

线性SVM

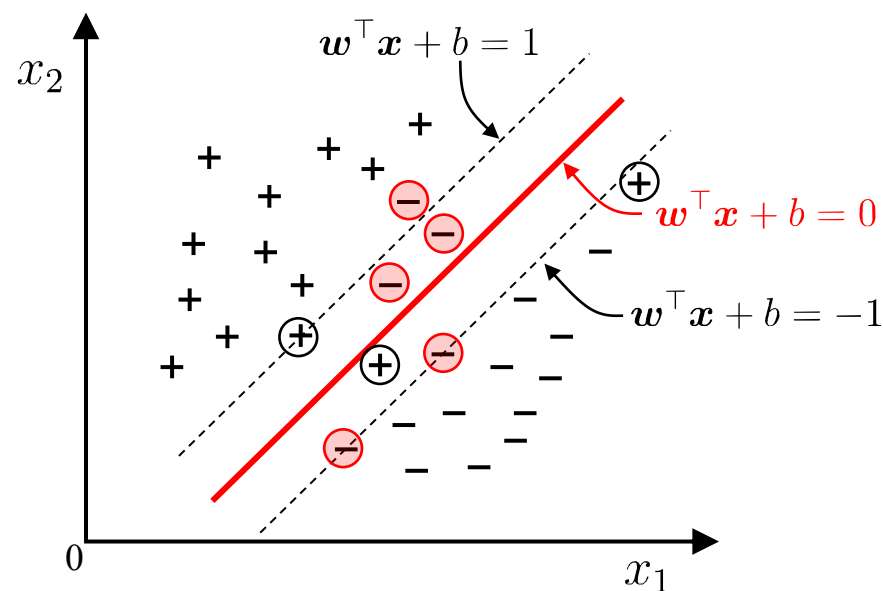
软间隔最大化

-Q: 对于线性可分问题，硬间隔最大化非常好。但如果样本中的噪声和特异点使训练数据线性不可分呢？

-A: 引入“软间隔”的概念，允许支持向量机在一些样本上不满足间隔 ≥ 1 的约束，即对每个样本点引入一个松弛变量 $\xi_i \geq 0$ 。因此，约束条件变为

$$y_i(w \cdot x_i + b) \geq 1 - \xi_i$$

$$\begin{aligned} \min_{w, b} \quad & \frac{1}{2} \|w\|^2 \\ \text{s.t.} \quad & y_i(w \cdot x_i + b) - 1 \geq 0, \quad i = 1, 2, \dots, N \end{aligned}$$



不满足约束的样本

$$\begin{aligned} \min_{w,b} \quad & \frac{1}{2} \|w\|^2 \\ \text{s.t.} \quad & y_i(w \cdot x_i + b) - 1 \geq 0, \quad i=1,2,\dots,N \end{aligned}$$

□ 线性不可分情况下的线性SVM的学习问题变为

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

□ 优化目标的两层含义：

1. 间隔尽量大
2. 不满足约束样本尽可能少

□ C是惩罚参数，超参数，起调和作用

□ 原始优化问题的拉格朗日函数是

$$L(w, b, \xi, \alpha, \mu) = \frac{1}{2} \|w\|^2 + C \sum_{i=1}^N \xi_i - \sum_{i=1}^N \alpha_i (y_i (w \cdot x_i + b) - 1 + \xi_i) - \sum_{i=1}^N \mu_i \xi_i$$
$$\alpha_i \geq 0, \mu_i \geq 0$$

□ 对偶问题是拉格朗日函数的极大极小问题

1. 求 $L(w, b, \xi, \alpha, \mu)$ 对 w, b, ξ 的极小

$$\nabla_w L(w, b, \xi, \alpha, \mu) = w - \sum_{i=1}^N \alpha_i y_i x_i = 0$$

$$\nabla_b L(w, b, \xi, \alpha, \mu) = -\sum_{i=1}^N \alpha_i y_i = 0$$

$$\nabla_{\xi_i} L(w, b, \xi, \alpha, \mu) = C - \alpha_i - \mu_i = 0$$



$$w = \sum_{i=1}^N \alpha_i y_i x_i$$
$$\sum_{i=1}^N \alpha_i y_i = 0$$
$$C - \alpha_i - \mu_i = 0$$

代入

拉格朗日对偶



□ 得到
$$\min_{w, b, \xi} L(w, b, \xi, \alpha, \mu) = -\frac{1}{2} \sum_{i=1}^N \sum_{j=1}^N \alpha_i \alpha_j y_i y_j (x_i \cdot x_j) + \sum_{i=1}^N \alpha_i$$

□ 再对 $\min_{w, b, \xi} L(w, b, \xi, \alpha, \mu)$ 求 α 的极大, 得到对偶问题:

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

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

$$C - \alpha_i - \mu_i = 0$$

$$\alpha_i \geq 0$$

$$\mu_i \geq 0, \quad i=1, 2, \dots, N$$

$$0 \leq \alpha_i \leq C$$

$$\min_{\alpha} \quad \frac{1}{2} \sum_{i=1}^N \sum_{j=1}^N \alpha_i \alpha_j y_i y_j (x_i \cdot x_j) - \sum_{i=1}^N \alpha_i$$

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

$$0 \leq \alpha_i \leq C, \quad i=1, 2, \dots, N$$

□ 设 $\alpha^* = (\alpha_1^*, \alpha_2^*, \dots, \alpha_N^*)^T$ 是对偶问题的一个解, 若存在 α^* 的一个分量 α_j^* , $0 < \alpha_j^* < C$, 则原始问题的解 w^*, b^* 可按下式求得

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

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

KKT条件



$$\nabla_w L(w^*, b^*, \xi^*, \alpha^*, \mu^*) = w^* - \sum_{i=1}^N \alpha_i^* y_i x_i = 0 \longrightarrow w^* = \sum_{i=1}^N \alpha_i^* y_i x_i$$

$$\nabla_b L(w^*, b^*, \xi^*, \alpha^*, \mu^*) = -\sum_{i=1}^N \alpha_i^* y_i = 0$$

$$\nabla_{\xi} L(w^*, b^*, \xi^*, \alpha^*, \mu^*) = C - \alpha^* - \mu^* = 0$$

$$\alpha_i^* (y_i (w^* \cdot x_i + b^*) - 1 + \xi_i^*) = 0$$

$$\mu_i^* \xi_i^* = 0$$

$$y_i (w^* \cdot x_i + b^*) - 1 + \xi_i^* \geq 0$$

$$\xi_i^* \geq 0$$

$$\alpha_i^* \geq 0$$

$$\mu_i^* \geq 0, \quad i=1, 2, \dots, N$$

若存在 α_j^* , $0 < \alpha_j^* < C$,

则 $y_j (w^* \cdot x_j + b^*) - 1 = 0$

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

分离超平面:

$$\sum_{i=1}^N \alpha_i^* y_i (x \cdot x_i) + b^* = 0$$

分类决策函数:

$$f(x) = \text{sign} \left(\sum_{i=1}^N \alpha_i^* y_i (x \cdot x_i) + b^* \right)$$

□ 输入：线性不可分训练数据集 $T = \{(x_1, y_1), (x_2, y_2), \dots, (x_N, y_N)\}$

$$x_i \in \mathcal{X} = \mathbf{R}^n \quad y_i \in \mathcal{Y} = \{-1, +1\}, \quad i = 1, 2, \dots, N$$

□ 输出：分离超平面和分类决策函数

1、构造并求解约束最优化问题

$$\min_{\alpha} \quad \frac{1}{2} \sum_{i=1}^N \sum_{j=1}^N \alpha_i \alpha_j y_i y_j (x_i \cdot x_j) - \sum_{i=1}^N \alpha_i$$

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

$$0 \leq \alpha_i \leq C, \quad i = 1, 2, \dots, N$$

求得最优解： $\alpha^* = (\alpha_1^*, \alpha_2^*, \dots, \alpha_N^*)^T$

2、计算 $w^* = \sum_{i=1}^N \alpha_i^* y_i x_i$

并选择 α^* 的一个分量 α_j^* ，适合条件 $0 < \alpha_j^* < C$ ，计算

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

3、求得分离超平面

$$w^* \cdot x + b^* = 0$$

分类决策函数

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

对任一适合条件 $0 < \alpha_j^* < C$ 的 α_j^* ，都可以求出 b^* ，理论上线性不可分情况的 b^* 不唯一。

软间隔的支持向量



$$\nabla_w L(w^*, b^*, \xi^*, \alpha^*, \mu^*) = w^* - \sum_{i=1}^N \alpha_i^* y_i x_i = 0$$

$$\nabla_b L(w^*, b^*, \xi^*, \alpha^*, \mu^*) = -\sum_{i=1}^N \alpha_i^* y_i = 0$$

$$\nabla_{\xi} L(w^*, b^*, \xi^*, \alpha^*, \mu^*) = C - \alpha^* - \mu^* = 0$$

$$\alpha_i^* (y_i (w^* \cdot x_i + b^*) - 1 + \xi_i^*) = 0$$

$$\mu_i^* \xi_i^* = 0$$

$$y_i (w^* \cdot x_i + b^*) - 1 + \xi_i^* \geq 0$$

$$\xi_i^* \geq 0$$

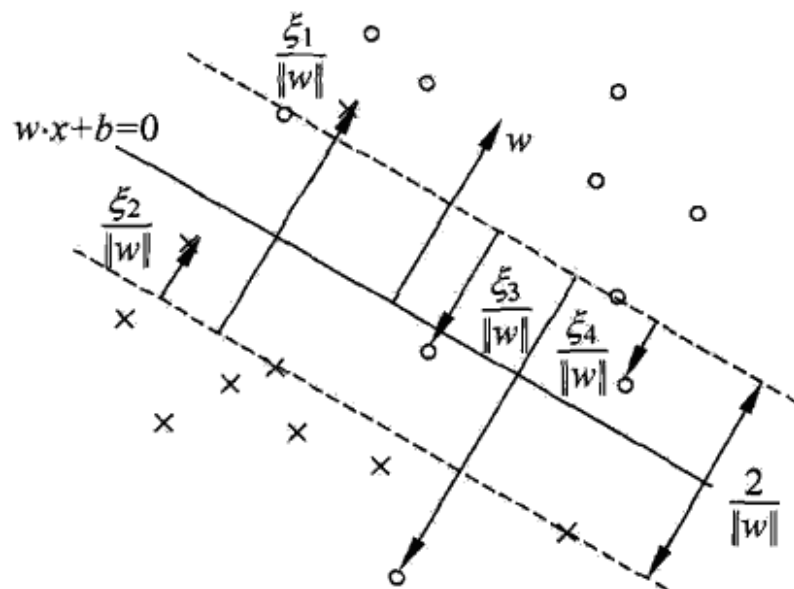
$$\alpha_i^* \geq 0$$

$$\mu_i^* \geq 0, \quad i=1, 2, \dots, N$$

□ 对任一训练样本，总有 $\alpha_i^* = 0$ 或 $y_i f(x_i) = 1 - \xi_i$

□ $\alpha_i^* = 0$ ，该样本对分类面无影响

□ $\alpha_i^* > 0$ ，支持向量



软间隔的支持向量



$$\nabla_w L(w^*, b^*, \xi^*, \alpha^*, \mu^*) = w^* - \sum_{i=1}^N \alpha_i^* y_i x_i = 0$$

$$\nabla_b L(w^*, b^*, \xi^*, \alpha^*, \mu^*) = -\sum_{i=1}^N \alpha_i^* y_i = 0$$

$$\nabla_{\xi} L(w^*, b^*, \xi^*, \alpha^*, \mu^*) = C - \alpha^* - \mu^* = 0$$

$$\alpha_i^* (y_i (w^* \cdot x_i + b^*) - 1 + \xi_i^*) = 0$$

$$\mu_i^* \xi_i^* = 0$$

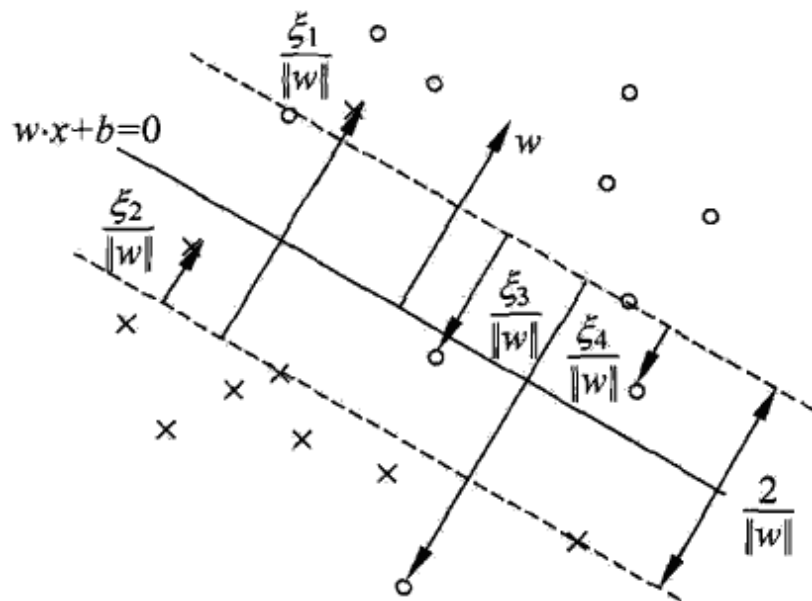
$$y_i (w^* \cdot x_i + b^*) - 1 + \xi_i^* \geq 0$$

$$\xi_i^* \geq 0$$

$$\alpha_i^* \geq 0$$

$$\mu_i^* \geq 0, \quad i=1, 2, \dots, N$$

- 若 $\alpha_i^* < C$ ，则 $\xi_i = 0$ ，间隔边界上
- 若 $\alpha_i^* = C$ ， $0 < \xi_i < 1$ ，分类正确，在间隔边界与超平面之间
- 若 $\alpha_i^* = C$ ， $\xi_i = 1$ ，在超平面上
- 若 $\alpha_i^* = C$ ， $\xi_i > 1$ ，误分类



- 最小化以下目标函数：最大化间隔的同时，让不满足约束的样本应尽可能少

$$\sum_{i=1}^N [1 - y_i (w \cdot x_i + b)]_+ + \lambda \|w\|^2$$

- 合页损失函数 (hinge loss)

$$[z]_+ = \begin{cases} z, & z > 0 \\ 0, & z \leq 0 \end{cases}$$

- 0/1损失不易优化求解

- 合页损失为“替代损失”

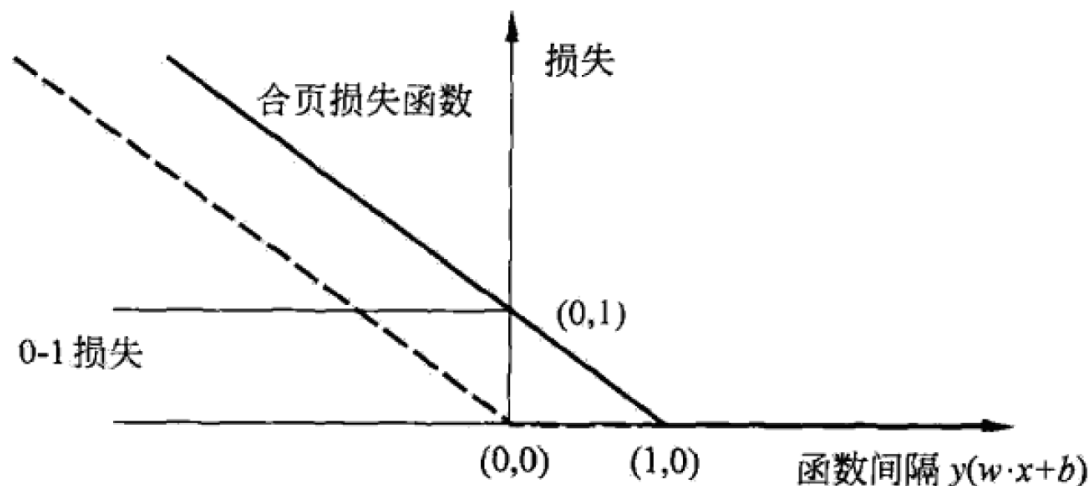


图 7.6 合页损失函数

□ 线性SVM原始最优化问题

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

□ 等价于最优化问题

$$\min_{w, b} \quad \sum_{i=1}^N [1 - y_i(w \cdot x_i + b)]_+ + \lambda \|w\|^2$$

$$\min_{w,b} \lambda \|w\|^2 + \sum_{i=1}^N [1 - y_i (w \cdot x_i + b)]_+$$

描述间隔大小

描述训练集上的误差

- 支持向量机学习模型的更一般形式

$$\min_f \Omega(f) + C \sum_{i=1}^m l(f(\mathbf{x}_i), y_i)$$

结构风险, 描述模型的某些性质,
正则化项

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)
-

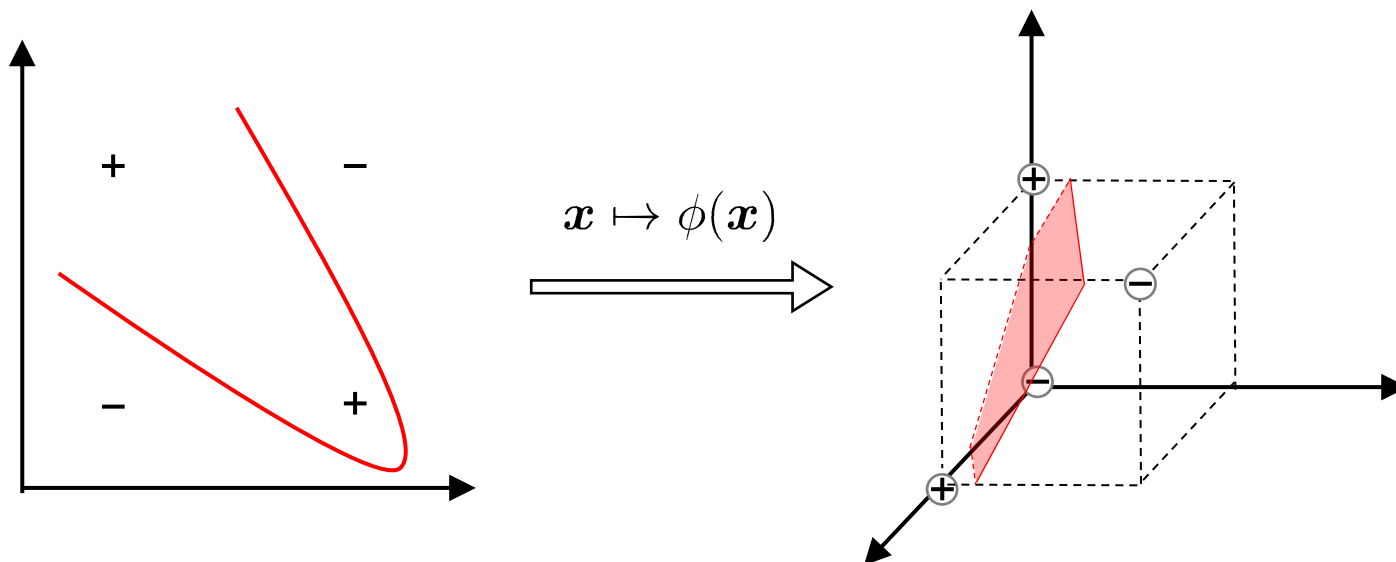
核函数与核方法

线性不可分



-Q: 若不存在一个能正确划分两类样本的线性超平面, 怎么办?

-A: 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.



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

原始问题

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

对偶问题

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

只以内积的形式出现

预测

$$f(\mathbf{x}) = \mathbf{w}^\top \phi(\mathbf{x}) + b = \sum_{i=1}^m \alpha_i y_i \phi(\mathbf{x}_i)^\top \phi(\mathbf{x}) + b$$

- $\phi(x)$ 一般是无限维，可以不知道 $\phi(x)$ 的显式表达，只要知道一个如下所示的核函数，则优化式依然可解。

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

- Mercer定理(充分非必要): 只要一个对称函数所对应的核矩阵半正定, 则它就能作为核函数来使用.

- 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, $\beta > 0, \theta < 0$

- 特征空间的选择对SVM的性能至关重要
- $\phi(x)$ 未知，难以确定什么样的核函数是合适的
- 文本数据通常采用线性核，情况不明时可先尝试高斯核

- 若 κ_1 和 κ_2 为核函数，
 - 对于任意正数 γ_1 、 γ_2 ， $\gamma_1\kappa_1 + \gamma_2\kappa_2$ 也是核函数
 - 直积 $\kappa_1 \otimes \kappa_2 = \kappa_1(x, z)\kappa_2(x, z)$ 也是核函数
 - 对于任意函数 $g(x)$ ， $\kappa(x, z) = g(x)\kappa_1(x, z)g(z)$ 也是核函数

SVM决策函数

$$f(\mathbf{x}) = \mathbf{w}^\top \phi(\mathbf{x}) + b = \sum_{i=1}^m \alpha_i y_i \kappa(\mathbf{x}_i, \mathbf{x}) + b$$

结论：不考虑偏移项**b**，SVM学得模型可以表示成核函数的线性组合。

□ 更一般的结论(表示定理)：对于任意单调增函数 Ω 和任意非负损失函数 l ，
优化问题

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

的解总可以写为 $h^* = \sum_{i=1}^m \alpha_i \kappa(\cdot, \mathbf{x}_i)$.

\mathbb{H} 空间中关于**h**的范数

□ 适用于一般的损失函数和正则项，显式出核函数的巨大威力。

□ 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

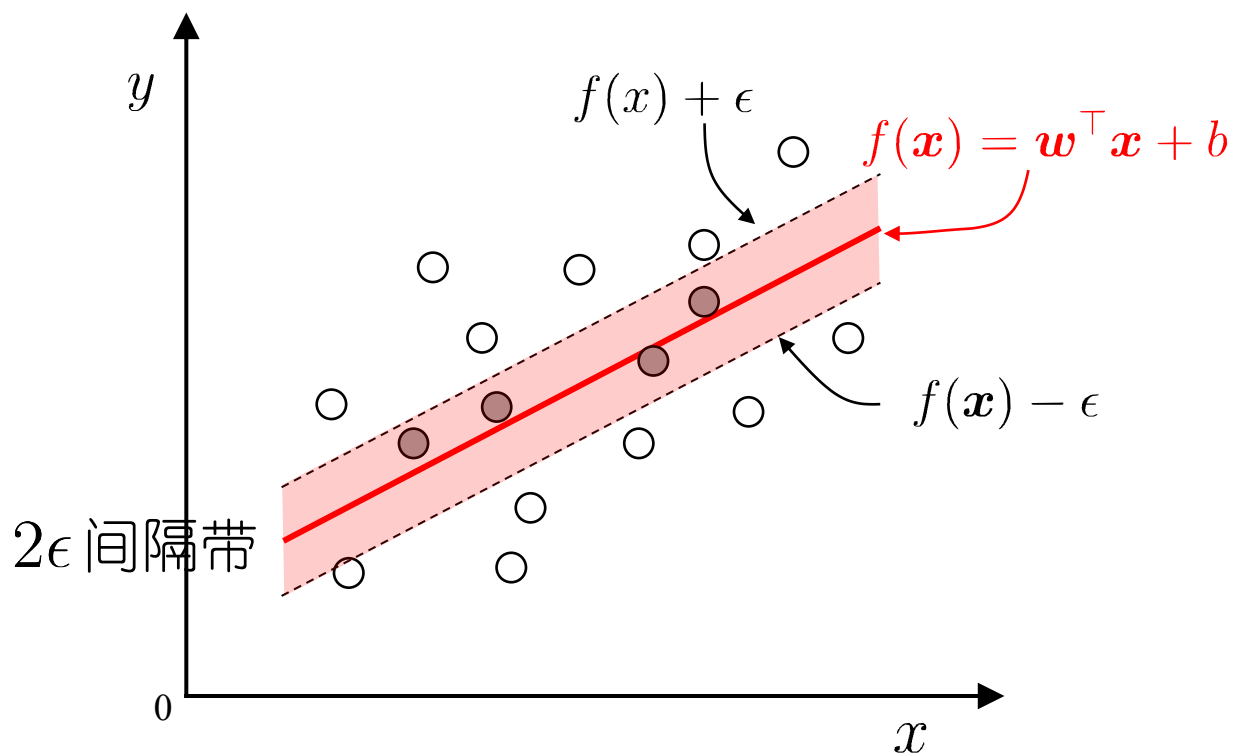
$$\begin{aligned} \max_{\mathbf{w}} J(\mathbf{w}) &= \frac{\mathbf{w}^\top \mathbf{S}_b^\phi \mathbf{w}}{\mathbf{w}^\top \mathbf{S}_w^\phi \mathbf{w}} \\ \downarrow \\ h(\mathbf{x}) &= \mathbf{w}^\top \phi(\mathbf{x}) = \sum_{i=1}^m \alpha_i \kappa(\mathbf{x}_i, \mathbf{x}) \\ \max_{\boldsymbol{\alpha}} J(\boldsymbol{\alpha}) &= \frac{\boldsymbol{\alpha}^\top \mathbf{M} \boldsymbol{\alpha}}{\boldsymbol{\alpha}^\top \mathbf{N} \boldsymbol{\alpha}} \end{aligned}$$

支持向量回归

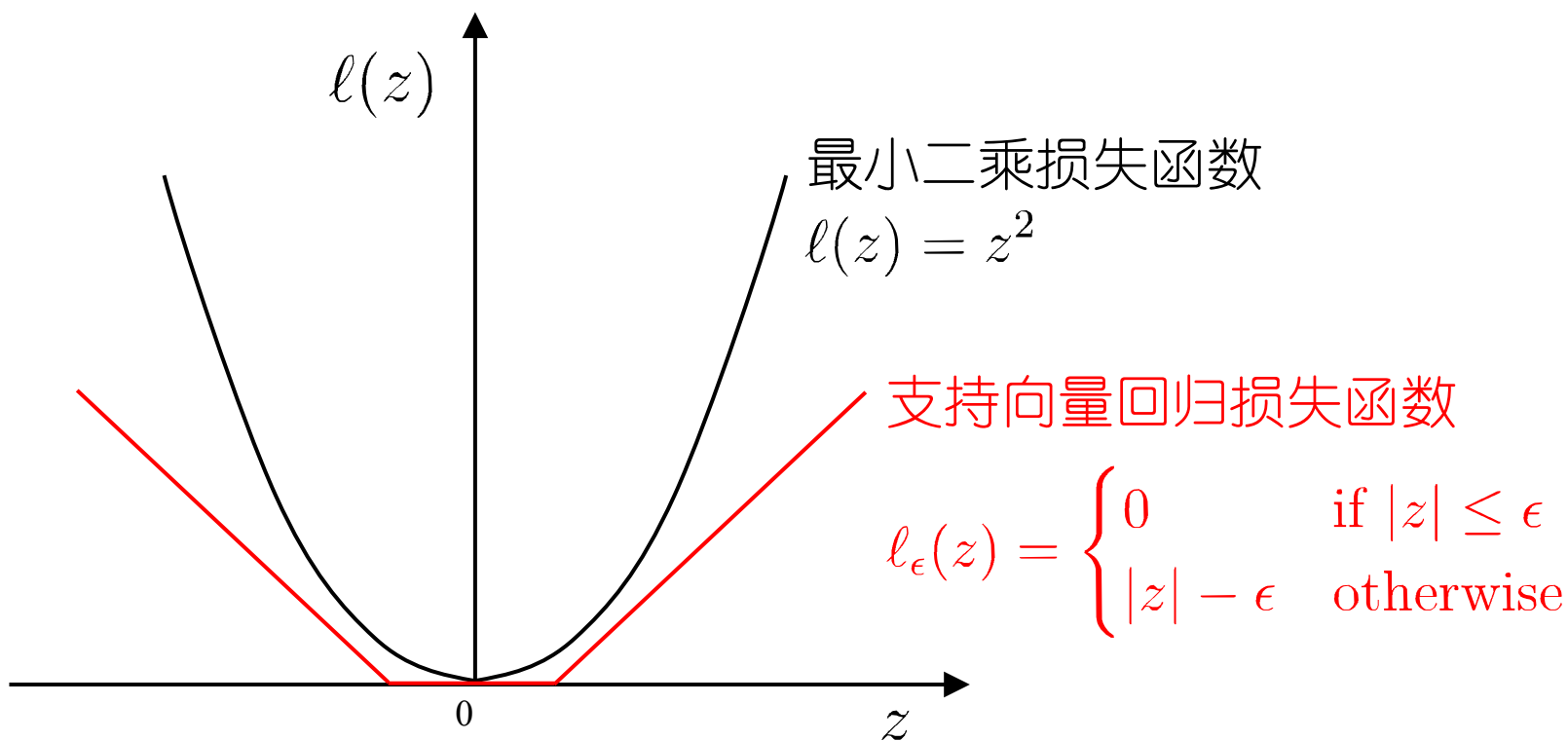
支持向量回归



特点：允许模型输出和实际输出间存在 2ϵ 的偏差。



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.



原始问题

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

对偶问题

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

预测

$$f(\mathbf{x}) = \mathbf{w}^\top \phi(\mathbf{x}) + b = \sum_{i=1}^m (\hat{\alpha}_i - \alpha_i) y_i \kappa(\mathbf{x}_i, \mathbf{x}) + b$$

Take Home Message



大连理工大学 人工智能学院
School of Artificial Intelligence, Dalian University of Technology

- 支持向量机的“最大间隔”思想
- 对偶问题及其解的稀疏性
- 引入“软间隔”缓解特征空间中线性不可分的问题
- 通过向高维空间映射解决线性不可分的问题
- 将核方法推广到其他学习模型
- 将支持向量的思想应用到回归问题上得到支持向量回归

成熟的SVM软件包



大连理工大学 人工智能学院
School of Artificial Intelligence, Dalian University of Technology

□ LIBSVM

<http://www.csie.ntu.edu.tw/~cjlin/libsvm/>

□ LIBLINEAR

<http://www.csie.ntu.edu.tw/~cjlin/liblinear/>

□ SVM^{light}、SVM^{perf}、SVM^{struct}

http://svmlight.joachims.org/svm_struct.html

□ Pegasos

<http://www.cs.huji.ac.il/~shais/code/index.html>