SVM

step(1) choose kernel function  $K(x, z)$  and penalty  $C$

$$\begin{aligned} \min_{\alpha} & \frac{1}{2} \sum_{i=1}^n \sum_{j=1}^n \alpha_i \alpha_j y_i y_j K(x_i, x_j) - \sum_{i=1}$$