

# 预备知识



# 凸集

---

对集合  $C$  内的任意两点  $x_1, x_2 \in C$ , 若它们之间连线上的所有点仍属于集合  $C$ , 即

$$\theta x_1 + (1 - \theta)x_2 \in C \quad (\forall 0 \leq \theta \leq 1), \quad (1.1)$$

则我们称集合  $C$  为 “凸” 的, 即  $C$  是一个凸集.

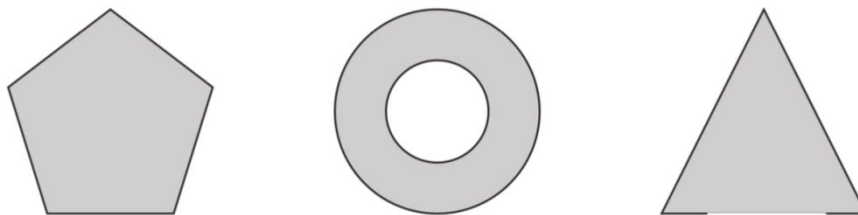


图 1.1 仅有第一个集合是凸的

# 凸函数和凹函数

对定义在凸集上的函数  $f: \mathbb{R}^d \mapsto \mathbb{R}$ , 令  $\text{dom}_f$  表示其定义域, 若  $\forall x, z \in \text{dom}_f$  均满足

$$f(\theta x + (1 - \theta)z) \leq \theta f(x) + (1 - \theta)f(z) \quad (\forall 0 \leq \theta \leq 1), \quad (1.2)$$

则我们称函数  $f(\cdot)$  为凸的, 即  $f(\cdot)$  是一个凸函数.

若将式1.2中的不等号反向, 则函数  $f(\cdot)$  是凹函数.

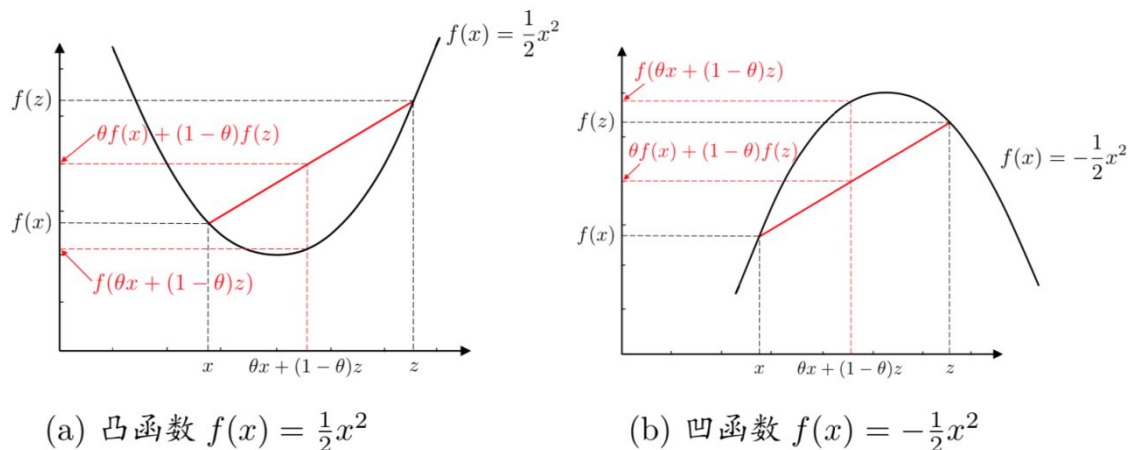


图 1.2 典型的凸函数和凹函数

# 常见的凸函数

表1.1列出了一些常见的凸函数.

表 1.1 常见凸函数

名称	函数形式	定义域	参数
1 维仿射函数	$ax + b$	$x \in \mathbb{R}$	$a, b \in \mathbb{R}$
1 维指数函数	$e^{ax}$	$x \in \mathbb{R}$	$a \in \mathbb{R}$
1 维幂函数	$x^a$	$x \in \mathbb{R}_+$	$a \geq 1$ 或 $a \leq 0$
1 维绝对值幂函数	$ x ^p$	$x \in \mathbb{R}$	$p \geq 1$
1 维负熵函数	$x \log x$	$x \in \mathbb{R}_+$	——
$d$ 维仿射函数	$\mathbf{a}^\top \mathbf{x} + b$	$x \in \mathbb{R}^d$	$\mathbf{a} \in \mathbb{R}^d, b \in \mathbb{R}$
$d$ 维范数	$\ \mathbf{x}\ _p = (\sum_{i=1}^d  x_i ^p)^{1/p}$	$x \in \mathbb{R}^d$	$p \geq 1$

# 梯度与凸函数

函数  $f: \mathbb{R}^d \mapsto \mathbb{R}$  的梯度(gradient)记为  $\nabla f(\mathbf{x}) = (\frac{\partial f(\mathbf{x})}{\partial x_1}; \cdots; \frac{\partial f(\mathbf{x})}{\partial x_d}) \in \mathbb{R}^d$  . 若函数  $f(\cdot)$  可微, 则它是凸函数当且仅当其定义域  $\text{dom}_f$  是凸集且  $\forall \mathbf{x}, \mathbf{z} \in \text{dom}_f$  都有

$$f(\mathbf{z}) \geq f(\mathbf{x}) + \nabla f(\mathbf{x})^T (\mathbf{z} - \mathbf{x}) . \quad (1.4)$$

上式意味着  $f(\cdot)$  在定义域中任意点的一阶泰勒展开是其下界. 例如, 图1.3显示的凸函数  $f(x) = \frac{1}{2}x^2$  及其在  $(1, f(1))$  处的一阶泰勒展开.

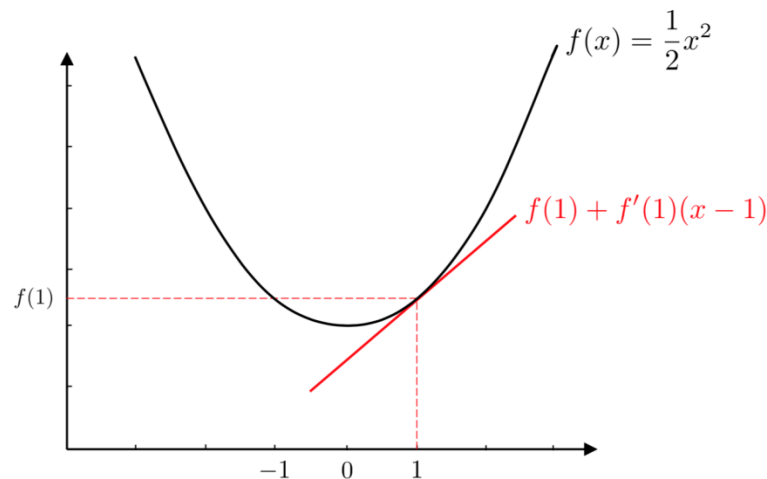


图 1.3 凸函数  $f(x) = \frac{1}{2}x^2$  及其在  $(1, f(1))$  处的一阶泰勒展开

# Hessian矩阵和凸函数

---

除了使用一阶信息，我们还可以基于二阶信息来判断函数的凸性.

函数  $f: \mathbb{R}^d \mapsto \mathbb{R}$  在定义域  $\text{dom}_f$  中  $x$  处的二阶导数矩阵(即Hessian矩阵)记为

$$\nabla^2 f(x) \in \mathbb{R}^{d \times d}, \text{ 其中 } \nabla^2 f(x)_{ij} = \frac{\partial^2 f(x)}{\partial x_i \partial x_j}.$$

若函数  $f(\cdot)$  二阶可微, 则它是凸函数当且仅当  $\text{dom}_f$  是凸集且  $\nabla^2 f(x) \succeq 0$ , 即  $\forall x \in \text{dom}_f$  的Hessian矩阵都是半正定的.

例如: 二次函数  $f(x) = \frac{1}{2}x^T \mathbf{A}x + b^T x + c$  是凸函数当且仅当  $\mathbf{A} \succeq 0$ .

# 数学变换与函数的凸性

---

一些数学变换能够保持函数的凸性, 例如

- $f$  是凸函数, 则  $g(x) = f(Ax + b)$  也是凸函数;
- $f_1, \dots, f_n$  是凸函数,  $w_1, \dots, w_n \geq 0$ , 则  $f(x) = \sum_{i=1}^n w_i f_i(x)$  也是凸函数;
- $f_1, \dots, f_n$  是凸函数, 则  $f(x) = \max \{f_1(x), \dots, f_n(x)\}$  也是凸函数;
- $\forall z \in \mathcal{X}$   $f(x, z)$  是关于  $x$  的凸函数, 则  $g(x) = \sup_{z \in \mathcal{X}} f(x, z)$  也是关于  $x$  的凸函数.

# 共轭函数

函数  $f: \mathbb{R}^d \mapsto \mathbb{R}$  的共轭函数定义为

$$f_*(z) = \sup_{x \in \text{dom } f} (z^T x - f(x)) , \quad (1.5)$$

其定义域

$$\text{dom } f_* = \left\{ z \mid \sup_{x \in \text{dom } f} (z^T x - f(x)) < \infty \right\} . \quad (1.6)$$

直观来看，共轭函数  $f_*(z)$  反映的是线性函数  $z^T x$  与  $f(x)$  之间的最大差值。图1.4 为我们显示了一个实例。

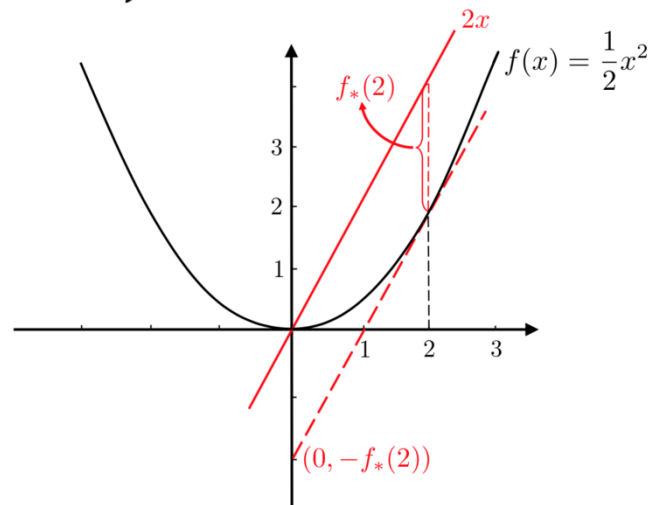


图 1.4 函数  $f(x) = \frac{1}{2}x^2$  的共轭函数在  $z = 2$  处的值的计算方法示意图



# 共轭函数

---

共轭函数有一些很好的性质：

- 无论原函数  $f$  是否是凸函数，共轭函数  $f_*$  一定是凸函数.
- 若函数  $f$  可微，则：

$$f_*(\nabla f(\mathbf{x})) = \nabla f(\mathbf{x})^T \mathbf{x} - f(\mathbf{x}) = - [f(\mathbf{x}) + \nabla f(\mathbf{x})^T (0 - \mathbf{x})] \quad . \quad (1.7)$$

# 重要不等式

---

【Jensen 不等式】 对任意凸函数  $f(\cdot)$  有

$$f(\mathbb{E}[X]) \leq \mathbb{E}[f(X)] .$$

由Jensen不等式可知  $(\mathbb{E}[X])^2 \leq \mathbb{E}[X^2]$ .

【Hölder 不等式】 对  $p, q \in \mathbb{R}_+$  且  $\frac{1}{p} + \frac{1}{q} = 1$ , 有

$$\mathbb{E}[|XY|] \leq (\mathbb{E}[|X|^p])^{\frac{1}{p}} (\mathbb{E}[|Y|^q])^{\frac{1}{q}} .$$

【Cauchy-Schwartz 不等式】

$$\mathbb{E}[|XY|] \leq \sqrt{\mathbb{E}[X^2] \mathbb{E}[Y^2]} .$$

# 重要不等式

---

【Liapounov 不等式】 对  $0 < r \leq s$  有

$$\sqrt[r]{\mathbb{E}[|X|^r]} \leq \sqrt[s]{\mathbb{E}[|X|^s]} .$$

【Minkowski 不等式】 对  $1 \leq p$  有

$$\sqrt[p]{\mathbb{E}[|X + Y|^p]} \leq \sqrt[p]{\mathbb{E}[|X|^p]} + \sqrt[p]{\mathbb{E}[|Y|^p]} .$$

【Bhatia - Davis 不等式】 对  $X \in [a, b]$  有

$$\mathbb{V}[X] \leq (b - \mathbb{E}[X])(\mathbb{E}[X] - a) \leq \frac{(b - a)^2}{4} .$$

# 重要不等式

---

【联合界 (Union Bound) 不等式】

$$P(X \cup Y) \leq P(X) + P(Y) .$$

【Markov 不等式】 对  $X \geq 0$ ,  $\forall \epsilon > 0$ , 有

$$P(X \geq \epsilon) \leq \frac{\mathbb{E}[X]}{\epsilon} .$$

【Chebyshev 不等式】  $\forall \epsilon > 0$  有

$$P(|X - \mathbb{E}[X]| \geq \epsilon) \leq \frac{\mathbb{V}[X]}{\epsilon^2} .$$

# 重要不等式

---

【Cantelli 不等式】  $\forall \epsilon > 0$  有

亦称单边 Chebyshev 不等式.

$$P(X - \mathbb{E}[X] \geq \epsilon) \leq \frac{\mathbb{V}[X]}{\mathbb{V}[X] + \epsilon^2}$$

$$P(X - \mathbb{E}[X] \leq -\epsilon) \leq \frac{\mathbb{V}[X]}{\mathbb{V}[X] + \epsilon^2} .$$

【Chernoff 不等式】 对  $m$  个独立随机变量  $X_i \in [a, b]$  ( $i = 1, \dots, m$ ), 令

$\bar{X} = \sum_{i=1}^m X_i / m$ , 有

$$P(\bar{X} - \mathbb{E}[\bar{X}] \geq \epsilon) \leq e^{-2m\epsilon^2 / (b-a)^2} ,$$

$$P(\bar{X} - \mathbb{E}[\bar{X}] \leq -\epsilon) \leq e^{-2m\epsilon^2 / (b-a)^2} .$$

# 重要不等式

---

【Cantelli 不等式】  $\forall \epsilon > 0$  有

亦称单边 Chebyshev 不等式.

$$P(X - \mathbb{E}[X] \geq \epsilon) \leq \frac{\mathbb{V}[X]}{\mathbb{V}[X] + \epsilon^2}$$

$$P(X - \mathbb{E}[X] \leq -\epsilon) \leq \frac{\mathbb{V}[X]}{\mathbb{V}[X] + \epsilon^2} .$$

【Chernoff 不等式】 对  $m$  个独立随机变量  $X_i \in [a, b]$  ( $i = 1, \dots, m$ ), 令

$\bar{X} = \sum_{i=1}^m X_i / m$ , 有

$$P(\bar{X} - \mathbb{E}[\bar{X}] \geq \epsilon) \leq e^{-2m\epsilon^2 / (b-a)^2} ,$$

$$P(\bar{X} - \mathbb{E}[\bar{X}] \leq -\epsilon) \leq e^{-2m\epsilon^2 / (b-a)^2} .$$

# 重要不等式

---

在机器学习研究中常用到Chernoff不等式的另一种表达形式，若

$$P(\bar{X} \geq \mathbb{E}[\bar{X}] + \epsilon) \leq e^{-2m\epsilon^2} = \delta ,$$

则下式至少以  $1 - \delta$  的概率成立.

$$\bar{X} \leq \mathbb{E}[\bar{X}] + \sqrt{\frac{1}{2m} \ln \frac{1}{\delta}} .$$

**【Hoeffding 不等式】** 对  $m$  个独立随机变量  $X_i \in [0, 1]$  ( $i = 1, \dots, m$ ), 令  $\bar{X} = \frac{1}{m} \sum_{i=1}^m X_i$ , 有

$$P(\bar{X} - \mathbb{E}[\bar{X}] \geq \epsilon) \leq e^{-2m\epsilon^2} .$$

# 重要不等式

---

**【McDiarmid 不等式】** 对  $m$  个独立随机变量  $X_i \in \mathcal{X}$  ( $i = 1, \dots, m$ ), 若  $f: \mathcal{X}^m \rightarrow \mathbb{R}$  是关于  $X_i$  的实值函数且  $\forall x_1, \dots, x_m, x'_i \in \mathcal{X}$  都有

$$|f(x_1, \dots, x_i, \dots, x_m) - f(x_1, \dots, x'_i, \dots, x_m)| \leq c_i ,$$

则  $\forall \epsilon > 0$  有

$$P(f(X_1, \dots, X_m) - \mathbb{E}[f(X_1, \dots, X_m)] \geq \epsilon) \leq e^{-2\epsilon^2 / \sum_{i=1}^m c_i^2} . \quad (1.25)$$

**【Bennett 不等式】** 对  $m$  个独立同分布的随机变量  $X_i$  ( $i = 1, \dots, m$ ), 令  $\bar{X} = \frac{1}{m} \sum_{i=1}^m X_i$ , 若  $X_i - \mathbb{E}[X_i] \leq 1$ , 则有

$$P(\bar{X} \geq \mathbb{E}[\bar{X}] + \epsilon) \leq \exp\left(\frac{-m\epsilon^2}{2\mathbb{V}[X_1] + 2\epsilon/3}\right) .$$



# 重要不等式

---

在机器学习研究中常用到Bennett不等式的另一种表达形式, 若

$$P(\bar{X} \geq \mathbb{E}[\bar{X}] + \epsilon) \leq \exp\left(\frac{-m\epsilon^2}{2\mathbb{V}[X_1] + 2\epsilon/3}\right) = \delta,$$

则下式至少以  $1 - \delta$  的概率成立.

$$\bar{X} \leq \mathbb{E}[\bar{X}] + \epsilon \leq \mathbb{E}[\bar{X}] + \frac{2 \ln 1/\delta}{3m} + \sqrt{\frac{2\mathbb{V}[X_1]}{m} \ln \frac{1}{\delta}}.$$

**【Bernstein 不等式】** 对  $m$  个独立同分布的随机变量  $X_i$  ( $i = 1, \dots, m$ ), 令  $\bar{X} = \frac{1}{m} \sum_{i=1}^m X_i$ , 若存在  $b > 0$  使得  $\forall k \geq 2$  有  $\mathbb{E}[|X_i|^k] \leq k! b^{k-2} \mathbb{V}[X_1]/2$  成立, 则有

$$P(\bar{X} \geq \mathbb{E}[\bar{X}] + \epsilon) \leq \exp\left(\frac{-m\epsilon^2}{2\mathbb{V}[X_1] + 2b\epsilon}\right).$$

# 重要不等式

---

【Azuma 不等式】 对于均值为  $\mu$  的鞅 (martingale)  $\{Z_m, m \geq 1\}$ , 令  $Z_0 = \mu$ , 若  $-c_i \leq Z_i - Z_{i-1} \leq c_i$ , 则  $\forall \epsilon > 0$  有

$$P\left(Z_m - \mu \geq \epsilon\right) \leq e^{-\epsilon^2/2 \sum_{i=1}^m c_i^2},$$

$$P\left(Z_m - \mu \leq -\epsilon\right) \leq e^{-\epsilon^2/2 \sum_{i=1}^m c_i^2}.$$

令  $X_i = Z_i - Z_{i-1}$  可以得到 鞅差序列 (martingale difference sequence)  $X_1, X_2, \dots, X_m$ , 于是有

$$P\left(\sum_{i=1}^m X_i \geq \epsilon\right) \leq e^{-\epsilon^2/2 \sum_{i=1}^m c_i^2},$$

$$P\left(\sum_{i=1}^m X_i \leq -\epsilon\right) \leq e^{-\epsilon^2/2 \sum_{i=1}^m c_i^2}.$$

# 最小二乘问题和线性规划问题

---

一种最简单的优化问题是最小二乘问题

$$\min_{\mathbf{x}} f(\mathbf{x}) = \sum_{i=1}^m (\mathbf{a}_i^T \mathbf{x} - b_i)^2 = \|\mathbf{A}\mathbf{x} - \mathbf{b}\|_2^2, \quad (1.28)$$

其中  $\mathbf{x} = (x_1; \dots; x_d) \in \mathbb{R}^d$  为  $d$  维优化变量,  $\mathbf{A} = (\mathbf{a}_1; \dots; \mathbf{a}_m) \in R^{m \times d}$ ,  $\mathbf{b} = (b_1; \dots; b_m) \in R^m$ .

该问题存在闭式最优解  $\mathbf{x}^* = (\mathbf{A}^T \mathbf{A})^{-1} \mathbf{A}^T \mathbf{b}$ . 计算复杂度为  $O(md^2)$ .

线性规划问题形如

$$\begin{aligned} \min_{\mathbf{x}} \quad & \mathbf{c}^T \mathbf{x} \\ \text{s. t.} \quad & \mathbf{a}_i^T \mathbf{x} \leq b_i \quad (i = 1, \dots, m), \end{aligned} \quad (1.29)$$

其中  $\mathbf{c}, \mathbf{a}_1, \dots, \mathbf{a}_m \in R^d$ ,  $b_1, \dots, b_m \in R$ .

该问题虽无闭式解, 但已有许多成熟的求解算法, 当  $m \geq d$  时计算复杂度仅为  $O(md^2)$ .

# 一般优化问题形式

---

一般的，一个优化问题可以表示为

$$\begin{aligned} \min_{\boldsymbol{x}} \quad & f(\boldsymbol{x}) \\ \text{s. t.} \quad & h_i(\boldsymbol{x}) \leq 0 \quad (i = 1, \dots, m), \end{aligned} \tag{1.30}$$

其中  $f : \mathbb{R}^d \mapsto \mathbb{R}$  称为优化目标函数,  $h_i : \mathbb{R}^d \mapsto \mathbb{R} \ (i = 1, \dots, m)$  称为约束函数.

该问题的最优解可以表达为  $\{\boldsymbol{x}^* \mid f(\boldsymbol{x}^*) \leq f(\boldsymbol{x}) \ (\forall \boldsymbol{x} \in \Omega)\}$ , 其中  $\Omega = \{\boldsymbol{x} \mid h_i(\boldsymbol{x}) \leq 0 \ (i = 1, \dots, m)\}$  称为可行域.

# 凸优化及其最优解

若式(1.30)中的目标函数和约束函数都是凸的，则该优化问题就是凸优化问题.

当目标函数  $f(\cdot)$  可微时,  $x^*$  是凸优化问题的最优解当且仅当  $x^* \in \Omega$  且  $\nabla f(x^*)^T(z - x^*) \geq 0 \ (\forall z \in \Omega)$ . 直观来看,  $-\nabla f(x)$  在  $x^*$  处定义了可行域  $\Omega$  的一个支撑面. 如下图1.5所示.

对于无约束的凸优化问题,  $x^*$  是最优解当且仅当  $x^* \in \Omega$  且  $\nabla f(x^*) = 0$ .

值得注意的是凸优化问题的一些性质:

- 任意一个局部最优解都是全局最优解
- 通常可以在多项式时间内求解

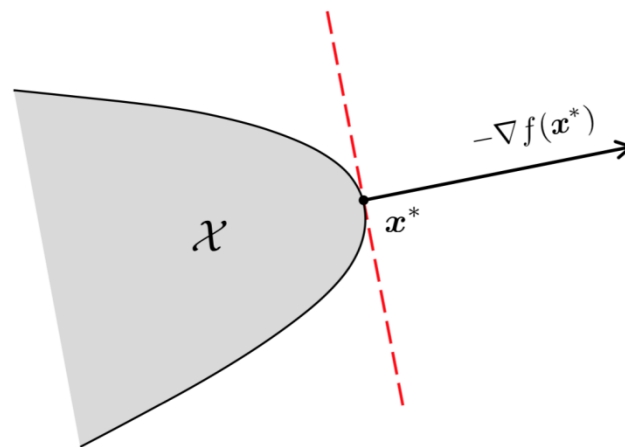


图 1.5  $-\nabla f(x)$  在  $x^*$  处定义了可行域  $\Omega$  的一个支撑面

# 主问题

---

一个优化问题可以从两个角度来考察, 即主问题和对偶问题.

主问题是(1.30)这样的原问题, 或显示列出 $m$  个不等式约束和  $n$  个等式约束写为

$$\begin{aligned} \min_{\boldsymbol{x}} \quad & f(\boldsymbol{x}) \\ \text{s.t.} \quad & h_i(\boldsymbol{x}) \leqslant 0 \quad (i = 1, \dots, m) , \\ & g_j(\boldsymbol{x}) = 0 \quad (j = 1, \dots, n) . \end{aligned} \tag{1.31}$$

我们假设其可行域  $\Omega \subset \mathbb{R}^d$  非空, 并将目标函数最优值记为  $p^*$ .

# 对偶问题

---

对优化问题(1.31), 引入拉格朗日乘子  $\lambda = (\lambda_1, \lambda_2, \dots, \lambda_m)^T$  和  $\mu = (\mu_1, \mu_2, \dots, \mu_n)^T$ , 相应的拉格朗日函数  $L : \mathbb{R}^d \times \mathbb{R}^m \times \mathbb{R}^n \mapsto \mathbb{R}$  为

$$L(\mathbf{x}, \lambda, \mu) = f(\mathbf{x}) + \sum_{i=1}^m \lambda_i h_i(\mathbf{x}) + \sum_{j=1}^n \mu_j g_j(\mathbf{x}), \quad (1.32)$$

其中  $\lambda_i$  和  $\mu_j$  是分别针对不等式约束  $h_i(\mathbf{x}) \leq 0$  和等式约束  $g_j(\mathbf{x}) = 0$  引入的拉格朗日乘子.

相应的拉格朗日对偶函数  $\Gamma : \mathbb{R}^m \times \mathbb{R}^n \mapsto \mathbb{R}$  为

$$\begin{aligned} \Gamma(\lambda, \mu) &= \inf_{\mathbf{x} \in \Omega} L(\mathbf{x}, \lambda, \mu) \\ &= \inf_{\mathbf{x} \in \Omega} \left( f(\mathbf{x}) + \sum_{i=1}^m \lambda_i h_i(\mathbf{x}) + \sum_{j=1}^n \mu_j g_j(\mathbf{x}) \right). \end{aligned} \quad (1.33)$$

# 对偶问题

---

由(1.31)可知, 对于任意  $\lambda \succeq 0$  都有

$$\sum_{i=1}^m \lambda_i h_i(\mathbf{x}) + \sum_{j=1}^n \mu_j g_j(\mathbf{x}) \leq 0, \quad (1.34)$$

对于任意  $\tilde{\mathbf{x}} \in \Omega$  有

$$\Gamma(\boldsymbol{\lambda}, \boldsymbol{\mu}) = \inf_{\mathbf{x} \in \Omega} L(\mathbf{x}, \boldsymbol{\lambda}, \boldsymbol{\mu}) \leq L(\tilde{\mathbf{x}}, \boldsymbol{\lambda}, \boldsymbol{\mu}) \leq f(\tilde{\mathbf{x}}), \quad (1.35)$$

于是, 对任意  $\lambda \succeq 0$  都有

$$\Gamma(\boldsymbol{\lambda}, \boldsymbol{\mu}) \leq p^*, \quad (1.36)$$

即对偶函数(1.33)给出了主问题(1.31)的目标函数最优值  $p^*$  的下界



# 对偶问题

---

基于对偶函数(1.33)可以定义

$$\max_{\lambda, \mu} \Gamma(\lambda, \mu) \quad \text{s.t. } \lambda \succeq 0, \quad (1.37)$$

这就是主问题(1.31)的对偶问题, 其中  $\lambda$  和  $\mu$  称为对偶变量.

由于  $\Gamma(\lambda, \mu)$  是一个凹函数, 而对偶问题(1.37)试图最大化一个凹函数, 因此它是凸优化问题, 且该问题的目标函数最优值  $d^*$  是主问题的目标函数最优值  $p^*$  的下界, 即

$$d^* \leq p^*, \quad (1.38)$$

这称为弱对偶性. 若

$$d^* = p^* \quad (1.39)$$

则称为强对偶性.

# 强对偶性的成立条件

---

对于一般的优化问题，强对偶性通常不成立.

但是，若主问题为凸优化问题，例如：

- (1.31)中的  $f(x)$  和  $h_i(x)$  均为凸函数,  $g_j(x)$  为仿射函数
- 可行域  $\Omega$  中至少有一处使不等式约束严格成立

则强对偶性成立.

此时，将拉格朗日函数(1.32)分别对原变量和对偶变量求导，再令导数等于零即可求解。

# KKT条件

---

KKT条件可以刻画主问题与对偶问题的最优解之间的关系.

令  $\mathbf{x}^*$  为主问题的(1.31)的最优解,  $(\lambda^*, \mu^*)$  为对偶问题(1.37)的最优解. 当强对偶性成立时,

$$\begin{aligned} f(\mathbf{x}^*) &= \Gamma(\lambda^*, \mu^*) \\ &= \inf_{\mathbf{x} \in \Omega} \left\{ f(\mathbf{x}) + \sum_{i=1}^m \lambda_i^* h_i(\mathbf{x}) + \sum_{j=1}^n \mu_j^* g_j(\mathbf{x}) \right\} \\ &\leq f(\mathbf{x}^*) + \sum_{i=1}^m \lambda_i^* h_i(\mathbf{x}^*) + \sum_{j=1}^n \mu_j^* h_j(\mathbf{x}^*) \\ &\leq f(\mathbf{x}^*) \end{aligned} \tag{1.40}$$

# KKT条件

---

显然, (1.40)中的不等式应该取等号. 于是有以下两个条件必定成立:

- 互补松弛条件:

$$\lambda_i^* h_i(\mathbf{x}^*) = 0 \quad (i = 1, \dots, m), \quad (1.41)$$

即  $\lambda_i^* > 0 \Rightarrow h_i(\mathbf{x}^*) = 0$  以及  $h_i(\mathbf{x}^*) < 0 \Rightarrow \lambda_i^* = 0$

- $\mathbf{x}^*$  是下面问题的最优解:

$$\begin{aligned} \mathbf{x}^* &= \arg \min_{\mathbf{x} \in \Omega} L(\mathbf{x}, \boldsymbol{\lambda}^*, \boldsymbol{\mu}^*) \\ &= \arg \min_{\mathbf{x} \in \Omega} \left\{ f(\mathbf{x}) + \sum_{i=1}^m \lambda_i^* h_i(\mathbf{x}) + \sum_{j=1}^n \mu_j^* g_j(\mathbf{x}) \right\}. \end{aligned} \quad (1.42)$$

通常  $\Omega$  为全集或  $\mathbf{x}^*$  位于  $\Omega$  内部, 因此拉格朗日函数  $L(\mathbf{x}, \boldsymbol{\lambda}^*, \boldsymbol{\mu}^*)$  在  $\mathbf{x}^*$  处的梯度为0.

# KKT条件

---

相应的，KKT条件由以下几部分组成：

(1) 主问题约束：

$$\begin{cases} h_i(\mathbf{x}^*) \leq 0 & (i = 1, \dots, m) \\ g_j(\mathbf{x}^*) = 0 & (j = 1, \dots, n) \end{cases}$$

(2) 对偶问题约束： $\lambda^* \succeq 0$

(3) 互补松弛条件： $\lambda_i^* h_i(\mathbf{x}^*) = 0 \quad (i = 1, \dots, m)$

(4) 拉格朗日函数在  $\mathbf{x}^*$  处的梯度为0：

$$\nabla f(\mathbf{x}^*) + \sum_{i=1}^m \lambda_i^* \nabla h_i(\mathbf{x}^*) + \sum_{j=1}^n \mu_j^* \nabla g_j(\mathbf{x}^*) = 0$$

# KKT条件

---

KKT条件具有如下重要性质:

- 强对偶性成立时, 对于任意优化问题, KKT条件是最优解的必要条件, 即主问题最优解和对偶问题最优解一定满足KKT条件;
- 对于凸优化问题, KKT条件是充分条件, 即满足KKT条件的解一定是最优解;
- 对于强对偶性成立的凸优化问题, KKT条件是充分必要条件, 即  $x^*$  是主问题最优解当且仅当存在  $(\lambda^*, \mu^*)$  满足KKT条件.

# 支持向量机

支持向量机是一类经典的机器学习方法.

给定训练样本集  $D = \{(\mathbf{x}_1, y_1), (\mathbf{x}_2, y_2), \dots, (\mathbf{x}_m, y_m)\}$ ,  $y_i \in \{-1, +1\}$ , 支持向量机试图找到恰好位于两类训练样本“正中间”的划分超平面.

如下图1.6所示, 距离超平面最近的这几个训练样本点被称为支持向量, 两个异类支持向量到超平面的距离之和  $\gamma = \frac{2}{\|\mathbf{w}\|}$  称为间隔.

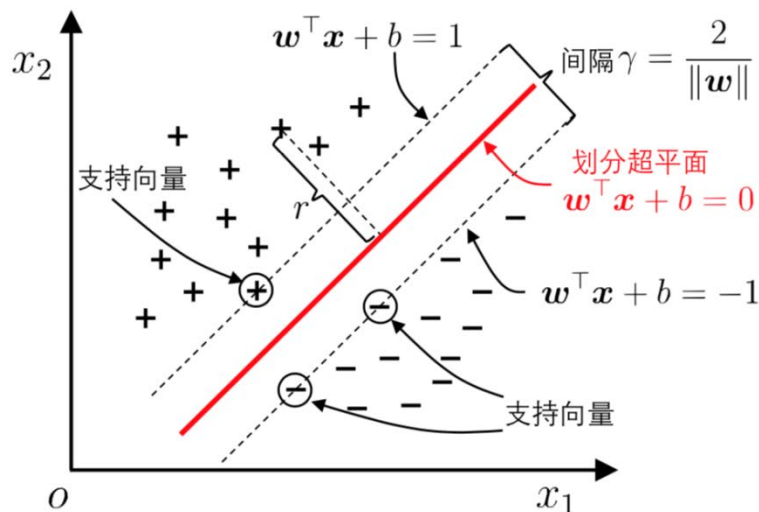


图 1.6 支持向量与间隔

# 支持向量机

---

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

$$\begin{cases} \mathbf{w}^T \mathbf{x}_i + b \geq +1, & y_i = +1; \\ \mathbf{w}^T \mathbf{x}_i + b \leq -1, & y_i = -1, \end{cases} \quad (1.44)$$

则求解最大间隔划分超平面对应于优化问题

$$\begin{aligned} \min_{\mathbf{w}, b} \quad & \frac{1}{2} \|\mathbf{w}\|^2 \\ \text{s.t.} \quad & y_i(\mathbf{w}^T \mathbf{x}_i + b) \geq 1, \quad i = 1, 2, \dots, m. \end{aligned} \quad (1.45)$$



# 支持向量机

---

由(1.32)可知, (1.45)的拉格朗日函数为:

$$L(\mathbf{w}, b, \boldsymbol{\alpha}) = \frac{1}{2} \|\mathbf{w}\|^2 + \sum_{i=1}^m \alpha_i (1 - y_i(\mathbf{w}^T \mathbf{x}_i + b)) , \quad (1.46)$$

其中拉格朗日乘子  $\alpha_i \geq 0$ ,  $\boldsymbol{\alpha} = (\alpha_1; \alpha_2; \dots; \alpha_m)$  .

令  $L(\mathbf{w}, b, \boldsymbol{\alpha})$  对  $\mathbf{w}$  和  $b$  的偏导为零可得

$$\mathbf{w} = \sum_{i=1}^m \alpha_i y_i \mathbf{x}_i , \quad (1.47)$$

$$0 = \sum_{i=1}^m \alpha_i y_i . \quad (1.48)$$

# 支持向量机

---

将(1.47)代入(1.46), 即可将  $L(w, b, \alpha)$  中的  $w$  和  $b$  消去, 再考虑(1.48)的约束, 就得到主问题(1.45)的对偶问题

$$\begin{aligned} \min_{\alpha} \quad & \frac{1}{2} \sum_{i=1}^m \sum_{j=1}^m \alpha_i \alpha_j y_i y_j \mathbf{x}_i^T \mathbf{x}_j - \sum_{i=1}^m \alpha_i \\ \text{s.t.} \quad & \sum_{i=1}^m \alpha_i y_i = 0, \\ & \alpha_i \geq 0, \quad i = 1, 2, \dots, m. \end{aligned} \tag{1.49}$$

上述过程需满足KKT条件

$$\begin{cases} \alpha_i \geq 0; \\ y_i f(\mathbf{x}_i) - 1 \geq 0; \\ \alpha_i (y_i f(\mathbf{x}_i) - 1) = 0. \end{cases} \tag{1.50}$$

# 线性不可分问题

对原始空间中线性不可分的问题，可将样本从原始空间映射到一个高维特征空间。

例如图1.7中所示的“异或问题”在原始二维空间中线性不可分，但若将原始二维空间映射到合适的三维空间则变得线性可分。

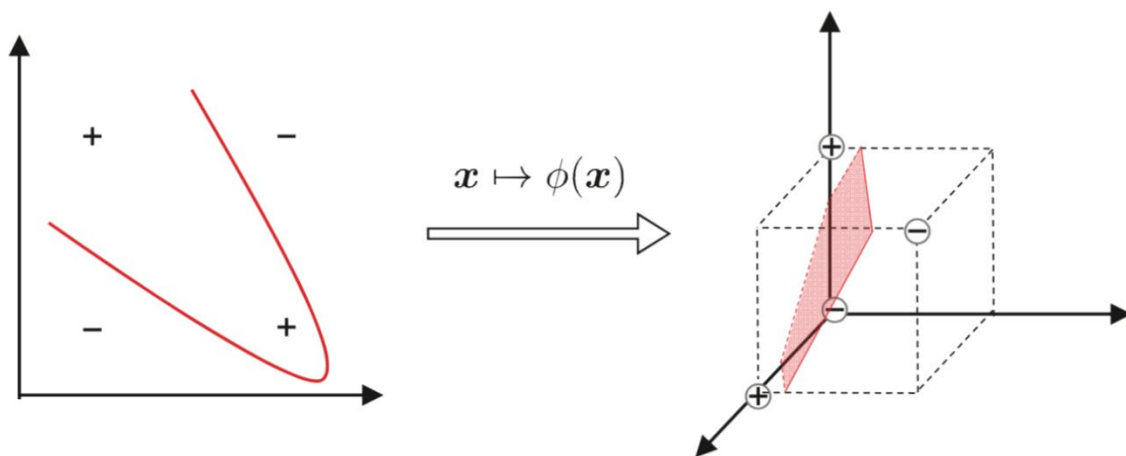


图 1.7 异或问题与非线性映射

# 线性不可分问题

---

引入非线性映射  $\phi(\mathbf{x})$  后, 支持向量机求解的主问题(1.45)变成

$$\begin{aligned} \min_{\mathbf{w}, b} \quad & \frac{1}{2} \|\mathbf{w}\|^2 \\ \text{s.t.} \quad & y_i(\mathbf{w}^T \phi(\mathbf{x}_i) + b) \geq 1, \quad i = 1, 2, \dots, m, \end{aligned} \tag{1.51}$$

相应的, 对偶问题(1.49)变成

$$\begin{aligned} \min_{\boldsymbol{\alpha}} \quad & \frac{1}{2} \sum_{i=1}^m \sum_{j=1}^m \alpha_i \alpha_j y_i y_j \phi(\mathbf{x}_i)^T \phi(\mathbf{x}_j) - \sum_{i=1}^m \alpha_i \\ \text{s.t.} \quad & \sum_{i=1}^m \alpha_i y_i = 0, \\ & \alpha_i \geq 0, \quad i = 1, 2, \dots, m. \end{aligned} \tag{1.52}$$

# 核函数

---

注意到, (1.52)涉及到计算  $\mathbf{x}_i$  与  $\mathbf{x}_j$  映射到特征空间之后的内积  $\phi(\mathbf{x}_i)^\top \phi(\mathbf{x}_j)$ , 当特征空间维数很高时, 直接计算  $\phi(\mathbf{x}_i)^\top \phi(\mathbf{x}_j)$  通常很困难. 为此, 考虑核函数

$$\kappa(\mathbf{x}_i, \mathbf{x}_j) = \langle \phi(\mathbf{x}_i), \phi(\mathbf{x}_j) \rangle = \phi(\mathbf{x}_i)^\top \phi(\mathbf{x}_j) , \quad (1.53)$$

即  $\mathbf{x}_i$  与  $\mathbf{x}_j$  在特征空间的内积等于它们在原始样本空间中通过  $\kappa(\cdot, \cdot)$  计算的结果.

此时, (1.52)可重写为

$$\begin{aligned} \min_{\boldsymbol{\alpha}} \quad & \frac{1}{2} \sum_{i=1}^m \sum_{j=1}^m \alpha_i \alpha_j y_i y_j \kappa(\mathbf{x}_i, \mathbf{x}_j) - \sum_{i=1}^m \alpha_i \\ \text{s.t.} \quad & \sum_{i=1}^m \alpha_i y_i = 0 , \\ & \alpha_i \geq 0 , \quad i = 1, 2, \dots, m . \end{aligned} \quad (1.54)$$

# 核函数

关于核函数有下面的定理：

**定理 1.1 核函数** 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  $\kappa$  是核函数当且仅当对于任意数据  $D = \{\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_m\}$ , 核矩阵(kernel matrix)  $\mathbf{K}$  总是半正定的. 这里  $\mathbf{K}$  是一个  $m$  阶方阵,  $\mathbf{K}_{ij} = \kappa(\mathbf{x}_i, \mathbf{x}_j)$ .

表1.2列出了几种常用的核函数

表 1.2 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, $\beta > 0, \theta < 0$

每个核函数都隐式地定义了一个特征空间, 称为再生希尔伯特空间, 其好坏对支持向量机的性能至关重要.

# 软间隔

---

目前为止，我们一直假设训练样本在样本空间或特征空间中线性可分。然而，在现实任务中往往很难确定合适的满足这一条件的核函数，有时貌似可分的结果甚至可能是由于过拟合而造成的。因此有必要允许支持向量机在少量样本上出错。

为此，引入软间隔的概念，允许某些样本不满足约束：

$$y_i(\mathbf{w}^T \phi(\mathbf{x}_i) + b) \geq 1. \quad (1.55)$$

在最大化间隔的同时，不满足约束的样本应尽可能少。于是得到如下优化目标

$$\min_{\mathbf{w}, b} \frac{1}{2} \|\mathbf{w}\|^2 + \beta \sum_{i=1}^m \ell_{0/1}(y_i(\mathbf{w}^T \phi(\mathbf{x}_i) + b) - 1), \quad (1.56)$$

其中  $\beta > 0$  是一个常数，

$$\ell_{0/1}(x) = \begin{cases} 1, & \text{if } x < 0; \\ 0, & \text{otherwise.} \end{cases} \quad (1.57)$$

# 软间隔

---

由于  $\ell_{0/1}$  非凸、不连续, (1.56)不易求解, 因此支持向量机用hinge损失函数  $\ell_{hinge}(x) = \max(0, 1 - x)$  作为替代损失得到(1.56)的等价形式(1.59)

$$\min_{\mathbf{w}, b} \quad \frac{1}{2} \|\mathbf{w}\|^2 + \beta \sum_{i=1}^m \max(0, 1 - y_i (\mathbf{w}^T \phi(\mathbf{x}_i) + b)) . \quad (1.59)$$

同时引入松弛变量  $\xi_i \geq 0$ , 可将(1.59)重写为

$$\min_{\mathbf{w}, b, \xi_i} \quad \frac{1}{2} \|\mathbf{w}\|^2 + \beta \sum_{i=1}^m \xi_i \quad (1.60)$$

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

$$\xi_i \geq 0, \quad i = 1, 2, \dots, m,$$

即软间隔支持向量机优化的主问题.



# 软间隔

---

相应的，软间隔支持向量机优化的对偶问题为

$$\begin{aligned} \min_{\alpha} \quad & \frac{1}{2} \sum_{i=1}^m \sum_{j=1}^m \alpha_i \alpha_j y_i y_j \kappa(\mathbf{x}_i, \mathbf{x}_j) - \sum_{i=1}^m \alpha_i \\ \text{s.t.} \quad & \sum_{i=1}^m \alpha_i y_i = 0, \\ & 0 \leq \alpha_i \leq \beta, \quad i = 1, 2, \dots, m. \end{aligned} \tag{1.61}$$

对比(1.61)和硬间隔下的对偶问题(1.54)可看出，两者的唯一差别就在于对偶变量的约束不同：前者是  $0 \leq \alpha_i \leq \beta$ ，后者是  $0 \leq \alpha_i$ 。