

机器学习

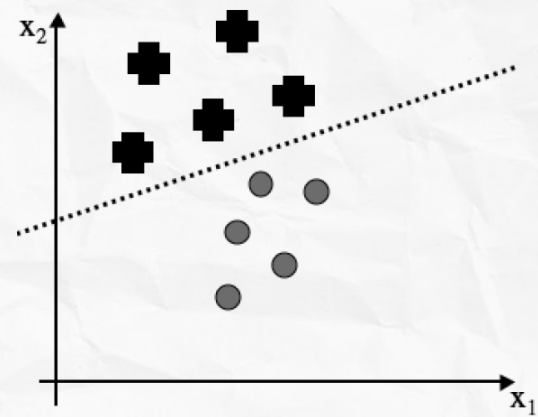
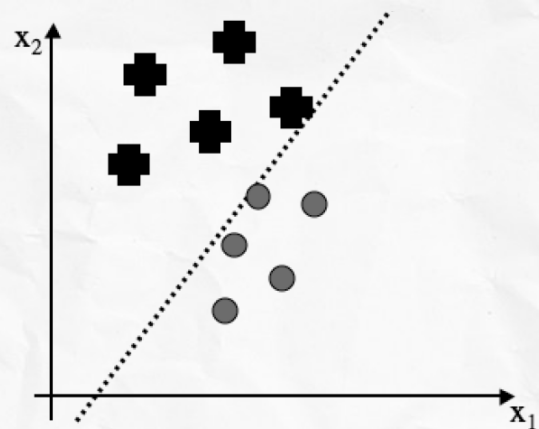
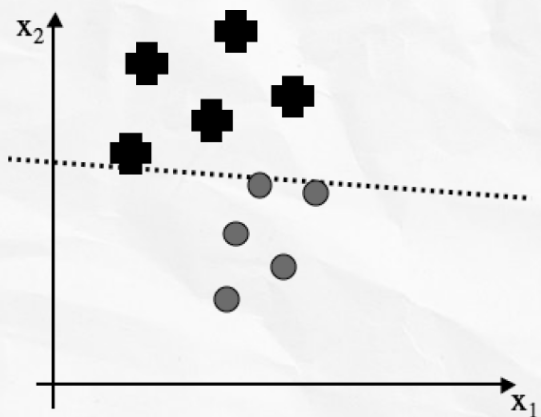
支持向量机

涂文婷

tu.wenting@mail.shufe.edu.cn

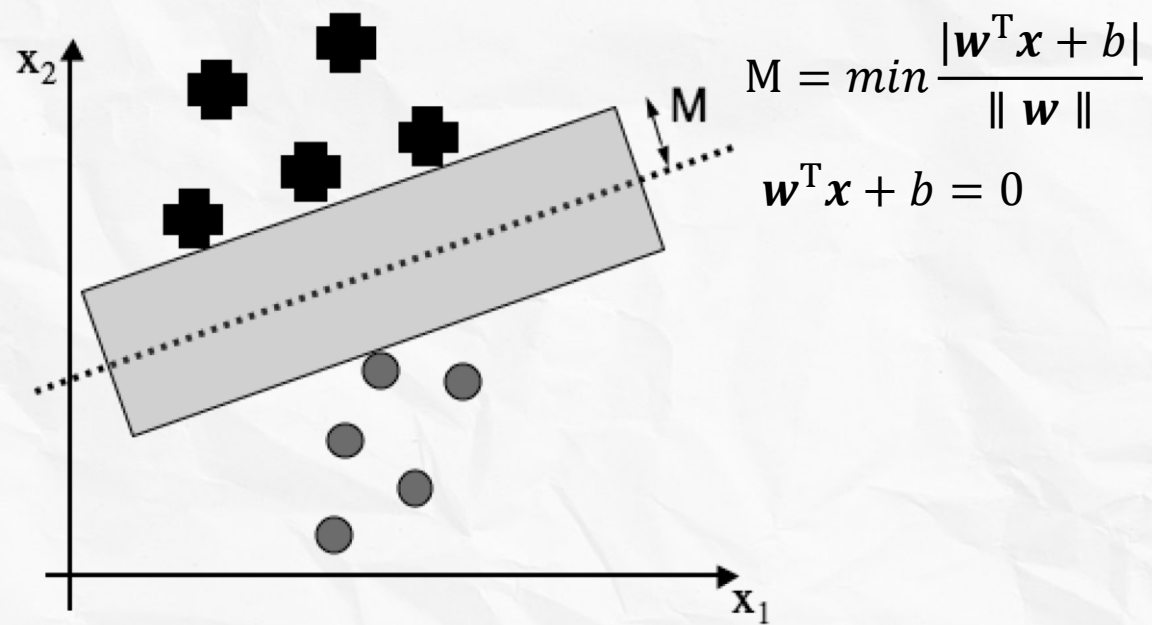
间隔

◦ 分类超平面的优劣



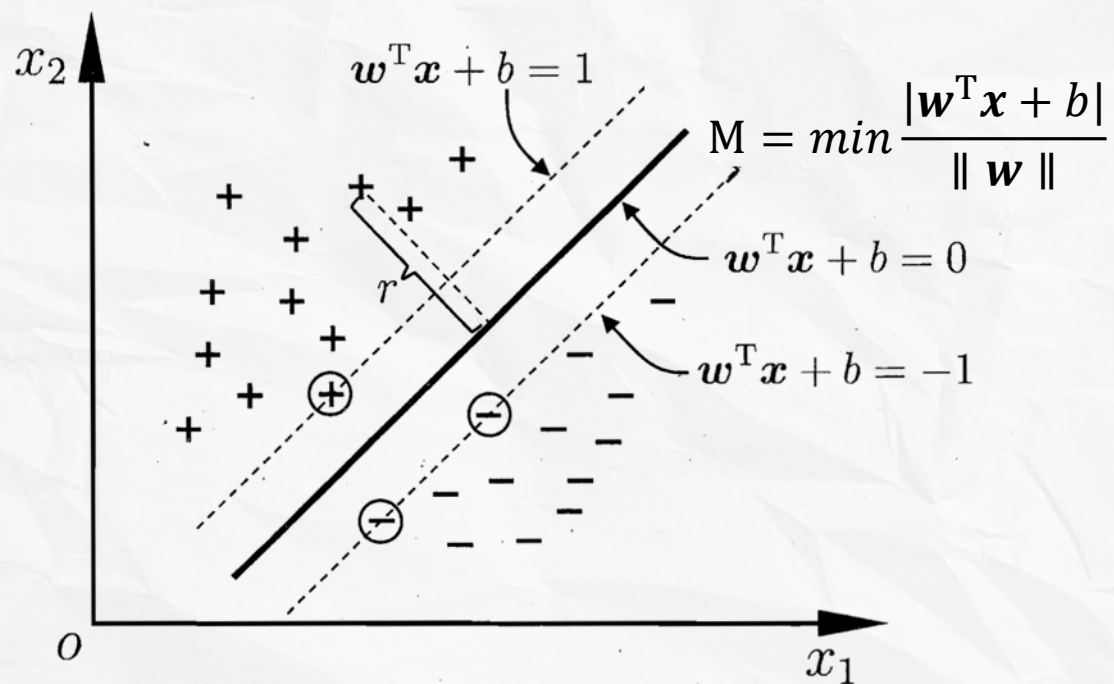
间隔

- 空间间隔



线性可分支持向量机

- 最大空间间隔

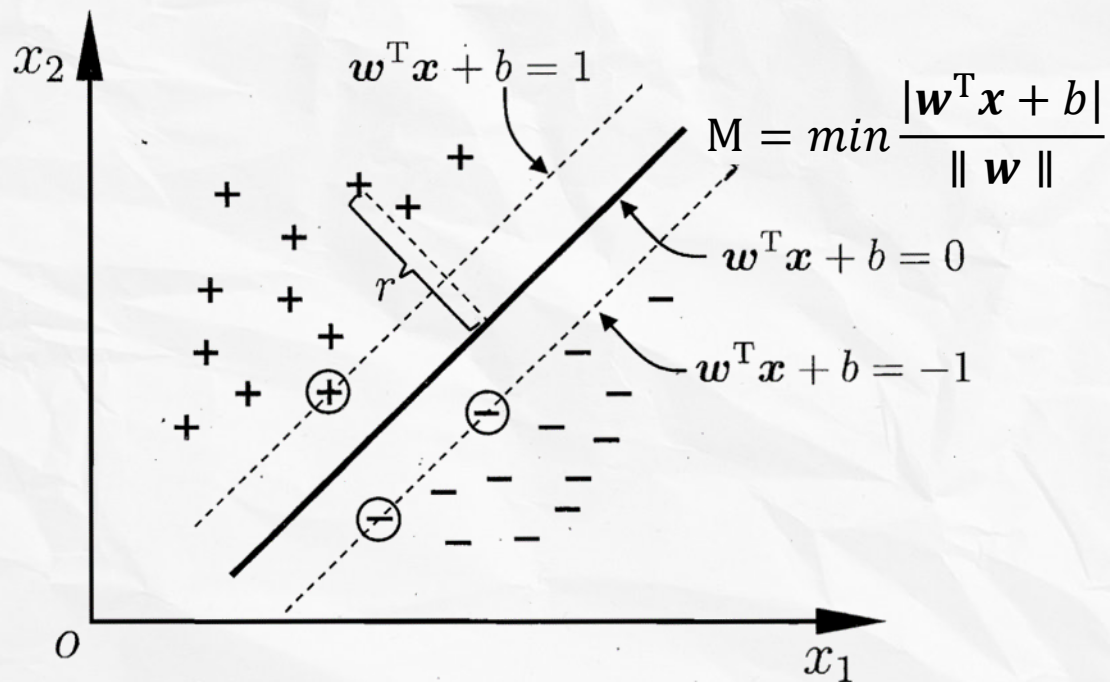


线性可分支持向量机

• 最大空间间隔

$$\max_{w,b} M$$

$$\text{s.t. } \begin{aligned} w^T x_i + b &\geq 0, & y_i &= +1 \\ w^T x_i + b &\leq 0, & y_i &= -1 \end{aligned}$$

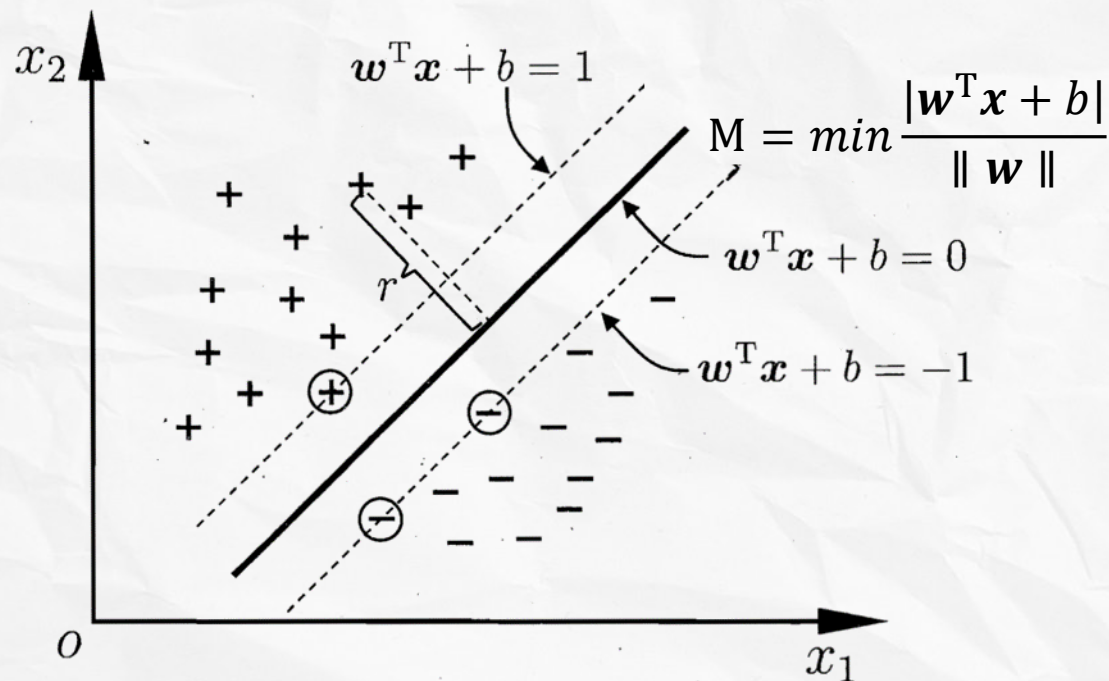


线性可分支持向量机

• 最大空间间隔

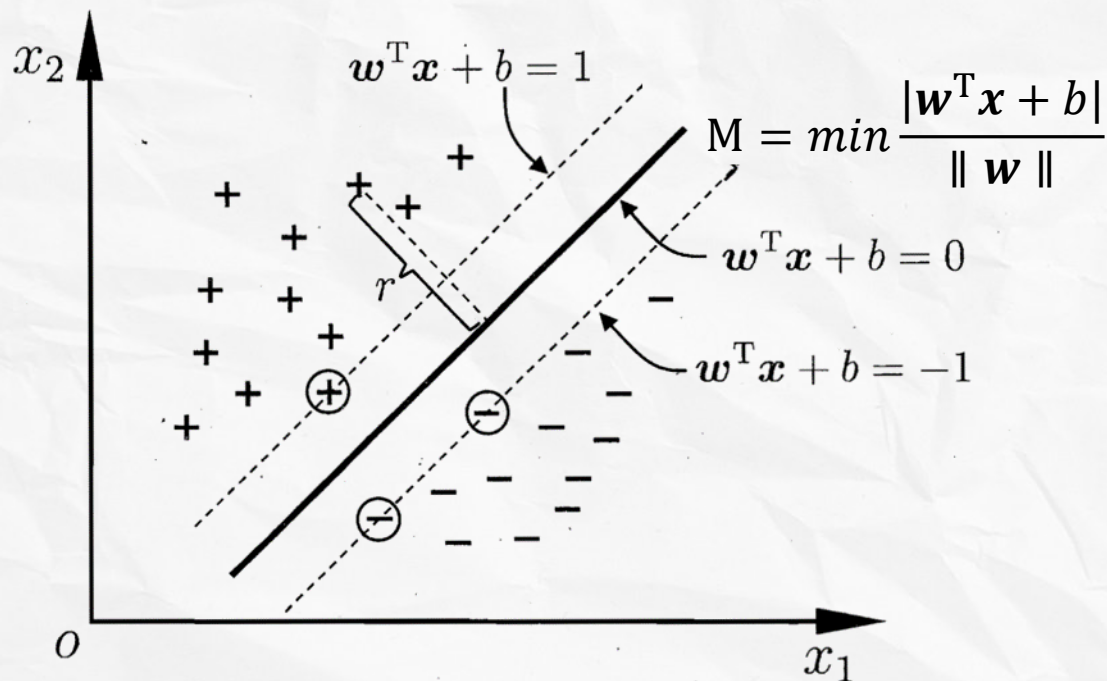
$$\max_{w,b} M$$

$$\text{s.t. } \begin{aligned} w^T x_i + b &\geq 0, & y_i &= +1 \\ w^T x_i + b &\leq 0, & y_i &= -1 \end{aligned} \quad \Leftrightarrow \quad y_i (w^T x_i + b) \geq 0$$



线性可分支持向量机

• 最大空间间隔



$$\max_{w,b} M$$

$$\text{s.t. } \begin{cases} \mathbf{w}^T \mathbf{x}_i + b \geq 0, & y_i = +1 \\ \mathbf{w}^T \mathbf{x}_i + b \leq 0, & y_i = -1 \end{cases} \Leftrightarrow y_i (\mathbf{w}^T \mathbf{x}_i + b) \geq 0$$

上述优化存在无穷多个最优解，加入限制

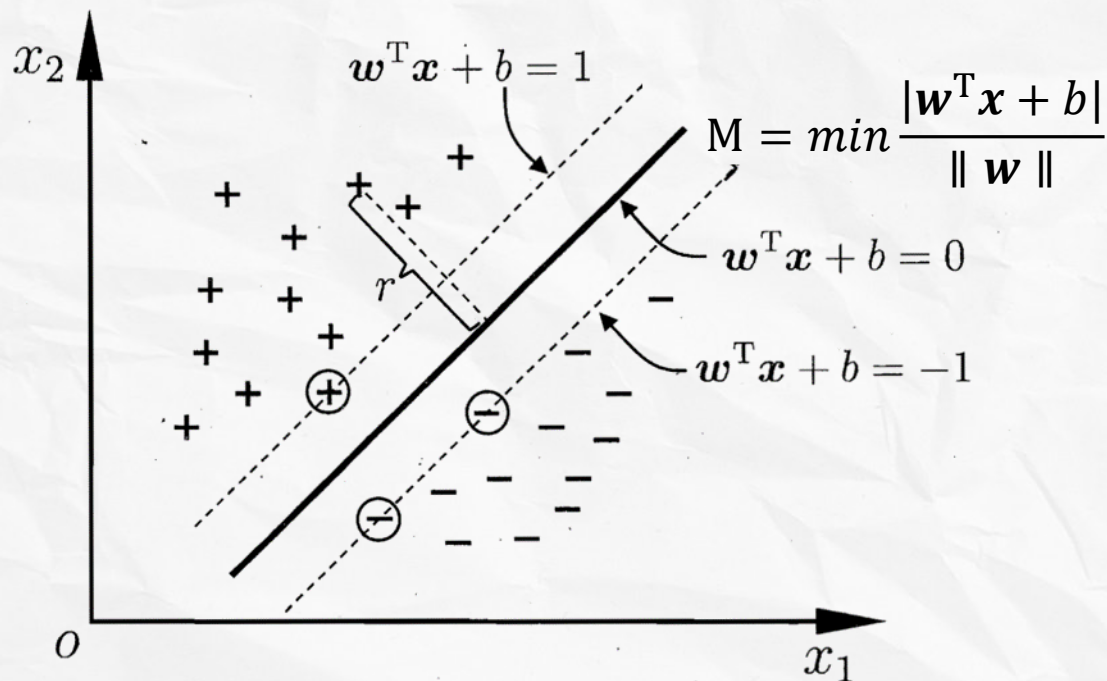
$$y_s (\mathbf{w}^T \mathbf{x}_s + b) = 1$$

离分类超平面最近的点

$$\Leftrightarrow y_i (\mathbf{w}^T \mathbf{x}_i + b) \geq 1$$

线性可分支持向量机

• 最大空间间隔



$$\max_{w,b} M$$

$$\text{s.t. } \begin{cases} \mathbf{w}^T \mathbf{x}_i + b \geq 0, & y_i = +1 \\ \mathbf{w}^T \mathbf{x}_i + b \leq 0, & y_i = -1 \end{cases} \Leftrightarrow y_i (\mathbf{w}^T \mathbf{x}_i + b) \geq 0$$

上述优化存在无穷多个最优解，加入限制

$$y_s (\mathbf{w}^T \mathbf{x}_s + b) = 1$$

离分类超平面最近的点

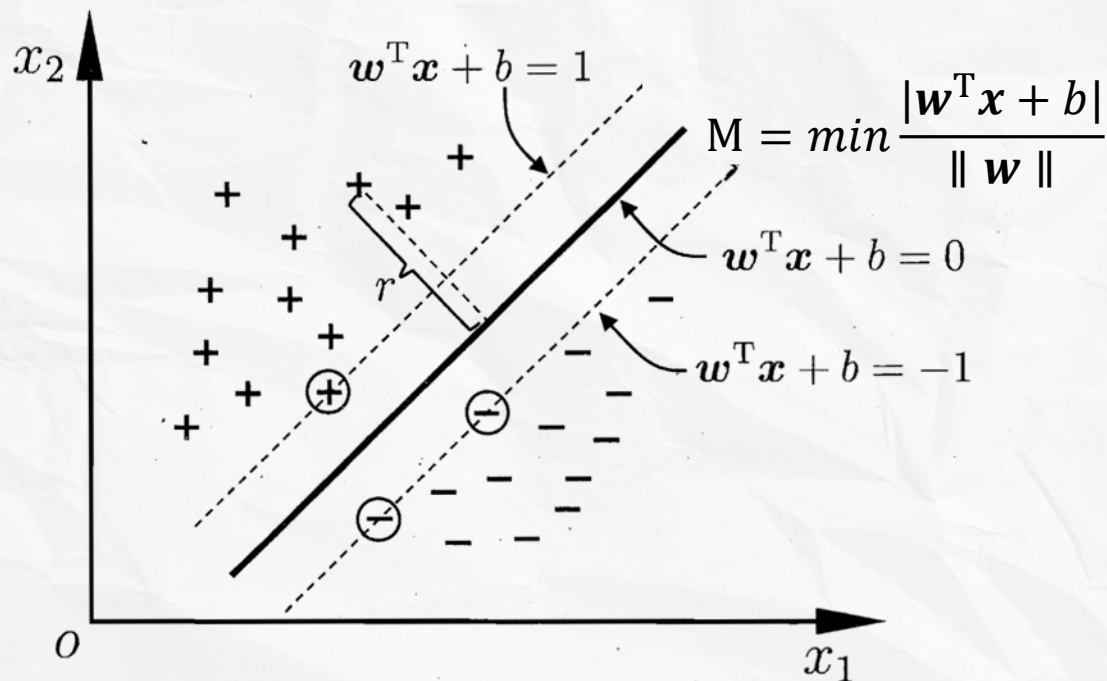
$$\Leftrightarrow y_i (\mathbf{w}^T \mathbf{x}_i + b) \geq 1$$

由于此时 $\min |\mathbf{w}^T \mathbf{x} + b| = \min y (\mathbf{w}^T \mathbf{x} + b) = 1$

$$M = \min \frac{|\mathbf{w}^T \mathbf{x} + b|}{\|\mathbf{w}\|} = \frac{1}{\|\mathbf{w}\|}$$

线性可分支持向量机

• 最大空间间隔



$$\max_{w,b} M$$

$$\text{s.t. } \begin{cases} \mathbf{w}^T \mathbf{x}_i + b \geq 0, & y_i = +1 \\ \mathbf{w}^T \mathbf{x}_i + b \leq 0, & y_i = -1 \end{cases} \Leftrightarrow y_i (\mathbf{w}^T \mathbf{x}_i + b) \geq 0$$

上述优化存在无穷多个最优解，加入限制

$$y_s (\mathbf{w}^T \mathbf{x}_s + b) = 1$$

离分类超平面最近的点

$$\Leftrightarrow y_i (\mathbf{w}^T \mathbf{x}_i + b) \geq 1$$

由于此时 $\min |\mathbf{w}^T \mathbf{x} + b| = \min y (\mathbf{w}^T \mathbf{x} + b) = 1$

$$M = \min \frac{|\mathbf{w}^T \mathbf{x} + b|}{\|\mathbf{w}\|} = \frac{1}{\|\mathbf{w}\|}$$

优化任务变为：

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

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

线性可分支持向量机

• 利用对偶问题求解

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

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

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

令 L 对 \mathbf{w} 和 b 的偏导为零: $\mathbf{w} = \sum_{i=1}^m \alpha_i y_i \mathbf{x}_i, 0 = \sum_{i=1}^m \alpha_i y_i$

$$\max_{\alpha} \sum_{i=1}^m \alpha_i - \frac{1}{2} \sum_{i=1}^m \sum_{j=1}^m \alpha_i \alpha_j y_i y_j \mathbf{x}_i^T \mathbf{x}_j$$

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

$$\alpha_i \geq 0, i = 1, 2, \dots, m$$

解出最优的 α 后

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

线性可分支持向量机

◦ Karush - Kuhn - Tucker条件与支持向量

对于非线性规划（Nonlinear Programming）问题能有最优化解法的一个必要和充分条件是KKT条件：

$$\alpha_i \geq 0$$

$$y_i f(\mathbf{x}_i) - 1 \geq 0$$

$$\alpha_i (y_i f(\mathbf{x}_i) - 1) = 0$$

可以发现，若样本的 α 满足 $\alpha_i > 0$ 则

(1), $y_i f(\mathbf{x}_i) = 1$ 即它是一个离分类超平面最近的点，即是 \mathbf{x}_s ，称为特征向量

(2), 可利用支持向量 $S = \{i \mid \alpha_i > 0, i = 1, 2, \dots, m\}$ 得到偏置项 b

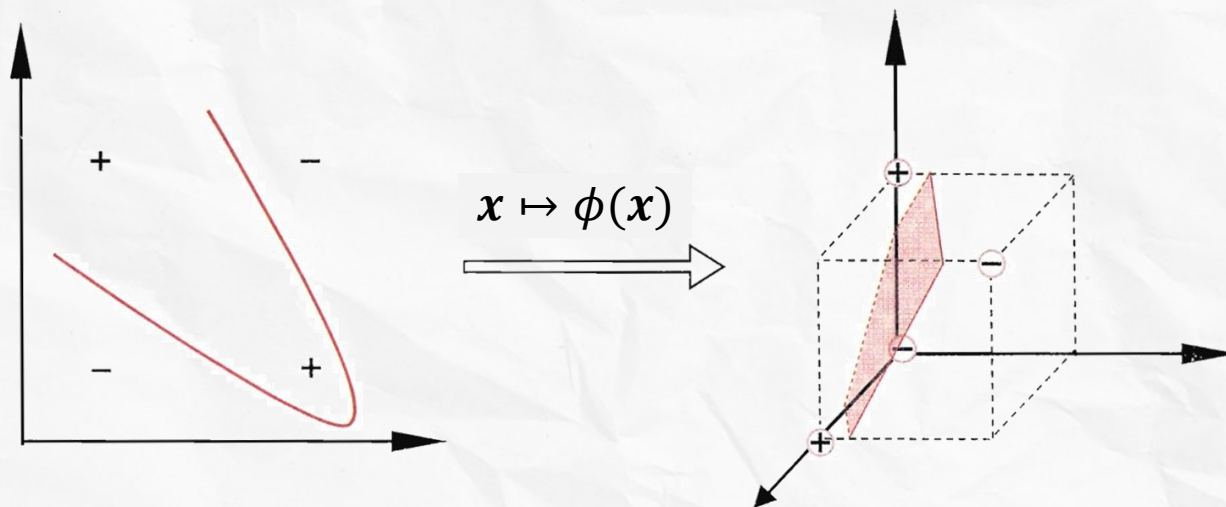
$$y_s \left(\sum_{i \in S} \alpha_i y_i \mathbf{x}_i^T \mathbf{x}_s + b \right) = 1 \rightarrow b = \frac{1}{|S|} \sum_{i \in S} (y_s - \sum_{i \in S} \alpha_i y_i \mathbf{x}_i^T \mathbf{x}_s)$$

(3), 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。支持向量机这个名字强调了此类学习器的关键是如何从支持向量构建出解；同时也暗示着其复杂度主要与支持向量的数目有关。

核支持向量机

◦ 核技术

• 空间变换



如果原始空间是有限维，即属性数有限，那么一定存在一个高维特征空间使样本可分.

核支持向量机

◦ 核技术

• 空间变换下的支持向量机

$$\mathbf{x} \mapsto \phi(\mathbf{x})$$

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

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

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

$$\max_{\alpha} \sum_{i=1}^m \alpha_i - \frac{1}{2} \sum_{i=1}^m \sum_{j=1}^m \alpha_i \alpha_j y_i y_j \phi(\mathbf{x}_i)^T \phi(\mathbf{x}_j)$$

整个过程的求解式总是涉及到 $\phi(\mathbf{x}_i)^T \phi(\mathbf{x}_j)$,
而直接计算高维空间下的样本内积可能是代价高昂的

$$f(\mathbf{x}) = \sum_{i=1}^m \alpha_i y_i \phi(\mathbf{x}_i)^T \phi(\mathbf{x}) + \frac{1}{|S|} \sum_{i \in S} (y_i - \sum_{i \in S} \alpha_i y_i \phi(\mathbf{x}_i)^T \phi(\mathbf{x}_i))$$

核支持向量机

◦ 核技术

• 核函数

$$\kappa(\mathbf{x}, \mathbf{y}) = \phi(\mathbf{x})^T \phi(\mathbf{y})$$

e. g. ,

$$\begin{aligned}\kappa(\mathbf{x}, \mathbf{y}) &= (\mathbf{x}^T \mathbf{y})^2 \\ &= x_1^2 y_1^2 + 2x_1 x_2 y_1 y_2 + x_2^2 y_2^2 \\ &= (x_1^2, \sqrt{2}x_1 x_2, x_2^2)(y_1^2, \sqrt{2}y_1 y_2, y_2^2) \\ &= \phi(\mathbf{x})^T \phi(\mathbf{y}), \phi(\mathbf{v}) = (v_1^2, \sqrt{2}v_1 v_2, v_2^2)\end{aligned}$$

核支持向量机

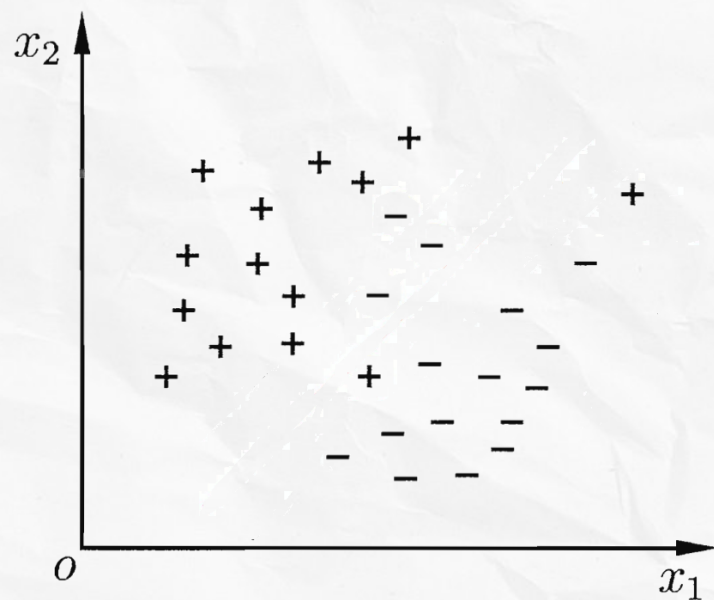
◦ 核技术

• 常用核函数

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

软间隔支持向量机

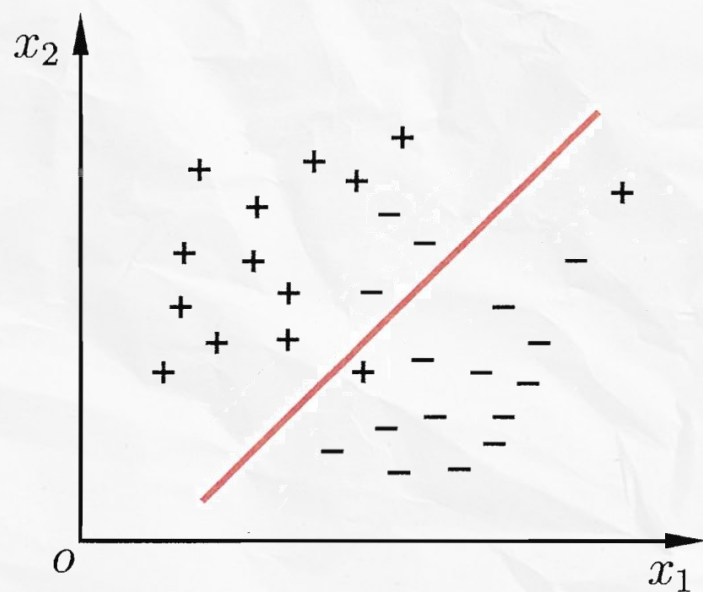
◦ 背景



一些离群点或噪音点可能会造成训练样本仅仅是“接近线性可分的”

软间隔支持向量机

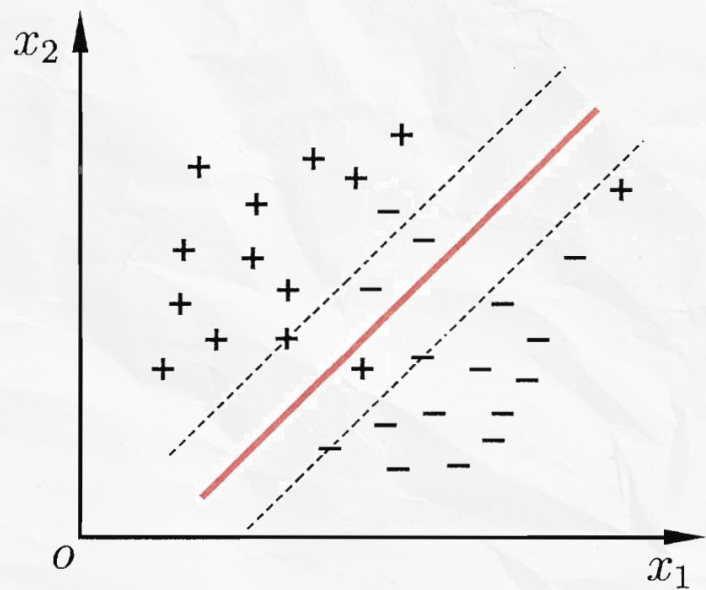
◦ 背景



当一些离群点或噪音点可能会造成训练样本仅仅是“接近线性可分的”，
如何学习到一个可以在一定程度接受破坏间隔约束的数据点存在的分类超平面？

软间隔支持向量机

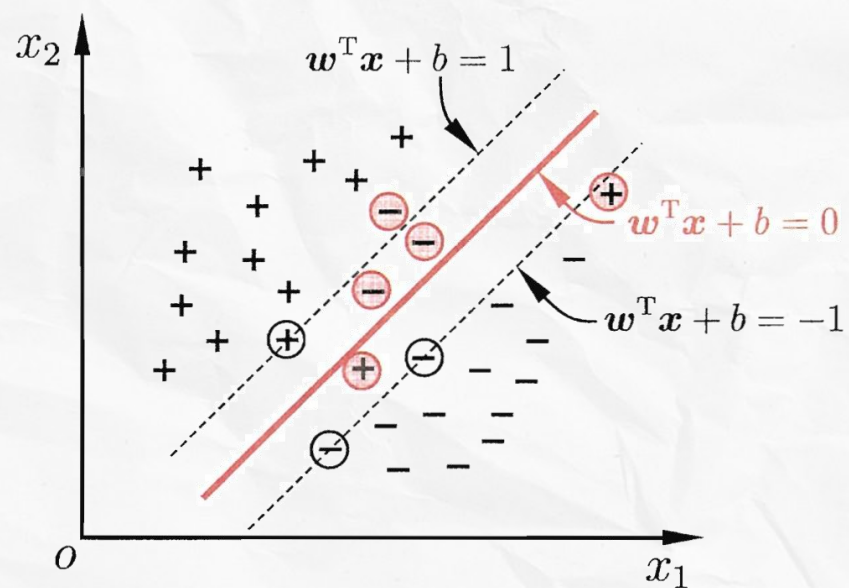
◦ 背景



当一些离群点或噪音点可能会造成训练样本仅仅是“接近线性可分的”，
如何学习到一个可以在一定程度接受破坏间隔约束的数据点存在的分类超平面？

软间隔支持向量机

• 软间隔



硬间隔要求所有样本满足 $y_i(w^T x_i + b) \geq 1$

软间隔允许某些样本不满足 $y_i(w^T x_i + b) \geq 1$

软间隔支持向量机

- 从正则化角度看支持向量机算法

- 硬间隔支持向量机

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

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

\Leftrightarrow

$$\min_{\mathbf{w}, b} \frac{1}{2} \|\mathbf{w}\|^2 + C \sum_{i=1}^m \ell_{0/1}(y_i(\mathbf{w}^T \mathbf{x}_i + b) - 1), C = +\infty, \ell_{0/1}(z) = \begin{cases} 1, & z < 0 \\ 0, & \text{otherwise} \end{cases}$$

软间隔支持向量机

• 从正则化角度看支持向量机算法

• 软间隔支持向量机

• hinge损失

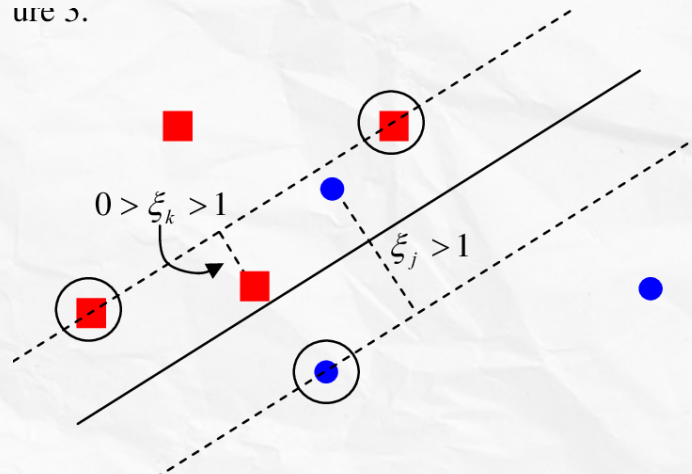
$$\ell_{\text{hinge}}(z) = \max(0, 1 - z)$$

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

ure 3.

• 松弛变量

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



软间隔支持向量机

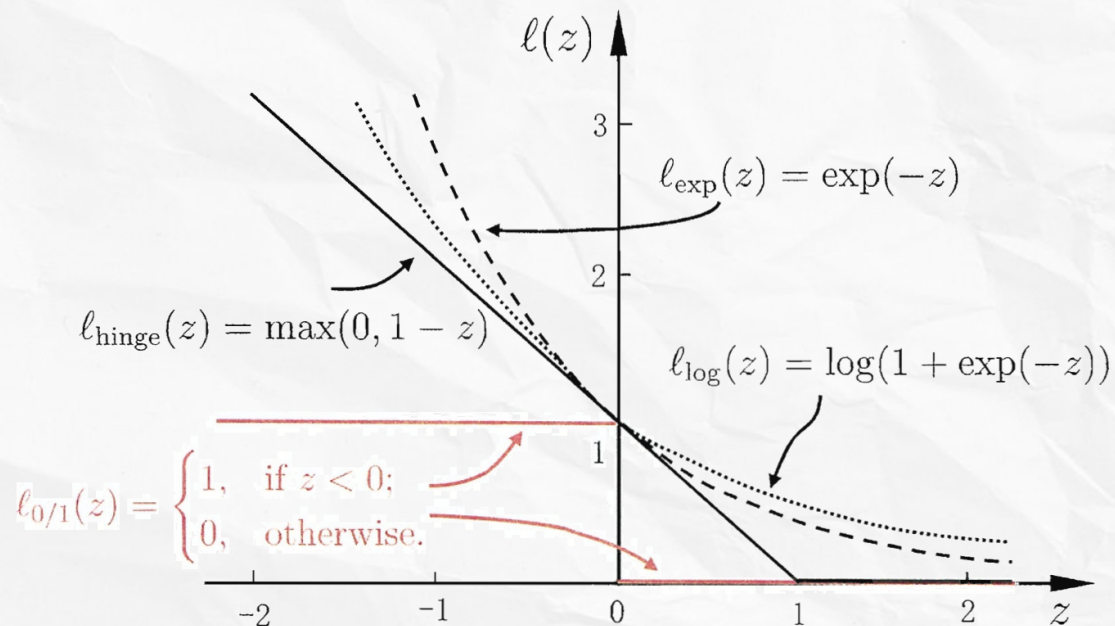
• 常用的分类损失

- 分类损失总与函数间隔 $z = y_i(\mathbf{w}^T \mathbf{x}_i + b)$ 相关

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

指数损失(exponential loss): $\ell_{\text{exp}}(z) = \exp(-z)$

对率损失(logistic loss): $\ell_{\text{log}}(z) = \log(1 + \exp(-z))$

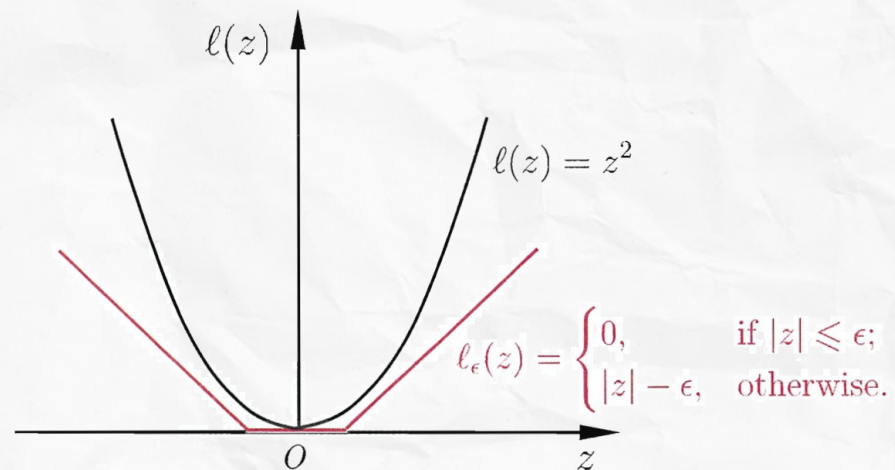


支持向量回归

- ϵ 不敏感损失

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

$$\ell_{\epsilon}(z) = \begin{cases} 0, & \text{if } |z| \leq \epsilon \\ |z| - \epsilon, & \text{otherwise} \end{cases}$$



支持向量机的实现

◦ 求解技术

- 二次规划 *quadratic programming (QP)*
- 序列最小化优化 *Platt's sequential minimal optimization (SMO)*
- 坐标下降法 *coordinate descent*

◦ 工具包

- `sklearn.svm.LinearSVC`
- `sklearn.svm.SVC`
- `sklearn.linear_model.SGDClassifier (loss='hinge')`