

支持向量机 (SVM)

翟婷婷

扬州大学
信息工程（人工智能）学院
zh tt@yzu.edu.cn

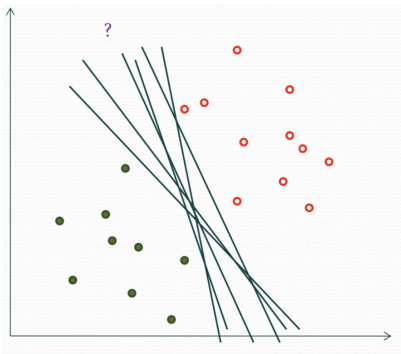
2023年春

课程目标

- 理解硬间隔最大化思想和线性可分SVM算法。
- 掌握拉格朗日对偶性的主要概念和结论，并灵活运用。
- 学会推导支持向量机的对偶优化问题。
- 理解支持向量的作用。

问题的提出

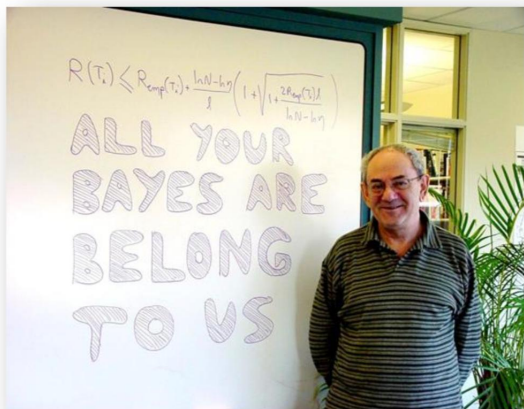
- 给定如图所示的一个线性可分的训练数据集，能将该数据集中的两类样本完全正确分开的直线有很多。
- 利用感知机算法可以找到一条直线，且找到的直线不是唯一的。图中每条直线都可能是感知机算法求得的直线。



- ✓ 哪一条直线最好？
- ✓ 感知机无法回答这个问题。
- ✓ SVM告诉我们，在所有能将训练样本正确分类的直线中**具有最大的间隔**的那条直线是最优的。

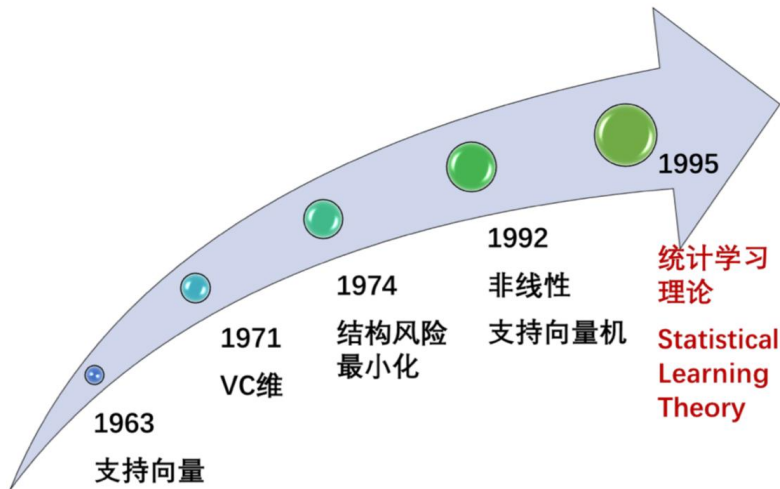
SVM简介

- 支持向量机算法是由俄罗斯著名的统计学家和数学家弗拉基米尔·瓦普尼克 (Vladimir Naumovich Vapnik) 于 1963 年提出的。



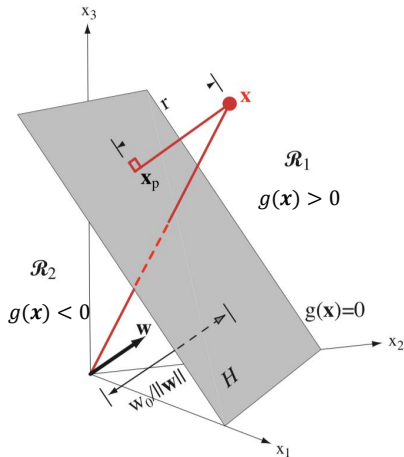
Vladimir N. Vapnik

SVM的发展历程



点到超平面的垂直距离

- 给定一个分割超平面 $\mathbf{w}^T \mathbf{x} + b = 0$ ，其对应的判别函数为 $g(\mathbf{x}) = \mathbf{w}^T \mathbf{x} + b$ ，分割超平面将整个特征空间分为两个决策域： $g(\mathbf{x}) > 0$ 和 $g(\mathbf{x}) < 0$ 。



- ✓ 点 \mathbf{x} 到超平面 $\mathbf{w}^T \mathbf{x} + b = 0$ 的垂直距离为：

$$r = \frac{|g(\mathbf{x})|}{||\mathbf{w}||} = \frac{|\mathbf{w}^T \mathbf{x} + b|}{||\mathbf{w}||}$$

- ✓ r 值的大小可以反应分类的确信程度。 r 值越小，点离决策超平面越近，对其分类的不确信程度就越高。极端情况 $r=0$ ，此时无法确定分类。

间隔

- 给定一个训练数据集 $\{(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_N, y_N)\}$, 其中 $\mathbf{x}_i \in \mathbb{R}^d$, $y_i \in \{+1, -1\}$ 。 $y_i = +1$ 对应的 \mathbf{x}_i 称为正类样本, $y_i = -1$ 对应的 \mathbf{x}_i 称为负类样本。
- (间隔): 超平面 $\mathbf{w}^\top \mathbf{x} + b = 0$ 关于样本点 (\mathbf{x}_i, y_i) 的间隔定义为:

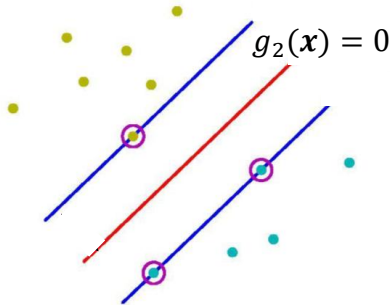
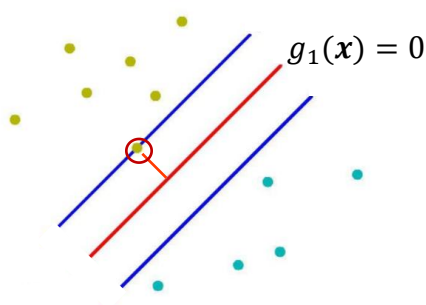
$$r_i = \frac{|\mathbf{w}^\top \mathbf{x}_i + b|}{\|\mathbf{w}\|}$$

超平面 $\mathbf{w}^\top \mathbf{x} + b = 0$ 关于训练数据集的间隔定义为:

超平面关于训练数据集中所有样本点的间隔中的最小间隔:

$$r_{\min} = \min_{i=1, \dots, N} r_i = \min_{i=1, \dots, N} \frac{|\mathbf{w}^\top \mathbf{x}_i + b|}{\|\mathbf{w}\|}$$

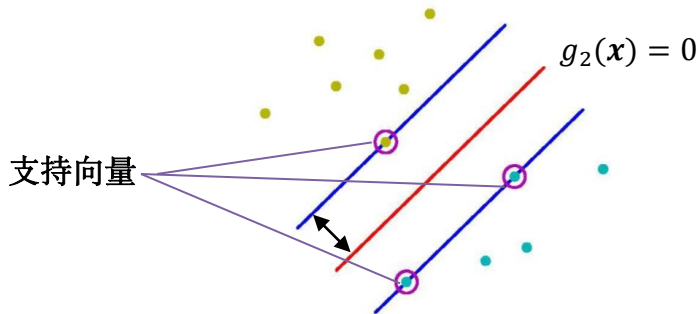
Support Vector Machine (SVM)的思想



- SVM的思想是：寻找能够正确划分训练数据集并且在训练数据集上具有最大间隔的超平面。
- ✓ 左图中的间隔小于右图中的间隔，所以右图的超平面更优。

Support Vector Machine (SVM)的思想

- **间隔最大化**的直观解释：不仅能将训练样本正确分开，而且对于最难分的样本点(离超平面最近)，也能够以足够大的置信度对其分类。满足这些条件的超平面一定具有很好的**泛化性能**。



- ✓ 训练集中具有最小间隔的样本点称为“**支持向量**”。

SVM思想的形式化描述

- 假设训练数据是线性可分的。
- 目标：找到分割超平面 $\mathbf{w}^\top \mathbf{x} + b = 0$ ，满足两个条件：
 - ① 正确划分训练样本： $\forall i = 1 \cdots N$, 都有 $y_i(\mathbf{w}^\top \mathbf{x}_i + b) > 0$
 - ② 最大化在训练集上的间隔：

$$\max_{\mathbf{w}, b} \left(\min_{i=1 \cdots N} \frac{|\mathbf{w}^\top \mathbf{x}_i + b|}{\|\mathbf{w}\|} \right)$$

- 满足条件①的超平面有无穷多个，同时满足①②的超平面只有一个。
- 公式简化：因为条件①， $|\mathbf{w}^\top \mathbf{x}_i + b| = y_i(\mathbf{w}^\top \mathbf{x}_i + b)$

$$\max_{\mathbf{w}, b} \left(\min_{i=1 \cdots N} \frac{y_i(\mathbf{w}^\top \mathbf{x}_i + b)}{\|\mathbf{w}\|} \right) = \max_{\mathbf{w}, b} \left(\frac{1}{\|\mathbf{w}\|} \min_{i=1 \cdots N} y_i(\mathbf{w}^\top \mathbf{x}_i + b) \right)$$

SVM思想的形式化描述

- 令 $\min_{i=1 \dots N} y_i(\mathbf{w}^\top \mathbf{x}_i + b) = 1$, 则有

$$\forall i, y_i(\mathbf{w}^\top \mathbf{x}_i + b) \geq 1$$

- 原优化问题等价于:

$$\max_{\mathbf{w}, b} \frac{1}{\|\mathbf{w}\|}$$

$$s. t. y_i(\mathbf{w}^\top \mathbf{x}_i + b) \geq 1, \forall i = 1 \dots N$$

- 因为:

$$\operatorname{argmax}_{\mathbf{w}} \frac{1}{\|\mathbf{w}\|} = \operatorname{argmin}_{\mathbf{w}} \frac{1}{2} \|\mathbf{w}\|^2$$

- 原优化问题等价于:

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

硬间隔最大化

凸二次规划

线性可分SVM的学习算法

- 输入：线性可分的训练数据集 $\{(\mathbf{x}_1, y_1) \cdots (\mathbf{x}_N, y_N)\}$, 其中 $\mathbf{x}_i \in \mathbb{R}^d$, $y_i \in \{+1, -1\}$.
- 输出：最大间隔的分离超平面和分类决策函数。

1. 求解如下优化问题

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

得到最优解 \mathbf{w}^* , b^* ;

2. 得到分割超平面: $\mathbf{w}^{*\top} \mathbf{x} + b^* = 0$

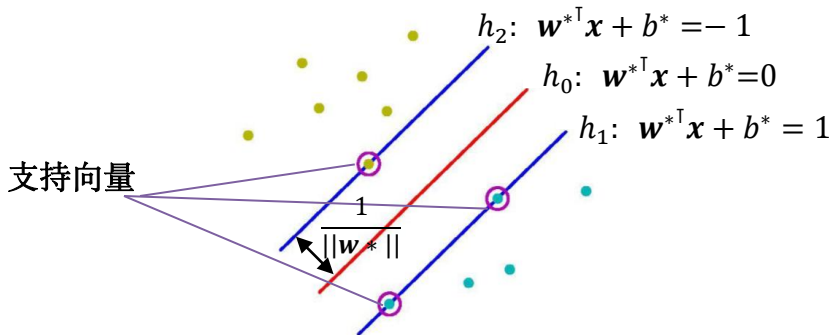
分类决策函数 $f(\mathbf{x}) = \text{sign}(\mathbf{w}^{*\top} \mathbf{x} + b^*)$

硬间隔最大化的SVM学习算法

支持向量

- 训练集中具有最小间隔的样本点称为“支持向量”。
- 支持向量 \mathbf{x}_i 一定满足: $y_i(\mathbf{w}^{*\top}\mathbf{x}_i + b^*) = 1$
- 支持向量到决策超平面 $\mathbf{w}^{*\top}\mathbf{x} + b^* = 0$ 的距离为: $\frac{1}{\|\mathbf{w}^*\|}$
- 对于正类($y_i = +1$)的支持向量在超平面
$$h_1: \mathbf{w}^{*\top}\mathbf{x} + b^* = 1$$
对于负类($y_i = -1$)的支持向量在超平面
$$h_2: \mathbf{w}^{*\top}\mathbf{x} + b^* = -1$$
超平面 h_1 和 h_2 称为间隔边界。

支持向量



- 在决定分类决策超平面的位置时，只有支持向量起作用，其它样本点不起作用，即使删除这些样本点，也不影响分类决策边界。所以这种分类模型称为“支持向量机”。
- 在线性可分的数据集上，支持向量的数目一般很少。

SVM的对偶优化问题

- SVM的原始优化问题:

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

$$s. t. \ y_i(\mathbf{w}^T \mathbf{x}_i + b) \geq 1, \ \forall i = 1 \cdots N$$

- 如何求解这个约束优化问题?
- 在约束最优化问题中, 常常利用拉格朗日对偶性 (Lagrange duality) 将原始问题转换为对偶问题, 通过解对偶问题得到原始问题的解。
- 这么做的好处是:
 - ① 对偶问题更容易求解;
 - ② 自然而然地引入核函数, 进而将算法推广到求解非线性分类问题。

拉格朗日对偶性的主要概念和结论

- 假设 $f_0(\mathbf{x})$, $f_i(\mathbf{x})$, $h_i(\mathbf{x})$ 是定义在 \mathbb{R}^d 上的连续可微函数, 一般的约束优化问题都可以表示为:

$$\begin{aligned} & \min_{\mathbf{x}} f_0(\mathbf{x}) \\ \text{subject to} \quad & f_i(\mathbf{x}) \leq 0, i = 1, 2 \cdots k \\ & h_i(\mathbf{x}) = 0, i = 1, 2 \cdots l \end{aligned} \quad (1)$$

称该问题为原始优化问题。

- (广义拉格朗日函数) 上述问题对应的拉格朗日函数为:

$$L(\mathbf{x}, \boldsymbol{\alpha}, \boldsymbol{\beta}) = f_0(\mathbf{x}) + \sum_{i=1}^k \alpha_i f_i(\mathbf{x}) + \sum_{i=1}^l \beta_i h_i(\mathbf{x})$$

其中 $\alpha_i \geq 0$ 是第 i 个不等式约束的拉格朗日乘子, β_i 是等式约束的乘子, $\boldsymbol{\alpha} = (\alpha_1, \alpha_2, \cdots, \alpha_k) \geq \mathbf{0}$, $\boldsymbol{\beta} = (\beta_1, \beta_2, \cdots, \beta_l)$ 。

原始问题

➤ 考虑如下的函数:

$$\theta_p(\mathbf{x}) = \max_{\alpha, \beta: \alpha \geq 0} L(\mathbf{x}, \alpha, \beta)$$

➤ 思考 $\theta_p(\mathbf{x})$ 的取值是什么?

① 当 \mathbf{x} 满足原始问题(1)的约束条件时,

$$\max_{\alpha, \beta: \alpha \geq 0} L(\mathbf{x}, \alpha, \beta) = f_0(\mathbf{x})$$

② 当 \mathbf{x} 违反原始问题(1)中至少一个约束条件时, 至少存在一个 i 使得 $f_i(\mathbf{x}) > 0$, 或者至少存在一个 i 使得 $h_i(\mathbf{x}) \neq 0$, 此时总是找到一个 $\alpha \geq 0$ 或 β , 使得

$$\max_{\alpha, \beta: \alpha \geq 0} L(\mathbf{x}, \alpha, \beta) = +\infty$$

➤ 综上所述:

$$\theta_p(\mathbf{x}) = \begin{cases} f_0(\mathbf{x}), & \text{如果}\mathbf{x}\text{满足原始问题(1)的约束条件} \\ +\infty, & \text{如果}\mathbf{x}\text{违反原始问题(1)的约束条件} \end{cases}$$

原始问题

- 因此, $\min_x \theta_p(x)$ 这个优化问题等价于原始优化问题(1), 即两个优化问题有相同的解, 又因为

$$\min_x \theta_p(x) = \min_x \max_{\alpha, \beta: \alpha \geq 0} L(x, \alpha, \beta)$$

- 观察到, 原始优化问题(1)等价于拉格朗日函数的极小极大问题。
- 定义如下的函数: $\theta_D(\alpha, \beta) = \min_x L(x, \alpha, \beta)$
- 原始优化问题(1)的对偶优化问题定义为:

$$\max_{\alpha, \beta: \alpha \geq 0} \theta_D(\alpha, \beta) = \max_{\alpha, \beta: \alpha \geq 0} \min_x L(x, \alpha, \beta)$$

可以看到, 原始优化问题(1)的对偶问题是拉格朗日函数的极大极小问题。

原始问题的对偶问题

➤ 换一种表述:

$$\begin{aligned} & \max_{\alpha, \beta} \theta_D(\alpha, \beta) \\ \text{s.t.} \quad & \alpha_i \geq 0, \quad i = 1 \cdots k \end{aligned} \quad (2)$$

➤ 对任意的一个优化问题，无论原始问题**是否是**凸优化问题，其对偶问题**一定是凸**优化问题。

➤ 记原始优化问题(1)的最优值为 p^* :

$$p^* = \min_x \theta_p(x)$$

➤ 记对偶问题(2)的最优值为 d^* :

$$d^* = \max_{\alpha, \beta: \alpha \geq 0} \theta_D(\alpha, \beta)$$

原始问题和对偶问题的关系(弱对偶性)

- **定理1:** 对任意的一个优化问题, 若原始问题和对偶问题都有最优值, 则**弱对偶性成立**, 即 $d^* \leq p^*$ 。

证明: 对于任意的 x, α, β , 有

$$\theta_D(\alpha, \beta) = \min_x L(x, \alpha, \beta) \leq L(x, \alpha, \beta) \leq \max_{\alpha, \beta: \alpha \geq 0} L(x, \alpha, \beta) = \theta_p(x)$$

也即是: $\theta_D(\alpha, \beta) \leq \theta_p(x)$ 。

因此, $\max_{\alpha, \beta: \alpha \geq 0} \theta_D(\alpha, \beta) \leq \min_x \theta_p(x)$, 即 $d^* \leq p^*$ 。

- 将 $p^* - d^*$ 称为对偶差距(**duality gap**)。
- 如果 $d^* = p^*$ 成立, 我们可以说, **强对偶性**成立。在什么情况下, 优化问题满足 $d^* = p^*$?

Slater条件 (强对偶性成立的条件)

- 定义1: 凸优化问题是形如:

$$\begin{array}{ll}\min_{\mathbf{x}} & f_0(\mathbf{x}) \\ \text{s.t.} & f_i(\mathbf{x}) \leq 0, \quad i = 1, 2, \dots, k \\ & \mathbf{a}_i^\top \mathbf{x} = b_i, \quad i = 1, 2, \dots, l\end{array}$$

的问题, 其中, $f_0(\mathbf{x}), f_1(\mathbf{x}), \dots, f_k(\mathbf{x})$ 都是凸函数。

- 定理2: 如果一个凸优化问题满足Slater条件: 存在 \mathbf{x} 满足: $f_i(\mathbf{x}) < 0, i = 1, 2, \dots, k, \mathbf{a}_i^\top \mathbf{x} = b_i, i = 1, 2, \dots, l$, 那么, 强对偶性成立, 即 $d^* = p^*$ 成立。
- 也就是说, 如果强对偶性成立, 则对原始问题的求解可以通过求解相应的对偶问题得到。

Slater条件 (强对偶性成立的条件)

➤ 定理2的语言描述:

对原始优化问题, 如果其目标函数 $f_0(\mathbf{x})$ 和所有的不等式约束函数 $f_i(\mathbf{x})$ 是凸函数, 所有的等式约束函数 $h_i(\mathbf{x})$ 是仿射函数, 并且所有的不等式约束是严格可行的,也即存在 \mathbf{x} 满足对 $i = 1, 2, \dots, k$ 都有 $f_i(\mathbf{x}) < 0$, 对所有的 $i = 1, 2, \dots, l$, 有 $h_i(\mathbf{x}) = 0$, 则存在 $\mathbf{x}^*, \boldsymbol{\alpha}^*, \boldsymbol{\beta}^*$, 使得 \mathbf{x}^* 是原始问题的最优解, $\boldsymbol{\alpha}^*, \boldsymbol{\beta}^*$ 是对偶问题的最优解, 并且

$$d^* = p^* = L(\mathbf{x}^*, \boldsymbol{\alpha}^*, \boldsymbol{\beta}^*)$$

改进的Slater条件

➤ 推论1: 如果一个优化问题具有如下形式:

$$\begin{aligned} \min_{\mathbf{x}} \quad & f_0(\mathbf{x}) \\ \text{s.t.} \quad & f_i(\mathbf{x}) \leq 0, \quad i = 1, 2, \dots, k \\ & \mathbf{a}_i^\top \mathbf{x} = b_i, \quad i = 1, 2, \dots, l \end{aligned}$$

其中, $f_0(\mathbf{x})$ 和前 r 个不等式约束函数 $f_1(\mathbf{x}), \dots, f_r(\mathbf{x})$ 是凸函数, $f_{r+1}(\mathbf{x}) \dots f_k(\mathbf{x})$ 是仿射函数, 且存在 \mathbf{x} 满足:

$$f_i(\mathbf{x}) < 0, \quad i = 1, 2, \dots, r,$$

$$f_i(\mathbf{x}) \leq 0, \quad i = r + 1, \dots, k, \quad \mathbf{a}_i^\top \mathbf{x} = b_i, \quad i = 1, 2, \dots, l$$

那么, 强对偶性成立, 即 $d^* = p^*$ 成立。

改进的Slater条件

- **推论1的语言描述:** 对原始问题和对偶问题, 如果目标函数 $f_0(\mathbf{x})$ 是凸函数, 不等式约束函数 $f_1(\mathbf{x})\cdots f_r(\mathbf{x})$ 是凸函数, $f_{r+1}(\mathbf{x})\cdots f_k(\mathbf{x})$ 是仿射函数, 所有的等式约束函数 $h_i(\mathbf{x})$ 是仿射函数, 并且所有的凸不等式约束是严格可行的, 则存在 $\mathbf{x}^*, \alpha^*, \beta^*$, 使得 \mathbf{x}^* 是原始问题的最优解, α^*, β^* 是对偶问题的最优解, 并且

$$d^* = p^* = L(\mathbf{x}^*, \alpha^*, \beta^*)$$

- 划线的部分称为“改进的Slater条件”。
- 上述定理说明: 如果原始问题的目标函数 $f_0(\mathbf{x})$ 是凸函数, 所有的等式和不等式约束都是仿射函数, 只要该问题的可行域非空, 则强对偶性成立。

KKT条件

- **定理3:** 如果一个**凸优化**问题具有**可微的**目标函数和约束函数，且满足Slater条件，则 \mathbf{x}^* 和 $\boldsymbol{\alpha}^*, \boldsymbol{\beta}^*$ 分别是原始问题和对偶问题的最优解的**充分必要条件**是 \mathbf{x}^* 和 $\boldsymbol{\alpha}^*, \boldsymbol{\beta}^*$ 满足**Karush-Kuhn-Tucker (KKT)**条件:

拉格朗日函数关于原问题的变量 \mathbf{x} 的梯度在最优解处为0

$$\left\{ \begin{array}{l} \nabla_{\mathbf{x}} L(\mathbf{x}^*, \boldsymbol{\alpha}^*, \boldsymbol{\beta}^*) = 0 \end{array} \right.$$

$$\alpha_i^* f_i(\mathbf{x}^*) = 0, i = 1, 2 \dots k$$

对偶互补条件

原始问题和对偶问题的约束条件

$$\left\{ \begin{array}{l} f_i(\mathbf{x}^*) \leq 0, i = 1, 2 \dots k \\ h_i(\mathbf{x}^*) = 0, i = 1, 2 \dots l \\ \alpha_i^* \geq 0, i = 1 \dots k \end{array} \right.$$

原始问题到对偶问题的转化

➤ 线性可分SVM的原始问题:

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

$$s. t. \quad y_i(\mathbf{w}^\top \mathbf{x}_i + b) \geq 1, \quad \forall i = 1 \dots N$$

➤ 原始问题满足slater条件吗?

① $\frac{1}{2} \|\mathbf{w}\|^2$ 是 \mathbf{w} 的凸函数, $1 - y_i(\mathbf{w}^\top \mathbf{x}_i + b)$ 是 \mathbf{w} 的仿射函数

② 原始问题的可行域非空。

综上, 原始问题满足改进的slater条件, 强对偶成立。

原问题是凸优化问题 + 满足slater条件

⇒ KKT条件是最优性的充分必要条件。

线性可分SVM的对偶优化问题

- 线性可分SVM的原始优化问题:

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

$$s. t. \ y_i(\mathbf{w}^\top \mathbf{x}_i + b) \geq 1, \ \forall i = 1 \cdots N$$

- 推导出其对偶优化问题:

- ① 定义拉格朗日函数为:

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

- ② 对偶优化目标函数为:

$$\theta_D(\boldsymbol{\alpha}) = \min_{\mathbf{w}, b} L(\mathbf{w}, b, \boldsymbol{\alpha})$$

如何求 $\min_{\mathbf{w}, b} L(\mathbf{w}, b, \boldsymbol{\alpha})$?

令 $L(\mathbf{w}, b, \boldsymbol{\alpha})$ 关于 \mathbf{w} 和 b 的偏导数分别为零:

SVM的对偶优化问题

$$\nabla_{\mathbf{w}} L(\mathbf{w}, b, \boldsymbol{\alpha}) = 0 \Rightarrow \mathbf{w} = \sum_{i=1}^N \alpha_i y_i \mathbf{x}_i$$

$$\nabla_b L(\mathbf{w}, b, \boldsymbol{\alpha}) = 0 \Rightarrow \sum_{i=1}^N \alpha_i y_i = 0$$

➤ 带入拉格朗日函数得到:

$$\begin{aligned} \theta_D(\boldsymbol{\alpha}) &= \frac{1}{2} \left(\sum_{i=1}^N \alpha_i y_i \mathbf{x}_i \right)^T \left(\sum_{i=1}^N \alpha_i y_i \mathbf{x}_i \right) + \sum_{i=1}^N \alpha_i \\ &\quad - \left(\sum_{i=1}^N \alpha_i y_i \mathbf{x}_i \right)^T \left(\sum_{i=1}^N \alpha_i y_i \mathbf{x}_i \right) - \sum_{i=1}^N \alpha_i y_i b \end{aligned}$$

SVM的对偶优化问题

➤ 整理得到:

$$\theta_D(\boldsymbol{\alpha}) = -\frac{1}{2} \sum_{i=1}^N \sum_{j=1}^N \alpha_i \alpha_j y_i y_j \mathbf{x}_i^\top \mathbf{x}_j + \sum_{i=1}^N \alpha_i$$

➤ 因为 $\max_{\boldsymbol{\alpha}} \theta_D(\boldsymbol{\alpha})$ 等价于 $\min_{\boldsymbol{\alpha}} -\theta_D(\boldsymbol{\alpha})$

➤ SVM的对偶优化问题为:

$$\begin{aligned} \min_{\boldsymbol{\alpha}} \quad & \frac{1}{2} \sum_{i=1}^N \sum_{j=1}^N \alpha_i \alpha_j y_i y_j \mathbf{x}_i^\top \mathbf{x}_j - \sum_{i=1}^N \alpha_i \\ \text{subject to} \quad & \alpha_i \geq 0, \quad i = 1 \cdots N \\ & \sum_{i=1}^N \alpha_i y_i = 0 \end{aligned}$$

由对偶解构造原始解

- 假设已经求出线性可分SVM的对偶问题的最优解为 $\alpha^* = (\alpha_1^* \cdots \alpha_N^*)$, 如何构造出原始问题的最优解 \mathbf{w}^*, b^* ?
- 思路: KKT条件
- 定理3告诉我们, 满足KKT条件的 \mathbf{w}^*, b^* 和 α^* 一定分别是原始问题和对偶问题的最优解。
- 由KKT条件可得到:

$$\nabla_{\mathbf{w}} L(\mathbf{w}^*, b^*, \alpha^*) = 0 \Rightarrow \mathbf{w}^* = \sum_{i=1}^N \alpha_i^* y_i \mathbf{x}_i$$

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

$$\text{对偶互补条件: } \alpha_i^* (1 - y_i (\mathbf{w}^{*\top} \mathbf{x}_i + b^*)) = 0, i = 1 \cdots N$$

由对偶解构造原始解

原问题约束: $y_i(\mathbf{w}^{*\top} \mathbf{x}_i + b^*) \geq 1, i = 1 \cdots N$

对偶约束: $\alpha_i^* \geq 0, i = 1 \cdots N$

- 根据KKT条件可得: 至少有一个拉格朗日乘子 $\alpha_j^* > 0$ 。
因为如果 $\forall i = 1 \cdots N, \alpha_i^* = 0$, 则 $\mathbf{w}^* = \mathbf{0}$, 而 $(\mathbf{0}, b)$ 不是问题的可行解, 产生矛盾。

由 $\alpha_j^* \neq 0$ 可得 $1 - y_j(\mathbf{w}^{*\top} \mathbf{x}_j + b^*) = 0$, 结合 $y_j^2 = 1$ 得:

$$b^* = y_j - \mathbf{w}^{*\top} \mathbf{x}_j = y_j - \sum_{i=1}^N \alpha_i^* y_i \mathbf{x}_i^\top \mathbf{x}_j$$

- ✓ **定理:** 若SVM对偶问题的最优解为 $\boldsymbol{\alpha}^* = (\alpha_1^* \cdots \alpha_N^*)$, 则存在下标 j , 使得 $\alpha_j^* > 0$, 可按如下公式得到原问题的最优解 (\mathbf{w}^*, b^*) :

$$\mathbf{w}^* = \sum_{i=1}^N \alpha_i^* y_i \mathbf{x}_i, \quad b^* = y_j - \sum_{i=1}^N \alpha_i^* y_i \mathbf{x}_i^\top \mathbf{x}_j$$

线性可分SVM的对偶学习算法

- 输入: 线性可分的训练数据集 $\{(\mathbf{x}_1, y_1) \cdots (\mathbf{x}_N, y_N)\}$, 其中 $\mathbf{x}_i \in \mathbb{R}^d$, $y_i \in \{+1, -1\}$.
- 输出: 最大间隔的分离超平面和分类决策函数.

1. 求解如下优化问题:

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

得到最优解 $\boldsymbol{\alpha}^* = (\alpha_1^* \cdots \alpha_N^*)$;

2. 计算原问题最优解: $\mathbf{w}^* = \sum_{i=1}^N \alpha_i^* y_i \mathbf{x}_i$, 选择 $\boldsymbol{\alpha}^*$ 的一个正分量 $\alpha_j^* > 0$, 计算 $b^* = y_j - \sum_{i=1}^N \alpha_i^* y_i \mathbf{x}_i^{\top} \mathbf{x}_j$

3. 得到分割超平面: $\mathbf{w}^{*\top} \mathbf{x} + b^* = 0$

分类决策函数 $f(\mathbf{x}) = \text{sign}(\mathbf{w}^{*\top} \mathbf{x} + b^*)$

支持向量

- 由对偶解和原始解之间的关系可得到： \mathbf{w}^* 和 b^* 仅依赖于 $\alpha_i^* > 0$ 的样本 (\mathbf{x}_i, y_i) ，其他样本对它们没有影响。
- 事实上，训练数据集中对应于 $\alpha_i^* > 0$ 的样本 (\mathbf{x}_i, y_i) 称为支持向量。
- 根据KKT条件中的对偶互补条件：

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

若 $\alpha_i^* > 0$ ，则一定有 $y_i(\mathbf{w}^{*\top}\mathbf{x}_i + b^*) = 1$ 。

所以支持向量在如下超平面上：

$$h_1: \mathbf{w}^{*\top}\mathbf{x} + b^* = 1$$

$$h_2: \mathbf{w}^{*\top}\mathbf{x} + b^* = -1$$

- 这与前面得到的结论是一致的。

小结

- 最简单的支持向量机是线性可分的支持向量机，也称为硬间隔的支持向量机，构建它的条件是训练数据集线性可分，其学习的策略是最大化间隔法，可以形式化为求解如下的凸二次规划问题：

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

$$s. t. \ y_i(\mathbf{w}^\top \mathbf{x}_i + b) \geq 1, \ \forall i = 1 \cdots N$$



求得最优解记为 \mathbf{w}^*, b^* ，得到的分类决策边界为 $\mathbf{w}^{*\top} \mathbf{x} + b^* = 0$ ，分类决策函数为 $f(\mathbf{x}) = \text{sign}(\mathbf{w}^{*\top} \mathbf{x} + b^*)$ 。

- 线性可分的支持向量机的最优解存在且唯一。
- 距离分类决策超平面最近的样本点称为支持向量。

小结

- 线性可分支持向量机的对偶优化问题为:

$$\begin{aligned} \min_{\alpha} \quad & \frac{1}{2} \sum_{i=1}^N \sum_{j=1}^N \alpha_i \alpha_j y_i y_j \mathbf{x}_i^{\top} \mathbf{x}_j - \sum_{i=1}^N \alpha_i \\ \text{subject to} \quad & \alpha_i \geq 0, \quad i = 1 \cdots N \\ & \sum_{i=1}^N \alpha_i y_i = 0 \end{aligned}$$

重难点

通过求解对偶问题的最优解 α^* ，可以求出原始问题的最优解 \mathbf{w}^*, b^* 。

- 支持向量对应的 $\alpha_i^* > 0$ ，且支持向量位于超平面

$$h_1: \mathbf{w}^{*\top} \mathbf{x} + b^* = 1$$

$$h_2: \mathbf{w}^{*\top} \mathbf{x} + b^* = -1$$