



2021年秋季 《机器学习概论》课程

第六章：支持向量机

主讲：连德富 特任教授 | 博士生导师

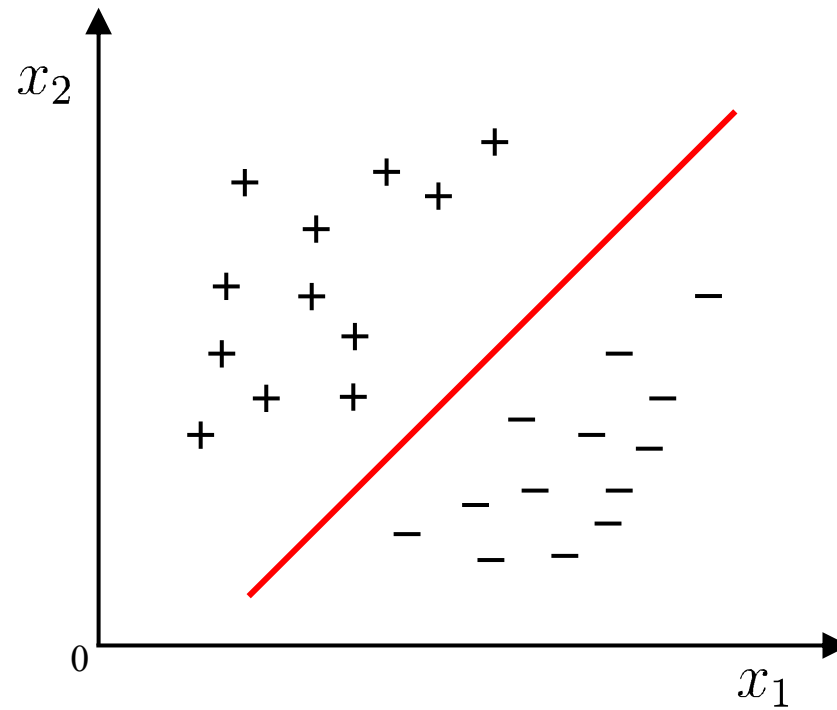
邮箱：liandefu@ustc.edu.cn

手机：13739227137

主页：<http://staff.ustc.edu.cn/~liandefu>

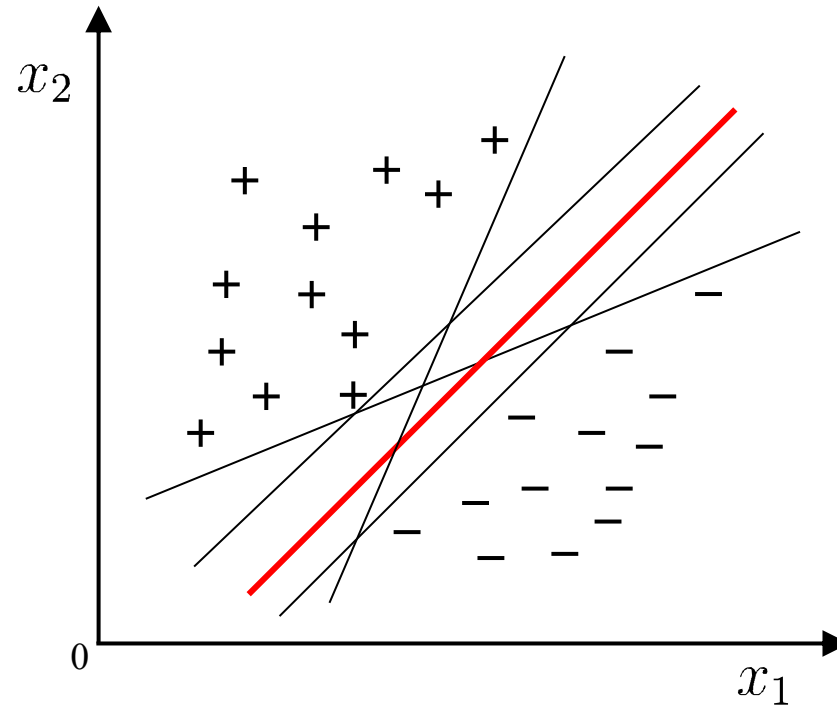
线性模型

在样本空间中寻找一个超平面, 将不同类别的样本分开



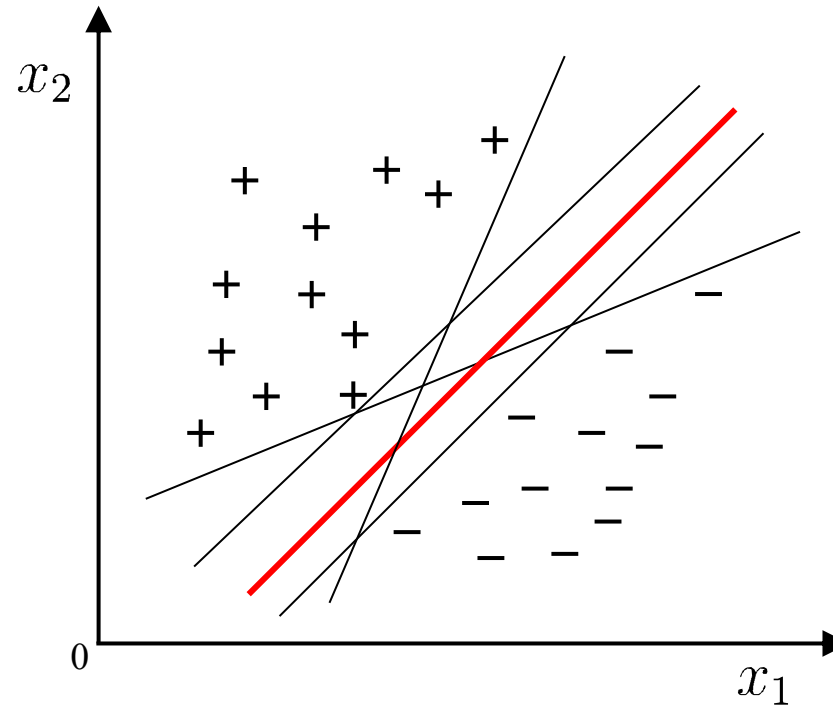
线性模型

- 将训练样本分开的超平面可能有很多, 哪一个好呢?



线性模型

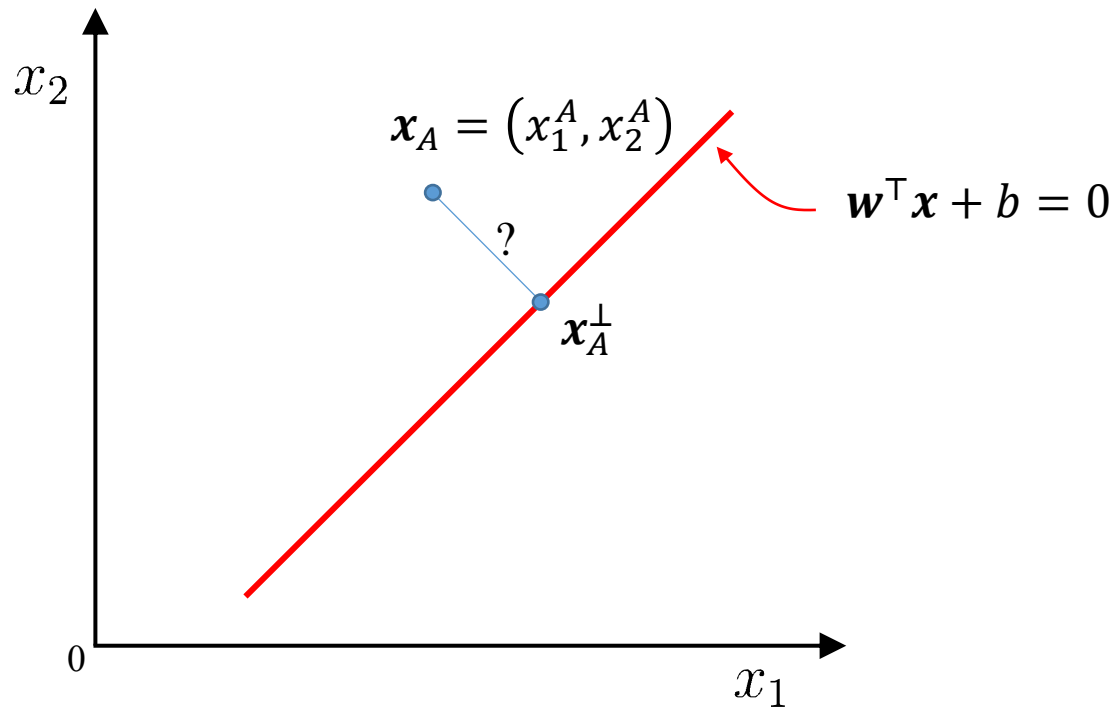
- 将训练样本分开的超平面可能有很多, 哪一个好呢?



应选择“ **正中间**”, 容忍性好, 鲁棒性高, 泛化能力最强

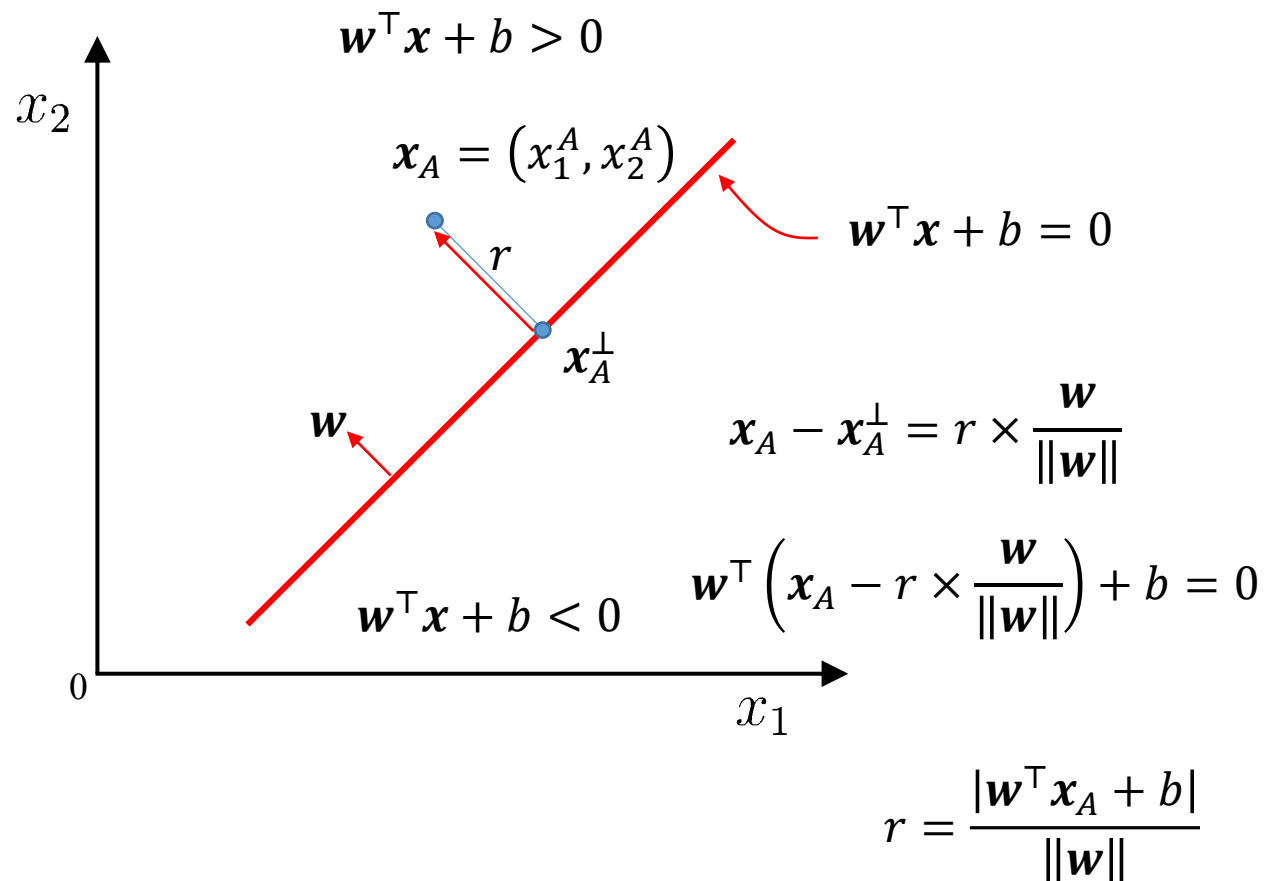
线性模型

- 超平面方程: $\mathbf{w}^\top \mathbf{x} + b = 0$



线性模型

- 超平面方程: $\mathbf{w}^\top \mathbf{x} + b = 0$



线性模型（线性可分）

- 假设超平面能将训练样本正确分类，即对于 $(x_i, y_i) \in D$
 - 若 $y_i = +1$ ，则有 $\mathbf{w}^\top \mathbf{x}_i + b > 0$
 - 若 $y_i = -1$ ，则有 $\mathbf{w}^\top \mathbf{x}_i + b < 0$

$$\Rightarrow y_i(\mathbf{w}^\top \mathbf{x}_i + b) > 0$$

- 分类平面对于向量 $\mathbf{w}' = [\mathbf{w}, b]$ 的长度具有不变性

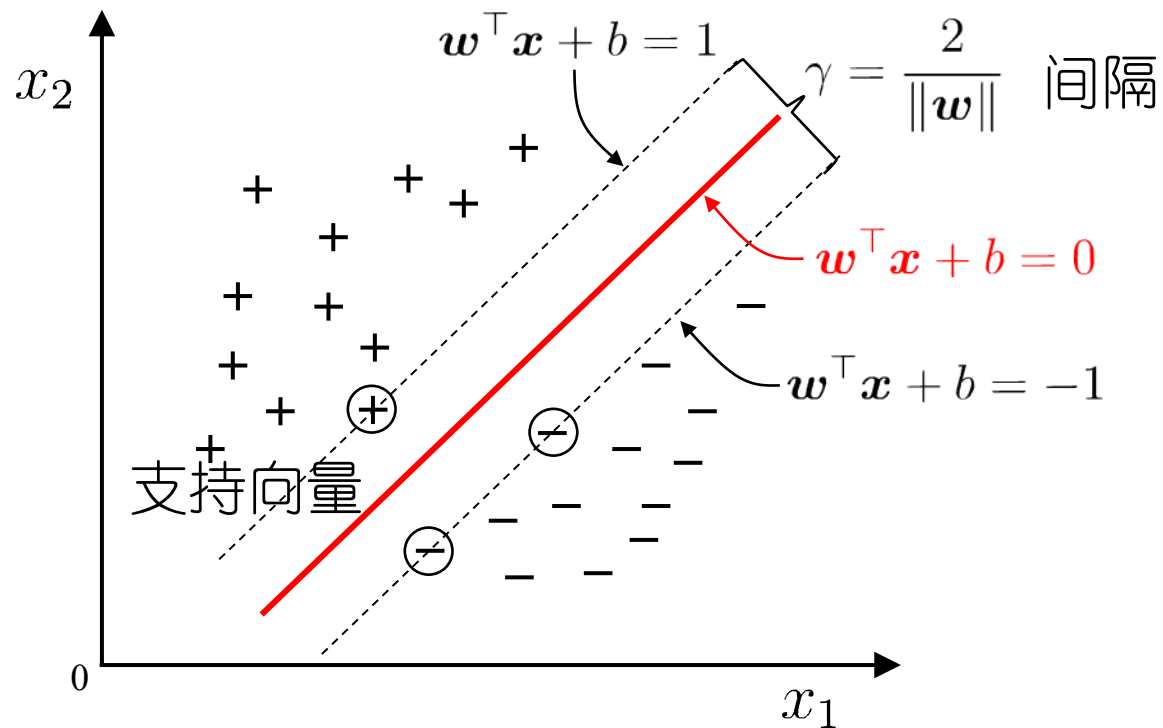
$$y_i(\mathbf{w}^\top \mathbf{x}_i + b) > 0 \quad \begin{array}{c} \mathbf{w}' \mapsto \alpha \mathbf{w}' \\ \alpha > 0 \end{array} \Rightarrow y_i(\alpha \mathbf{w}^\top \mathbf{x}_i + \alpha b) > 0$$

- 假设 $y_i(\mathbf{w}^\top \mathbf{x}_i + b) \geq \epsilon$ ，则 $y_i \left(\frac{\mathbf{w}^\top}{\epsilon} \mathbf{x}_i + \frac{b}{\epsilon} \right) \geq 1$

$$y_i \left(\frac{\mathbf{w}^\top}{\epsilon} \mathbf{x}_i + \frac{b}{\epsilon} \right) \geq 1 \quad \begin{array}{c} \frac{1}{\epsilon}(\mathbf{w}, b) \mapsto (\mathbf{w}', b') \end{array} \Leftrightarrow y_i(\mathbf{w}'^\top \mathbf{x}_i + b') \geq 1$$

线性模型（线性可分）

- $y_i(\mathbf{w}^\top \mathbf{x}_i + b) \geq 1$ \rightarrow 若 $y_i = +1$, 则有 $\mathbf{w}^\top \mathbf{x}_i + b \geq +1$
若 $y_i = -1$, 则有 $\mathbf{w}^\top \mathbf{x}_i + b \leq -1$



支持向量机基本型

- 最大间隔:寻找参数 \mathbf{w} 和 b , 使得间隔 r 最大

$$\begin{aligned} & \max_{\mathbf{w}, b} \frac{2}{\|\mathbf{w}\|} \\ & \text{s. t. } y_i(\mathbf{w}^\top \mathbf{x}_i + b) \geq 1 \quad i = 1, \dots, m \end{aligned}$$



$$\begin{aligned} & \min_{\mathbf{w}, b} \frac{1}{2} \|\mathbf{w}\|^2 \\ & \text{s. t. } y_i(\mathbf{w}^\top \mathbf{x}_i + b) \geq 1 \quad i = 1, \dots, m \end{aligned}$$

对偶问题

$$\min_{\mathbf{w}, b} \frac{1}{2} \|\mathbf{w}\|^2 \quad \text{s.t.} \quad y_i(\mathbf{w}^\top \mathbf{x}_i + b) \geq 1 \quad i = 1, \dots, m$$

凸函数

凸函数

凸优化问题

$$\begin{aligned} \min_x & f(\mathbf{x}) \\ \text{s.t.} & g_i(\mathbf{x}) \leq 0, i = 1, 2, \dots, m \\ & \mathbf{a}_i^\top \mathbf{x} = b, j = 1, 2, \dots, n \end{aligned}$$

其中 $f(\mathbf{x}), g_i(\mathbf{x})$ 是凸函数

对偶问题

原问题 凸优化

$$\begin{aligned} \min_x & f(x) \\ \text{s.t.} & g_i(x) \leq 0, i = 1, 2, \dots, m \\ & h_j(x) = 0, j = 1, 2, \dots, n \end{aligned}$$

广义拉格朗日函数

$$L(x, u, v) = f(x) + \sum_i u_i g_i(x) + \sum_j v_j h_j(x)$$

其中 $u_i \geq 0$

对偶问题 凸优化

$$\begin{aligned} \max_{u, v} & g(u, v) \text{ s.t. } u \geq 0 \\ & \text{其中 } g(u, v) = \min_x L(x, u, v) \end{aligned}$$

对偶问题

原问题

凸优化

$$\begin{aligned} \min_{\mathbf{w}, b} & \|\mathbf{w}\|^2 / 2 \\ \text{s.t. } & y_i(\mathbf{w}^\top \mathbf{x}_i + b) \geq 1, i = 1, 2, \dots, m \end{aligned}$$

广义拉格朗日函数

$$L(\mathbf{w}, b, \boldsymbol{\alpha}) = \frac{1}{2} \|\mathbf{w}\|^2 - \sum_{i=1}^m \alpha_i (y_i(\mathbf{w}^\top \mathbf{x}_i + b) - 1)$$

其中 $\alpha_i \geq 0$

对偶问题

凸优化

$$\begin{aligned} \max_{\boldsymbol{\alpha}} & g(\boldsymbol{\alpha}) \text{ s.t. } \boldsymbol{\alpha} \geq 0 \\ \text{其中 } & g(\boldsymbol{\alpha}) = \min_{\mathbf{w}, b} L(\mathbf{w}, b, \boldsymbol{\alpha}) \end{aligned}$$

对偶问题

- $g(\alpha) = \min_{\mathbf{w}, b} L(\mathbf{w}, b, \alpha)$

$$L(\mathbf{w}, b, \alpha) = \frac{1}{2} \|\mathbf{w}\|^2 - \sum_{i=1}^m \alpha_i (y_i (\mathbf{w}^\top \mathbf{x}_i + b) - 1)$$

$L(\mathbf{w}, b, \alpha)$ 是关于 \mathbf{w} 和 b 的凸函数, 在其梯度等于0时取得最优值

$$\begin{aligned} \frac{\partial L(\mathbf{w}, b, \alpha)}{\partial \mathbf{w}} = 0 & \quad \Rightarrow \quad \mathbf{w} = \sum_{i=1}^m \alpha_i y_i \mathbf{x}_i \\ \frac{\partial L(\mathbf{w}, b, \alpha)}{\partial b} = 0 & \quad \Rightarrow \quad \sum_{i=1}^m \alpha_i y_i = 0 \end{aligned}$$

对偶问题

$$L(\mathbf{w}, b, \alpha) = \frac{1}{2} \|\mathbf{w}\|^2 - \sum_{i=1}^m \alpha_i (y_i (\mathbf{w}^\top \mathbf{x}_i + b) - 1)$$

$$\mathbf{w} = \sum_{i=1}^m \alpha_i y_i \mathbf{x}_i$$



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

$$\begin{aligned} g(\alpha) &= \frac{1}{2} \left\| \sum_{i=1}^m \alpha_i y_i \mathbf{x}_i \right\|^2 - \sum_{i=1}^m \alpha_i \left(y_i \left(\left(\sum_{j=1}^m \alpha_j y_j \mathbf{x}_j \right)^\top \mathbf{x}_i \right) \right) + \sum_i \alpha_i \\ &= \sum_i \alpha_i - \frac{1}{2} \sum_{i=1}^m \sum_{j=1}^m \alpha_i \alpha_j y_i y_j \mathbf{x}_i^\top \mathbf{x}_j \end{aligned}$$

对偶
问题

$$\begin{aligned} \max_{\alpha} g(\alpha) &= \sum_{i=1}^m \alpha_i - \frac{1}{2} \sum_{i=1}^m \sum_{j=1}^m \alpha_i \alpha_j y_i y_j \mathbf{x}_i^\top \mathbf{x}_j \\ \text{s.t. } \alpha &\geq 0 \text{ and } \sum_{i=1}^m \alpha_i y_i = 0 \end{aligned}$$

对偶问题

原问题 凸优化

$$\begin{aligned} & \min_x f(x) \\ \text{s.t. } & g_i(x) \leq 0, i = 1, 2, \dots, m \\ & h_j(x) = 0, j = 1, 2, \dots, n \end{aligned}$$

对偶问题 凸优化

$$\max_{u,v} g(u, v) \text{ s.t. } u \geq 0$$

$$\text{其中 } g(u, v) = \min_x L(x, u, v)$$

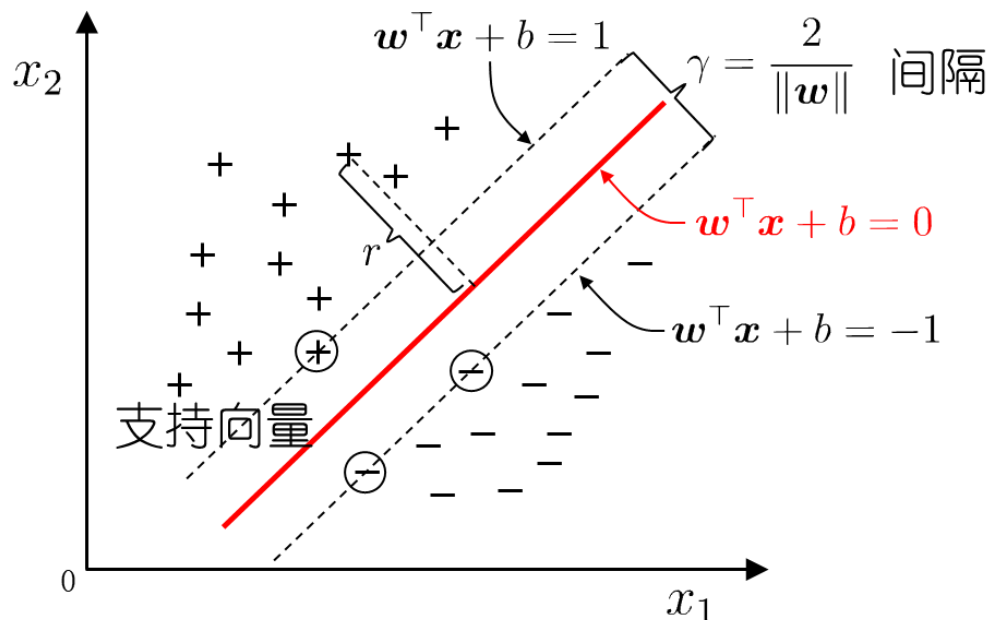
强对偶性

$$f^* = g^* = \max_{u,v} g(u, v)$$

Slater条件

原问题为凸优化问题，且可行域中至少有一个点使得不等式约束严格成立

对偶问题



- 对于线性可分问题, 一定存在 (w, b) 使得 $y_i(w^T x_i + b) \geq 1$, $i = 1, 2, \dots, m$
- 设 $\epsilon < 1$, 则 $y_i(w^T x_i + b) > \epsilon \implies y_i \left(\frac{w^T}{\epsilon} x_i + \frac{b}{\epsilon} \right) > 1$

则 $(\frac{w}{\epsilon}, \frac{b}{\epsilon})$ 严格满足不等式。因此, 该优化问题满足强对偶性条件

对偶问题

对偶
问题

$$\begin{aligned} \max_{\alpha} g(\alpha) &= \sum_{i=1}^m \alpha_i - \frac{1}{2} \sum_{i=1}^m \sum_{j=1}^m \alpha_i \alpha_j y_i y_j \underbrace{\mathbf{x}_i^T \mathbf{x}_j}_{K_{ij}} \\ \text{s.t. } \alpha &\geq 0 \text{ and } \sum_{i=1}^m \alpha_i y_i = 0 \end{aligned}$$

通过序列最小优化（SMO）求解最优的 α

基本思路：不断执行如下两个步骤直至收敛

第一步：选取一对需更新的变量 α_i 和 α_j

第二步：固定 α_i 和 α_j 以外的参数，求解对偶问题更新 α_i 和 α_j

$$g(\alpha_i, \alpha_j) = \alpha_i + \alpha_j - \frac{1}{2} \left(\alpha_i^2 K_{ii} + \alpha_j^2 K_{jj} + 2\alpha_i \alpha_j y_i y_j K_{ij} + \alpha_i y_i \sum_{k \neq i, j} \alpha_k y_k K_{ik} + \alpha_j y_j \sum_{k \neq i, j} \alpha_k y_k K_{jk} \right)$$

仅考虑 α_i 和 α_j 时，对偶问题的约束变为

$$\alpha_i y_i + \alpha_j y_j = - \sum_{k \neq i, j} \alpha_k y_k = \zeta \quad \Rightarrow \quad \alpha_i = (\zeta - \alpha_j y_j) y_i$$

对偶问题

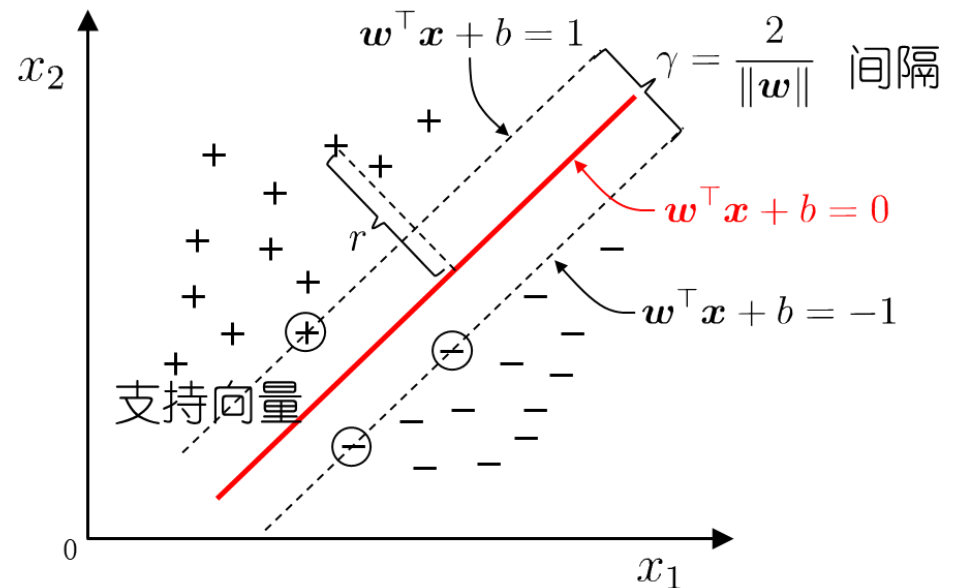
- 终止条件：KKT条件

$$\begin{cases} \alpha_i \geq 0 \\ y_i(\mathbf{w}^\top \mathbf{x}_i + b) \geq 1 \\ \alpha_i(y_i(\mathbf{w}^\top \mathbf{x}_i + b) - 1) = 0 \end{cases}$$

$$y_i(\mathbf{w}^\top \mathbf{x}_i + b) > 1 \quad \Rightarrow \quad \alpha_i = 0$$

$$\alpha_i > 0 \quad \Rightarrow \quad y_i(\mathbf{w}^\top \mathbf{x}_i + b) = 1$$

支持向量机解的稀疏性：
训练完成后，大部分的训练样本都不需保留，最终模型仅与支持向量有关



对偶问题

- 如何选择两个变量?

$$\begin{cases} \alpha_i \geq 0 \\ y_i(\mathbf{w}^\top \mathbf{x}_i + b) \geq 1 \\ \alpha_i(y_i(\mathbf{w}^\top \mathbf{x}_i + b) - 1) = 0 \end{cases}$$

若 α_i 和 α_j 有一个违反了KKT条件，目标函数就会在迭代后增大

Osuna et al. 1997

- KKT条件违背的程度越大，则变量更新后可能导致的目标函数值增幅越大

第一个变量：选取违背KKT条件程度最大的变量
第二个变量：与第一个变量的间隔最大的变量

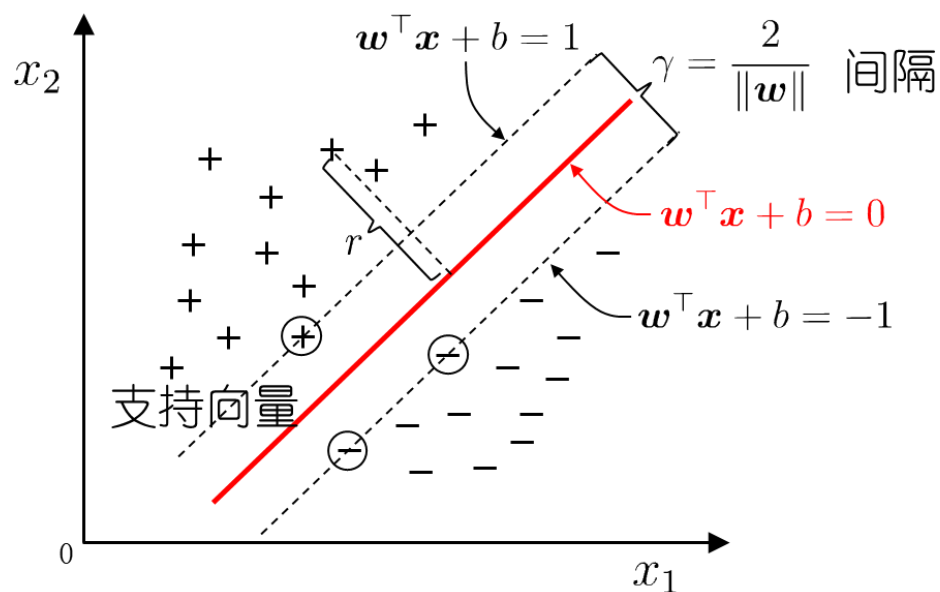
对偶问题

给定最优 α , 求解 w 和 b

$$\alpha^* = \max_{\alpha} g(\alpha)$$



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



对于任意的支持向量 $x_i \in S$, 均满足 $w^T x_i + b = y_i$, 则

$$b^* = \frac{1}{|S|} \sum_{i \in S} (y_i - x_i^T w^*)$$

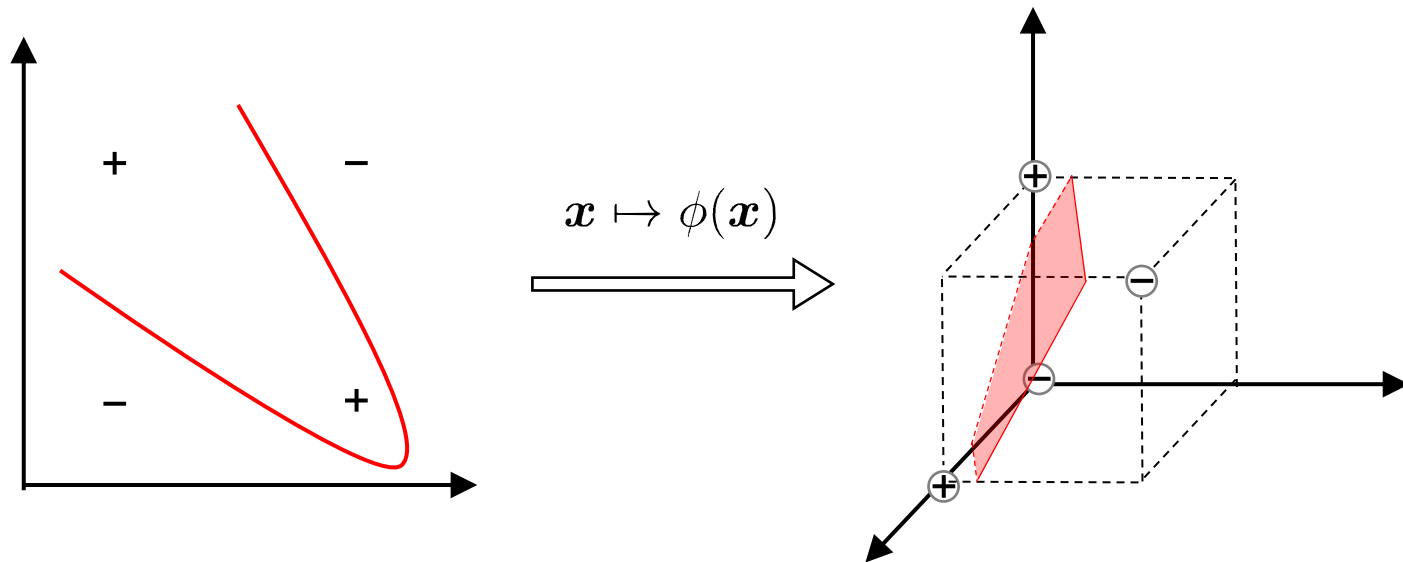
支持向量机解的稀疏性: 训练完成后, 大部分的训练样本都不需保留, 最终模型仅与支持向量有关

$$y = x^T w + b = \sum_i \alpha_i^* y_i x_i^T x_i + b^*$$

线性模型（线性不可分）

- 若不存在一个能正确划分两类样本的超平面, 怎么办

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



核支持向量机

- 设样本 x 映射后的向量为 $\phi(x)$, 划分超平面为 $\mathbf{w}^\top \phi(x) + b = 0$

原问题

$$\begin{aligned} \min_{\mathbf{w}} & \|\mathbf{w}\|^2 / 2 \\ \text{s.t.} & y_i(\mathbf{w}^\top \phi(\mathbf{x}_i) + b) \geq 1, i = 1, 2, \dots, m \end{aligned}$$

对偶问题

$$\begin{aligned} \max_{\alpha} g(\alpha) &= \sum_{i=1}^m \alpha_i - \frac{1}{2} \sum_{i=1}^m \sum_{j=1}^m \alpha_i \alpha_j y_i y_j \phi(\mathbf{x}_i)^\top \phi(\mathbf{x}_j) \\ \text{s.t.} & \alpha \geq 0 \text{ and } \sum_{i=1}^m \alpha_i y_i = 0 \end{aligned}$$

只以内积的形式出现

预测函数

$$y = \mathbf{x}^\top \mathbf{w} + b = \sum_i^m \alpha_i^* y_i \phi(\mathbf{x})^\top \phi(\mathbf{x}_i) + b^*$$

核支持向量机

- 由于特征空间维数可能很高，甚至是无穷维，因此直接计算 $\phi(\mathbf{x})^\top \phi(\mathbf{x}_i)$ 通常是困难的。
- 可以设计函数 $\kappa(\mathbf{x}_i, \mathbf{x}_j) = \phi(\mathbf{x}_i)^\top \phi(\mathbf{x}_j)$

核函数

对偶
问题

$$\begin{aligned} \max_{\alpha} g(\alpha) &= \sum_{i=1}^m \alpha_i - \frac{1}{2} \sum_{i=1}^m \sum_{j=1}^m \alpha_i \alpha_j y_i y_j \kappa(\mathbf{x}_i, \mathbf{x}_j) \\ \text{s.t. } \alpha &\geq 0 \text{ and } \sum_{i=1}^m \alpha_i y_i = 0 \end{aligned}$$

预测
函数

$$y = \mathbf{x}^\top \mathbf{w} + b = \sum_i^m \alpha_i^* y_i \kappa(\mathbf{x}_i, \mathbf{x}) + b^*$$

核支持向量机

- 若已知 $\phi(\cdot)$, 则可写出核函数 $\kappa(\cdot, \cdot)$; 但现实任务中通常不知道 $\phi(\cdot)$ 的形式
 - 核函数是否存在?
 - 什么样的函数可以作为核函数?

令 \mathcal{X} 为输入空间, \mathcal{H} 为特征空间, 如果存在 $\phi(\cdot): \mathcal{X} \mapsto \mathcal{H}$, 使得对所有 $\mathbf{x}, \mathbf{z} \in \mathcal{X}$, 函数 $\kappa(\cdot, \cdot)$ 满足条件

$$\kappa(\mathbf{x}, \mathbf{z}) = \phi(\mathbf{x})^\top \phi(\mathbf{z})$$

则称 $\kappa(\cdot, \cdot)$ 为核函数

Mercer定理: 令 \mathcal{X} 为输入空间, $\kappa(\cdot, \cdot)$ 是定义在 $\mathcal{X} \times \mathcal{X}$ 上的对称函数, 则 κ 是核函数当且仅当对于任意数据集 $D = \{\mathbf{x}_1, \dots, \mathbf{x}_m\}$, 核矩阵 K 总是半正定的

$$K = \begin{bmatrix} \kappa(\mathbf{x}_1, \mathbf{x}_1) & \cdots & \kappa(\mathbf{x}_1, \mathbf{x}_m) \\ \vdots & \ddots & \vdots \\ \kappa(\mathbf{x}_m, \mathbf{x}_1) & \cdots & \kappa(\mathbf{x}_m, \mathbf{x}_m) \end{bmatrix}$$

核支持向量机

- 常用核函数:

名称	表达式	参数
线性核	$\kappa(\mathbf{x}_i, \mathbf{x}_j) = \mathbf{x}_i^\top \mathbf{x}_j$	
多项式核	$\kappa(\mathbf{x}_i, \mathbf{x}_j) = (\mathbf{x}_i^\top \mathbf{x}_j)^d$	$d \geq 1$ 为多项式的次数
高斯核	$\kappa(\mathbf{x}_i, \mathbf{x}_j) = \exp\left(-\frac{\ \mathbf{x}_i - \mathbf{x}_j\ ^2}{2\delta^2}\right)$	$\delta > 0$ 为高斯核的带宽(width)
拉普拉斯核	$\kappa(\mathbf{x}_i, \mathbf{x}_j) = \exp\left(-\frac{\ \mathbf{x}_i - \mathbf{x}_j\ }{\delta}\right)$	$\delta > 0$
Sigmoid核	$\kappa(\mathbf{x}_i, \mathbf{x}_j) = \tanh(\beta \mathbf{x}_i^\top \mathbf{x}_j + \theta)$	\tanh 为双曲正切函数, $\beta > 0, \theta < 0$

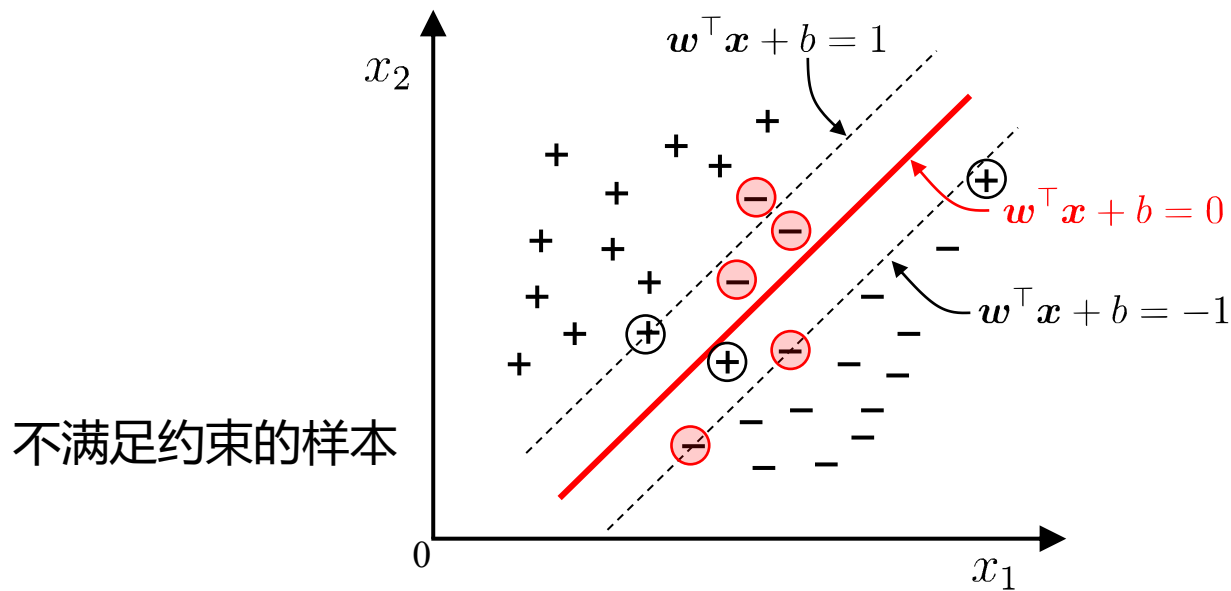
- 函数组合得到

- 核函数线性组合 $\gamma_1 \kappa_1 + \gamma_2 \kappa_2$
- 核函数直积 $\kappa_1 \otimes \kappa_2(\mathbf{x}, \mathbf{z}) = \kappa_1(\mathbf{x}, \mathbf{z}) \otimes \kappa_1(\mathbf{x}, \mathbf{z})$
- $g(\mathbf{x})\kappa(\mathbf{x}, \mathbf{z})g(\mathbf{z})$

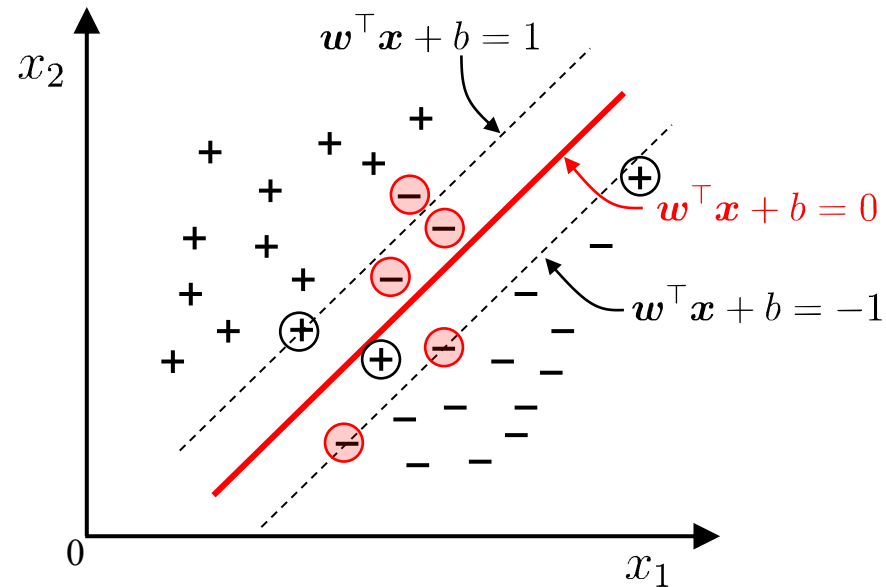
线性模型（线性不可分）

- 很难确定合适的核函数使得训练样本在特征空间中线性可分
- 一个线性可分的结果也很难断定是否是有过拟合造成的

引入**软间隔**的概念, 允许支持向量机在一些样本上不满足约束



软间隔支持向量机



$$\min_{\mathbf{w}, b, \xi} \frac{1}{2} \|\mathbf{w}\|^2 + C \sum_{i=1}^m \xi_i$$

正则化

$$s.t. \quad y_i(\mathbf{w}^T \phi(\mathbf{x}_i) + b) \geq 1 - \xi_i, i = 1, 2, \dots, m$$

为每个样本 (\mathbf{x}_i, y_i) 引入松弛变量 ξ_i $\xi_i \geq 0, i = 1, 2, \dots, m$

软间隔支持向量机

原问题 凸优化

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



广义拉格朗日函数

$$\begin{aligned} L(\mathbf{w}, b, \xi, \alpha, \beta) &= \frac{1}{2} \|\mathbf{w}\|^2 + C \sum_{i=1}^m \xi_i - \sum_{i=1}^m \alpha_i (y_i(\mathbf{w}^\top \phi(\mathbf{x}_i) + b) - 1 + \xi_i) - \sum_i \mu_i \xi_i \\ & \alpha_i, \mu_i \geq 0 \end{aligned}$$



对偶问题

$$\begin{aligned} \max_{\alpha} g(\alpha) &= \sum_{i=1}^m \alpha_i - \frac{1}{2} \sum_{i=1}^m \sum_{j=1}^m \alpha_i \alpha_j y_i y_j \phi(\mathbf{x}_i) \phi(\mathbf{x}_j) \\ \text{s.t. } & C \succcurlyeq \alpha \succcurlyeq 0 \text{ and } \sum_{i=1}^m \alpha_i y_i = 0 \end{aligned}$$

软间隔支持向量机

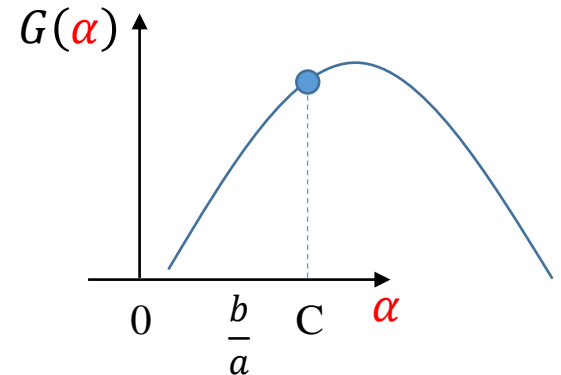
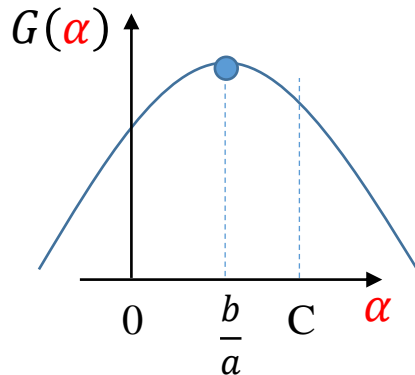
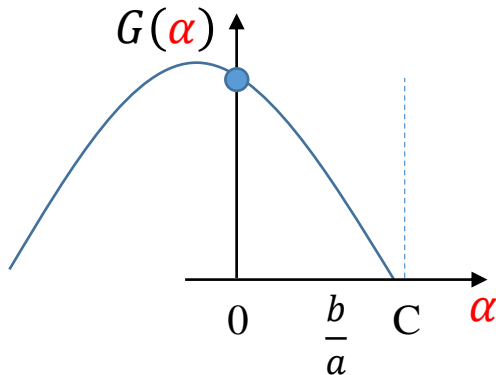
对偶问题

$$\begin{aligned} \max_{\alpha} g(\alpha) &= \sum_{i=1}^m \alpha_i - \frac{1}{2} \sum_{i=1}^m \sum_{j=1}^m \alpha_i \alpha_j y_i y_j \mathbf{x}_i^T \mathbf{x}_j \\ \text{s.t. } C &\geq \alpha \geq 0 \text{ and } \sum_{i=1}^m \alpha_i y_i = 0 \end{aligned}$$



SMO

$$\begin{aligned} \max_{\alpha} G(\alpha) &= -a\alpha^2 + 2b\alpha + c \\ \text{s.t. } 0 &\leq \alpha \leq C \end{aligned}$$



软间隔支持向量机

- 终止条件: KKT条件

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

对于样本 (\mathbf{x}_i, y_i) , 若 $\alpha_i > 0$, 则 $y_i(\mathbf{w}^\top \mathbf{x}_i + b) = 1 - \xi_i$. 该样本是支持向量

若 $\alpha_i < C$, 由 $C = \alpha_i + \mu_i$ 得出 $\mu_i > 0$, 从而得出 $\xi_i = 0$. 该样本在最大间隔边界上

若 $\alpha_i = C$, 由 $C = \alpha_i + \mu_i$ 得出 $\mu_i = 0$

若 $\xi_i \leq 1$, 该样本在最大间隔内

若 $\xi_i > 1$, 该样本被错误分类

软间隔支持向量机

- 目标函数视角

基本想法：最大化间隔的同时, 让不满足约束的样本应尽可能少.

$$\min_{\mathbf{w}, b} \frac{1}{2} \|\mathbf{w}\|^2 + C \sum_{i=1}^m \ell_{0/1}(y_i(\mathbf{w}^\top \phi(\mathbf{x}_i) + b) - 1)$$

其中 $\ell_{0/1}$ 是0/1损失函数

$$\ell_{0/1}(z) = \begin{cases} 1, & z < 0 \\ 0, & \text{otherwise} \end{cases}$$

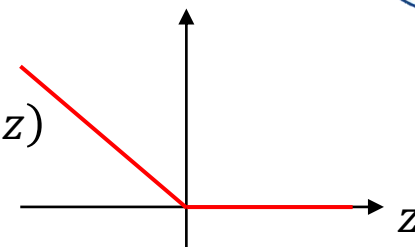
存在的问题：0/1损失函数非凸、非连续, 不易优化!

软间隔支持向量机

$$\ell_{0/1}(z) = \begin{cases} 1, & z < 0 \\ 0, & \text{otherwise} \end{cases}$$



$$\max(0, -z)$$



$$\min_{\mathbf{w}, b} \frac{1}{2} \|\mathbf{w}\|^2 + C \sum_{i=1}^m \max(0, 1 - y_i(\mathbf{w}^\top \phi(\mathbf{x}_i) + b))$$



$$\xi_i = \max(0, 1 - y_i(\mathbf{w}^\top \phi(\mathbf{x}_i) + b))$$

$$\min_{\mathbf{w}, b, \xi} \frac{1}{2} \|\mathbf{w}\|^2 + C \sum_{i=1}^m \xi_i$$

$$\text{s.t. } 1 - y_i(\mathbf{w}^\top \phi(\mathbf{x}_i) + b) \leq \xi_i$$



$$y_i(\mathbf{w}^\top \phi(\mathbf{x}_i) + b) \geq 1 - \xi_i$$

$$0 \leq \xi_i$$

替代损失函数

$$y_i \in \{\pm 1\}$$

$$\text{令 } z = y_i(\mathbf{w}^\top \phi(\mathbf{x}_i) + b)$$

函数间隔

软间隔SVM

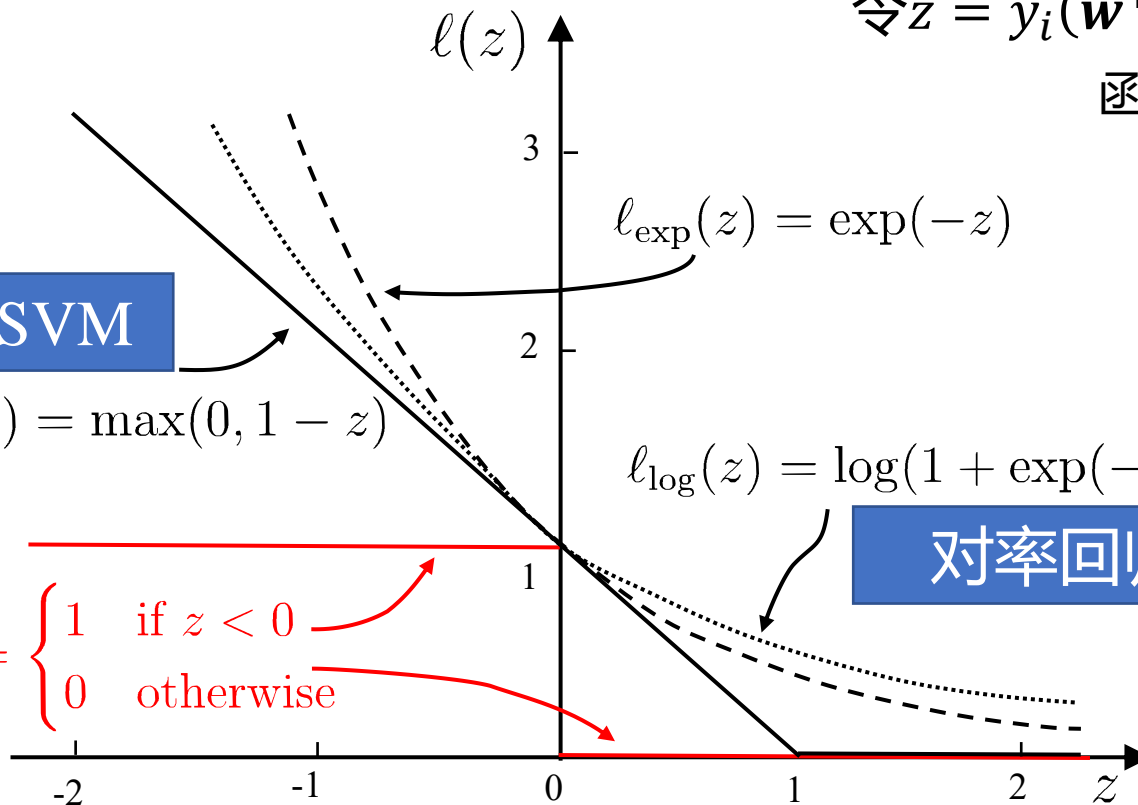
$$\ell_{\text{hinge}}(z) = \max(0, 1 - z)$$

$$\ell_{\text{exp}}(z) = \exp(-z)$$

$$\ell_{\log}(z) = \log(1 + \exp(-z))$$

对率回归


$$\ell_{0/1}(z) = \begin{cases} 1 & \text{if } z < 0 \\ 0 & \text{otherwise} \end{cases}$$



替代损失函数数学性质较好, 一般是0/1损失函数的上界

正则化

- 支持向量机学习模型的更一般形式

$$\min_f \Omega(f) + C \sum_i^m \ell(f(x_i), y_i)$$


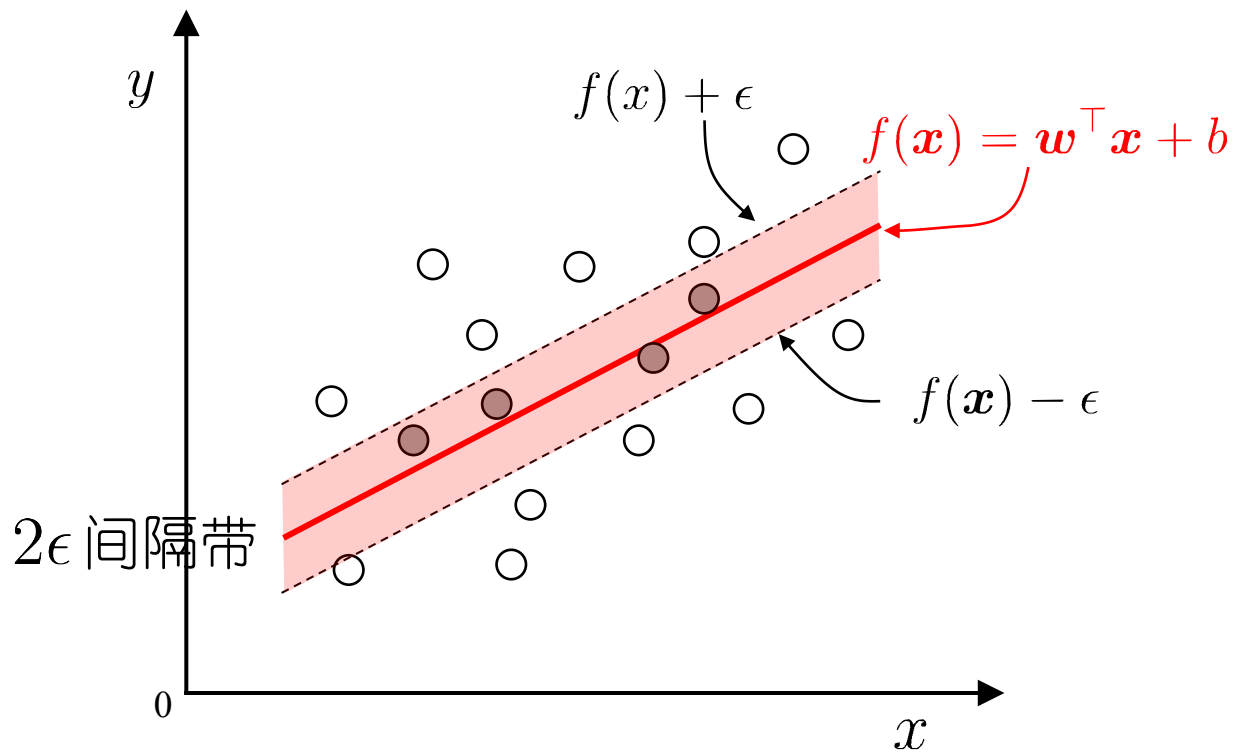
结构风险, 描述
模型的某些性质

经验风险, 描述模型与
训练数据的契合程度

- 通过替换上面两个部分, 可以得到许多其他学习模型
 - 对数几率回归(Logistic Regression)
 - 最小绝对收缩选择算子(LASSO)
 -

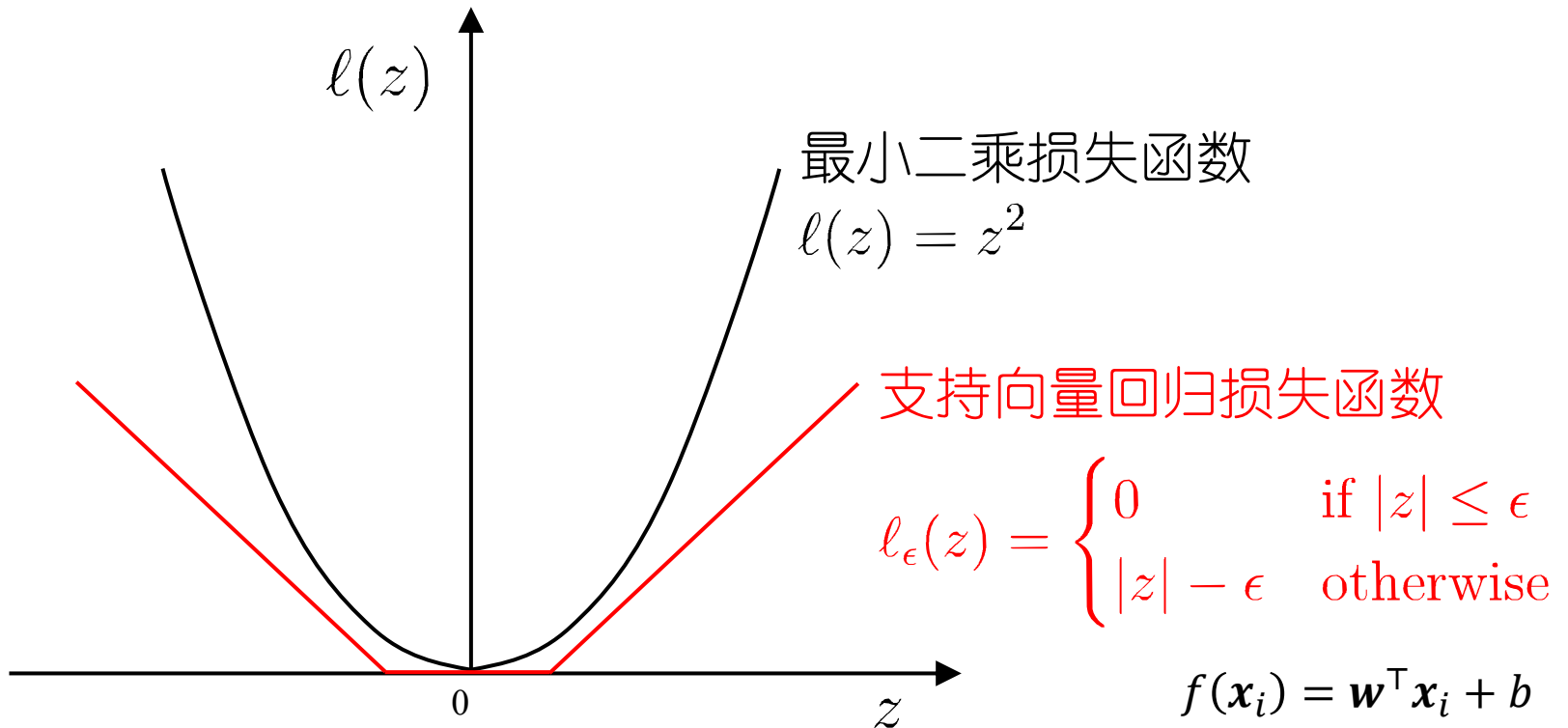
支持向量回归

- 允许模型输出和实际输出间存在 2ϵ 的偏差.



损失函数

- 落入中间 2ϵ 间隔带的样本不计算损失, 从而使得模型获得稀疏性



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

支持向量回归

原问题

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



对偶问题

$$\begin{aligned} \max_{\alpha} g(\alpha, \hat{\alpha}) &= -\frac{1}{2} \sum_{i=1}^m \sum_{j=1}^m (\alpha_i - \hat{\alpha}_i)(\alpha_j - \hat{\alpha}_j) \kappa(\mathbf{x}_i, \mathbf{x}_j) \\ &\quad + \sum_{i=1}^m (y_i(\hat{\alpha}_i - \alpha_i) - \epsilon(\hat{\alpha}_i + \alpha_i)) \\ \text{s.t.} & C \succcurlyeq \alpha, \hat{\alpha} \succcurlyeq 0 \text{ and } \sum_{i=1}^m (\alpha_i - \hat{\alpha}_i) = 0 \end{aligned}$$

预测
函数

$$y = \mathbf{w}^\top \phi(\mathbf{x}) + b = \sum_i^m (\hat{\alpha}_i^* - \alpha_i^*) y_i \kappa(\mathbf{x}_i, \mathbf{x}) + b^*$$

核方法

- 支持向量机
$$y = \mathbf{w}^\top \phi(\mathbf{x}) + b = \sum_i^m \alpha_i^* y_i \kappa(\mathbf{x}_i, \mathbf{x}) + b^*$$
- 支持向量回归
$$y = \mathbf{w}^\top \phi(\mathbf{x}) + b = \sum_i^m (\hat{\alpha}_i^* - \alpha_i^*) y_i \kappa(\mathbf{x}_i, \mathbf{x}) + b^*$$

无论是支持向量机还是支持向量回归, 学得模型总可以表示成
核函数的线性组合

更一般的结论 (表示定理): 令 \mathbb{H} 为核函数 κ 对应的再生核希尔伯特空间, $\|h\|_{\mathbb{H}}$ 表示在 \mathbb{H} 空间中关于 h 的范数, 对于任意单调增函数 $\Omega: [0, \infty) \mapsto \mathbb{R}$ 和任意非负损失函数 $\ell: \mathbb{R}^m \mapsto [0, \infty)$, 优化问题

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

的解总可以写成 $h^* = \sum_i^m \alpha_i \kappa(\cdot, \mathbf{x}_i)$

核线性判别分析

- 通过表示定理可以得到很多线性模型的”核化”版本
 - 核SVM
 - 核LDA
 - 核PCA
 -
- 核LDA: 先将样本映射到高维特征空间, 然后在此特征空间中做线性判别分析

$$\begin{aligned}
 \max_{\mathbf{w}} J(\mathbf{w}) &= \frac{\mathbf{w}^\top \mathbf{S}_b^\phi \mathbf{w}}{\mathbf{w}^\top \mathbf{S}_w^\phi \mathbf{w}} & \mathbf{S}_b^\phi &= (\boldsymbol{\mu}_1^\phi - \boldsymbol{\mu}_0^\phi)(\boldsymbol{\mu}_1^\phi - \boldsymbol{\mu}_0^\phi)^\top \\
 & & \mathbf{S}_w^\phi &= \sum_{i=0}^1 \sum_{x \in X_i} (\phi(x) - \boldsymbol{\mu}_i^\phi)(\phi(x) - \boldsymbol{\mu}_i^\phi)^\top \\
 & \Downarrow & h(x) &= \mathbf{w}^\top \phi(x) = \sum_i^m \alpha_i \kappa(x_i, x) \\
 \max_{\boldsymbol{\alpha}} J(\boldsymbol{\alpha}) &= \frac{\boldsymbol{\alpha}^\top \mathbf{M} \boldsymbol{\alpha}}{\boldsymbol{\alpha}^\top \mathbf{N} \boldsymbol{\alpha}} & \mathbf{M} &= (\hat{\boldsymbol{\mu}}_1 - \hat{\boldsymbol{\mu}}_0)(\hat{\boldsymbol{\mu}}_1 - \hat{\boldsymbol{\mu}}_0)^\top \\
 & & \mathbf{N} &= \mathbf{K} \mathbf{K}^\top - \sum_{i=0}^1 m_i \hat{\boldsymbol{\mu}}_i \hat{\boldsymbol{\mu}}_i^\top
 \end{aligned}$$

核线性判别分析

由表示定理可得 $w = \sum_{i=1}^m \alpha_i \phi(x_i) \quad \mu_i^\phi = \frac{1}{m_i} \sum_{x \in X_i} \phi(x)$

$$S_b^\phi = (\mu_1^\phi - \mu_0^\phi)(\mu_1^\phi - \mu_0^\phi)^\top \quad w^\top \phi(x) = \sum_{i=1}^m \alpha_i \kappa(x_i, x) = \alpha^\top K(:, x)$$

$$w^\top \mu_i^\phi = \frac{1}{m_i} \sum_{x \in X_i} w^\top \phi(x) = \frac{1}{m_i} \sum_{x \in X_i} \alpha^\top K(:, x) = \alpha^\top \underbrace{\frac{1}{m_i} K \mathbf{1}_i}_{\hat{\mu}_i}$$

$$w^\top S_b^\phi w = w^\top (\mu_1^\phi - \mu_0^\phi)(\mu_1^\phi - \mu_0^\phi)^\top w = \alpha^\top \underbrace{(\hat{\mu}_1 - \hat{\mu}_0)(\hat{\mu}_1 - \hat{\mu}_0)^\top}_M \alpha$$

核线性判别分析

$$\mathbf{w} = \sum_{i=1}^m \alpha_i \phi(\mathbf{x}_i) \quad \mathbf{w}^\top \phi(\mathbf{x}) = \sum_{i=1}^m \alpha_i \kappa(\mathbf{x}_i, \mathbf{x}) = \boldsymbol{\alpha}^\top \mathbf{K}(:, \mathbf{x})$$

$$\begin{aligned} \mathbf{w}^\top \mathbf{S}_w^\phi \mathbf{w} &= \sum_{i=0}^1 \sum_{\mathbf{x} \in X_i} \mathbf{w}^\top (\phi(\mathbf{x}) - \boldsymbol{\mu}_i^\phi) (\phi(\mathbf{x}) - \boldsymbol{\mu}_i^\phi)^\top \mathbf{w} \\ &= \sum_{i=1}^m \mathbf{w}^\top \phi(\mathbf{x}_i) \phi(\mathbf{x}_i)^\top \mathbf{w} - \sum_{i=0}^1 m_i \mathbf{w}^\top \boldsymbol{\mu}_i^\phi (\boldsymbol{\mu}_i^\phi)^\top \mathbf{w} \\ &= \boldsymbol{\alpha}^\top \left(\mathbf{K} \mathbf{K}^\top - \sum_{i=0}^1 m_i \hat{\boldsymbol{\mu}}_i \hat{\boldsymbol{\mu}}_i^\top \right) \boldsymbol{\alpha} \end{aligned}$$

$$\sum_{i=1}^m \mathbf{w}^\top \phi(\mathbf{x}_i) \phi(\mathbf{x}_i)^\top \mathbf{w} = \sum_{i=1}^m \boldsymbol{\alpha}^\top \mathbf{K}(:, \mathbf{x}_i) \mathbf{K}(:, \mathbf{x}_i)^\top \boldsymbol{\alpha} = \boldsymbol{\alpha}^\top \mathbf{K} \mathbf{K}^\top \boldsymbol{\alpha}$$

$$\mathbf{w}^\top \boldsymbol{\mu}_i^\phi = \frac{1}{m_i} \sum_{\mathbf{x} \in X_i} \mathbf{w}^\top \phi(\mathbf{x}) = \frac{1}{m_i} \sum_{\mathbf{x} \in X_i} \boldsymbol{\alpha}^\top \mathbf{K}(:, \mathbf{x}) = \boldsymbol{\alpha}^\top \underbrace{\frac{1}{m_i} \mathbf{K} \mathbf{1}_i}_{\hat{\boldsymbol{\mu}}_i}$$

$$\sum_{i=0}^1 m_i \mathbf{w}^\top \boldsymbol{\mu}_i^\phi (\boldsymbol{\mu}_i^\phi)^\top \mathbf{w} = \sum_{i=0}^1 m_i \boldsymbol{\alpha}^\top \hat{\boldsymbol{\mu}}_i \hat{\boldsymbol{\mu}}_i^\top \boldsymbol{\alpha} = \boldsymbol{\alpha}^\top \left(\sum_{i=0}^1 m_i \hat{\boldsymbol{\mu}}_i \hat{\boldsymbol{\mu}}_i^\top \right) \boldsymbol{\alpha}$$

作业

- 6.4
- 6.6
- 6.9

支持向量回归的对偶问题如下,

$$\begin{aligned} \max_{\alpha, \hat{\alpha}} g(\alpha, \hat{\alpha}) = & -\frac{1}{2} \sum_{i=1}^m \sum_{j=1}^m (\alpha_i - \hat{\alpha}_i)(\alpha_j - \hat{\alpha}_j) \kappa(\mathbf{x}_i, \mathbf{x}_j) + \sum_{i=1}^m (y_i(\hat{\alpha}_i - \alpha_i) - \epsilon(\hat{\alpha}_i + \alpha_i)) \\ \text{s.t. } & C \geq \alpha, \hat{\alpha} \geq 0 \text{ and } \sum_{i=1}^m (\alpha_i - \hat{\alpha}_i) = 0 \end{aligned}$$

请将该问题转化为类似于如下标准型的形式 ($\mathbf{u}, \mathbf{v}, \mathbf{K}$ 均已知) ,

$$\begin{aligned} \max_{\alpha} g(\alpha) = & \alpha^T \mathbf{v} - \frac{1}{2} \alpha^T \mathbf{K} \alpha \\ \text{s.t. } & C \geq \alpha \geq 0 \text{ and } \alpha^T \mathbf{u} = 0 \end{aligned}$$

例如在软间隔SVM中 $\mathbf{v} = \mathbf{1}, \mathbf{u} = \mathbf{y}, \mathbf{K}[i, j] = y_i y_j \kappa(\mathbf{x}_i, \mathbf{x}_j)$.