



# 第五章 支持向量机

§ 5.1 背景知识

§ 5.2 理论推导

§ 5.3 方法扩展



## § 5.1 背景知识

一、问题引入

二、间隔与支持向量



## □ 感知器方法回顾

### ■ 数据集:

$$\{(\mathbf{x}_1, y_1), (\mathbf{x}_2, y_2), \dots, (\mathbf{x}_N, y_N)\}$$

### ■ 模型:

$$\hat{y} = \text{sgn}(\mathbf{w}^T \mathbf{x} + w_0)$$

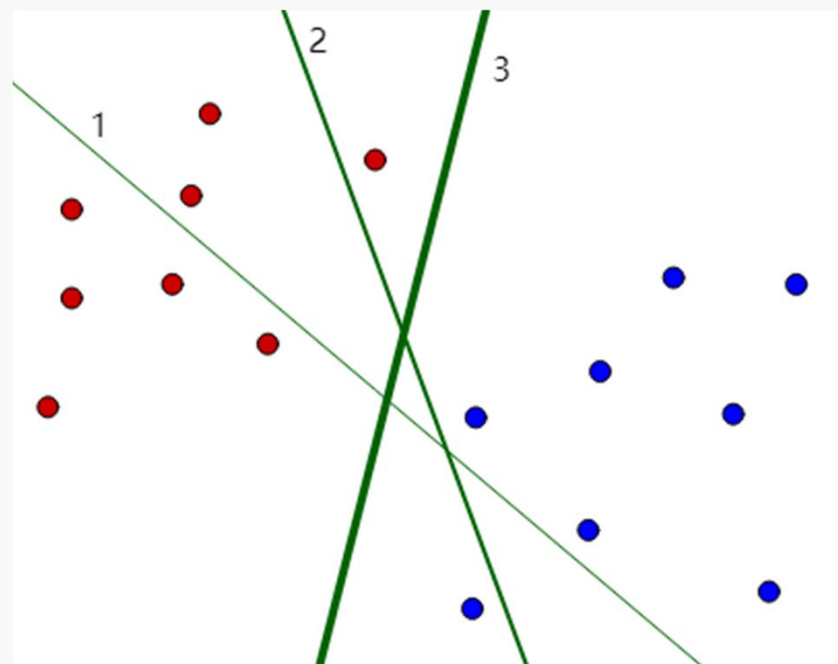
### ■ 惩罚函数:

$$J(\mathbf{w}, w_0) = \sum_{y_i(\mathbf{w}^T \mathbf{x}_i + w_0) \leq 0} -y_i(\mathbf{w}^T \mathbf{x}_i + w_0)$$

### ■ 优化方法: 梯度下降法

$$\mathbf{w}^{(t+1)} = \mathbf{w}^{(t)} - \rho \nabla_{\mathbf{w}} J(\mathbf{w}, w_0)$$

$$w_0^{(t+1)} = w_0^{(t)} - \rho \nabla_{w_0} J(\mathbf{w}, w_0)$$



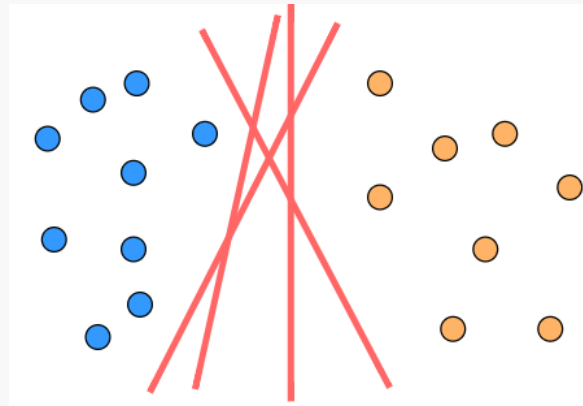
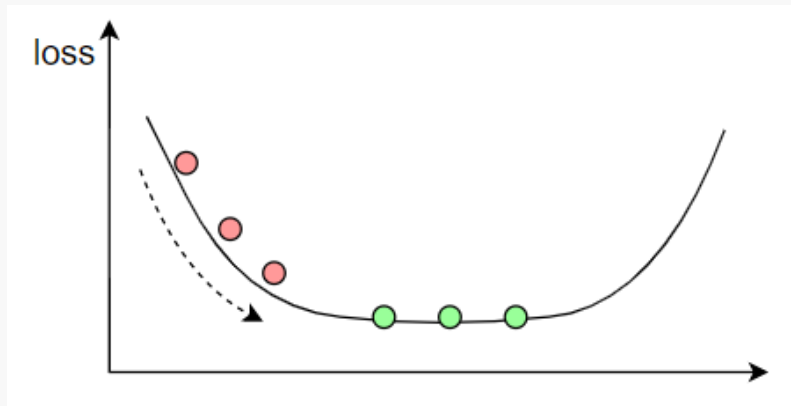
感知器收敛过程示意图



# 问题引入

## □ 感知器分类面特点

- 当训练数据线性可分时，能得到线性分类面
- **理论上：**满足优化目标的线性分类面通常有无数个
- **实际上：**受初始化、学习率等因素影响，会收敛到不同解

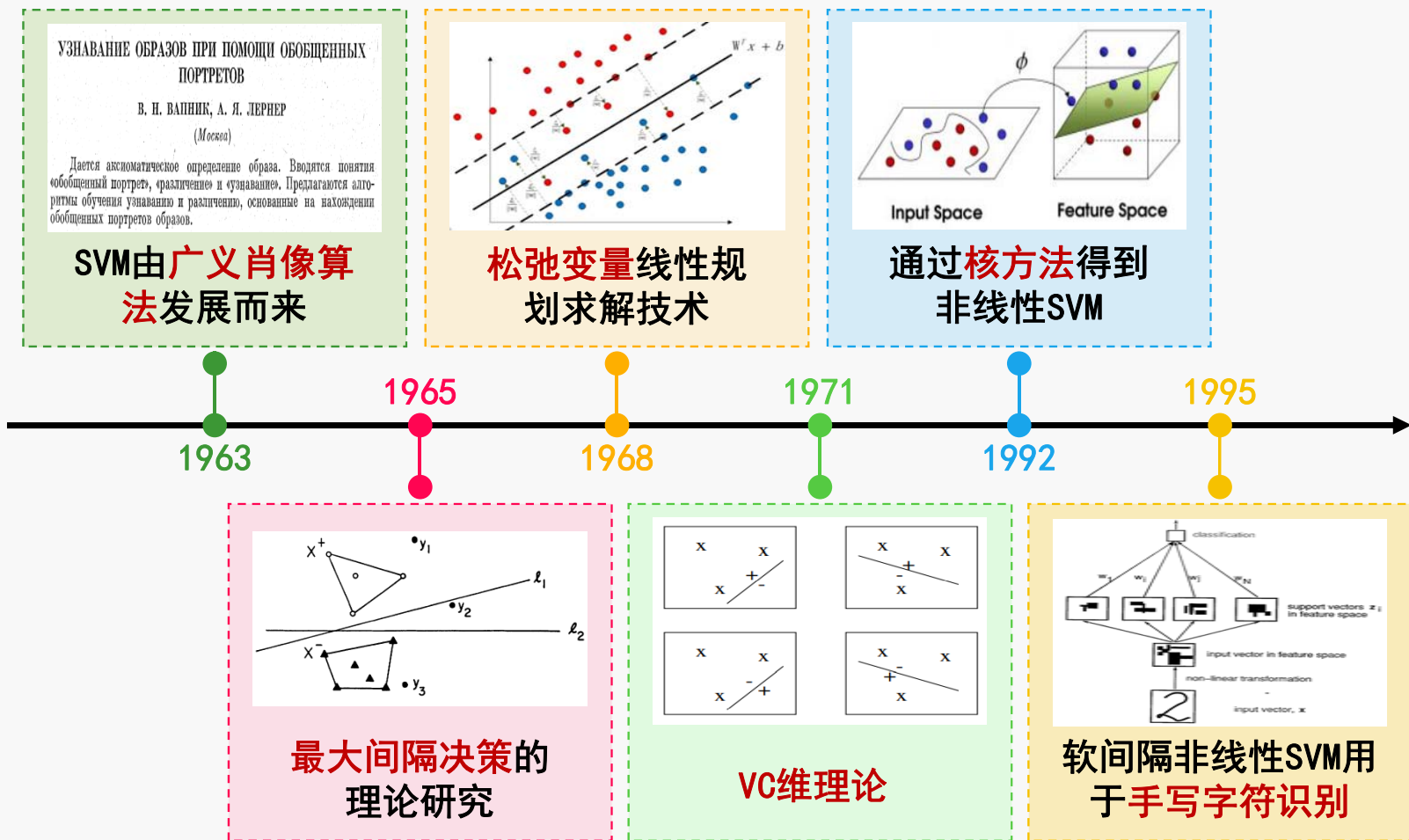


## □ 感知器是否存在最优分类面？

**不一定存在，多数情况下不存在，最优解不能保证**



# 问题引入



V. N. Vapnik



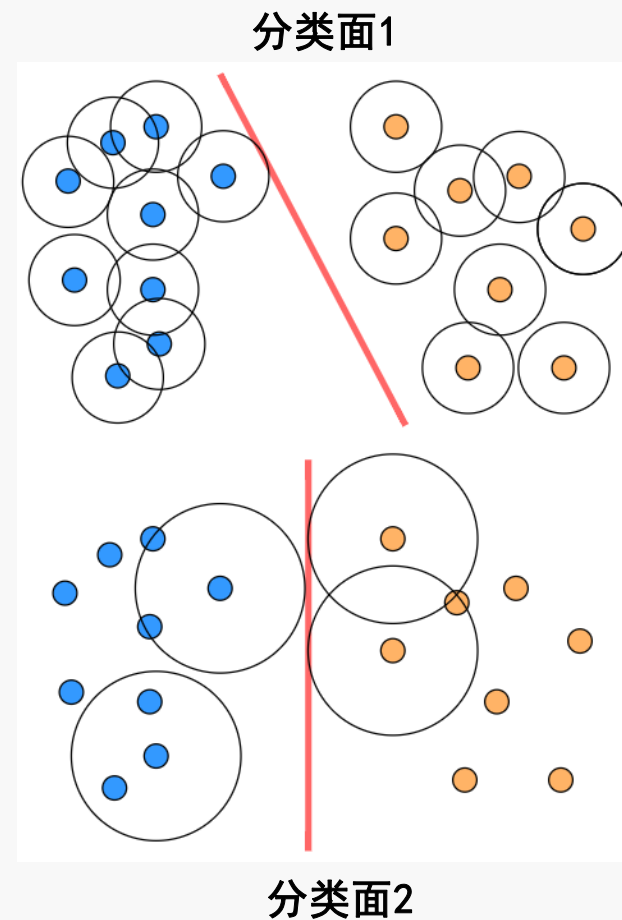
# 间隔与支持向量

## □ 线性分类面的优劣

- 一个好的分类面应该对局部扰动不敏感
- 右图中，以超球的半径代表扰动的大小。很明显，分类面2比分类面1能“容忍”更大的扰动
- 在线性可分的情况下，随着超球半径增大，最后可以得到唯一的线性分类面，能“容忍”最大程度的扰动

## □ 最优超平面

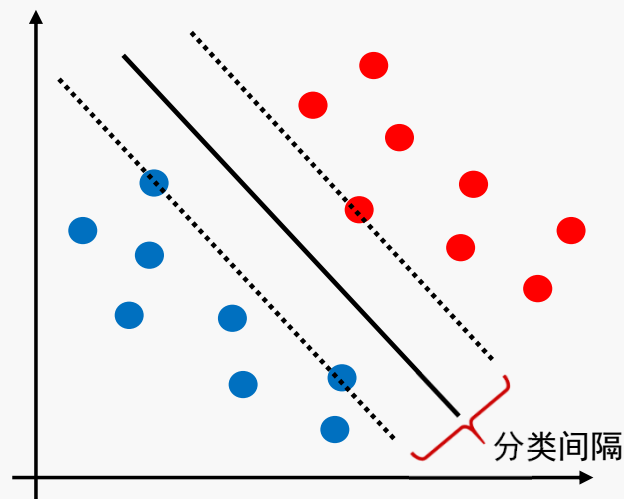
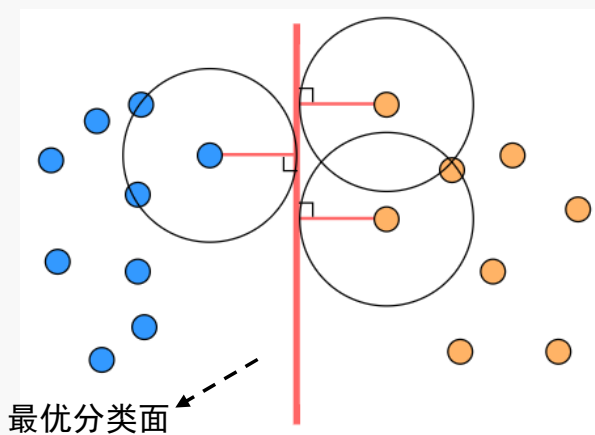
- 定义：一个超平面，如果它能够将训练样本没有错误地分开，并且**两类训练样本中离超平面最近的样本与超平面之间的距离是最大的**，则把这个超平面**称为最优分类超平面**





# 间隔与支持向量

- 支持向量：训练样本中离超平面最近的样本
- 分类间隔：两个异类支持向量到超平面的距离之和
- 最优超平面：一个超平面，如果它能够将训练样本没有错误地分开，并且两类训练样本中离超平面最近的样本与超平面之间的距离是最大，则把这个超平面称作**最优（分类）超平面**，或**最大间隔超平面**





# 间隔与支持向量

## □ 如何求解最优超平面呢？

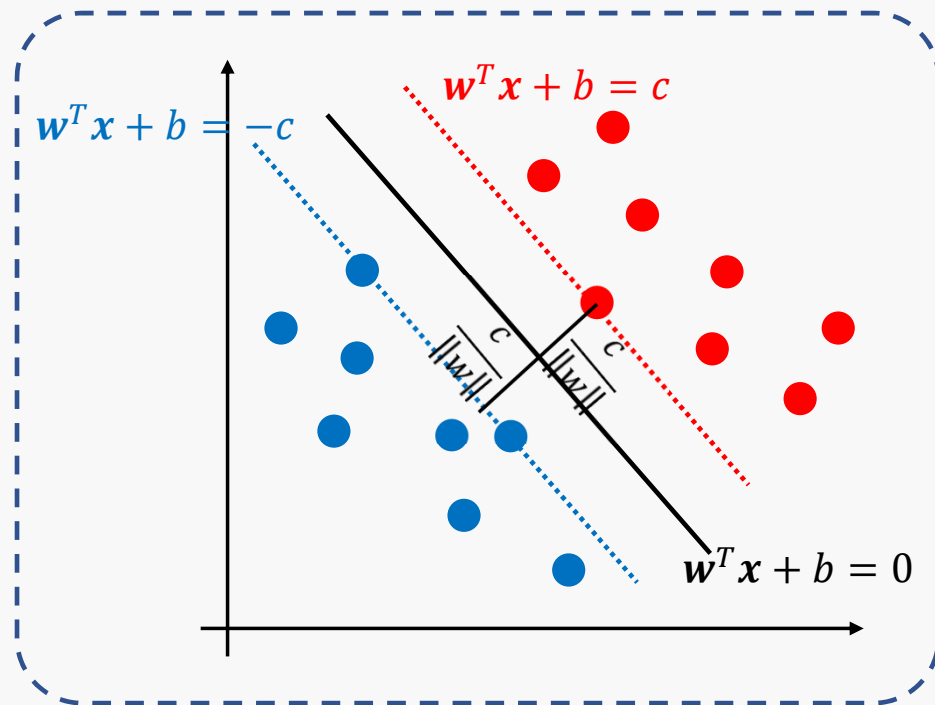
### ■ 准备工作

1. 最优超平面的必要条件：  
异类支持向量到超平面的距离相等
2. 同一个超平面可对应无数组  $(w, b, c)$   
为确定它们的尺度，不妨固定  $c = 1$

### ■ 数学形式：有约束的优化问题

$$\max_{w, b} \text{Margin}(w, b) = \frac{2c}{\|w\|_2} = \frac{2}{\|w\|_2}$$

$$\text{s.t. } y_i(w \cdot x_i + b) \geq c = 1, 1 \leq i \leq n$$



分类间隔为  $\frac{2c}{\|w\|}$ ，为什么？





## □ 求解最优超平面

$$\max_{w,b} \text{Margin}(w, b) = \frac{2}{\|w\|_2}$$

$$\text{s. t. } y_i(w \cdot x_i + b) \geq 1, 1 \leq i \leq n$$

■ 上述优化问题等价于下面的形式：

支持向量机的原问题：

$$\min_{w,b} \frac{1}{2} \|w\|_2^2$$

$$\text{s. t. } y_i(w \cdot x_i + b) \geq 1, 1 \leq i \leq n$$

注：  $\|w\|_2^2$  相比  $\|w\|_2^{-1}$  更容易优化，系数  $1/2$  仅是为了计算方便。



# SVM应用举例

## □ 利用线性SVM求解分类问题

$$\omega_1: \{(1, 1)^T, (0, 1)^T, (2, 0)^T, (0, 0)^T\}$$

$$\omega_2: \{(2, 2)^T, (2, 3)^T, (0, 3)^T\}$$

解答：通过观察法获取支持向量为

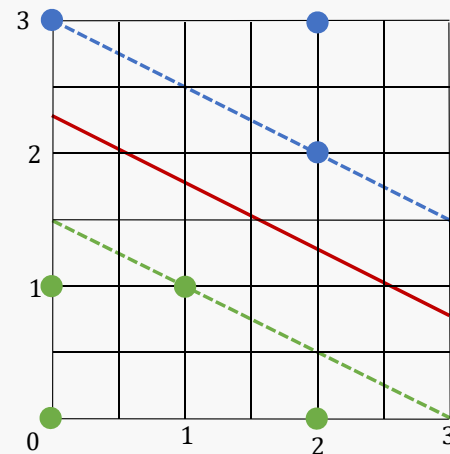
$$\omega_1: (1, 1)^T \quad \omega_2: (2, 2)^T, (0, 3)^T$$

列式计算

$$w^T \begin{pmatrix} 1 \\ 1 \end{pmatrix} + b = 1 \quad w^T \begin{pmatrix} 0 \\ 3 \end{pmatrix} + b = -1 \quad w^T \begin{pmatrix} 2 \\ 2 \end{pmatrix} + b = -1$$

解得

$$w = \begin{pmatrix} -\frac{4}{3} \\ \frac{2}{3} \end{pmatrix}, \quad b = 3, \quad -\frac{4}{3}x_1 - \frac{2}{3}x_2 + 3 = 0$$





## § 5.2 理论推导

一、对偶问题

二、核函数

三、软间隔与正则化



□ 考虑如下有约束优化问题：

$$\begin{aligned} \min_{\mathbf{w}} \quad & f(\mathbf{w}) \\ \text{s. t.} \quad & g_i(\mathbf{w}) \leq 0, i = 1, \dots, k \\ & h_i(\mathbf{w}) = 0, i = 1, \dots, l \end{aligned}$$

可行域 $\mathcal{R}$

□ 该问题的拉格朗日函数如下：

$$\begin{aligned} \mathcal{L}(\mathbf{w}, \boldsymbol{\alpha}, \boldsymbol{\beta}) = f(\mathbf{w}) + \sum_{i=1}^k \alpha_i g_i(\mathbf{w}) + \sum_{i=1}^l \beta_i h_i(\mathbf{w}) \\ \text{s. t.} \quad \alpha_i \geq 0, \quad i = 1, \dots, k \end{aligned}$$

其中 $\alpha$ 和 $\beta$ 为拉格朗日乘子

□ 进一步考虑如下函数：

$$\theta(\mathbf{w}) = \max_{\boldsymbol{\alpha}, \boldsymbol{\beta}: \alpha_i \geq 0} \mathcal{L}(\mathbf{w}, \boldsymbol{\alpha}, \boldsymbol{\beta})$$

说明 $\theta(\mathbf{w})$ 在 $\mathcal{R}$ 上的取值等于 $f(\mathbf{w})$ ,  
在 $\mathcal{R}$ 外取值为 $+\infty$



□ 如果 $w$ 的某个取值 $\hat{w}$ 不满足对 $g(w)$ 或 $h(w)$ 的约束条件, 即:

$$\exists i, \text{ s.t. } g_i(\hat{w}) > 0 \text{ or } h_i(\hat{w}) \neq 0$$

则 $\theta(w)$ 的取值满足:

$$\theta(\hat{w}) = \max_{\alpha, \beta: \alpha_i \geq 0} \left[ f(\hat{w}) + \sum_{i=1}^k \alpha_i g_i(\hat{w}) + \sum_{i=1}^l \beta_i h_i(\hat{w}) \right] = +\infty$$

由此易知:

$$\theta(w) = \begin{cases} f(w), & \text{如果 } w \text{ 满足约束} \\ +\infty, & \text{如果 } w \text{ 不满足约束} \end{cases}$$

□ 由以上讨论可知, 我们仅需求解  $\min_w \theta(w) = \min_w \max_{\alpha, \beta: \alpha_i \geq 0} \mathcal{L}(w, \alpha, \beta)$  即可。



## □ 上面两页总结：

$$\begin{aligned} \min_w f(w) \\ \text{s.t. } g_i(w) \leq 0, i = 1, \dots, k \\ h_i(w) = 0, i = 1, \dots, l \end{aligned} \quad \text{等价于求解} \quad \min_w \max_{\alpha, \beta: \alpha_i \geq 0} \mathcal{L}(w, \alpha, \beta)$$

## □ 在一定条件下，可以进一步得到：

$$\min_w \max_{\alpha, \beta: \alpha_i \geq 0} \mathcal{L}(w, \alpha, \beta) = \max_{\alpha, \beta: \alpha_i \geq 0} \min_w \mathcal{L}(w, \alpha, \beta)$$

且存在解 $(w^*, \alpha^*, \beta^*)$ 满足：

1. 使等式两边同时取得最值
2. 满足Karush-Kuhn-Tucker (KKT) 条件

很重要：想要的不是目标函数的最值，而是取最值时的参数值



□ 上页的“一定条件”包括：

- $f(\cdot)$ 和 $g_i(\cdot)$ 为凸函数
- $h_i(\cdot)$ 为线性函数
- $\forall i, \exists w, \text{ s.t. } g_i(w) < 0$

□ 在SVM原问题中：

- $f(w, b) = \frac{1}{2} \|w\|_2^2$ 为凸函数
- $g_i(w, b) = 1 - y_i(w \cdot x_i + b)$ 为凸函数
- 不存在等式约束条件
- 对任意样本 $(x_i, y_i)$ ，显然存在 $(w, b)$ 满足 $y_i(w \cdot x_i + b) > 1$

□ KKT条件：

- $\frac{\partial}{\partial w} \mathcal{L}(w^*, \alpha^*, \beta^*) = 0$
- $\frac{\partial}{\partial \beta} \mathcal{L}(w^*, \alpha^*, \beta^*) = 0$
- $\alpha_i^* g_i(w^*) = 0, \quad i = 1, \dots, k$
- $g_i(w^*) \leq 0, \quad i = 1, \dots, k$
- $\alpha_i^* \geq 0, \quad i = 1, \dots, k$

利用这个结论可进一步推导SVM原问题！



□ 支持向量机的原问题：

$$\begin{aligned} \min_{\mathbf{w}, b} \quad & \frac{1}{2} \|\mathbf{w}\|_2^2 \\ \text{s. t.} \quad & y_i(\mathbf{w} \cdot \mathbf{x}_i + b) \geq 1, 1 \leq i \leq n \end{aligned}$$

□ 由上页的结论，原问题与下面的问题同解：

$$\max_{\alpha: \alpha_i \geq 0} \min_{\mathbf{w}, b} \mathcal{L}(\mathbf{w}, b, \alpha) = \frac{1}{2} \|\mathbf{w}\|_2^2 + \sum_{i=1}^n \alpha_i (1 - y_i(\mathbf{w} \cdot \mathbf{x}_i + b))$$

求解上式分为三个步骤：

1. 求解  $\min_{\mathbf{w}, b} \mathcal{L}(\mathbf{w}, b, \alpha)$ ，得到  $\mathbf{w}, b$  关于  $\alpha$  的表达式，带入  $\mathcal{L}(\mathbf{w}, b, \alpha)$
2. 继续求解  $\max_{\alpha: \alpha_i \geq 0} \mathcal{L}(\mathbf{w}(\alpha), b(\alpha), \alpha)$
3. 利用  $\alpha^*$  回代求解  $(\mathbf{w}^*, b^*)$





# 对偶问题

- 第一步，令  $\nabla_w \mathcal{L}$  和  $\nabla_b \mathcal{L}$  等于零，得到： 这一步是无约束优化，直接求导即可

$$\nabla_w \mathcal{L} = w - \sum_{i=1}^n \alpha_i y_i x_i = 0 \quad \nabla_b \mathcal{L} = \sum_{i=1}^n \alpha_i y_i = 0$$

- 第二步，将上面两式带入  $\max_{\alpha: \alpha_i \geq 0} \min_{w, b} \mathcal{L}(w, b, \alpha)$  得到：

支持向量机的  
对偶问题：

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



- 第三步，利用 $\alpha^*$ 回代求解 $w^*$ 和判别函数（ $b^*$ 在后面给出求解方法）：

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

$$\hat{y} = \text{sgn}(w^* \cdot x + b^*) = \text{sgn}\left(\sum_{i=1}^n \alpha_i^* y_i (x_i \cdot x) + b^*\right)$$

- 回顾 $(w^*, b^*, \alpha^*)$ 满足的KKT条件：

$$\begin{aligned} \alpha_i^* (1 - y_i (w^* \cdot x_i + b^*)) &= 0, & i &= 1, \dots, n \\ 1 - y_i (w \cdot x_i + b) &\leq 0, & i &= 1, \dots, n \\ \alpha_i &\geq 0, & i &= 1, \dots, n \end{aligned}$$

- 若 $\alpha_i > 0$ ，必有  $y_i (w \cdot x_i + b) = 1$ ，即 $(x_i, y_i)$ 为支持向量
- 若  $y_i (w \cdot x_i + b) > 1$ ，必有 $\alpha_i = 0$ ，即非支持向量对应的 $\alpha_i$ 等于0



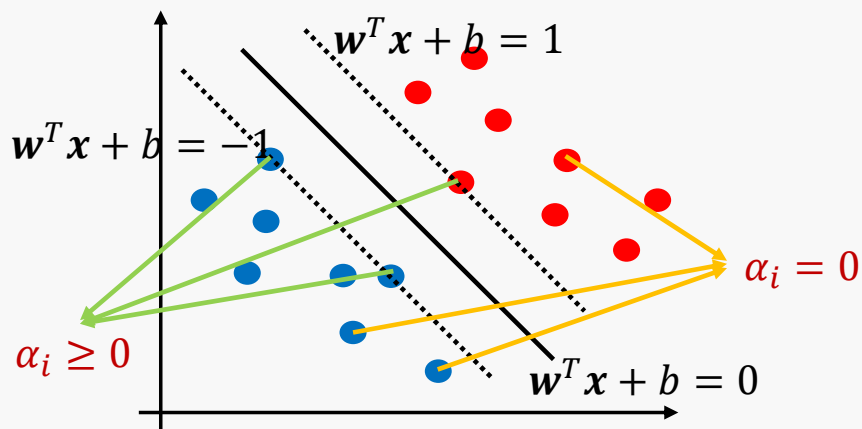
□ 回顾 $(w^*, b^*, \alpha^*)$ 满足的KKT条件:

$$\alpha_i^*(1 - y_i(w^* \cdot x_i + b^*)) = 0, \quad i = 1, \dots, n$$

$$1 - y_i(w \cdot x_i + b) \leq 0, \quad i = 1, \dots, n$$

$$\alpha_i \geq 0, \quad i = 1, \dots, n$$

- 若 $\alpha_i > 0$ , 必有  $y_i(w \cdot x_i + b) = 1$ , 即 $(x_i, y_i)$ 为支持向量
- 若  $y_i(w \cdot x_i + b) > 1$ , 必有 $\alpha_i = 0$ , 即非支持向量对应的 $\alpha_i$ 等于0





□ 从支持向量的角度重新审视 $(w^*, b^*)$ 以及判别函数，以 $\mathcal{S}$ 表示支持向量的集合

■  $w^*$ 的表达式：

$$w^* = \sum_{i=1}^n \alpha_i^* y_i x_i = \sum_{i \in \mathcal{S}} \alpha_i^* y_i x_i$$

■ 因为对于支持向量有  $y_i(w \cdot x_i + b) = 1$ ，所以 $b^*$ 的表达式：

$$b^* = \frac{1}{|\mathcal{S}|} \sum_{i \in \mathcal{S}} \left( y_i - \sum_{s \in \mathcal{S}} \alpha_s y_s (x_s \cdot x_i) \right)$$

■ 判别函数：

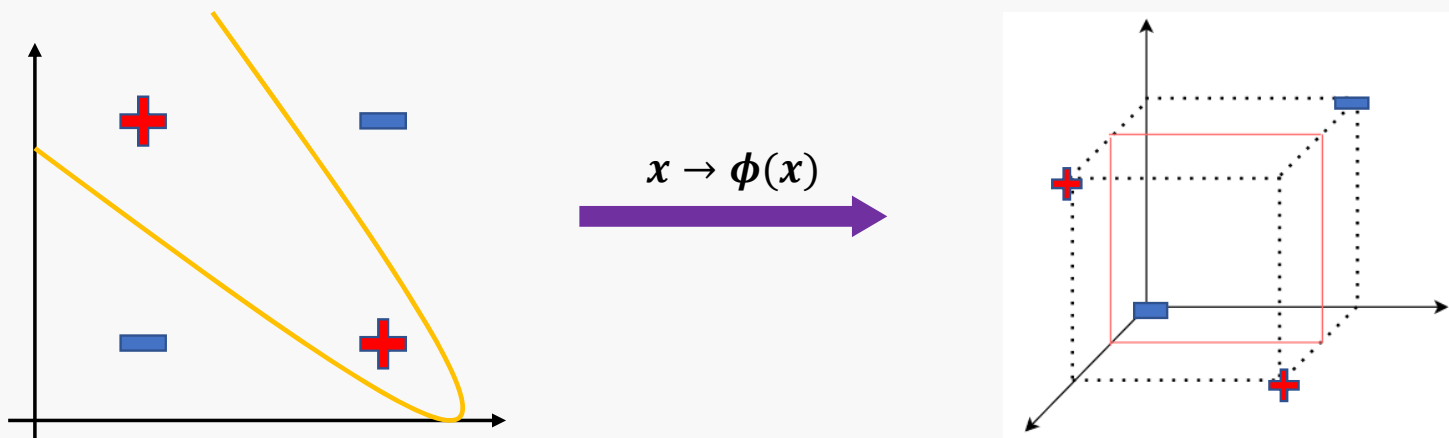
$$\hat{y} = \text{sgn}(w^* \cdot x + b^*) = \text{sgn} \left( \sum_{i \in \mathcal{S}} \alpha_i^* y_i (x_i \cdot x) + b^* \right)$$

□ 以上三式中，只有支持向量发挥作用！最终模型只与支持向量有关！



# 核函数

□ 前面讨论了训练样本线性可分的情况，如果不满足线性可分条件呢？



- 可将样本从原始空间映射到一个更高维的特征空间，使得样本在这个特征空间内线性可分
- 如果原始空间是有限维，那么一定存在一个高维特征空间使样本可分



- 令 $\phi(x)$ 表示高维特征向量，则高维特征空间中的分类面可表示为：

$$f(x) = \mathbf{w} \cdot \phi(x) + b$$

- 支持向量机的原问题可表示如下：

$$\min_{\mathbf{w}, b} \frac{1}{2} \|\mathbf{w}\|_2^2$$

$$\text{s. t. } y_i(\mathbf{w} \cdot \phi(x_i) + b) \geq 1, \quad i = 1, \dots, n$$

- 其对偶问题为：

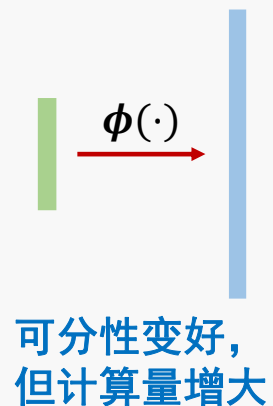
$$\max_{\alpha} \sum_{i=1}^n \alpha_i - \frac{1}{2} \sum_{i,j=1}^n \alpha_i \alpha_j y_i y_j (\phi(x_i) \cdot \phi(x_j))$$

$$\text{s. t. } \sum_{i=1}^n \alpha_i y_i = 0,$$

$$\alpha_i \geq 0, \quad i = 1, \dots, n$$



- 函数  $\phi(\cdot)$  又被称为基函数，它将  $x$  映射到更高维度的特征空间，这将导致内积运算  $\phi(x_i) \cdot \phi(x_j)$  的计算量增加
- 多项式基函数是一种常用的基函数
  - $(x_1, x_2, x_3)$  对应的2次基函数为  $(x_1, x_2, x_3, x_1x_2, x_1x_3, x_2x_3, x_1^2, x_2^2, x_3^2)$
  - $x \in \mathbb{R}^d$  对应的  $p$  次基函数的维度为  $\sum_{i=1}^p \binom{i+d-1}{i}$
  - 当  $d$  或  $p$  较大时，多项式基函数将引发维度灾难
- 先计算基函数，然后计算内积不可取：计算量过大！
- 有什么解决方案呢？





## □ 回顾线性可分情况的对偶问题和判别函数

$$\max_{\alpha} \sum_{i=1}^n \alpha_i - \frac{1}{2} \sum_{i,j=1}^n y_i y_j \alpha_i \alpha_j (x_i \cdot x_j)$$
$$\hat{y} = \text{sgn} \left( \sum_{i \in \mathcal{S}} \alpha_i^* y_i (x_i \cdot x) + b^* \right)$$

对偶问题的求解以及判别函数都只涉及内积 $(x_i \cdot x_j)$ !

可以在原始空间中计算高维空间的内积, 以降低计算量!

## □ 核函数的定义:

$$K(x_i, x_j) = \phi(x_i) \cdot \phi(x_j)$$



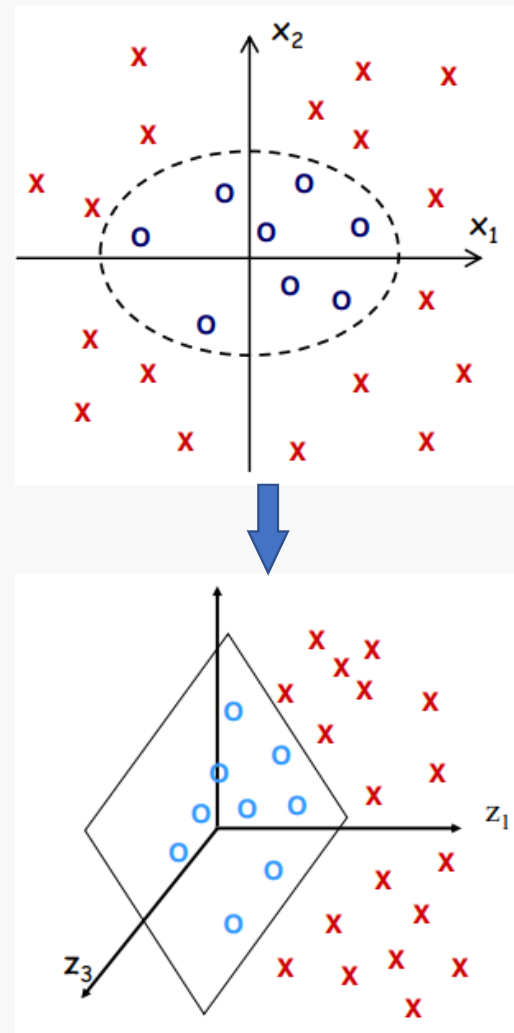


## □ 核函数举例

- 原始特征  $x = (x_1, x_2)$ , 基函数  $\phi(x) = (x_1^2, \sqrt{2}x_1x_2, x_2^2)$

$$\begin{aligned}\phi(x) \cdot \phi(y) &= (x_1^2, \sqrt{2}x_1x_2, x_2^2) \cdot (y_1^2, \sqrt{2}y_1y_2, y_2^2) \\ &= x_1^2y_1^2 + 2x_1x_2y_1y_2 + x_2^2y_2^2 \\ &= (x_1y_1 + x_2y_2)^2 \\ &= ((x_1, x_2) \cdot (y_1, y_2))^2 \\ &= (x \cdot y)^2 \\ &:= K(x, y)\end{aligned}$$

- 通过核函数  $K(x, y) = (x \cdot y)^2$ , 我们能够在原始特征空间 (2维) 计算高维特征空间 (3维) 的内积





## □ 引入核函数后的支持向量机

### ■ 对偶问题：

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

### ■ 求解后可以得到：

$$\begin{aligned} \hat{y} &= \text{sgn}(w \cdot \phi(x) + b) \\ &= \text{sgn} \left( \sum_{i=1}^n \alpha_i y_i (\phi(x_i) \cdot \phi(x)) + b \right) \\ &= \text{sgn} \left( \sum_{i=1}^n \alpha_i y_i K(x_i, x_j) + b \right) \end{aligned}$$

### ■ 求解过程和判别函数均不涉及高维基函数！



## □ 常用核函数列表

### ■ 线性核:

$$K(x_i, x_j) = x_i \cdot x_j$$

### ■ 多项式核:

$$K(x_i, x_j) = (x_i \cdot x_j)^d$$

$$K(x_i, x_j) = (x_i \cdot x_j + 1)^d$$

### ■ 径向基函数 (RBF) 核:

$$K(x_i, x_j) = \exp \left\{ -\frac{\|x_i - x_j\|_2^2}{\sigma^2} \right\}$$

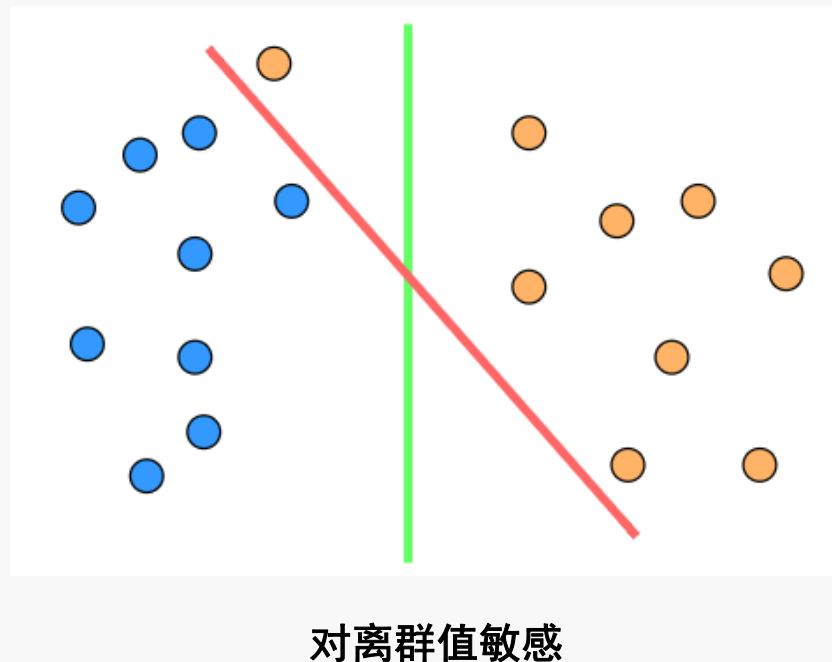
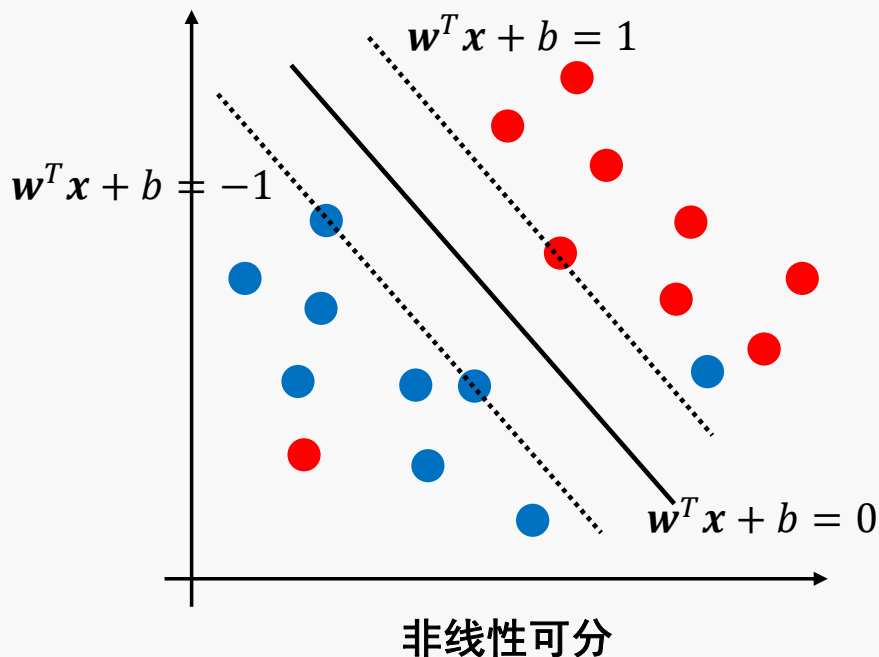
### ■ Sigmoid核:

$$K(x_i, x_j) = \tanh(\beta x_i \cdot x_j + \theta)$$



# 软间隔与正则化

- 在前面的讨论中，我们一直假定训练样本在样本空间或者特征空间中是线性可分的，实际中可能存在一些问题：
  - 很难确定合适的核函数使得训练样本在特征空间线性可分
  - 离群值很可能导致分类面过拟合





## □ 软间隔支持向量机

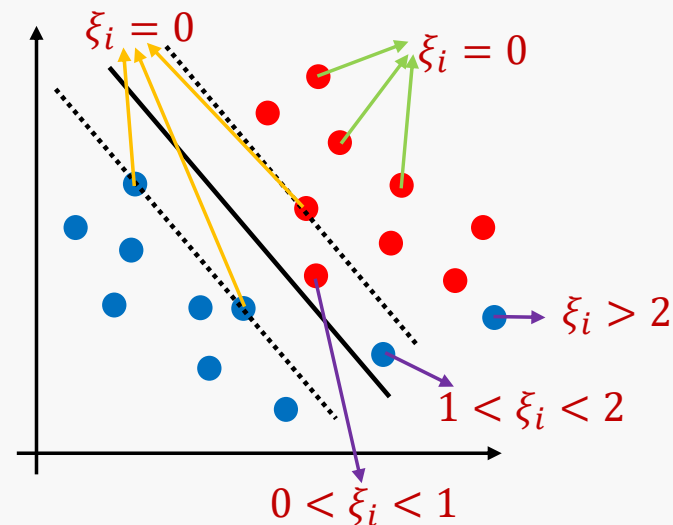
- 目标：让支持向量机应对非线性可分数据以及缓解其对离群值的敏感性

- 原问题：

$$\begin{aligned} \min_{\mathbf{w}, b, \xi} \quad & \frac{1}{2} \|\mathbf{w}\|_2^2 + C \sum_{i=1}^n \xi_i \\ \text{s. t.} \quad & y_i(\mathbf{w} \cdot \mathbf{x}_i + b) \geq 1 - \xi_i, \quad i = 1, \dots, n \\ & \xi_i \geq 0, \quad i = 1, \dots, n \end{aligned}$$

- 解释：

- 对于部分样本，允许函数间隔为  $1 - \xi_i$  ( $\xi_i > 0$ )，此时优化目标会增加  $C\xi_i$  的惩罚项
- 超参数  $C$  用来权衡“大间隔”与“保证训练样本均在间隔带以外”





## □ 软间隔支持向量机的对偶问题

$$\begin{aligned} & \max_{\substack{\alpha: \alpha_i \geq 0 \\ \mu: \mu_i \geq 0}} \min_{\mathbf{w}, b, \xi} \mathcal{L}(\mathbf{w}, b, \xi, \alpha, \mu) \\ &= \frac{1}{2} \|\mathbf{w}\|_2^2 + C \sum_{i=1}^n \xi_i + \sum_{i=1}^n \alpha_i (1 - \xi_i - y_i(\mathbf{w} \cdot \mathbf{x}_i + b)) - \sum_{i=1}^n \mu_i \xi_i \end{aligned}$$

令 $\mathcal{L}(\mathbf{w}, b, \xi, \alpha, \mu)$ 对 $\mathbf{w}, b, \xi$ 的偏导数为0, 可将上式化简为:

软间隔支持向量  
机的对偶问题

$$\begin{aligned} & \max_{\alpha} \sum_{i=1}^n \alpha_i - \frac{1}{2} \sum_{i,j=1}^n y_i y_j \alpha_i \alpha_j (\mathbf{x}_i \cdot \mathbf{x}_j) \\ & \text{s. t. } 0 \leq \alpha_i \leq C, \quad i = 1, \dots, n \\ & \quad \sum_{i=1}^n \alpha_i y_i = 0 \end{aligned}$$



# 软间隔与正则化

□ 与软间隔对应，线性可分条件下的支持向量机也可称为硬间隔支持向量机

□ 软间隔与硬间隔支持向量机的相同之处：

- 两者对偶问题的优化目标函数一致，解法一致
- 两者的 $w$ 和判别函数的表达式一致

$$w^* = \sum_{i=1}^n \alpha_i^* y_i x_i$$
$$\hat{y} = \text{sgn} \left( \sum_{i=1}^n \alpha_i^* y_i (x_i \cdot x) + b^* \right)$$

□ 软间隔与硬间隔支持向量机的不同之处：

- 硬间隔版本要求 $\alpha_i \geq 0$
- 软间隔版本要求 $0 \leq \alpha_i \leq C$



# 软间隔与正则化

□ 从KKT条件的角度看  $0 \leq \alpha_i \leq C$ :

$$\begin{cases} \alpha_i \geq 0, & \mu_i \geq 0, \\ y_i(\mathbf{w} \cdot \mathbf{x}_i + b) \geq 1 - \xi_i \\ \alpha_i(y_i(\mathbf{w} \cdot \mathbf{x}_i + b) - 1 + \xi_i) = 0, \\ \xi_i \geq 0, & \mu_i \xi_i = 0 \end{cases}$$

$$\alpha_i + \mu_i = C \leftarrow \text{----- } \nabla_{\xi_i} \mathcal{L} = 0 \text{ 要求}$$

- 当 $\alpha_i = 0$ 时,  $\mu_i = C$ ,  $\xi_i = 0$ ,  $y_i(\mathbf{w} \cdot \mathbf{x}_i + b) \geq 1$ , 非支持向量
- 当 $0 < \alpha_i < C$ 时,  $0 < \mu_i < C$ ,  $\xi_i = 0$ ,  $y_i(\mathbf{w} \cdot \mathbf{x}_i + b) = 1$ , 支持向量, 且在最大间隔边界上
- 当 $\alpha_i = C$ 时,  $\mu_i = 0$ ,  $\xi_i \geq 0$ ,  $y_i(\mathbf{w} \cdot \mathbf{x}_i + b) = 1 - \xi_i \leq 1$ , 支持向量, 且越过最大间隔边界

□ 软间隔支持向量机的最终模型也仅与支持向量有关!

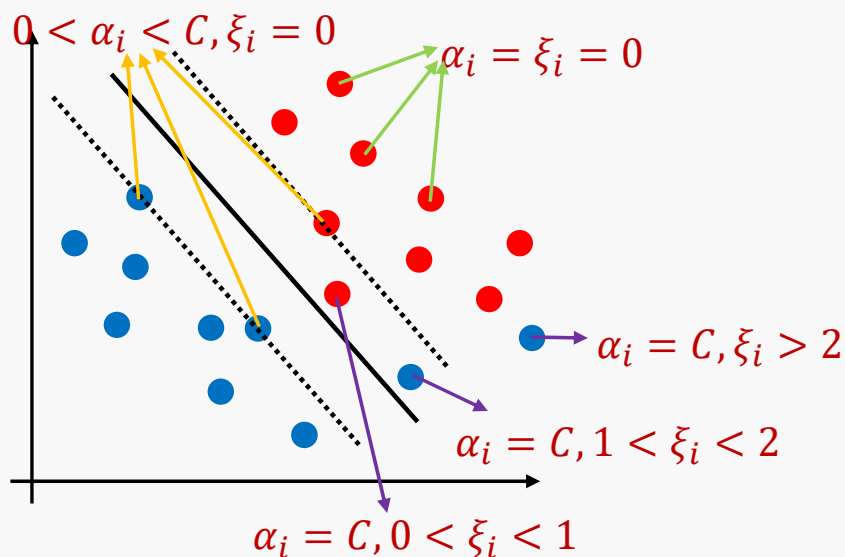




# 软间隔与正则化

## □ 上页结论:

- 当 $\alpha_i = 0$ 时,  $\mu_i = C$ ,  $\xi_i = 0$ ,  $y_i(w \cdot x_i + b) \geq 1$ , 非支持向量
- 当 $0 < \alpha_i < C$ 时,  $0 < \mu_i < C$ ,  $\xi_i = 0$ ,  $y_i(w \cdot x_i + b) = 1$ , 支持向量, 且在最大间隔边界上
- 当 $\alpha_i = C$ 时,  $\mu_i = 0$ ,  $\xi_i \geq 0$ ,  $y_i(w \cdot x_i + b) = 1 - \xi_i \leq 1$ , 支持向量, 且越过最大间隔边界





## § 5.3 方法扩展

- 一、支持向量回归
- 二、核方法
- 三、多分类问题



## □ 回顾：线性回归

在于如何衡量 $f(x_i)$ 和 $y_i$ 之间的差异，均方误差是回归任务中最常用的性能度量，因此我们可试图让均方误差最小化，即

$$\begin{aligned}(w^*, b^*) &= \operatorname{argmin}_{(w, b)} \sum_{i=1}^m (f(x_i) - y_i)^2 \\ &= \operatorname{argmin}_{(w, b)} \sum_{i=1}^m (y_i - wx_i - b)^2\end{aligned}$$

□ 均方误差有非常好的几何意义，它对应了常用的欧几里得距离或简称“**欧氏距离**” (Euclidean distance)



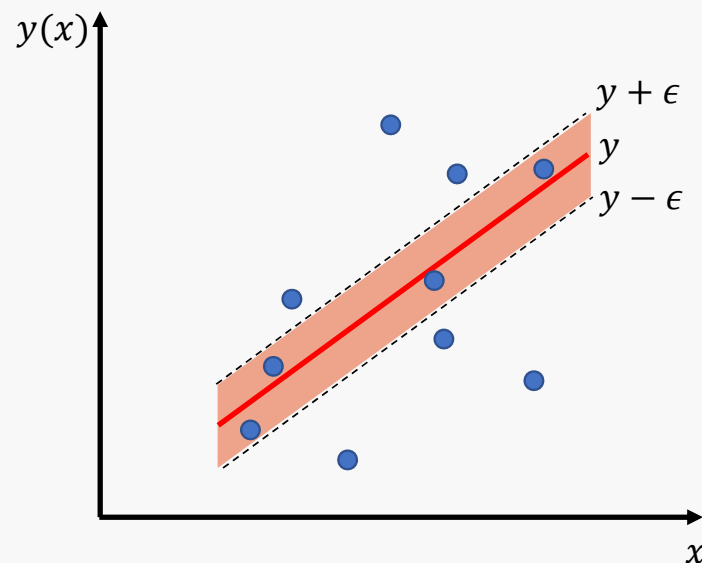
# 支持向量回归

□ 目标：学习回归模型  $f(x) = w\phi(x) + b$

- **传统回归模型**：当且仅当  $f(x)$  与  $y$  完全相同时，损失才为零
- **支持向量回归**：对  $f(x)$  与  $y$  的偏差大于容忍度  $\epsilon$  才计算损失

□ 支持向量的对比

- 分类（SVM）：当数据点远离分类平面时损失为零，支持向量为离分类平面最近的点
- 回归（SVR）：当数据点与回归平面足够近时损失为0，支持向量为  $|f(x) - y| > \epsilon$  的点



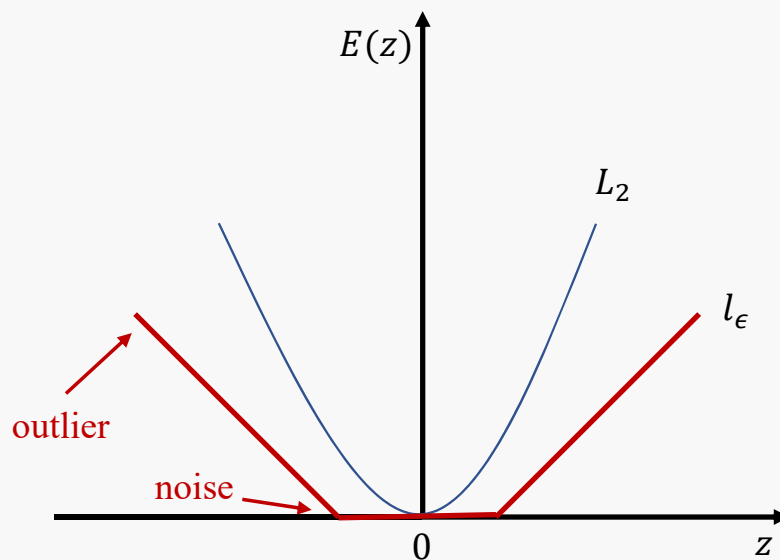


## □ 优化目标

- SVR问题可形式化为：

$$\min_{w,b} \frac{1}{2} \|w\|^2 + C \sum_{i=1}^m l_{\epsilon}(f(x_i) - y_i)$$

- 不敏感损失函数  $l_{\epsilon}(z) = \max(0, |z| - \epsilon)$
- 对噪声（较小偏差）、异常（较大偏差）的鲁棒性均优于  $L_2$  损失函数





## □ 优化目标

- 间隔带两侧损失函数可有所不同
- 引入松弛变量 $\xi_i$ 和 $\hat{\xi}_i$ ，重写为：

$$\begin{aligned} \min_{\mathbf{w}, b, \xi_i, \hat{\xi}_i} \quad & \frac{1}{2} \|\mathbf{w}\|^2 + C \sum_{i=1}^m (\xi_i + \hat{\xi}_i) \\ \text{s. t.} \quad & f(\mathbf{x}_i) - y_i \leq \epsilon + \xi_i, \\ & y_i - f(\mathbf{x}_i) \leq \epsilon + \hat{\xi}_i, \\ & \xi_i \geq 0, \hat{\xi}_i \geq 0, i = 1, 2, \dots, m. \end{aligned}$$

- 拉格朗日函数：

$$\begin{aligned} L(\mathbf{w}, b, \alpha, \hat{\alpha}, \xi, \hat{\xi}, \mu, \hat{\mu}) \\ = \frac{1}{2} \|\mathbf{w}\|^2 + C \sum_{i=1}^m (\xi_i + \hat{\xi}_i) - \sum_{i=1}^m \mu_i \xi_i - \sum_{i=1}^m \hat{\mu}_i \hat{\xi}_i \\ + \sum_{i=1}^m \alpha_i (f(\mathbf{x}_i) - y_i - \epsilon - \xi_i) + \sum_{i=1}^m \hat{\alpha}_i (y_i - f(\mathbf{x}_i) - \epsilon - \hat{\xi}_i) \end{aligned}$$



## □ 优化目标

- 令 $L$ 对变量的偏导为零可得：

$$w = \sum_{i=1}^m (\hat{\alpha}_i - \alpha_i) x_i$$

$$0 = \sum_{i=1}^m (\hat{\alpha}_i - \alpha_i)$$

$$C = \alpha_i + \mu_i$$

$$C = \hat{\alpha}_i + \hat{\mu}_i$$

- 得到对偶问题：

$$\min_{\alpha, \hat{\alpha}} \sum_{i=1}^m y_i (\hat{\alpha}_i - \alpha_i) - \epsilon (\hat{\alpha}_i + \alpha_i) - \frac{1}{2} \sum_{i=1}^m \sum_{j=1}^m (\hat{\alpha}_i - \alpha_i) (\hat{\alpha}_j - \alpha_j) x_i^T x_j$$

$$\begin{aligned} \text{s. t. } & \sum_{i=1}^m (\hat{\alpha}_i - \alpha_i) = 0 \\ & 0 \leq \alpha_i, \hat{\alpha}_i \leq C \end{aligned}$$



## □ 优化目标

- 上述过程需满足KKT条件：

$$\begin{cases} \alpha_i(f(x_i) - y_i - \epsilon - \xi_i) = 0 \\ \widehat{\alpha}_i(y_i - f(x_i) - \epsilon - \widehat{\xi}_i) = 0 \\ \alpha_i \widehat{\alpha}_i = 0, \xi_i \widehat{\xi}_i = 0 \\ (C - \alpha_i)\xi_i = 0, (C - \widehat{\alpha}_i)\widehat{\xi}_i = 0 \end{cases}$$

- 样本落在 $\epsilon$  - 间隔带之外时,  $\alpha_i$ 和 $\widehat{\alpha}_i$ 才能取非零值, 此时样本为支持向量

- 将 $w$ 代入回归模型有：

$$f(x) = \sum_{i=1}^m (\widehat{\alpha}_i - \alpha_i) x_i^T x + b$$

- SVR的支持向量满足 $(\alpha_i - \widehat{\alpha}_i) \neq 0$

- 求解 $b$ ：先求解对偶问题得到 $\alpha_i$ , 由 $\alpha_i(f(x_i) - y_i - \epsilon - \xi_i) = 0$ 和 $(C - \alpha_i)\xi_i = 0$ , 选取 $0 < \alpha_i < C$ , 则可以得到：

$$b = y_i + \epsilon - \sum_{j=1}^m (\widehat{\alpha}_j - \alpha_j) x_j^T x_i$$





## □ 考虑特征映射

### ■ SVM:

$$f(x) = \sum_{i=1}^m \alpha_i y_i \kappa(x, x_i) + b$$

### ■ SVR:

$$f(x) = \sum_{i=1}^m (\hat{\alpha}_i - \alpha_i) \kappa(x, x_i) + b$$

### ■ 若不考虑偏移项，则SVM和SVR学到的模型总是能表示成核函数 $\kappa(x, x_i)$ 的线性组合

### ■ 更一般的结论（表示定理）：

**定理 6.2 (表示定理)** 令  $\mathbb{H}$  为核函数  $\kappa$  对应的再生核希尔伯特空间,  $\|h\|_{\mathbb{H}}$  表示  $\mathbb{H}$  空间中关于  $h$  的范数, 对于任意单调递增函数  $\Omega: [0, \infty] \mapsto \mathbb{R}$  和任意非负损失函数  $\ell: \mathbb{R}^m \mapsto [0, \infty]$ , 优化问题

$$\min_{h \in \mathbb{H}} F(h) = \Omega(\|h\|_{\mathbb{H}}) + \ell(h(\mathbf{x}_1), h(\mathbf{x}_2), \dots, h(\mathbf{x}_m)) \quad (6.57)$$

的解总可写为

$$h^*(\mathbf{x}) = \sum_{i=1}^m \alpha_i \kappa(\mathbf{x}, \mathbf{x}_i) . \quad (6.58)$$



□ **核方法：**一系列基于核函数的学习方法

□ **核线性判别分析（KLDA）**

- 假设通过映射 $\phi$ ，将样本映射到特征空间 $\mathbb{F}$ 执行线性判别分析：

$$h(x) = w^T \phi(x)$$

- KLDA的学习目标为：

$$\min_w J(w) = \frac{w^T S_b^\phi w}{w^T S_w^\phi w}$$

- 记两类样本集为 $X_0, X_1$ ，总样本数为 $m = m_0 + m_1$ ，则第 $i$ 类样本在特征空间 $\mathbb{F}$ 中的均值为：

$$\mu_i^\phi = \frac{1}{m_i} \sum_{x \in X_i} \phi(x)$$

- 两个散度矩阵分别为：

$$S_b^\phi = (\mu_1^\phi - \mu_0^\phi)(\mu_1^\phi - \mu_0^\phi)^T, S_w^\phi = \sum_{i=0}^1 \sum_{x \in X_i} (\phi(x) - \mu_i^\phi)(\phi(x) - \mu_i^\phi)^T$$



## □ 核线性判别分析 (KLDA)

- 将 $J(w)$ 作为损失函数 $l$ , 令 $\Omega \equiv 0$ , 则由表示定理,  $h(x)$ 可以表示为:

$$h(x) = \sum_{i=1}^m \alpha_i \kappa(x, x_i)$$

- 得到:

$$w = \sum_{i=1}^m \alpha_i \phi(x_i)$$

- 通常难以直接写出映射 $\phi$ 的具体形式, 可通过核函数对该问题进行求解。令 $K \in \mathbb{R}^{m \times m}$ 为核矩阵, 令 $\mathbf{1}_i \in \{1, 0\}^{m \times 1}$ 为第 $i$ 类样本的指示向量, 再令:

$$\hat{\mu}_0 = \frac{1}{m_0} K \mathbf{1}_0, \hat{\mu}_1 = \frac{1}{m_1} K \mathbf{1}_1, M = (\hat{\mu}_0 - \hat{\mu}_1)(\hat{\mu}_0 - \hat{\mu}_1)^T, N = K K^T - \sum_{i=0}^1 m_i \hat{\mu}_i \hat{\mu}_i^T$$

- 原问题等价于:

$$\max_{\alpha} J(\alpha) = \frac{\alpha^T M \alpha}{\alpha^T N \alpha} \longrightarrow \text{同LDA解法得到 } \alpha \longrightarrow \text{即可写出 } h(x) \text{ 的形式}$$



## □ 成对分类方法 (one-against-one)

- 每两个类之间都构造一个SVM进行判别:

$$\begin{aligned} \min_{w^{ij}, b^{ij}, \xi^{ij}} & \frac{1}{2} \|w^{ij}\|^2 + C \sum_{t=1}^m \xi_t^{ij} \\ \text{s. t. } & (w^{ij})^T x_t + b^{ij} \geq 1 - \xi_t^{ij}, \text{ if } y_t = i \\ & (w^{ij})^T x_t + b^{ij} \leq -1 + \xi_t^{ij}, \text{ if } y_t = j \\ & \xi_t^{ij} \geq 0, t = 1, 2, \dots, m. \end{aligned}$$

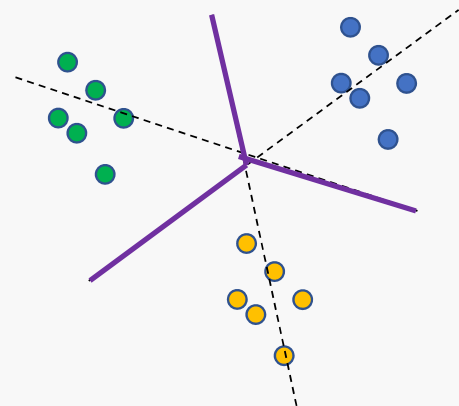
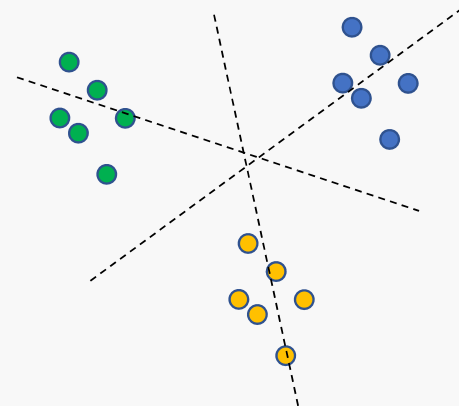
- 若一共有M个类，则需要训练 $\frac{1}{2} M(M-1)$ 个SVM分类器

- 训练完成后，对于新数据采用投票策略进行分类:

- 每个分类器对判定的类别投一票，决策函数为:

$$y_{new}^{ij} = \text{sign} \left[ (w^{ij})^T x_t + b^{ij} \right]$$

- 票数最多的类别为最终预测
- 票数相同时简单选择索引最小的类别





## □ 一对多分类方法 (one-against-all)

- 对每个类都构造一个SVM进行判别，区分第*i*类和其余*M*-1类：

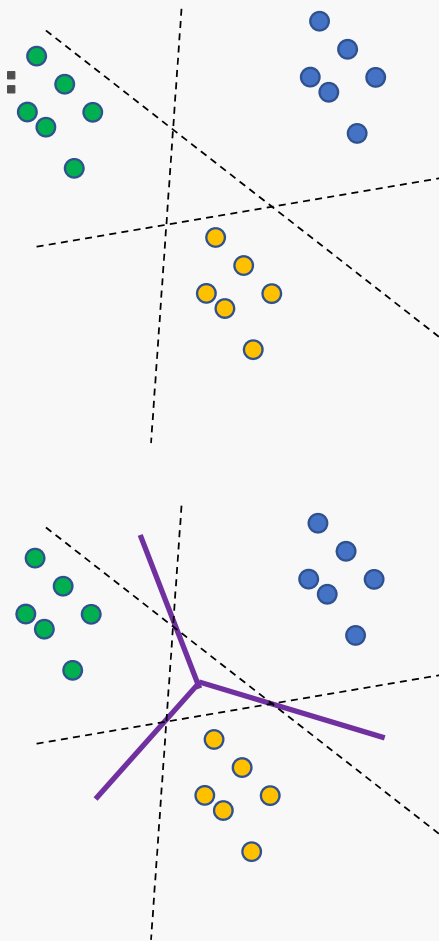
$$\begin{aligned} \min_{w^i, b^i, \xi^i} & \frac{1}{2} \|w^i\|^2 + C \sum_{t=1}^m \xi_t^i \\ \text{s. t. } & (w^i)^T x_t + b^i \geq 1 - \xi_t^i, \text{ if } y_t = i \\ & (w^i)^T x_t + b^i \leq -1 + \xi_t^i, \text{ if } y_t \neq i \\ & \xi_t^i \geq 0, t = 1, 2, \dots, m. \end{aligned}$$

- 训练完成后，对于新数据根据预测值的大小进行分类：

- 修改决策函数，去掉符号函数：

$$d_{new}^i = (w^i)^T x_t + b^i$$

- 选取预测值最大的一类作为最终预测
- **问题：**不同的分类器并非同时训练，且不同的分类器训练时正负样本的比例各有不同，因此去掉符号函数以后预测值可能不具有很好的可比较性

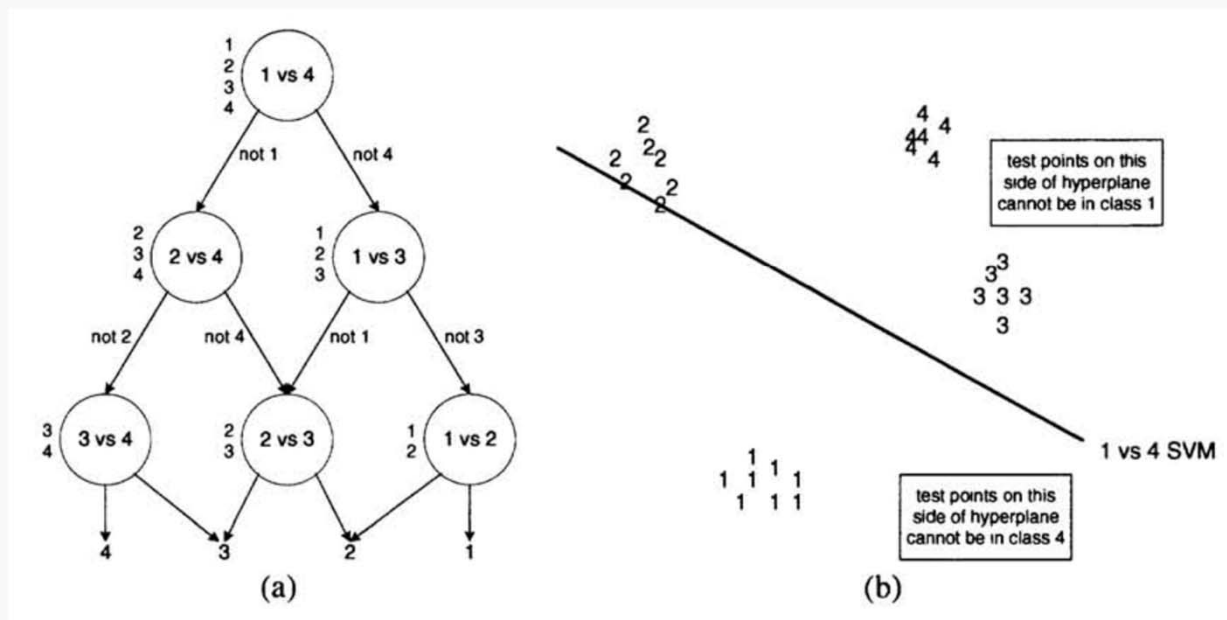




# 多分类问题

## □ 有向无环图支持向量机 (DAGSVM)

- 训练阶段与成对分类方法相同
- 测试时通过有根节点（共有 $\frac{1}{2}M(M-1)$ 个节点）的有向无环图进行判断





# 思考题

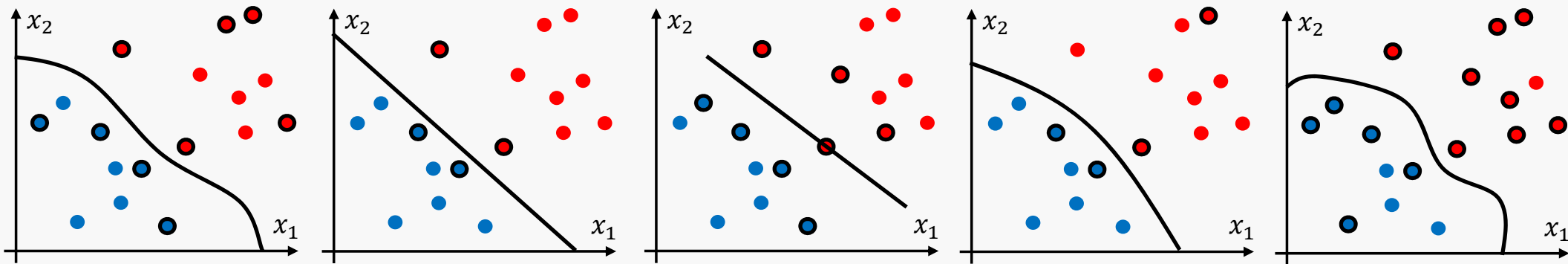
□ 为下面的每个SVM模型标出正确的图形

$$\min \frac{1}{2} \|\mathbf{w}\|^2 + C \sum_{i=1}^N \xi_i$$

s. t.  $\xi_i \geq 0, y_i(\mathbf{w}^T \mathbf{x} + w_0) \geq 1 - \xi_i, C = 0.1$

$$\min \frac{1}{2} \|\mathbf{w}\|^2 + C \sum_{i=1}^N \xi_i$$

s. t.  $\xi_i \geq 0, y_i(\mathbf{w}^T \mathbf{x} + w_0) \geq 1 - \xi_i, C = 1$

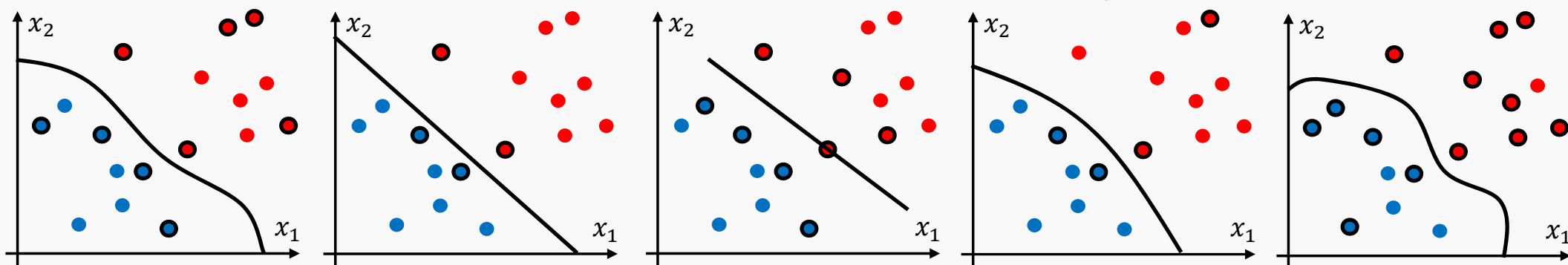




# 思考题

□ 为下面的每个SVM模型标出正确的图形

$$\max \left( \sum_{i=1}^N \alpha_i - \frac{1}{2} \sum_{i=1}^N \sum_{j=1}^N \alpha_i \alpha_j y_i y_j k(x_i, x_j) \right)$$
$$\text{s. t. } \alpha_i \geq 0, \sum_{i=1}^N \alpha_i y_i = 0, k(x, x') = x^T x' + (x^T x')^2$$







# 思考题

□ 为下面的每个SVM模型标出正确的图形

$$\max \left( \sum_{i=1}^N \alpha_i - \frac{1}{2} \sum_{i=1}^N \sum_{j=1}^N \alpha_i \alpha_j y_i y_j k(x_i, x_j) \right)$$
$$\text{s.t. } \alpha_i \geq 0, \sum_{i=1}^N \alpha_i y_i = 0, k(x, x') = \exp \left( -\frac{1}{2} \|x - x'\|^2 \right)$$

$$\max \left( \sum_{i=1}^N \alpha_i - \frac{1}{2} \sum_{i=1}^N \sum_{j=1}^N \alpha_i \alpha_j y_i y_j k(x_i, x_j) \right)$$
$$\text{s.t. } \alpha_i \geq 0, \sum_{i=1}^N \alpha_i y_i = 0, k(x, x') = \exp \left( -\|x - x'\|^2 \right)$$

