

机器学习

于元隆、朱丹红

福州大学计算机与大数据学院
Email: yu.yuanlong@fzu.edu.cn

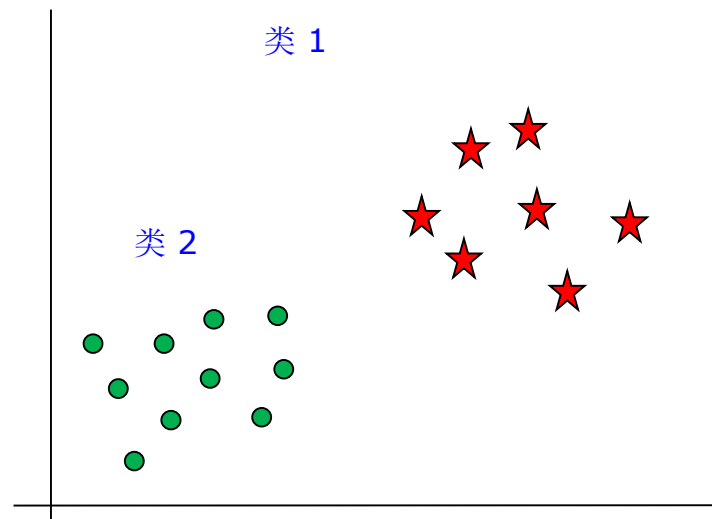


线性判据

——支持向量机

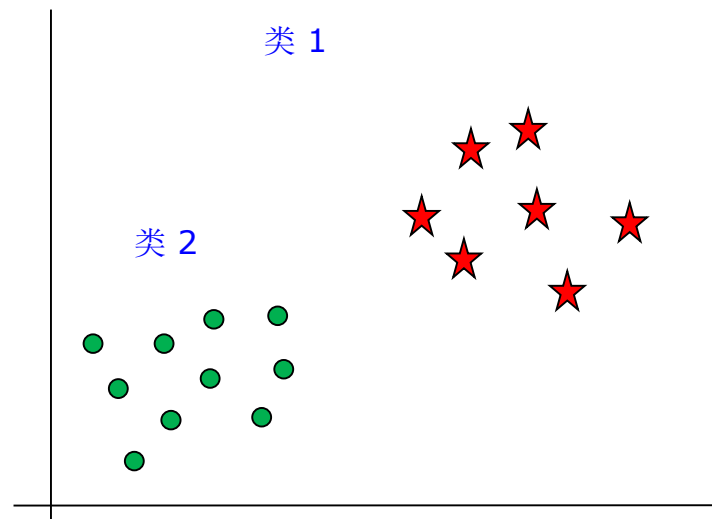
支持向量机：设计动机

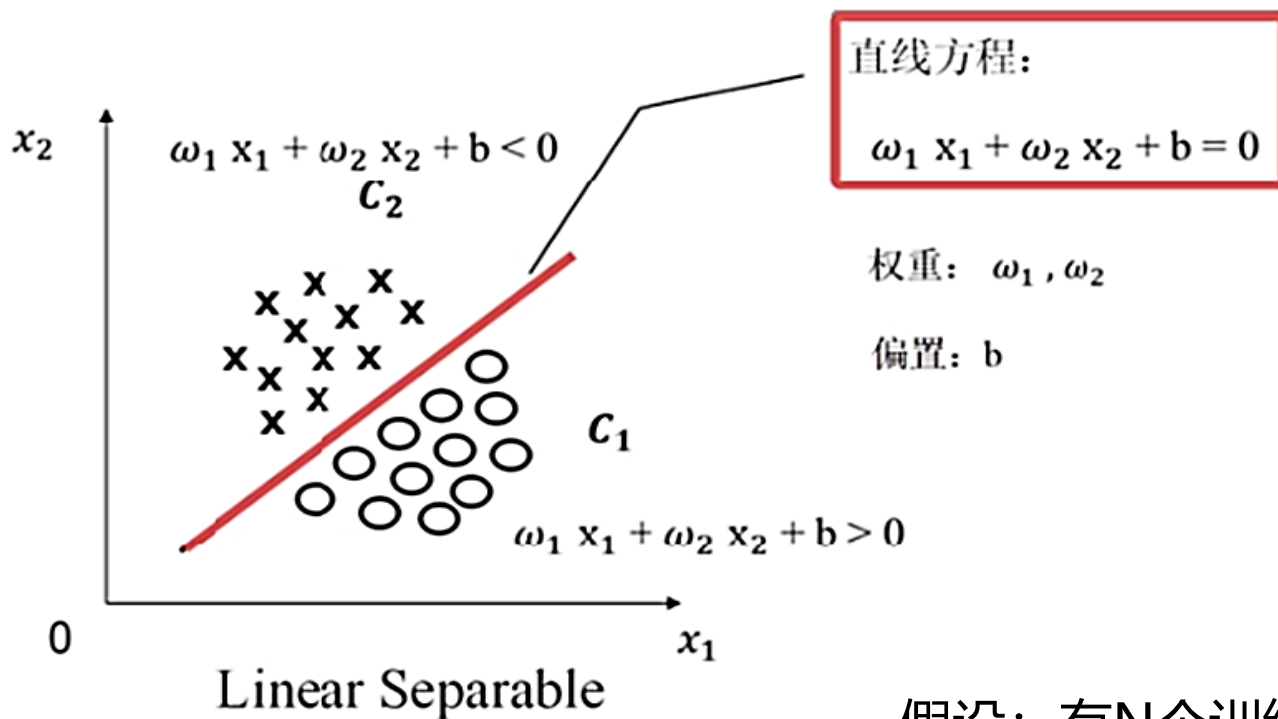
- 如果两个类是线性可分的，则存在着多个线性决策平面用作决策边界。
- 感知机从最小化分类误差角度来设计。
- Fisher线性判据首先降维到一维空间，从最大化类间距离同时最小化类内散度的角度来设计。
- 如何拥有更好的泛化性能？



支持向量机：设计动机

- 如果两个类是线性可分的，则存在着多个线性决策平面用作决策边界。
- 感知机从最小化分类误差角度来设计。
- Fisher线性判据首先降维到一维空间，从最大化类间距离同时最小化类内散度的角度来设计。
- 如何拥有更好的泛化性能？





- 假设: 有N个训练样本和标签:

$$\{(X_1, y_1), (X_2, y_2), \dots, (X_N, y_N)\}$$

其中 $X_i = [x_{i1}, x_{i2}]^T$

$$y_i = \{+1, -1\}$$

\uparrow \uparrow
 x_i 属于C1 x_i 属于C2

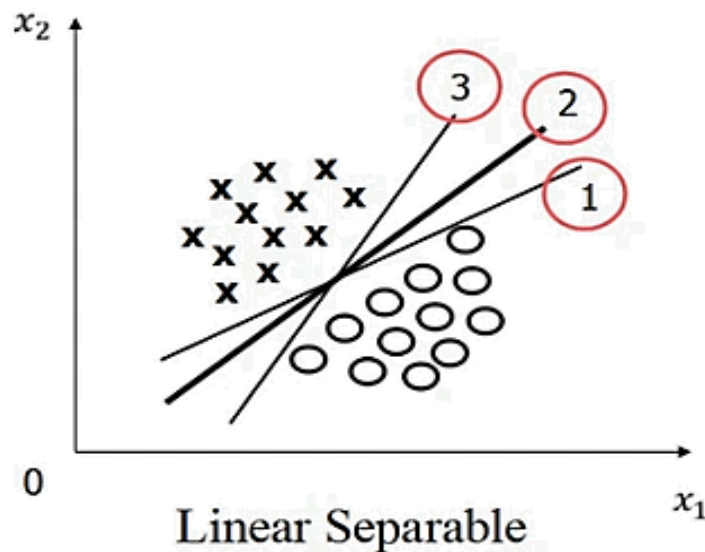
- 线性可分的向量形式定义:

(1) 若 $y_i = +1$, 则 $\omega^T X_i + b > 0$

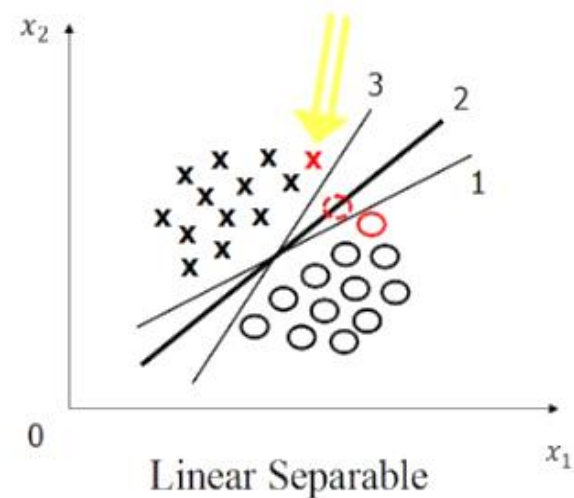
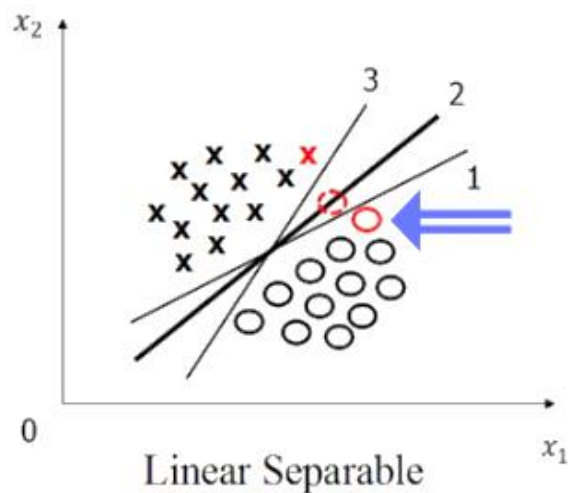
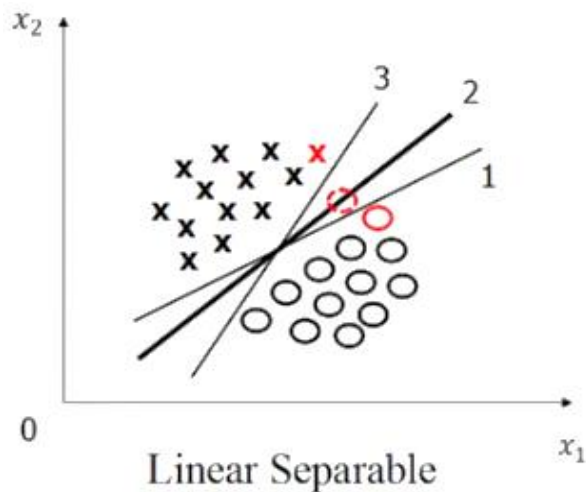
(2) 若 $y_i = -1$, 则 $\omega^T X_i + b < 0$

➔ $y_i(\omega^T X_i + b) > 0$

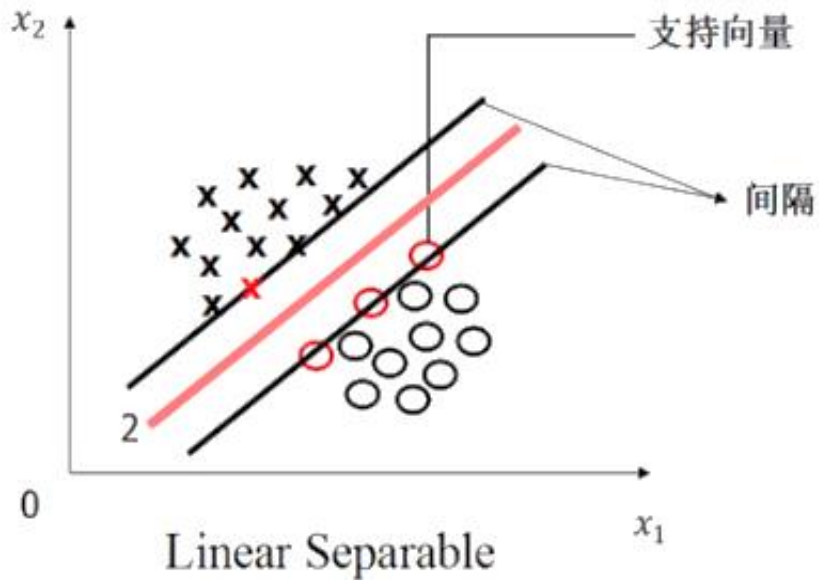
- 如果一个数据集是线性可分的, 存在多少个超平面将各个类别分开?



- 假设训练样本的位置在特征空间上有测量误差：



- 2号线更能抵御训练样本位置的测量误差。



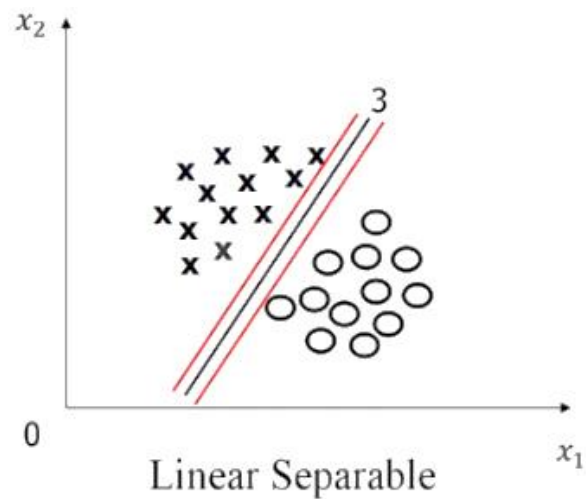
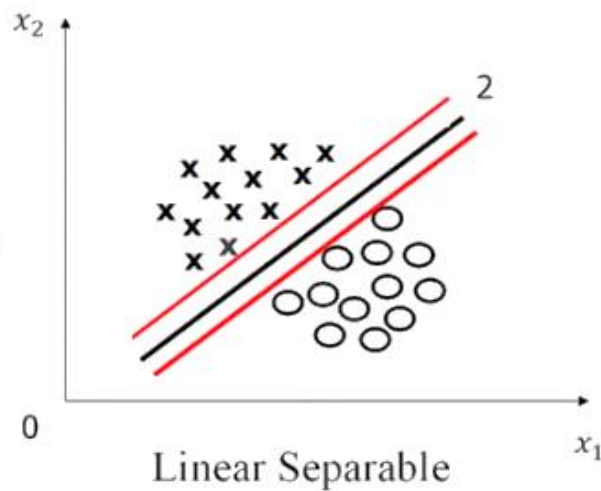
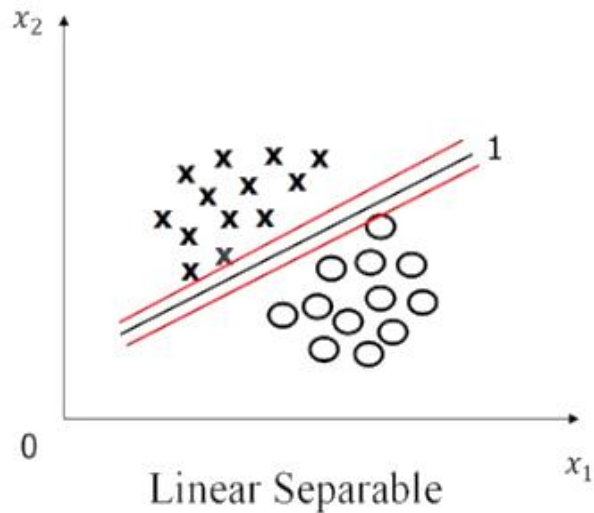
- 怎样画出2号线？基于最优化理论。



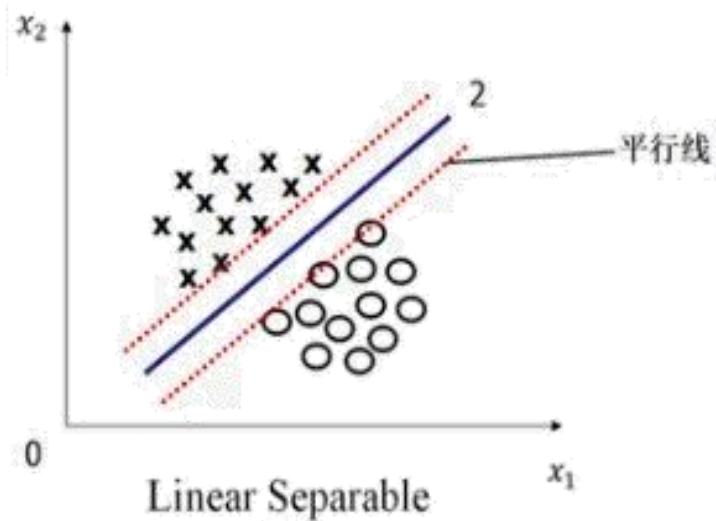
Vladimir Vapnik

- 间隔 (*Margin*) 最大的是2号线。

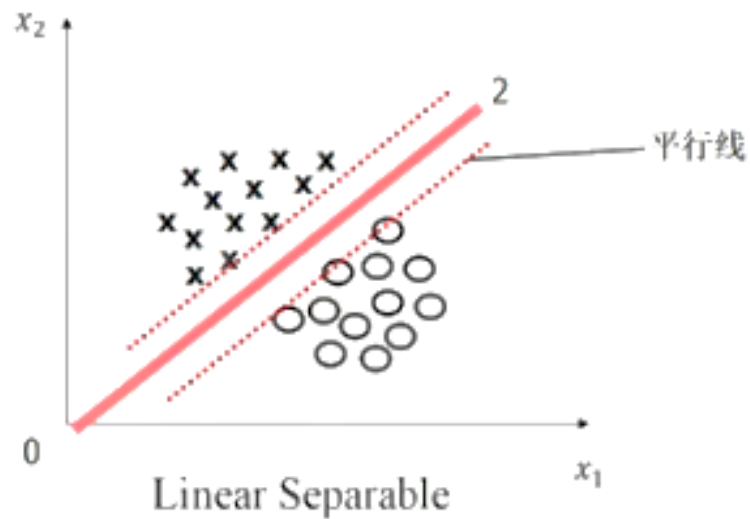
- 哪种方式得到的间隔最大？



- 使用 $Margin$ 最大还不能唯一确定一条直线。

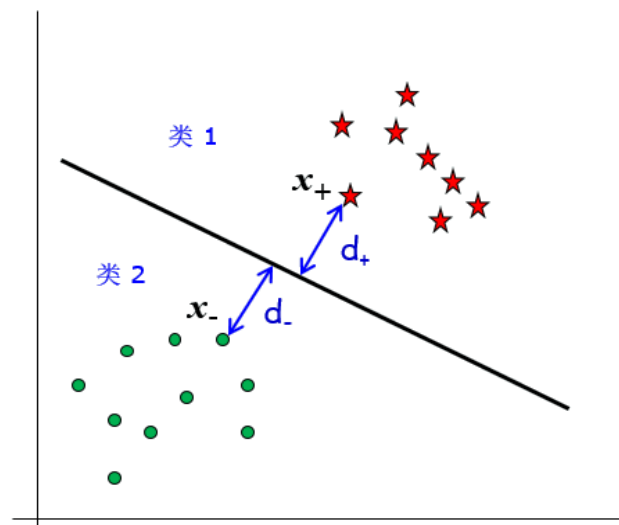


- 这条线应在平行线的中间。



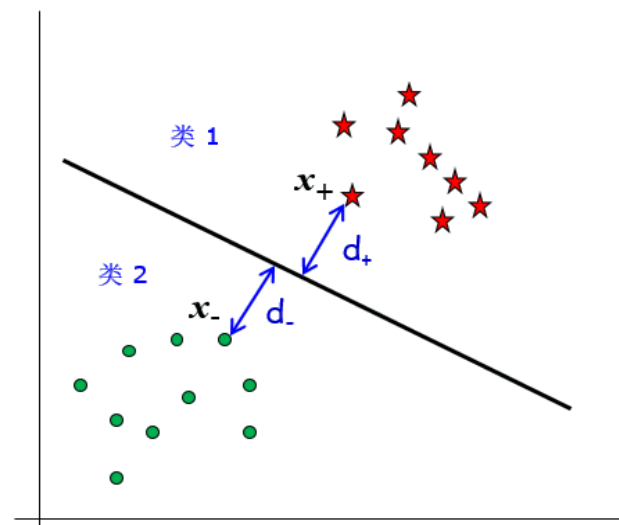
支持向量机：设计思想

- 支持向量机寻找的最优分类直线（决策边界）应满足：
 - (1) 该直线分开了两类；
 - (2) 该直线最大号间隔；
 - (3) 该直线处于间隔的中间，两个类中与决策边界最近的训练样本到决策边界之间的间隔最大。
- 高维空间中，直线 \Rightarrow 超平面
- 寻找最优分类超平面 \Rightarrow 最优化问题



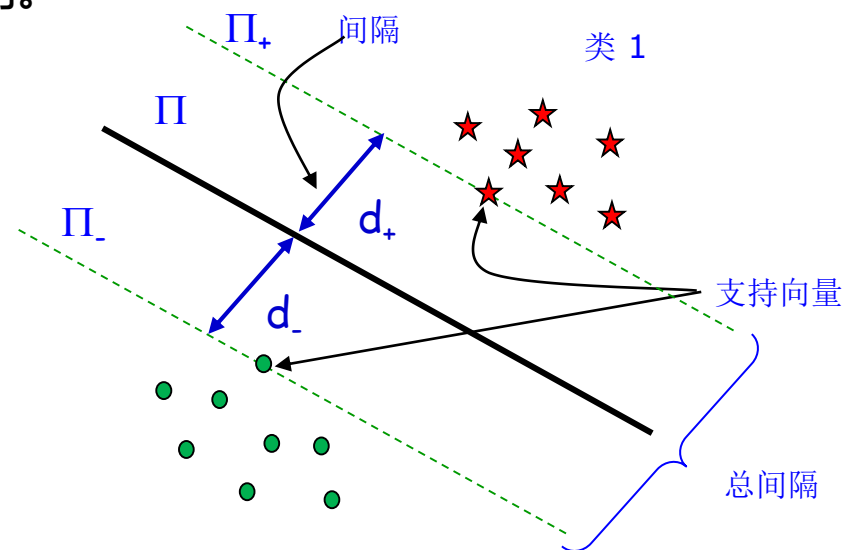
间隔的数学定义

- 间隔的数学定义：
- 在两个类的训练样本中，分别找到与决策边界最近的两个训练样本，记作 x_+ 和 x_- 。
- x_+ 和 x_- 到决策边界的垂直距离叫作间隔，记作 d_+ 和 d_- 。



支持向量的概念

- 决策边界记作 Π ，平行于 Π 且分别通过 x_+ 和 x_- 的两个超平面记作 Π_+ 和 Π_- ，称为**间隔边界**。
- 没有任何训练样本落在这两个超平面中间的间隔区域。
- 位于超平面 Π_+ 和 Π_- 上的样本被称为**支持向量 (Support vector)**。
- 支持向量在确定决策边界 Π 中起到核心作用。



分类器重新表达

- 在支持向量机中，正负类训练样本输出真值分别用+1和-1来表达。
- 给定标记过的训练样本 $\{(\mathbf{x}_n, t_n)\}$ ，线性分类器可以表达为：

$$\mathbf{w}^T \mathbf{x}_n + w_0 > 0 \quad t_n = +1, \quad \text{Class } C_1$$

$$\mathbf{w}^T \mathbf{x}_n + w_0 < 0 \quad t_n = -1, \quad \text{Class } C_2$$

- 加入间隔的概念，引入一个正常数 Δ ，分类器进一步表达为：

$$\mathbf{w}^T \mathbf{x}_n + w_0 \geq +\Delta \quad t_n = +1$$

$$\mathbf{w}^T \mathbf{x}_n + w_0 \leq -\Delta \quad t_n = -1$$

意味着没有训练样本落在 $\pm\Delta$ 间隔范围内。

- 加入间隔的概念，引入一个正常数 Δ ，分类器进一步表达为：

$$\mathbf{w}^T \mathbf{x}_n + w_0 \geq +1 \quad t_n = +1$$

$$\mathbf{w}^T \mathbf{x}_n + w_0 \leq -1 \quad t_n = -1$$

$$\longrightarrow t_n (\mathbf{w}^T \mathbf{x}_n + w_0) - 1 \geq 0, \quad \forall n$$

支持向量机：支持向量

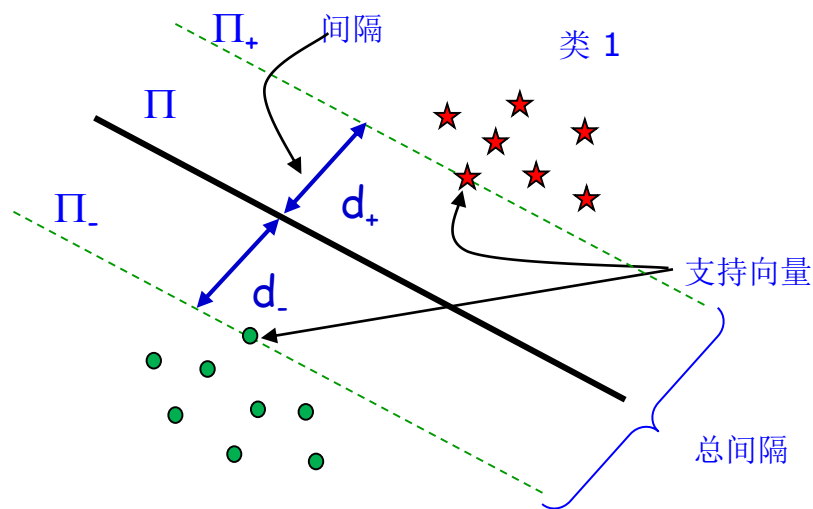
- 分类器新的表达中，什么时候等式成立呢？

$$t_n (\mathbf{w}^T \mathbf{x}_n + w_0) - 1 \geq 0, \quad \forall n$$

- 当 \mathbf{x}_n 是支持向量时，等式成立，例如， \mathbf{x}_+ 和 \mathbf{x}_- 。

$$\mathbf{w}^T \mathbf{x}_+ + w_0 = +1$$

$$\mathbf{w}^T \mathbf{x}_- + w_0 = -1$$



支持向量机：间隔计算

- 根据线性判据的几何含义，点 x_+ 到决策边界 Π 的距离为：

$$r = \frac{f(\mathbf{x}_+)}{\|\mathbf{w}\|}$$

$$\Rightarrow d_+ = |r| = \frac{|f(\mathbf{x}_+)|}{\|\mathbf{w}\|}$$

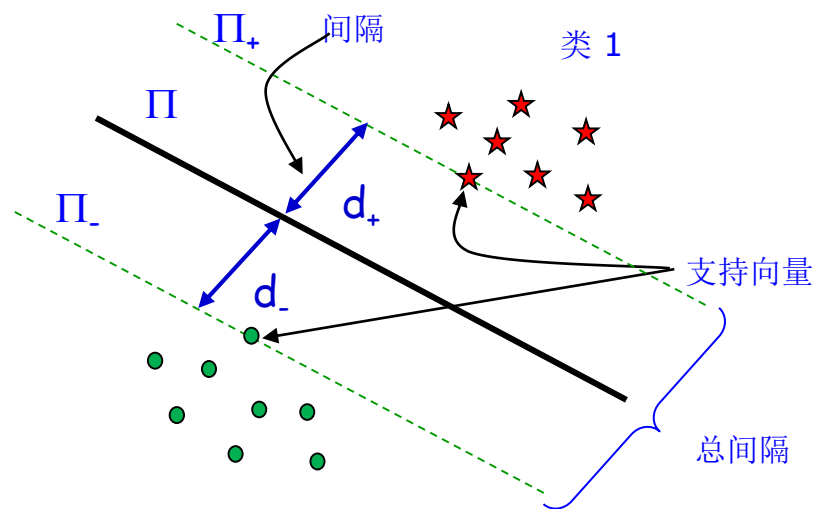
$$\Rightarrow d_- = |r| = \frac{|f(\mathbf{x}_-)|}{\|\mathbf{w}\|}$$

- 所以，总间隔为：

$$M = d_+ + d_- = \frac{|f(\mathbf{x}_+)| + |f(\mathbf{x}_-)|}{\|\mathbf{w}\|}$$

- 由于 x_+ 和 x_- 是位于间隔边界的支持向量，可以得到：

$$\begin{aligned} M = d_+ + d_- &= \frac{|f(\mathbf{x}_+)| + |f(\mathbf{x}_-)|}{\|\mathbf{w}\|} \\ &= \frac{2}{\|\mathbf{w}\|} \end{aligned}$$



支持向量机：目标函数

- 支持向量机 (support vector machine, SVM) 的目标: **最大化总间隔。**
- 最大化间隔, 等价于最小化 $\|\mathbf{w}\|$, 所以目标函数设计为:

$$\min_{\mathbf{w}} \frac{1}{2} \|\mathbf{w}\|_2^2 \quad \longrightarrow \quad \min_{\mathbf{w}} \frac{1}{2} \mathbf{w}^T \mathbf{w}$$

- 同时满足如下约束条件:

$$s.t. \quad t_n (\mathbf{w}^T \mathbf{x}_n + w_0) - 1 \geq 0 \quad \forall n$$

当 \mathbf{x}_n 位于超平面 Π_+ 和 Π_- 上时, 该约束条件等于0, 其余情况都为大于0。

- 该目标函数是条件优化问题(Constrained Optimization)。
- 目标函数自变量为 \mathbf{w} , 是关于 \mathbf{w} 的二次型函数。约束条件是关于 \mathbf{w} 的仿射函数 (线性函数) 。

拉格朗日乘数法

条件优化问题

- 支持向量机的目标函数是一个条件优化问题(Constrained Optimization)。
- 拉格朗日乘数法 (Lagrange Multiplier) 是常用的解决该类问题的方法。

不等式约束优化问题

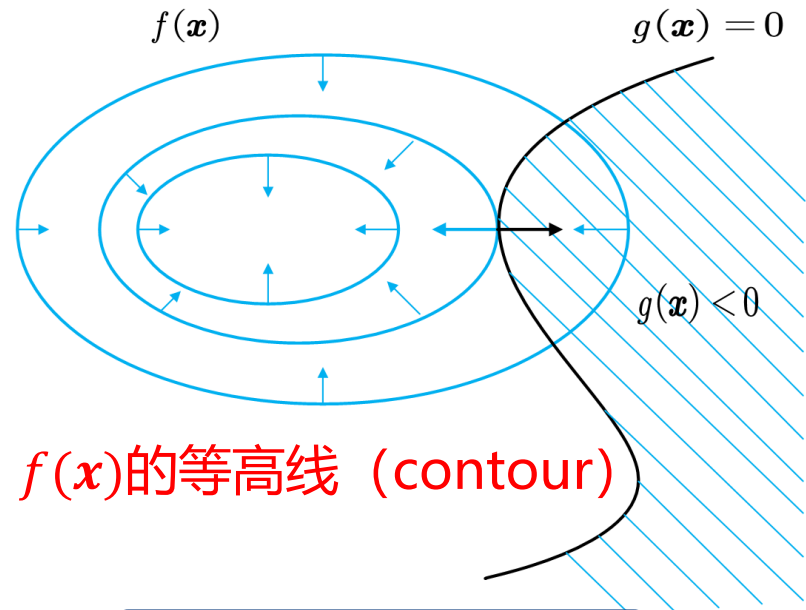
$$\begin{aligned} \min f(\mathbf{x}) \\ \text{s.t. } g(\mathbf{x}) \leq 0 \end{aligned}$$

- 可行域 (Feasible region) : $g(x) \leq 0$ 的区域。

等式约束优化问题

$$\begin{aligned} \min f(\mathbf{x}) \\ \text{s.t. } g(\mathbf{x}) = 0 \end{aligned}$$

- 可行域: $g(x) = 0$ 的区域。



条件优化求解

- 在可行域内寻找 $f(x)$ 的最小值。

拉格朗日乘数法：等式约束

等式约束优化求解思路

- 函数在等高面上任意一点的梯度方向与其等高面（切线方向）正交，且朝向（即正方向）函数值较高方向。
- $f(x)$ 的极值点 x^* 必须位于曲线 $g(x) = 0$ 上。
- 搜寻极值点 x^* ：沿着 $g(x) = 0$ 的切线方向、向着 $f(x)$ 负梯度方向移动。当出现沿着切线方向、无法再向 $f(x)$ 负梯度方向移动时停止。

- 函数在等高面上任意一点的梯度方向与其等高面（切线方向）正交，且朝向（即正方向）函数值较高方向。
- $f(x)$ 的极值点 x^* 必须位于曲线 $g(x) = 0$ 上。
- 搜寻极值点 x^* ：沿着 $g(x) = 0$ 的切线方向、向着 $f(x)$ 负梯度方向移动。当出现沿着切线方向、无法再向 $f(x)$ 负梯度方向移动时停止。

蓝色箭头： $f(x)$ 负梯度方向，即梯度下降最快方向， $f(x)$ 取值沿着该方向不断减小。

等式约束：拉格朗日函数

拉格朗日函数

- 因此，存在一个 $\lambda \neq 0$ ，使得： $f(\mathbf{x})$ 与 $g(\mathbf{x})$ 的梯度记作 $\nabla f(\mathbf{x})$ 和 $\nabla g(\mathbf{x})$ 。

$$\nabla f(\mathbf{x}^*) + \lambda \nabla g(\mathbf{x}^*) = 0,$$

- λ ：拉格朗日乘子（可正可负），**无符号限制**。
- 由此，可以定义一个**拉格朗日函数**：

$$L(\mathbf{x}, \lambda) = f(\mathbf{x}) + \lambda g(\mathbf{x})$$

- 拉格朗日函数**满足驻点(stationary point)条件和约束条件**：

$$\frac{\partial L}{\partial \mathbf{x}} = 0 \Rightarrow \nabla f(\mathbf{x}^*) + \lambda \nabla g(\mathbf{x}^*) = 0$$

$$\frac{\partial L}{\partial \lambda} = 0 \Rightarrow g(\mathbf{x}) = 0$$

等式约束：等价优化问题

等价优化问题

- 因此，因等式约束问题可以转换为等价的不带约束的优化问题：

$$\begin{aligned} \min_{\mathbf{x}} \quad & f(\mathbf{x}) \\ \text{s.t.} \quad & g(\mathbf{x}) = 0 \end{aligned}$$



$$L(\mathbf{x}, \lambda) = f(\mathbf{x}) + \lambda g(\mathbf{x})$$

拉格朗日乘数法：不等式约束

情况1：极值点落在可行域内

- 假设极值点 x^* 落在了可行域内（不含边界），即极值点位于区域 $g(x) < 0$ 范围内。
- 这种情况下，约束条件不起作用。
- 直接通过 $\nabla f(x) = 0$ 获得极值点。
- 此时，在极值点 x^* 上：

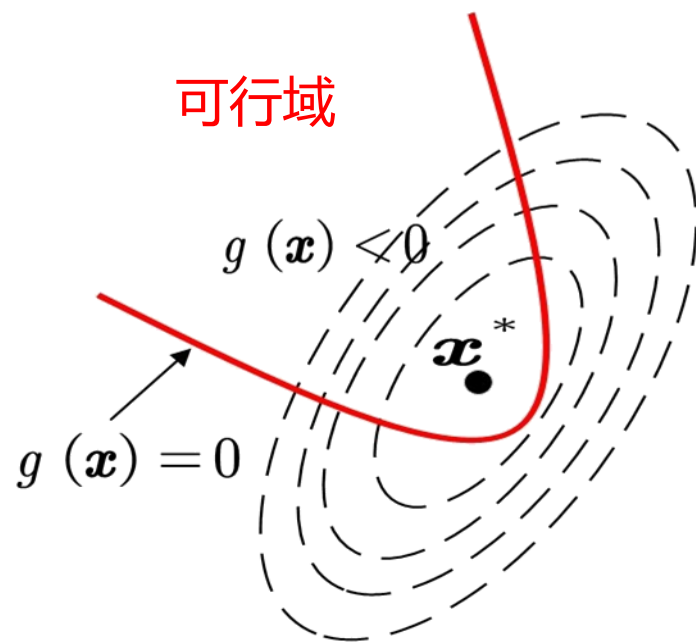
$$\begin{cases} \nabla f(\mathbf{x}^*) = 0 \\ g(\mathbf{x}^*) < 0 \end{cases}$$

- 这种情况相当于在拉格朗日函数中设置 $\lambda = 0$ ：

$$L(\mathbf{x}, \lambda) = f(\mathbf{x}) + \lambda g(\mathbf{x})$$

- 从而满足驻点条件且 $\lambda = 0$ ：

$$\frac{\partial L}{\partial \mathbf{x}} = 0 \Rightarrow \nabla f(\mathbf{x}^*) + \lambda \nabla g(\mathbf{x}^*) = 0$$



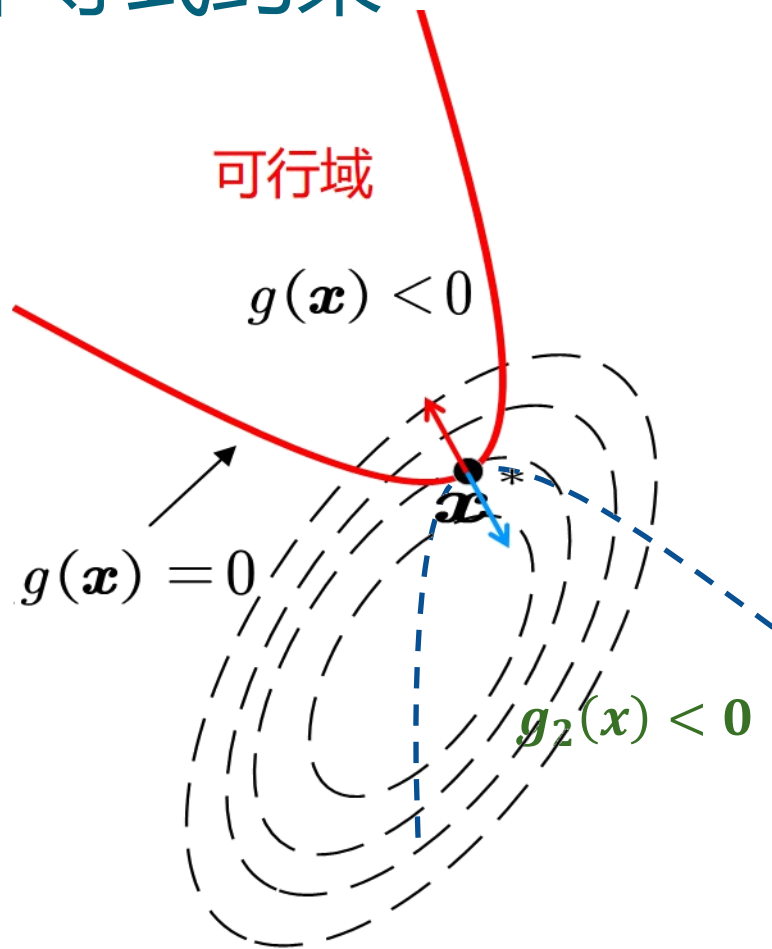
拉格朗日乘数法：不等式约束

情况2：极值点落在可行域边界

- 假设极值点落在了可行域边界，即极值点位于区域 $g(x) = 0$ 区域。
- 搜寻极值点 x^* ：当出现沿着 $g(x) = 0$ 切线方向、无法再向 $f(x)$ 负梯度方向移动时停止。在该点， $f(x)$ 等高线与 $g(x) = 0$ 相切，该点为 $f(x)$ 的极值点 x^* 。
- 对于不等式约束，在极值点 x^* ， $f(x)$ 与 $g(x)$ 的负梯度方向平行且相反。
- 梯度的幅值可能不同。
- 此种情况相当于存在一个 $\lambda > 0$ 满足驻点条件：

$$\nabla f(x^*) + \lambda \nabla g(x^*) = 0$$

- 注意： λ 有符号限制。



如果在极值点 x^* ， $f(x)$ 与 $g(x)$ 的梯度方向相同，则说明在可行域内有更小的 $f(x)$ 值存在。

拉格朗日乘数法：不等式约束

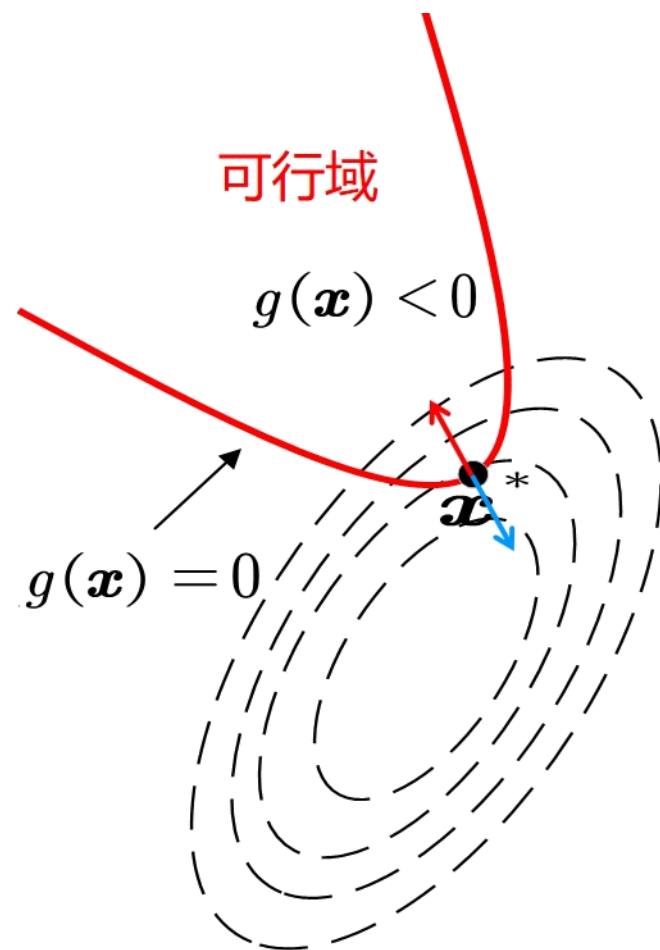
综合两种情况

- 无论是 $g(\mathbf{x}) < 0$ ($\lambda = 0$) 还是 $g(\mathbf{x}) = 0$ ($\lambda > 0$) 的约束情况, 始终存在一个 $\lambda \geq 0$ (对偶可行性), 满足:

$$\lambda g(\mathbf{x}) = 0$$

- 同时, 始终存在一个 $\lambda \geq 0$, 满足驻点条件:

$$\nabla f(\mathbf{x}^*) + \lambda \nabla g(\mathbf{x}^*) = 0$$



不等式约束：KKT条件

KKT条件&等价优化问题

- 在 $g(\mathbf{x}) \leq 0$ 约束条件下最小化 $f(\mathbf{x})$ 的问题，可以转化为如下约束条件（KKT条件）下的拉格朗日函数优化问题：

$$\min_{\mathbf{x}} f(\mathbf{x})$$

$$s.t. \quad g(\mathbf{x}) \leq 0$$



$$L(\mathbf{x}, \lambda) = f(\mathbf{x}) + \lambda g(\mathbf{x})$$

$$s.t. \quad \begin{cases} g(\mathbf{x}) \leq 0 \\ \lambda \geq 0 \\ \lambda g(\mathbf{x}) = 0 \end{cases}$$

Karush-Kuhn-Tucker
(KKT) 条件

拉格朗日乘数法：多个约束

KKT条件&等价优化问题

- 在多个约束条件下最小化 $f(\mathbf{x})$ 的问题，转化为KKT条件下的拉格朗日函数优化问题：

$$\begin{aligned} \min_{\mathbf{x}} \quad & f(\mathbf{x}) \\ \text{s.t.} \quad & h_j(\mathbf{x}) = 0 \\ & g_i(\mathbf{x}) \leq 0 \end{aligned}$$



$$\begin{aligned} L(\mathbf{x}, \mathbf{\Lambda}, \mathbf{\Gamma}) &= f(\mathbf{x}) + \sum_i \lambda_i g_i(\mathbf{x}) + \sum_j \gamma_j h_j(\mathbf{x}) \\ \text{s.t.} \quad & \begin{cases} g_i(\mathbf{x}) \leq 0 \\ \lambda_i \geq 0 \\ \lambda_i g_i(\mathbf{x}) = 0 \\ h_j(\mathbf{x}) = 0 \end{cases} \quad \begin{aligned} \mathbf{\Lambda} &= [\lambda_1, \lambda_2, \dots]^T \\ \mathbf{\Gamma} &= [\gamma_1, \gamma_2, \dots]^T \end{aligned} \end{aligned}$$

拉格朗日对偶问题

拉格朗日乘数法：主问题

主问题 (Primal Problem)

- 根据原问题的约束条件 $h_j(\mathbf{x}) = 0$, 对于 γ_j 的任意取值, 拉格朗日函数第三项都可以消去。
- 根据原问题的约束条件 $g_i(\mathbf{x}) \leq 0$, 对于 $\lambda_i > 0$ 的任意取值, 使得 $\lambda_i g_i(\mathbf{x}) \leq 0$
- 可见, 拉格朗日函数 $L(\mathbf{x}, \Lambda, \Gamma)$ 关于 Λ, Γ 的最大值就是 $f(\mathbf{x})$ 。
- **主问题**: 带约束的原问题等价于如下 (关于 \mathbf{x} 的) 无约束问题。

$$\min_{\mathbf{x}} \max_{\lambda_i \geq 0, \gamma_j} L(\mathbf{x}, \Lambda, \Gamma)$$

注意: λ 的约束依然存在。

$$\lambda_i \geq 0$$

主问题如何求解？

针对不等式约束，主问题难以求解

$$\min_{\mathbf{x}} \max_{\lambda_i \geq 0, \gamma_j} f(\mathbf{x}) + \sum_i \lambda_i g_i(\mathbf{x}) + \sum_j \gamma_j h_j(\mathbf{x})$$

- λ_i 求L偏导并设偏导等于0:

$$\frac{\partial L}{\partial \lambda_i} = g_i(\mathbf{x}) \quad \frac{\partial L}{\partial \lambda_i} = 0 \Rightarrow g_i(\mathbf{x}) = 0$$

- 偏导中不含有 λ_i ，无法解析得到最优的 λ_i^*
- 根据KKT条件，如果要求 $g_i(x) = 0$ ，则 $\lambda_i > 0$ ，但无法确定最优值。

$g_i(x) < 0$	$g_i(x) = 0$
$\lambda_i = 0$	$\lambda_i > 0$

主问题难以求解或者是NP难问题，如何解决？

拉格朗日对偶函数

- 取拉格朗日函数关于 \boldsymbol{x} 在其可行域内的最小值，记作 L_D ：

$$\begin{aligned} L_D(\boldsymbol{\Lambda}, \boldsymbol{\Gamma}) &= \min_{\boldsymbol{x}} L(\boldsymbol{x}, \boldsymbol{\Lambda}, \boldsymbol{\Gamma}) \\ &= \min_{\boldsymbol{x}} \left(f(\boldsymbol{x}) + \sum_i \lambda_i h_i(\boldsymbol{x}) + \sum_j \gamma_j g_j(\boldsymbol{x}) \right) \end{aligned}$$

- 对偶函数 L_D 是关于 $\boldsymbol{\Lambda} > 0$ 和 $\boldsymbol{\Gamma}$ 的函数，与 \boldsymbol{x} 无关。
 - 对偶变量： $\boldsymbol{\Lambda}$ 和 $\boldsymbol{\Gamma}$ 。
 - 主变量： \boldsymbol{x} 。

主问题最优值的下界

对偶函数是主问题的最优值下界

- 针对可行域 R 内的任意 \mathbf{x} , 对任意 $\Lambda > 0$ 和 Γ , 都存在

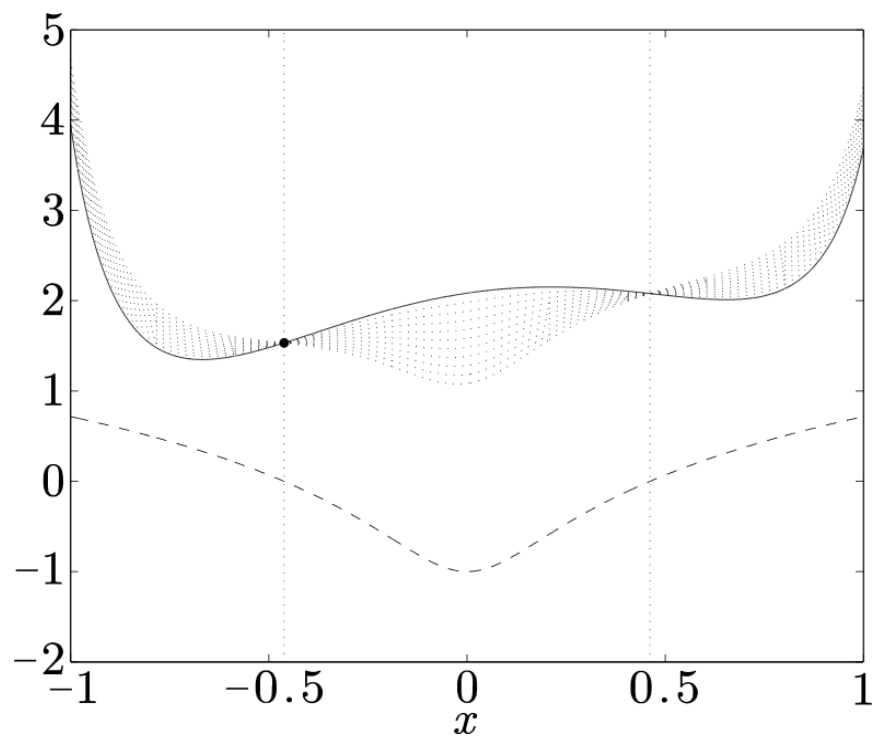
$$\sum_i \lambda_i g_i(\mathbf{x}) + \sum_j \gamma_j h_j(\mathbf{x}) \leq 0$$

- 因此可以得到: $L_D(\Lambda, \Gamma) \leq L(\mathbf{x}, \Lambda, \Gamma) \leq f(\mathbf{x})$
- 设主问题的最小值是 $p^* = f(\mathbf{x}^*)$, 则得到

$$L_D(\Lambda, \Gamma) \leq L(\mathbf{x}^*, \Lambda, \Gamma) \leq f(\mathbf{x}^*)$$

- 可见, 对偶函数 L_D 给出了主问题最优值的下界。
- 该下界只跟对偶变量 Λ 和 Γ 有关, 与 \mathbf{x} 无关。

主问题最优值的下界



- 实线：目标函数 $f(x)$ 。
- 虚线：约束函数 $g(x)$ ， x 的可行域为 $[-0.46, 0.46]$ ，最优点为 $x^* = -0.46$ ， $p^* = 1.54$
- 点线：取不同 λ 值($\lambda = 0.1, 0.2, \dots, 1.0$)的拉格朗日函数 $L(x, \lambda)$ 。由于 $L(x, \lambda) \leq f(x^*)$ ，每条曲线都有一个最小值小于 p^* 。

对偶问题

对偶问题 (Dual Problem)

- 针对 $\Lambda > 0$ 和 Γ ，最大化对偶函数 L_D ，得到主问题的对偶问题：

对偶问题

$$\max_{\Lambda, \Gamma} L_D(\Lambda, \Gamma)$$

$$\max_{\Lambda, \Gamma} \min_{\mathbf{x}} L(\mathbf{x}, \Lambda, \Gamma)$$

$$\text{s.t. } \Lambda \geq \mathbf{0}$$

主问题

$$\min_{\mathbf{x}} \max_{\Lambda, \Gamma} L(\mathbf{x}, \Lambda, \Gamma)$$

$$\text{s.t. } \Lambda \geq \mathbf{0}$$

- 首先求取 L 关于 \mathbf{x} 的最小值（下界），再求取下界关于 Λ, Γ 的最大值。

对偶函数的凹凸性

对偶函数：分析

- 为什么要建立对偶问题？

$$L_D(\Lambda, \Gamma) = \min_x \left(f(\mathbf{x}) + \sum_i \lambda_i h_i(\mathbf{x}) + \sum_j \gamma_j g_j(\mathbf{x}) \right)$$

- 对偶函数是以 Λ 和 Γ 为自变量的、与 x 无关。
- 因此，里面的拉格朗日函数 L 看做关于对偶变量 Λ 和 Γ 的仿射组合。
- 对偶函数 L_D 则是拉格朗日函数 L 的逐点（pointwise）最小值函数。

对偶函数的凹凸性

最小值函数的凹凸性

Given a min function $f(x) = \min_i x_i$, for $0 \leq \theta \leq 1$,

$$f(\theta x + (1 - \theta)y)$$

$$= \min_i (\theta x_i + (1 - \theta)y_i)$$

$$\geq \theta \min_i x_i + (1 - \theta) \min_i y_i$$

$$= \theta f(x) + (1 - \theta)f(y)$$

- 可见，逐点 (pointwise) 最小值函数min是凹函数。

对偶函数的凹凸性

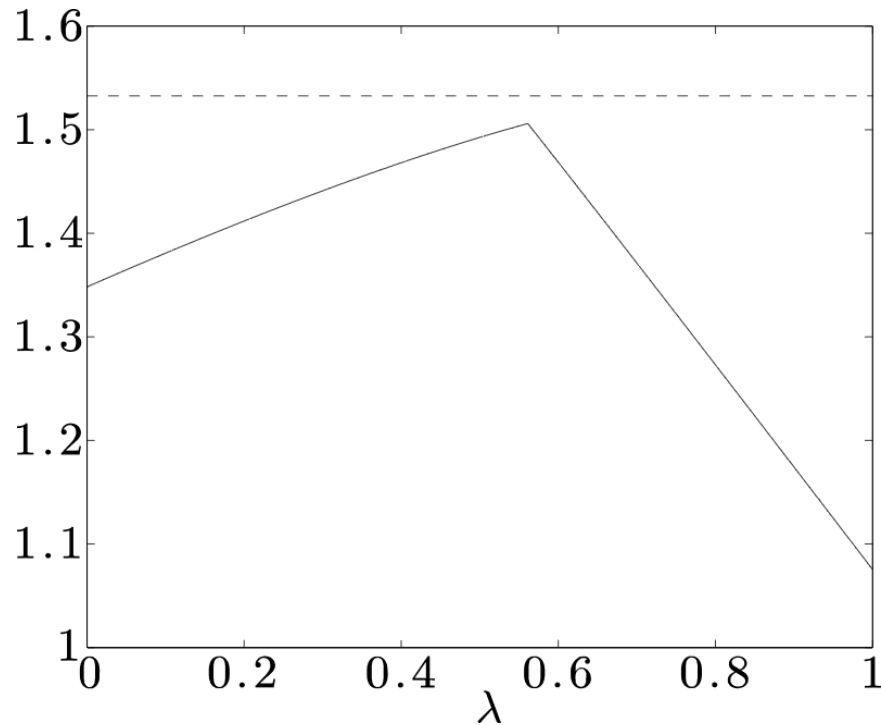
对偶函数是凹函数

Given $g(\mathbf{x}) = f(\mathbf{Ax} + \mathbf{b})$, if the function f is concave, then g is concave.

$$L_D(\mathbf{\Lambda}, \mathbf{\Gamma}) = \min_{\mathbf{x}} \left(f(\mathbf{x}) + \sum_i \lambda_i h_i(\mathbf{x}) + \sum_j \gamma_j g_j(\mathbf{x}) \right)$$

- 由于逐点最大化函数是凹函数，且拉格朗日函数可以看做关于对偶变量的仿射组合，所以对偶函数 L_D 是凹函数。

对偶函数是凹函数



- 实线：对偶函数 $L_D(\lambda)$ 是凹函数。由前图可知， $f(x)$ 和 $g(x)$ 都是非凸的。
- 虚线：最优点 $p^* = 1.54$ 。

对偶问题：凸优化

对偶问题是凸还是非凸？

$$\max_{\Lambda, \Gamma} L_D(\Lambda, \Gamma)$$

$$\text{s.t. } \Lambda \geq \mathbf{0}$$

- 由于目标函数 L_D 是凹函数，约束条件是凸函数，所以对偶问题是凸优化问题。
- 无论主问题的凸性如何，对偶问题始终是一个凸优化问题。
- 凸优化的性质：局部极值点就是全局极值点。
- 所以，对偶问题的极值是唯一的全局极值点。
- 因此，对于难以求解的主问题（例如，非凸问题或者NP难问题），可以通过求解其对偶问题，得到原问题的一个下界估计。

对偶问题与主问题的差异？

弱对偶性 (weak duality)

- 设对偶问题的最优值为 d^* 、主问题的最优值为 p^* 无论主问题的凸性如何，对偶问题始终是一个凸优化问题。
- 对于所有的优化问题都存在： $d^* \leq p^*$

强对偶性 (strong duality)

- 强对偶性： $d^* = p^*$
- 如果强对偶性成立，则对偶问题获得主问题的最优下界。

强对偶性成立的条件 (Slater条件)

$f(x)$ 是凸函数

$g_i(x)$ 是凸函数

$h_j(x)$ 是仿射函数

在可行域至少有一点使得不等式约束严格成立。

- 如果强对偶性成立, 则对偶问题获得主问题的最优下界。

如何使用拉格朗日对偶法求解支持向量机呢?

支持向量机学习算法

支持向量机

线性判据

$$f(\mathbf{x}) = \mathbf{w}^T \mathbf{x} + w_0$$

支持向量机目标函数

给定标记过的训练样本 $\{(\mathbf{x}_n, t_n)\}_{n=1, \dots, N}$:

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

$$s.t. \quad t_n (\mathbf{w}^T \mathbf{x}_n + w_0) - 1 \geq 0, \quad \forall n$$

目标函数如何求解

- 带不等式约束的优化问题，使用拉格朗日对偶法求解。

构建拉格朗日函数

拉格朗日函数

$$L(\mathbf{w}, w_0, \mathbf{\Lambda}) = \frac{1}{2} \|\mathbf{w}\|_2^2 - \sum_{n=1}^N \lambda_n [t_n(\mathbf{w}^T \mathbf{x}_n + w_0) - 1]$$

■ 拉格朗日乘数向量: $\mathbf{\Lambda} = [\lambda_1, \lambda_2, \dots, \lambda_N]^T$

■ KKT条件: $s.t. \begin{cases} -t_n(\mathbf{w}^T \mathbf{x}_n + w_0) + 1 \leq 0 \\ \lambda_n \geq 0 \\ \lambda_n [t_n(\mathbf{w}^T \mathbf{x}_n + w_0) - 1] = 0 \end{cases}$

构建对偶函数

对偶函数

$$\begin{aligned} L_D(\Lambda) &= \min_{\mathbf{w}, w_0} L(\mathbf{w}, w_0, \Lambda) \\ &= \min_{\mathbf{w}, w_0} \frac{1}{2} \|\mathbf{w}\|_2^2 - \sum_{n=1}^N \lambda_n [t_n (\mathbf{w}^T \mathbf{x}_n + w_0) - 1] \end{aligned}$$

- L 对参数 w 和 w_0 求导，并设偏导等于0:

$$\frac{\partial L(\mathbf{w}, w_0, \Lambda)}{\partial \mathbf{w}} = \mathbf{w} - \sum_{n=1}^N t_n \lambda_n \mathbf{x}_n = 0 \quad \rightarrow \quad \mathbf{w} = \sum_{n=1}^N t_n \lambda_n \mathbf{x}_n$$

$$\frac{\partial L(\mathbf{w}, w_0, \Lambda)}{\partial w_0} = \sum_{n=1}^N t_n \lambda_n = 0 \quad \rightarrow \quad \sum_{n=1}^N t_n \lambda_n = 0$$

构建对偶函数

在极值点计算 $\|\mathbf{w}\|_2^2$

$$L_D(\Lambda) = \min_{\mathbf{w}, w_0} \frac{1}{2} \|\mathbf{w}\|_2^2 - \sum_{n=1}^N \lambda_n [t_n (\mathbf{w}^T \mathbf{x}_n + w_0) - 1]$$

$$\mathbf{w} = \sum_{n=1}^N t_n \lambda_n \mathbf{x}_n$$



$$\begin{aligned} \|\mathbf{w}\|_2^2 &= \mathbf{w}^T \mathbf{w} = \left(\sum_{n=1}^N t_n \lambda_n \mathbf{x}_n \right)^T \left(\sum_{n=1}^N t_n \lambda_n \mathbf{x}_n \right) \\ &= \sum_{n=1}^N \sum_{m=1}^M t_n t_m \lambda_n \lambda_m \mathbf{x}_n^T \mathbf{x}_m \end{aligned}$$

构建对偶函数

在极值点计算 $\|\mathbf{w}\|_2^2$

$$L_D(\Lambda) = \min_{\mathbf{w}, w_0} \frac{1}{2} \|\mathbf{w}\|_2^2 - \sum_{n=1}^N \lambda_n [t_n (\mathbf{w}^T \mathbf{x}_n + w_0) - 1]$$

$$\begin{aligned} \mathbf{w} = \sum_{n=1}^N t_n \lambda_n \mathbf{x}_n \\ \sum_{n=1}^N t_n \lambda_n = 0 \end{aligned} \left. \vphantom{\begin{aligned} \mathbf{w} = \sum_{n=1}^N t_n \lambda_n \mathbf{x}_n \\ \sum_{n=1}^N t_n \lambda_n = 0 \end{aligned}} \right\} \begin{aligned} & \sum_{n=1}^N \lambda_n [t_n (\mathbf{w}^T \mathbf{x}_n + w_0) - 1] \\ & = \sum_{n=1}^N \sum_{m=1}^M t_n t_m \lambda_n \lambda_m \mathbf{x}_n^T \mathbf{x}_m - \sum_{n=1}^N \lambda_n \end{aligned}$$

构建对偶函数

在极值点得到 L_D

$$L_D(\Lambda) = \min_{\mathbf{w}, w_0} \frac{1}{2} \|\mathbf{w}\|_2^2 - \sum_{n=1}^N \lambda_n [t_n (\mathbf{w}^T \mathbf{x}_n + w_0) - 1]$$

- 将上述偏导等于0的结果带入拉格朗日函数，消去 \mathbf{w} 和 w_0 ：

$$\begin{aligned} L_D(\Lambda) &= -\frac{1}{2} \sum_{n=1}^N \sum_{m=1}^M t_n t_m \lambda_n \lambda_m \mathbf{x}_n^T \mathbf{x}_m + \sum_{n=1}^N \lambda_n \\ &= -\frac{1}{2} \Lambda^T \mathbf{H} \Lambda + \Lambda^T \mathbf{1} \end{aligned}$$

where \mathbf{H} is Hessian matrix, $h_{nm} = t_n t_m \mathbf{x}_n^T \mathbf{x}_m$

构建对偶函数

对偶函数的约束条件

- 对偶函数是关于 Λ 的函数，所以约束条件只需考虑 Λ 的相关项。

- 对偶可行性：对于不等式约束，构建拉格朗日函数必须满足对偶可行性条件。

$$\Lambda \geq 0$$

- 极值点上关于 λ 的约束项：

$$\Lambda^T \mathbf{t} = 0$$

$$\text{where } \mathbf{t} = [t_1, t_2, \dots, t_N]^T$$

对偶问题

对偶问题

- 对偶函数是关于 Λ 的函数，所以约束条件只需考虑 Λ 的相关项。

$$L_D(\Lambda) = -\frac{1}{2}\Lambda^T H \Lambda + \Lambda^T \mathbf{1}$$

$$\text{s.t.} \quad \begin{cases} \Lambda \geq 0 \\ \Lambda^T \mathbf{t} = 0 \end{cases}$$

- 得到最优的 Λ 。即可得到最优的参数 w 和 w_0 。

对偶问题的求解

求解对偶问题

- 这是标准的关于 λ 的二次规划(quadratic programming)问题。
- 可以调用Matlab提供的quadprog函数来求解。

quadprog函数

$$\min_x \frac{1}{2} \mathbf{x}^T \mathbf{P} \mathbf{x} + \mathbf{q}^T \mathbf{x} + r$$

$$s.t. \quad \mathbf{G} \mathbf{x} \leq \mathbf{h},$$

$$\mathbf{A} \mathbf{x} = \mathbf{b}$$

支持向量

求解支持向量

- 用二次规划求解得到最优的 Λ^* ，包含 N 个最优的拉格朗日乘数。

$$\Lambda^* = [\lambda_1^*, \lambda_2^*, \dots, \lambda_N^*]^T$$

- 根据KKT条件可知：

$\lambda_n^* = 0$	$\lambda_n^* > 0$
$t_n(\mathbf{w}^T \mathbf{x}_n + w_0) > 1$	$t_n(\mathbf{w}^T \mathbf{x}_n + w_0) = 1$
该样本位于超平面 Π_+ 和 Π_- 之外	该样本位于超平面 Π_+ 或 Π_- 上
样本 x_n 非支持向量	样本 x_n 是支持向量

参数最优解： w

w 最优解

- 根据找到的支持向量 x_n 以及对应的拉格朗日乘子 λ_n^* 构建 w^* :

$$w^* = \sum_{n=1}^{N_s} t_n \lambda_n^* x_n$$

- 其中, N_s 表示支持向量的个数。

参数最优解： w

w 最优解

- 根据找到的支持向量 x_n 以及对应的拉格朗日乘子 λ_n^* 构建 w^* :

$$w^* = \sum_{n=1}^{N_s} t_n \lambda_n^* x_n$$

- 其中, N_s 表示支持向量的个数。

参数最优解： w_0

w_0 最优解

- 根据支持向量机定义的约束条件，针对任意一个支持向量 \mathbf{x}_s ：

$$\lambda_s [t_s (\mathbf{w}^T \mathbf{x}_s + w_0) - 1] = 0$$

- 由此得到： $w_0^* = \frac{1}{t_s} - \mathbf{w}^{*T} \mathbf{x}_s$

- w_0 通常由所有支持向量取均值得到：

$$w_0^* = \frac{1}{N_s} \sum_{n=1}^{N_s} \left(\frac{1}{t_n} - \mathbf{w}^{*T} \mathbf{x}_n \right)$$

支持向量机：决策过程

如何用于识别决策过程？

- 给定一个测试模式 \mathbf{x}_{test} ，支持向量机分类器可表达为：

$$y(\mathbf{x}_{test}) = \text{sign}\left(\sum_{k=1}^{N_s} t_k \lambda_k \mathbf{x}_{test}^T \mathbf{x}_k + w_0^*\right)$$

$$\begin{cases} \mathbf{x}_{test} \in C_1 & \text{if } y(\mathbf{x}_{test}) = 1 \\ \mathbf{x}_{test} \in C_2 & \text{if } y(\mathbf{x}_{test}) = -1 \end{cases}$$

- \mathbf{w} 和 w_0 的学习过程实际上是从训练样本中选择一组支持向量，并将这些支持向量存储下来，用作线性分类器。