

第十一章 最优性条件和对偶理论

第 32 讲 最优性条件和对偶理论

黄定江

DaSE @ ECNU

djhuang@dase.ecnu.edu.cn

- 1 32.1 无约束优化最优性条件
- 2 32.2 拉格朗日对偶函数
- 3 32.3 拉格朗日对偶问题
- 4 32.4 Slater 约束准则与强对偶原理
- 5 32.5 KKT 最优性条件
- 6 32.6 对偶问题在数据科学中优化问题的应用

- 1 32.1 无约束优化最优性条件
- 2 32.2 拉格朗日对偶函数
- 3 32.3 拉格朗日对偶问题
- 4 32.4 Slater 约束准则与强对偶原理
- 5 32.5 KKT 最优性条件
- 6 32.6 对偶问题在数据科学中优化问题的应用

关于最优性条件

为什么需要考虑优化问题的最优性条件？

- 方便验证一个点是否为极小值点
- 通过最优性条件获得解析解求解的表达式
- 通过最优性条件获得数值解求解的迭代方法

本小节主要考虑如下问题的最优性理论：

- 无约束可微问题
- 无约束不可微问题

32.1.1 无约束可微问题的最优性理论

先考虑如下无约束可微优化问题：

$$\min_{\mathbf{x} \in \mathbb{R}^n} f(\mathbf{x}), \quad (1)$$

其中 $f(\mathbf{x})$ 是连续可微函数.

1. 一阶最优性条件

根据多元微积分的知识：

定理 1

假设 f 在全空间 \mathbb{R}^n 可微。若 \mathbf{x}^* 是一个局部极小解，那么

$$\nabla f(\mathbf{x}^*) = 0 \quad (2)$$

Proof.

任取 $\boldsymbol{v} \in \mathbb{R}^n$, 考虑 f 在点 $\boldsymbol{x} = \boldsymbol{x}^*$ 处的泰勒展开

$$f(\boldsymbol{x}^* + t\boldsymbol{v}) = f(\boldsymbol{x}^*) + t\boldsymbol{v}^T \nabla f(\boldsymbol{x}^*) + o(t),$$

整理得

$$\frac{f(\boldsymbol{x}^* + t\boldsymbol{v}) - f(\boldsymbol{x}^*)}{t} = \boldsymbol{v}^T \nabla f(\boldsymbol{x}^*) + o(1).$$

根据 \boldsymbol{x}^* 的最优性, 在上式中分别对 t 取点 0 处的左、右极限可知

$$\begin{aligned} \lim_{t \rightarrow 0^+} \frac{f(\boldsymbol{x}^* + t\boldsymbol{v}) - f(\boldsymbol{x}^*)}{t} &= \boldsymbol{v}^T \nabla f(\boldsymbol{x}^*) \geq 0, \\ \lim_{t \rightarrow 0^-} \frac{f(\boldsymbol{x}^* + t\boldsymbol{v}) - f(\boldsymbol{x}^*)}{t} &= \boldsymbol{v}^T \nabla f(\boldsymbol{x}^*) \leq 0, \end{aligned}$$

即对任意的 \boldsymbol{v} 有 $\boldsymbol{v}^T \nabla f(\boldsymbol{x}^*) = 0$, 由 \boldsymbol{v} 的任意性知 $\nabla f(\boldsymbol{x}^*) = 0$.

2. 二阶最优性条件

如果一阶必要条件满足，我们仍然不能确定当前点是否是一个局部极小点。这里考虑使用二阶信息来进一步判断给定点的最优性。

定理 2

假设 f 在点 \mathbf{x}^* 的一个开邻域内是二阶连续可微的，则以下最优性条件成立：

- 二阶必要条件：如果 \mathbf{x}^* 是 f 的一个局部极小点，那么

$$\nabla f(\mathbf{x}^*) = 0, \quad \nabla^2 f(\mathbf{x}^*) \succeq 0.$$

- 二阶充分条件：如果在点 \mathbf{x}^* 处，有

$$\nabla f(\mathbf{x}^*) = 0, \quad \nabla^2 f(\mathbf{x}^*) \succ 0$$

成立，那么 \mathbf{x}^* 是 f 的一个局部极小点。

Proof.

考虑 $f(\mathbf{x})$ 在点 \mathbf{x}^* 处的二阶泰勒展开,

$$f(\mathbf{x} + \mathbf{d}) = f(\mathbf{x}) + \frac{1}{2} \mathbf{d}^T \nabla^2 f(\mathbf{x}) \mathbf{d} + o(\|\mathbf{d}\|^2)$$

这里因为一阶必要条件成立, 所以 $\nabla f(\mathbf{x}^*) = \mathbf{0}$. 反设 $\nabla^2 f(\mathbf{x}^*) \succeq \mathbf{0}$ 不成立, 即 $\nabla^2 f(\mathbf{x}^*)$ 有负的特征值. 取 \mathbf{d} 为其负特征值 λ_- 对应的特征向量, 通过对上式变形得到

$$\frac{f(\mathbf{x}^* + \mathbf{d}) - f(\mathbf{x}^*)}{\|\mathbf{d}\|^2} = \frac{1}{2} \frac{\mathbf{d}^T}{\|\mathbf{d}\|} \nabla^2 f(\mathbf{x}^*) \frac{\mathbf{d}}{\|\mathbf{d}\|} + o(1).$$

这里注意 $\frac{\mathbf{d}}{\|\mathbf{d}\|}$ 是 \mathbf{d} 的单位化, □

Proof.

(续) 因此

$$\frac{f(\mathbf{x}^* + \mathbf{d}) - f(\mathbf{x}^*)}{\|\mathbf{d}\|^2} = \frac{1}{2}\lambda_- + o(1).$$

当 $\|\mathbf{d}\|$ 充分小时, $f(\mathbf{x}^* + \mathbf{d}) < f(\mathbf{x}^*)$, 这和点 \mathbf{x}^* 的最优性矛盾. 因此二阶必要条件成立. 当 $\nabla^2 f(\mathbf{x}^*) \succ \mathbf{0}$ 时, 对任意的 $\mathbf{d} \neq \mathbf{0}$ 有 $\mathbf{d}^T \nabla^2 f(\mathbf{x}^*) \mathbf{d} \geq \lambda_{\min} \|\mathbf{d}\|^2 > 0$, 这里 $\lambda_{\min} > 0$ 是 $\nabla^2 f(\mathbf{x}^*)$ 的最小特征值. 因此我们有

$$\frac{f(\mathbf{x}^* + \mathbf{d}) - f(\mathbf{x}^*)}{\|\mathbf{d}\|^2} \geq \frac{1}{2}\lambda_{\min} + o(1).$$

当 $\|\mathbf{d}\|$ 充分小时有 $f(\mathbf{x}^* + \mathbf{d}) \geq f(\mathbf{x}^*)$, 即二阶充分条件成立. □

我们以线性最小二乘问题为例来说明其最优性条件的具体形式.

例 1

线性最小二乘问题可以表示为

$$\min_{\mathbf{x} \in \mathbb{R}^n} f(\mathbf{x}) \stackrel{\text{def}}{=} \frac{1}{2} \|\mathbf{b} - \mathbf{A}\mathbf{x}\|_2^2,$$

其中 $\mathbf{A} \in \mathbb{R}^{m \times n}$, $\mathbf{b} \in \mathbb{R}^m$ 分别是给定的矩阵和向量. 易知 $f(\mathbf{x})$ 是可微且凸的, 因此, \mathbf{x}^* 为一个全局最优解当且仅当

$$\nabla f(\mathbf{x}^*) = \mathbf{A}^T (\mathbf{A}\mathbf{x}^* - \mathbf{b}) = 0.$$

因此, 线性最小二乘问题本质上等于求解线性方程组, 可以利用数值代数知识对其有效求解.

32.1.2 无约束不可微问题的最优性条件

仍考虑问题:

$$\min_{\mathbf{x} \in \mathbb{R}^n} f(\mathbf{x}),$$

但其中 $f(\mathbf{x})$ 为不可微函数. 很多实际问题的目标函数不是光滑的.

例如 $f(\mathbf{x}) = \|\mathbf{x}\|_1$. 对于此类问题, 由于目标函数可能不存在梯度和海瑟矩阵, 因此上一小节中的一阶和二阶条件不适用. 此时我们必须使用其他最优性条件来判断不可微问题的最优点。

1. 凸优化问题一阶最优性条件

对于目标函数是凸函数的情形, 我们已经引入了次梯度的概念并给出了其计算法则. 一个自然的问题是: 可以利用次梯度代替梯度来构造最优性条件吗? 实际上有如下定理:

定理 3

假设 f 是适当且凸的函数, 则 \mathbf{x}^* 为无约束优化问题的一个全局极小点当且仅当

$$0 \in \partial f(\mathbf{x}^*)$$

Proof.

先证必要性. 因为 \mathbf{x}^* 为全局极小点, 所以

$$f(\mathbf{y}) \geq f(\mathbf{x}^*) = f(\mathbf{x}^*) + 0^T (\mathbf{y} - \mathbf{x}^*), \quad \forall \mathbf{y} \in \mathbb{R}^n.$$

因此, $0 \in \partial f(\mathbf{x}^*)$. 再证充分性. 如果 $0 \in \partial f(\mathbf{x}^*)$, 那么根据次梯度的定义

$$f(\mathbf{y}) \geq f(\mathbf{x}^*) + 0^T (\mathbf{y} - \mathbf{x}^*) = f(\mathbf{x}^*), \quad \forall \mathbf{y} \in \mathbb{R}^n.$$

因而 \mathbf{x}^* 为一个全局极小点. □

这说明条件 $0 \in \partial f(\mathbf{x}^*)$ 是 \mathbf{x}^* 为全局最优解的充要条件. 这个结论比前面的一阶条件要强, 其原因是凸问题有非常好的性质。

2. 复合优化问题的一阶最优性条件

在实际问题中, 目标函数不一定是凸函数, 但它可以写成一个光滑函数与一个非光滑凸函数的和。在压缩感知中, 我们使用 ℓ_1 范数来获得信号的稀疏性; 再比如在机器学习中使用 ℓ_1 正则化; 还有经典的 LASSO 回归问题。这时我们需要考虑复合优化问题

$$\min_{\mathbf{x} \in \mathbb{R}^n} \psi(\mathbf{x}) \stackrel{\text{def}}{=} f(\mathbf{x}) + h(\mathbf{x}), \quad (3)$$

其中 f 为光滑函数 (可能非凸), h 为凸函数 (可能非光滑).

不可微复合优化一阶最优性条件

对于这种情形, 我们给出如下一阶必要条件:

定理 4

令 \mathbf{x}^* 为问题 (3) 的一个局部极小点, 那么

$$-\nabla f(\mathbf{x}^*) \in \partial h(\mathbf{x}^*),$$

其中 $\partial h(\mathbf{x}^*)$ 为凸函数 h 在点 \mathbf{x}^* 处的次梯度集合.

Proof.

因为 \mathbf{x}^* 为一个局部极小点, 所以对于任意单位向量 $\mathbf{d} \in \mathbb{R}^n$ 和足够小的 $t > 0$,

$$f(\mathbf{x}^* + t\mathbf{d}) + h(\mathbf{x}^* + t\mathbf{d}) \geq f(\mathbf{x}^*) + h(\mathbf{x}^*).$$

给定任一方向 $\mathbf{d} \in \mathbb{R}^n$, 其中 $\|\mathbf{d}\| = 1$. 因为对光滑函数和凸函数都可以考虑方向导数, 根据方向导数的定义,

$$\begin{aligned}\psi'(\mathbf{x}^*; \mathbf{d}) &= \lim_{t \rightarrow 0^+} \frac{\psi(\mathbf{x}^* + t\mathbf{d}) - \psi(\mathbf{x}^*)}{t} \\ &= \nabla f(\mathbf{x}^*)^T \mathbf{d} + \partial h(\mathbf{x}^*; \mathbf{d}) \\ &= \nabla f(\mathbf{x}^*)^T \mathbf{d} + \sup_{\theta \in \partial h(\mathbf{x}^*)} \theta^T \mathbf{d},\end{aligned}$$

其中 $\partial h(\mathbf{x}^*; \mathbf{d})$ 表示凸函数 $h(\mathbf{x})$ 在点 \mathbf{x}^* 处的方向导数, 最后一个等式利用了凸函数方向导数和次梯度的关系。



Proof.

(续) 现在用反证法证明我们所需要的结论. 反设 $-\nabla f(\mathbf{x}^*) \notin \partial h(\mathbf{x}^*)$, 根据次梯度的性质可知 $\partial h(\mathbf{x}^*)$ 是有界闭凸集, 又根据严格分离定理, 存在 $\mathbf{d} \in \mathbb{R}^n$ 以及常数 b 使得

$$\theta^T \mathbf{d} < b < -\nabla f(\mathbf{x}^*)^T \mathbf{d}, \quad \forall \theta \in \partial h(\mathbf{x}^*).$$

根据 $\partial h(\mathbf{x}^*)$ 是有界闭集可知对此方向 \mathbf{d} ,

$$\psi'(\mathbf{x}^*; \mathbf{d}) = \nabla f(\mathbf{x}^*)^T \mathbf{d} + \sup_{\theta \in \partial h(\mathbf{x}^*)} \theta^T \mathbf{d} < 0.$$

这说明对充分小的非负实数 t ,

$$\psi(\mathbf{x}^* + t\mathbf{d}) < \psi(\mathbf{x}^*).$$

这与 \mathbf{x}^* 的局部极小性矛盾. 因此 $-\nabla f(\mathbf{x}^*) \in \partial h(\mathbf{x}^*)$. □

例 2

我们以 ℓ_1 范数正则化的优化问题为例, 给出其最优解的最优性条件。前面我们已经介绍其一般形式可以写成

$$\min_{\mathbf{x} \in \mathbb{R}^n} \psi(\mathbf{x}) \stackrel{\text{def}}{=} f(\mathbf{x}) + \mu \|\mathbf{x}\|_1,$$

其中 $f(\mathbf{x}) : \mathbb{R}^n \rightarrow \mathbb{R}$ 为光滑函数, 正则系数 $\mu > 0$ 用来调节解的稀疏度。

尽管 $\|\mathbf{x}\|_1$ 不是可微的, 但我们可以计算其次微分, 在次梯度计算的例子中, 我们已经计算出

$$\partial_i \|\mathbf{x}\|_1 = \begin{cases} \{1\}, & x_i > 0, \\ [-1, 1], & x_i = 0, \\ \{-1\}, & x_i < 0. \end{cases}$$

例2 续

因此, 如果 \mathbf{x}^* 是优化问题的一个局部最优解, 那么其满足

$$-\nabla f(\mathbf{x}^*) \in \mu \partial \|\mathbf{x}^*\|_1,$$

即

$$\nabla_{if}(\mathbf{x}^*) = \begin{cases} -\mu, & x_i^* > 0, \\ a \in [-\mu, \mu], & x_i^* = 0, \\ \mu, & x_i^* < 0. \end{cases}$$

进一步地, 如果 $f(\mathbf{x})$ 是凸的 (比如在 LASSO 回归中 $f(\mathbf{x}) = \frac{1}{2}\|\mathbf{A}\mathbf{x} - \mathbf{b}\|^2$), 那么满足上式的 \mathbf{x}^* 就是全局最优解。

- 1 32.1 无约束优化最优性条件
- 2 32.2 拉格朗日对偶函数
- 3 32.3 拉格朗日对偶问题
- 4 32.4 Slater 约束准则与强对偶原理
- 5 32.5 KKT 最优性条件
- 6 32.6 对偶问题在数据科学中优化问题的应用

32.2.1 约束优化问题

现考虑标准形式的约束优化问题 (不一定是凸的)

$$\begin{aligned} \min \quad & f_0(\mathbf{x}) \\ \text{s.t.} \quad & f_i(\mathbf{x}) \leq 0, \quad i = 1, \dots, m \\ & h_j(\mathbf{x}) = 0, \quad j = 1, \dots, p \end{aligned} \tag{4}$$

其中自变量 $\mathbf{x} \in \mathbb{R}^n$, 假设定义域 $\mathcal{D} = \bigcap_{i=0}^m \mathbf{dom} f_i \cap \bigcap_{j=1}^p \mathbf{dom} h_j$ 非空。

约束优化问题

例 3

最大割问题：

$$\begin{aligned} \min \quad & \mathbf{x}^T \mathbf{W} \mathbf{x} \\ \text{s.t.} \quad & x_i^2 = 1, \quad i = 1, \dots, n \end{aligned}$$

其中， $\mathbf{W} \in S^n$ 。

约束优化问题

例 4

支持向量机:

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

32.2.2 引入：为什么需要对偶

- 约束优化问题的最优性理论分析
- 非凸的优化问题可以转化为凸优化问题便于求解
- 凸的计算复杂度高的问题转化为凸的计算复杂度低的问题
- 帮助证明原问题无解
- 便于敏感性分析

拉格朗日函数

拉格朗日函数: $L: \mathcal{D} \times \mathbb{R}_+^m \times \mathbb{R}^p \rightarrow \mathbb{R}$ 为

$$L(\mathbf{x}, \boldsymbol{\lambda}, \boldsymbol{\nu}) = f_0(\mathbf{x}) + \sum_{i=1}^m \lambda_i f_i(\mathbf{x}) + \sum_{j=1}^p \nu_j h_j(\mathbf{x})$$

其中定义域为 $\text{dom } L = \mathcal{D} \times \mathbb{R}_+^m \times \mathbb{R}^p$

- 拉格朗日函数: 添加约束条件的加权和, 得到增广的目标函数
- λ_i 称为第 i 个不等式约束 $f_i(\mathbf{x}) \leq 0$ 对应的拉格朗日乘子
- ν_j 称为第 j 个等式约束 $h_j(\mathbf{x}) = 0$ 对应的拉格朗日乘子
- 向量 $\boldsymbol{\lambda}$ 和 $\boldsymbol{\nu}$ 称为对偶变量或者问题 (4) 的拉格朗日乘子向量

拉格朗日对偶函数

定义 1

拉格朗日对偶函数 (或简称对偶函数) $g: \mathbb{R}_+^m \times \mathbb{R}^p \rightarrow \mathbb{R}$ 是拉格朗日函数关于 \mathbf{x} 取得的下确界: 即对 $\boldsymbol{\lambda} \in \mathbb{R}_+^m, \boldsymbol{\nu} \in \mathbb{R}^p$, 有

$$g(\boldsymbol{\lambda}, \boldsymbol{\nu}) = \inf_{\mathbf{x} \in \mathcal{D}} L(\mathbf{x}, \boldsymbol{\lambda}, \boldsymbol{\nu}) = \inf_{\mathbf{x} \in \mathcal{D}} \left(f_0(\mathbf{x}) + \sum_{i=1}^m \lambda_i f_i(\mathbf{x}) + \sum_{j=1}^p \nu_j h_j(\mathbf{x}) \right)$$

- 如果函数关于 \mathbf{x} 无下界, 则对偶函数取值为 $-\infty$
- 对偶函数是一族关于 $(\boldsymbol{\lambda}, \boldsymbol{\nu})$ 的仿射函数的逐点下确界, 所以即使原问题 (4) 非凸, 对偶函数也是凹函数

32.2.3 最优值的下界

引理 1

对偶函数构成了原问题 (4) 最优值 p^* 的下界：即对任意 $\lambda \geq 0$ 和 ν ，下式成立

$$g(\lambda, \nu) \leq p^* \quad (5)$$

Proof.

设 $\tilde{\mathbf{x}}$ 是原问题 (4) 的一个可行点，即 $f_i(\tilde{\mathbf{x}}) \leq 0$ 且 $h_j(\tilde{\mathbf{x}}) = 0$ 。根据假设， $\lambda \geq 0$ ，有

$$\sum_{i=1}^m \lambda_i f_i(\tilde{\mathbf{x}}) + \sum_{j=1}^p \nu_j h_j(\tilde{\mathbf{x}}) \leq 0$$

左边第一项非正，第二项为零。 □

最优值的下界

证明 (续) .

根据上述不等式, 有

$$L(\tilde{\mathbf{x}}, \boldsymbol{\lambda}, \boldsymbol{\nu}) = f_0(\tilde{\mathbf{x}}) + \sum_{i=1}^m \lambda_i f_i(\tilde{\mathbf{x}}) + \sum_{j=1}^p \nu_j h_j(\tilde{\mathbf{x}}) \leq f_0(\tilde{\mathbf{x}}).$$

因此

$$g(\boldsymbol{\lambda}, \boldsymbol{\nu}) = \inf_{\mathbf{x}} L(\mathbf{x}, \boldsymbol{\lambda}, \boldsymbol{\nu}) \leq L(\tilde{\mathbf{x}}, \boldsymbol{\lambda}, \boldsymbol{\nu}) \leq f_0(\tilde{\mathbf{x}})$$

由于每一个可行点 $\tilde{\mathbf{x}}$ 都满足 $g(\boldsymbol{\lambda}, \boldsymbol{\nu}) \leq f_0(\tilde{\mathbf{x}})$, 因此 $g(\boldsymbol{\lambda}, \boldsymbol{\nu}) \leq p^*$ 成立 □

32.2.4 求拉格朗日对偶函数的解析表达式

线性方程组的最小二乘解

例 5

考虑问题

$$\begin{aligned} \min \quad & \mathbf{x}^T \mathbf{x} \\ \text{s.t.} \quad & \mathbf{A} \mathbf{x} = \mathbf{b} \end{aligned} \tag{6}$$

其中 $\mathbf{A} \in \mathbb{R}^{p \times n}$, 求其对偶函数。

解

此问题无不等式约束，有 p 个 (线性) 等式约束，其拉格朗日函数为

$$L(\mathbf{x}, \boldsymbol{\nu}) = \mathbf{x}^T \mathbf{x} + \boldsymbol{\nu}^T (\mathbf{A}\mathbf{x} - \mathbf{b})$$

定义域为 $\mathbb{R}^n \times \mathbb{R}^p$ 。对偶函数为

$$g(\boldsymbol{\nu}) = \inf_{\mathbf{x}} L(\mathbf{x}, \boldsymbol{\nu})$$

解

(续) 因为 $L(\boldsymbol{x}, \boldsymbol{\nu})$ 是 \boldsymbol{x} 的二次凸函数, 可通过求解最优性条件得到函数的最小值

- L 关于 \boldsymbol{x} 取极小值, 梯度为零:

$$\nabla_{\boldsymbol{x}} L(\boldsymbol{x}, \boldsymbol{\nu}) = 2\boldsymbol{x} + \boldsymbol{A}^T \boldsymbol{\nu} = \mathbf{0} \implies \boldsymbol{x} = -(1/2)\boldsymbol{A}^T \boldsymbol{\nu}$$

- 代入到 L 得到对偶函数 g :

$$g(\boldsymbol{\nu}) = L((-1/2)\boldsymbol{A}^T \boldsymbol{\nu}, \boldsymbol{\nu}) = -(1/4)\boldsymbol{\nu}^T \boldsymbol{A} \boldsymbol{A}^T \boldsymbol{\nu} - \boldsymbol{b}^T \boldsymbol{\nu}$$

是关于 $\boldsymbol{\nu}$ 的二次凹函数, 定义域为 \mathbb{R}^p

最优值下界: 根据引理 1, 对 $\forall \boldsymbol{\nu} \in \mathbb{R}^p$

$$p^* \geq -(1/4)\boldsymbol{\nu}^T \boldsymbol{A} \boldsymbol{A}^T \boldsymbol{\nu} - \boldsymbol{b}^T \boldsymbol{\nu}$$

标准形式的线性规划

例 6

考虑问题

$$\begin{aligned} \min \quad & \mathbf{c}^T \mathbf{x} \\ \text{s.t.} \quad & \mathbf{Ax} = \mathbf{b} \\ & \mathbf{x} \geq \mathbf{0} \end{aligned} \tag{7}$$

求其对偶函数。

解

构造拉格朗日函数:

$$\begin{aligned} L(\mathbf{x}, \boldsymbol{\lambda}, \boldsymbol{\nu}) &= \mathbf{c}^T \mathbf{x} - \sum_{i=1}^n \lambda_i x_i + \boldsymbol{\nu}^T (\mathbf{A}\mathbf{x} - \mathbf{b}) \\ &= -\mathbf{b}^T \boldsymbol{\nu} + (\mathbf{c} + \mathbf{A}^T \boldsymbol{\nu} - \boldsymbol{\lambda})^T \mathbf{x} \end{aligned}$$

 L 关于 \mathbf{x} 线性, 因此对偶函数为

$$g(\boldsymbol{\lambda}, \boldsymbol{\nu}) = \inf_{\mathbf{x}} L(\mathbf{x}, \boldsymbol{\lambda}, \boldsymbol{\nu}) = \begin{cases} -\mathbf{b}^T \boldsymbol{\nu} & \mathbf{A}^T \boldsymbol{\nu} - \boldsymbol{\lambda} + \mathbf{c} = \mathbf{0} \\ -\infty & \text{其他情况} \end{cases}$$

最优值下界: 如果 $\mathbf{A}^T \boldsymbol{\nu} + \mathbf{c} \geq \mathbf{0}$, $p^* \geq -\mathbf{b}^T \boldsymbol{\nu}$

最大割问题（双向划分问题）

例 7

考虑问题

$$\begin{aligned} \min \quad & \mathbf{x}^T \mathbf{W} \mathbf{x} \\ \text{s.t.} \quad & x_i^2 = 1, \quad i = 1, \dots, n \end{aligned} \tag{8}$$

其中， $\mathbf{W} \in S^n$ ，求其对偶函数。

解

拉格朗日函数:

$$\begin{aligned} L(\mathbf{x}, \boldsymbol{\nu}) &= \mathbf{x}^T \mathbf{W} \mathbf{x} + \sum_{i=1}^n \nu_i (x_i^2 - 1) \\ &= \mathbf{x}^T (\mathbf{W} + \text{diag}(\boldsymbol{\nu})) \mathbf{x} - \mathbf{1}^T \boldsymbol{\nu} \end{aligned}$$

对偶函数:

$$\begin{aligned} g(\boldsymbol{\nu}) &= \inf_{\mathbf{x}} (\mathbf{x}^T \mathbf{W} \mathbf{x} + \sum_i \nu_i (x_i^2 - 1)) = \inf_{\mathbf{x}} \mathbf{x}^T (\mathbf{W} + \text{diag}(\boldsymbol{\nu})) \mathbf{x} - \mathbf{1}^T \boldsymbol{\nu} \\ &= \begin{cases} -\mathbf{1}^T \boldsymbol{\nu} & \mathbf{W} + \text{diag}(\boldsymbol{\nu}) \succeq \mathbf{0} \\ -\infty & \text{其他情况} \end{cases} \end{aligned}$$

最优值下界: 若取 $\boldsymbol{\nu} = -\lambda_{\min}(\mathbf{W})\mathbf{1}$, 则 $\mathbf{W} + \text{diag}(\boldsymbol{\nu}) \succeq \mathbf{0}$, 可得 p^* 的一个下界:

$$p^* \geq -\mathbf{1}^T \boldsymbol{\nu} = n\lambda_{\min}(\mathbf{W})$$

32.2.5 拉格朗日对偶函数与共轭函数的关系

回忆函数 $f: \mathbb{R}^n \rightarrow \mathbb{R}$ 的共轭函数 f^* 为

$$f^*(\mathbf{y}) = \sup_{\mathbf{x} \in \text{dom } f} (\mathbf{y}^T \mathbf{x} - f(\mathbf{x}))$$

例 8

考虑问题

$$\begin{aligned} \min \quad & f_0(\mathbf{x}) \\ \text{s.t.} \quad & \mathbf{x} = \mathbf{0} \end{aligned} \tag{9}$$

该问题的拉格朗日函数为 $L(\mathbf{x}, \boldsymbol{\nu}) = f_0(\mathbf{x}) + \boldsymbol{\nu}^T \mathbf{x}$, 对偶函数为

$$g(\boldsymbol{\nu}) = \inf_{\mathbf{x}} (f_0(\mathbf{x}) + \boldsymbol{\nu}^T \mathbf{x}) = -\sup_{\mathbf{x}} ((-\boldsymbol{\nu})^T \mathbf{x} - f_0(\mathbf{x})) = -f_0^*(-\boldsymbol{\nu})$$

例 9

考虑问题

$$\begin{aligned}
 \min \quad & f_0(\mathbf{x}) \\
 \text{s.t.} \quad & \mathbf{Ax} \leq \mathbf{b} \\
 & \mathbf{Cx} = \mathbf{d}
 \end{aligned} \tag{10}$$

利用函数 f_0 的共轭函数，其对偶函数为

$$\begin{aligned}
 g(\lambda, \nu) &= \inf_{\mathbf{x}} (f_0(\mathbf{x}) + \lambda^{\text{T}}(\mathbf{Ax} - \mathbf{b}) + \nu^{\text{T}}(\mathbf{Cx} - \mathbf{d})) \\
 &= -\mathbf{b}^{\text{T}}\lambda - \mathbf{d}^{\text{T}}\nu + \inf_{\mathbf{x}} (f_0(\mathbf{x}) + (\mathbf{A}^{\text{T}}\lambda + \mathbf{C}^{\text{T}}\nu)^{\text{T}}\mathbf{x}) \\
 &= -\mathbf{b}^{\text{T}}\lambda - \mathbf{d}^{\text{T}}\nu - f_0^*(-\mathbf{A}^{\text{T}}\lambda - \mathbf{C}^{\text{T}}\nu)
 \end{aligned} \tag{11}$$

函数 g 的定义域 $\text{dom } g = \{(\lambda, \nu) \mid -\mathbf{A}^{\text{T}}\lambda - \mathbf{C}^{\text{T}}\nu \in \text{dom } f_0^*\}$

- 这些例子告诉我们对偶函数与共轭函数之间的转换关系
- 在第 10 章的内容中，已经计算过许多函数的共轭函数
- 因此，可以利用前面章节共轭函数的结论，以及它与对偶函数的关系，直接求得对偶函数

32.2.6 用共轭函数求解拉格朗日对偶函数

等式约束条件下的范数极小化

例 10

考虑问题

$$\begin{aligned} \min \quad & f_0(\mathbf{x}) = \|\mathbf{x}\| \\ \text{s.t.} \quad & \mathbf{Ax} = \mathbf{b} \end{aligned} \tag{12}$$

其中 $\|\cdot\|$ 是任意范数，求其对偶函数。

等式约束条件下的范数极小化

解

 $f_0 = \|\cdot\|$ 的共轭函数为

$$f_0^*(\mathbf{y}) = \begin{cases} 0 & \|\mathbf{y}\|_* \leq 1 \\ \infty & \text{其他情况} \end{cases}$$

可以看出此函数是对偶范数单位球的示性函数。利用结论 (11) 得到问题 (12) 的对偶函数为

$$g(\boldsymbol{\nu}) = -\mathbf{b}^T \boldsymbol{\nu} - f_0^*(-\mathbf{A}^T \boldsymbol{\nu}) = \begin{cases} -\mathbf{b}^T \boldsymbol{\nu} & \|\mathbf{A}^T \boldsymbol{\nu}\|_* \leq 1 \\ -\infty & \text{其他情况} \end{cases}$$

熵的最大化

例 11

考虑问题

$$\begin{aligned} \min \quad & f_0(\mathbf{x}) = \sum_{i=1}^n x_i \log x_i \\ \text{s.t.} \quad & \mathbf{A}\mathbf{x} \leq \mathbf{b} \\ & \mathbf{1}^T \mathbf{x} = 1 \end{aligned} \tag{13}$$

其中 $\text{dom } f_0 = \mathbb{R}_{++}^n$, 求其对偶函数。

熵的最大化

解

关于实变量 x 的负熵函数 $x \log x$ 的共轭函数是 e^{y-1} 。由于函数 f_0 是不同变量的负熵函数的和，其共轭函数为

$$f_0^*(\mathbf{y}) = \sum_{i=1}^n e^{y_i-1}$$

其定义域为 $\text{dom } f_0^* = \mathbb{R}^n$ 。根据前述结论 (11)，问题 (13) 的对偶函数为

$$g(\boldsymbol{\lambda}, \nu) = -\mathbf{b}^T \boldsymbol{\lambda} - \nu - \sum_{i=1}^n e^{-\mathbf{a}_i^T \boldsymbol{\lambda} - \nu - 1} = -\mathbf{b}^T \boldsymbol{\lambda} - \nu - e^{-\nu-1} \sum_{i=1}^n e^{-\mathbf{a}_i^T \boldsymbol{\lambda}}$$

其中 \mathbf{a}_i 是矩阵 A 的第 i 列向量。

- 1 32.1 无约束优化最优性条件
- 2 32.2 拉格朗日对偶函数
- 3 32.3 拉格朗日对偶问题**
- 4 32.4 Slater 约束准则与强对偶原理
- 5 32.5 KKT 最优性条件
- 6 32.6 对偶问题在数据科学中优化问题的应用

32.3.1 显式表达对偶约束

拉格朗日对偶问题

$$\begin{aligned} \max \quad & g(\boldsymbol{\lambda}, \boldsymbol{\nu}) \\ \text{s.t.} \quad & \boldsymbol{\lambda} \geq \mathbf{0} \end{aligned} \tag{14}$$

- 从拉格朗日函数能够得到的最好下界
- 凸优化问题；最优值记为 d^* ；如果 $\boldsymbol{\lambda} \geq \mathbf{0}, (\boldsymbol{\lambda}, \boldsymbol{\nu}) \in \text{dom } g$ ，则 $(\boldsymbol{\lambda}, \boldsymbol{\nu})$ 是对偶可行的，称 $(\boldsymbol{\lambda}^*, \boldsymbol{\nu}^*)$ 是对偶最优解或者是最优拉格朗日乘子
- 经常，通过使隐含约束变为显式约束，问题能够得到简化

标准形式线性规划的拉格朗日对偶

例 12

考虑问题

$$\begin{aligned} \min \quad & \mathbf{c}^T \mathbf{x} \\ \text{s.t.} \quad & \mathbf{Ax} = \mathbf{b}, \quad \mathbf{x} \geq \mathbf{0} \end{aligned} \tag{15}$$

给出其对偶问题。

标准形式线性规划的拉格朗日对偶

解

对偶函数：

$$g(\boldsymbol{\lambda}, \boldsymbol{\nu}) = \begin{cases} -\mathbf{b}^T \boldsymbol{\nu} & \mathbf{A}^T \boldsymbol{\nu} - \boldsymbol{\lambda} + \mathbf{c} = \mathbf{0} \\ -\infty & \text{其他情况} \end{cases}$$

- 其对偶问题是在满足约束 $\boldsymbol{\lambda} \geq \mathbf{0}$ 的条件下极大化对偶函数 g

对偶问题：

$$\begin{aligned} \max \quad & g(\boldsymbol{\lambda}, \boldsymbol{\nu}) = \begin{cases} -\mathbf{b}^T \boldsymbol{\nu} & \mathbf{A}^T \boldsymbol{\nu} - \boldsymbol{\lambda} + \mathbf{c} = \mathbf{0} \\ -\infty & \text{其他情况} \end{cases} \\ \text{s.t.} \quad & \boldsymbol{\lambda} \geq \mathbf{0} \end{aligned} \quad (16)$$

标准形式线性规划的拉格朗日对偶

解

(续) 当且仅当 $A^T \nu - \lambda + c = 0$ 时对偶函数 g 有界, 通过隐式约束变为显式约束, 得到等价问题

$$\begin{aligned} \max \quad & -b^T \nu \\ \text{s.t.} \quad & A^T \nu - \lambda + c = 0 \\ & \lambda \geq 0 \end{aligned} \tag{17}$$

进一步, 此问题可以表述为

$$\begin{aligned} \max \quad & -b^T \nu \\ \text{s.t.} \quad & A^T \nu + c \geq 0 \end{aligned} \tag{18}$$

这个是一个不等式形式的线性规划。

不等式形式线性规划 (LP) 的拉格朗日对偶

例 13

考虑问题

$$\begin{aligned} \min \quad & \mathbf{c}^T \mathbf{x} \\ \text{s.t.} \quad & \mathbf{A} \mathbf{x} \leq \mathbf{b} \end{aligned} \tag{19}$$

给出其对偶问题。

不等式形式线性规划 (LP) 的拉格朗日对偶

解

拉格朗日函数:

$$L(\mathbf{x}, \boldsymbol{\lambda}) = \mathbf{c}^T \mathbf{x} + \boldsymbol{\lambda}^T (\mathbf{A}\mathbf{x} - \mathbf{b}) = -\mathbf{b}^T \boldsymbol{\lambda} + (\mathbf{A}^T \boldsymbol{\lambda} + \mathbf{c})^T \mathbf{x}$$

对偶函数:

$$g(\boldsymbol{\lambda}) = \inf_{\mathbf{x}} L(\mathbf{x}, \boldsymbol{\lambda}) = -\mathbf{b}^T \boldsymbol{\lambda} + \inf_{\mathbf{x}} (\mathbf{A}^T \boldsymbol{\lambda} + \mathbf{c})^T \mathbf{x} = \begin{cases} -\mathbf{b}^T \boldsymbol{\lambda} & \mathbf{A}^T \boldsymbol{\lambda} + \mathbf{c} = \mathbf{0} \\ -\infty & \text{其他情况} \end{cases} \quad (20)$$

称对偶变量 $\boldsymbol{\lambda}$ 是对偶可行的, 如果 $\boldsymbol{\lambda} \geq \mathbf{0}$ 且 $\mathbf{A}^T \boldsymbol{\lambda} + \mathbf{c} = \mathbf{0}$

不等式形式线性规划 (LP) 的拉格朗日对偶

解

(续) 对偶问题

$$\begin{aligned} \max \quad & -\mathbf{b}^T \boldsymbol{\lambda} \\ \text{s.t.} \quad & \mathbf{A}^T \boldsymbol{\lambda} + \mathbf{c} = \mathbf{0} \\ & \boldsymbol{\lambda} \geq \mathbf{0} \end{aligned} \tag{21}$$

- 可发现：不等式形式线性规划的对偶问题是一个标准形式的线性规划
- 可自行证明对偶问题的对偶问题等价于原问题

二次规划的拉格朗日对偶

例 14

考虑问题

$$\begin{aligned} \min \quad & \mathbf{x}^T \mathbf{P} \mathbf{x} \\ \text{s.t.} \quad & \mathbf{A} \mathbf{x} \leq \mathbf{b} \end{aligned} \tag{22}$$

其中 $\mathbf{P} \in S_+^m$ ，求其对偶问题。

二次规划的拉格朗日对偶

解

对偶函数:

$$g(\lambda) = \inf_x (\mathbf{x}^T \mathbf{P} \mathbf{x} + \lambda^T (\mathbf{A} \mathbf{x} - \mathbf{b})) = -\frac{1}{4} \lambda^T \mathbf{A} \mathbf{P}^{-1} \mathbf{A}^T \lambda - \mathbf{b}^T \lambda$$

对偶问题:

$$\begin{aligned} \max \quad & -\frac{1}{4} \lambda^T \mathbf{A} \mathbf{P}^{-1} \mathbf{A}^T \lambda - \mathbf{b}^T \lambda \\ \text{s.t.} \quad & \lambda \geq \mathbf{0} \end{aligned} \tag{23}$$

- 1 32.1 无约束优化最优性条件
- 2 32.2 拉格朗日对偶函数
- 3 32.3 拉格朗日对偶问题
- 4 32.4 Slater 约束准则与强对偶原理**
- 5 32.5 KKT 最优性条件
- 6 32.6 对偶问题在数据科学中优化问题的应用

32.4.1 弱对偶性

用 p^* 标记原问题的最优值， d^* 标记 Lagrange 对偶问题的最优值。则如下定理成立：

定理 5

不等式

$$d^* \leq p^* \quad (24)$$

成立。这个性质称为弱对偶性。

- 不等式总是成立的（不论是凸问题还是非凸问题）
- $(p^* - d^*)$ 被称为最优对偶间隙
- 当原问题很难求解时，弱对偶不等式可以给出原问题最优值的一个下界

弱对偶性

例如，考虑双向划分问题，其对偶问题是一个半定规划问题

$$\begin{aligned} \max \quad & -\mathbf{1}^T \boldsymbol{\nu} \\ \text{s.t.} \quad & \mathbf{W} + \text{diag}(\boldsymbol{\nu}) \succeq \mathbf{0} \end{aligned}$$

其中， $\boldsymbol{\nu} \in \mathbb{R}^n$ 。对偶问题能够被有效地求解，而且对偶问题的最优值给出了双向划分问题最优值的一个下界。

32.4.2 强对偶性和 Slater 约束准则

一个很自然的问题： p^* 和 d^* 在什么条件下相等？也即，所谓的强对偶性：

定义 2

如果原问题和对偶问题的最优值相等，即等式

$$p^* = d^* \quad (25)$$

成立，最优对偶间隙为零，那么，称它们满足强对偶性。

- 对于一般情况，强对偶性不成立
- (一般) 只对凸问题成立
- 保证强对偶性在凸问题中成立的条件称为约束准则

强对偶性和 Slater 约束准则

现在，考虑如下形式的凸优化问题：

$$\begin{aligned} \min \quad & f_0(\mathbf{x}) \\ \text{s.t.} \quad & f_i(\mathbf{x}) \leq 0, \quad i = 1, \dots, m, \\ & \mathbf{Ax} = \mathbf{b}, \end{aligned} \tag{26}$$

其中，函数 f_0, \dots, f_m 是凸函数。我们将证明当它的约束满足一定准则时，强对偶性将成立。

强对偶性和 Slater 约束准则

一个简单的约束准则是 **Slater 条件**: 存在一点 $\mathbf{x} \in \text{relint } \mathcal{D}$ ¹使得

$$f_i(\mathbf{x}) < 0, \quad i = 1, \dots, m, \quad \mathbf{Ax} = \mathbf{b}. \quad (27)$$

- 满足上述条件的点有时称为**严格可行**

¹给定集合 \mathcal{D} , 记其仿射包 $\text{affine } \mathcal{D}$ 。则:

$$\text{relint } \mathcal{D} = \{\mathbf{x} \in \mathcal{D} | \exists r > 0, \text{使得 } B(\mathbf{x}, r) \cap \text{affine } \mathcal{D} \subset \mathcal{D}\}$$

强对偶性和 Slater 约束准则

当 Slater 条件成立（且原问题是凸问题）时，强对偶性成立：

定理 6

（强对偶性定理）假设函数 f_0, f_1, \dots, f_m 以及 h_1, \dots, h_p 均为凸函数，而且满足 Slater 条件，那么

$$\sup_{\lambda \geq \mathbf{0}, \nu} g(\lambda, \nu) = \inf_{x \in K} f(x).$$

即对偶间隙为零。

强对偶性和 Slater 约束准则

改进的 Slater 条件：当不等式约束函数 f_i 中有一些仿射函数时，仿射不等式不需要严格成立。该条件为：存在一点 $\mathbf{x} \in \text{relint } \mathcal{D}$ 使得

$$f_i(\mathbf{x}) \leq 0, \quad i = 1, \dots, k, \quad f_i(\mathbf{x}) < 0, \quad i = k + 1, \dots, m, \quad \mathbf{Ax} = \mathbf{b} \quad (28)$$

- 当所有约束条件都是线性等式或不等式且 $\text{dom } f_0$ 是开集时，改进的 Slater 条件 (28) 就是可行性条件

32.4.3 强对偶性成立的一些示例

线性方程组的最小二乘解

例 15

考虑问题

$$\begin{aligned} \min \quad & \mathbf{x}^T \mathbf{x} \\ \text{s.t.} \quad & \mathbf{A}\mathbf{x} = \mathbf{b} \end{aligned}$$

求其对偶问题及讨论强对偶性。

线性方程组的最小二乘解

解

原问题：

$$\begin{aligned} \min \quad & \mathbf{x}^T \mathbf{x} \\ \text{s.t.} \quad & \mathbf{A} \mathbf{x} = \mathbf{b} \end{aligned}$$

对偶问题：

$$\max \quad -(1/4) \boldsymbol{\nu}^T \mathbf{A} \mathbf{A}^T \boldsymbol{\nu} - \mathbf{b}^T \boldsymbol{\nu}$$

- Slater 条件此时是原问题的可行性条件，强对偶性通常成立。

二次约束二次规划 (QCQP)

例 16

考虑问题

$$\begin{aligned} \min \quad & (1/2)\mathbf{x}^T \mathbf{P}_0 \mathbf{x} + \mathbf{q}_0^T \mathbf{x} + r_0 \\ \text{s.t.} \quad & (1/2)\mathbf{x}^T \mathbf{P}_i \mathbf{x} + \mathbf{q}_i^T \mathbf{x} + r_i \leq 0, \quad i = 1, \dots, m \end{aligned} \quad (29)$$

其中, $\mathbf{P}_0 \in S_{++}^m$, $\mathbf{P}_i \in S_{++}^m$, $i = 1, \dots, m$ 。求其对偶问题及讨论强对偶性。

二次约束二次规划 (QCQP)

解

原问题

$$\begin{aligned} \min \quad & (1/2)\mathbf{x}^T \mathbf{P}_0 \mathbf{x} + \mathbf{q}_0^T \mathbf{x} + r_0 \\ \text{s.t.} \quad & (1/2)\mathbf{x}^T \mathbf{P}_i \mathbf{x} + \mathbf{q}_i^T \mathbf{x} + r_i \leq 0, \quad i = 1, \dots, m \end{aligned}$$

其中, $\mathbf{P}_0 \in S_{++}^m$, $\mathbf{P}_i \in S_{++}^m$, $i = 1, \dots, m$

Lagrange 函数:

$$L(\mathbf{x}, \boldsymbol{\lambda}) = (1/2)\mathbf{x}^T \mathbf{P}(\boldsymbol{\lambda}) \mathbf{x} + \mathbf{q}(\boldsymbol{\lambda})^T \mathbf{x} + r(\boldsymbol{\lambda})$$

$$\text{其中, } \mathbf{P}(\boldsymbol{\lambda}) = \mathbf{P}_0 + \sum_{i=1}^m \lambda_i \mathbf{P}_i, \quad \mathbf{q}(\boldsymbol{\lambda}) = \mathbf{q}_0 + \sum_{i=1}^m \lambda_i \mathbf{q}_i, \quad r(\boldsymbol{\lambda}) = r_0 + \sum_{i=1}^m \lambda_i r_i$$

二次约束二次规划 (QCQP)

解

(续) 如果 $\lambda \geq 0$, 有 $P(\lambda) \succ \mathbf{0}$ 及

$$g(\lambda) = \inf_x L(\mathbf{x}, \lambda) = -(1/2)\mathbf{q}(\lambda)^T P(\lambda)^{-1} \mathbf{q}(\lambda) + r(\lambda)$$

对偶问题:

$$\begin{aligned} \max \quad & -(1/2)\mathbf{q}(\lambda)^T P(\lambda)^{-1} \mathbf{q}(\lambda) + r(\lambda) \\ \text{s.t.} \quad & \lambda \geq 0 \end{aligned} \tag{30}$$

- 根据 Slater 条件, 当二次不等式约束严格成立时, 即存在一点 \mathbf{x} 使得

$$(1/2)\mathbf{x}^T P_i \mathbf{x} + \mathbf{q}_i^T \mathbf{x} + r_i < 0, \quad i = 1, \dots, m$$

对偶问题 (30) 和原问题 (29) 之间强对偶性成立。

- 1 32.1 无约束优化最优性条件
- 2 32.2 拉格朗日对偶函数
- 3 32.3 拉格朗日对偶问题
- 4 32.4 Slater 约束准则与强对偶原理
- 5 32.5 KKT 最优性条件**
- 6 32.6 对偶问题在数据科学中优化问题的应用

32.5.1 互补松弛性

设强对偶性成立。令 \mathbf{x}^* 是原问题的最优解, $(\boldsymbol{\lambda}^*, \boldsymbol{\nu}^*)$ 是对偶问题的最优解, 这表明

$$\begin{aligned} f_0(\mathbf{x}^*) &= g(\boldsymbol{\lambda}^*, \boldsymbol{\nu}^*) = \inf_{\mathbf{x}} (f_0(\mathbf{x}) + \sum_{i=1}^m \lambda_i^* f_i(\mathbf{x}) + \sum_{j=1}^p \nu_j^* h_j(\mathbf{x})) \\ &\leq f_0(\mathbf{x}^*) + \sum_{i=1}^m \lambda_i^* f_i(\mathbf{x}^*) + \sum_{j=1}^p \nu_j^* h_j(\mathbf{x}^*) \\ &\leq f_0(\mathbf{x}^*) \end{aligned}$$

因此, 两个不等式取等号。这样便得出一些有意义的结论。

互补松弛性

$L(\mathbf{x}, \boldsymbol{\lambda}^*, \boldsymbol{\nu}^*)$ 在 \mathbf{x}^* 处取得最小值

$$\sum_{i=1}^m \lambda_i^* f_i(\mathbf{x}^*) = 0$$

求和项的每一项都非正，因此有

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

被称为互补松弛性，即对每个 i ， λ_i^* 与 $f_i(\mathbf{x}^*)$ 至少有一个为 0。

32.5.2 Karush-Kuhn-Tucker(KKT) 最优性条件

将前面的讨论汇在一起，便形成以下 KKT 条件

定义 3

令 \mathbf{x}^* 和 (λ^*, ν^*) 分别是原问题和对偶问题的最优解，其对偶间隙为零。因此有

- 原始约束：

$$f_i(\mathbf{x}^*) \leq 0, \quad i = 1, \dots, m$$

$$h_j(\mathbf{x}^*) = 0, \quad j = 1, \dots, p$$

- 对偶约束： $\lambda_i^* \geq 0, \quad i = 1, \dots, m$

- 互补松弛： $\lambda_i^* f_i(\mathbf{x}^*) = 0, \quad i = 1, \dots, m$

- 稳定性条件： $\nabla f_0(\mathbf{x}^*) + \sum_{i=1}^m \lambda_i^* \nabla f_i(\mathbf{x}^*) + \sum_{j=1}^p \nu_j^* \nabla h_j(\mathbf{x}^*) = 0$

这些公式被称为 **KKT (Karush-Kuhn-Tucker)** 条件 (f_i, h_j 可微)。

Karush-Kuhn-Tucker(KKT) 最优性条件

凸问题的最优性条件

定理 7

对于凸优化问题 (26), 如果 Slater 条件成立, 那么 $\mathbf{x}^*, \boldsymbol{\lambda}^*, \boldsymbol{\nu}^*$ 分别是原始、对偶全局最优解当且仅当它们满足 KKT 条件。

- 对于非凸优化问题, 通常 KKT 条件只是局部最优解的一个必要条件。
- 它是无约束优化问题一阶最优性条件 $\nabla f_0(\mathbf{x}) = 0$ 的一个推广。
- 有时可以直接利用 KKT 条件, 求得优化问题的解。

等式约束二次凸问题求极小

例 17

考虑问题

$$\begin{aligned} \min \quad & (1/2)\mathbf{x}^T \mathbf{P}\mathbf{x} + \mathbf{q}^T \mathbf{x} + r \\ \text{s.t.} \quad & \mathbf{A}\mathbf{x} = \mathbf{b} \end{aligned}$$

其中 $\mathbf{P} \in S_+^m$ 。求其 KKT 条件。

等式约束二次凸问题求极小

解

此问题的 KKT 条件为

$$Ax^* = b, \quad Px^* + q + A^T \nu^* = 0$$

可将其写成

$$H_x = \begin{bmatrix} P & A^T \\ A & 0 \end{bmatrix} \begin{bmatrix} x^* \\ \nu^* \end{bmatrix} = \begin{bmatrix} -q \\ b \end{bmatrix}$$

求解变量 x^*, ν^* 的 $m+n$ 个方程, 其中变量的维数为 $m+n$, 可以得到原问题和对偶问题的最优解。

32.5.3 通过解对偶问题求解原问题

在原问题不易求解，且强对偶性成立时，可先将其转化为对偶问题再求解。具体方式：

- 求得对偶最优解 $(\boldsymbol{\lambda}^*, \boldsymbol{\nu}^*)$
- 然后求下列问题的解

$$\min \quad f_0(\mathbf{x}) + \sum_{i=1}^m \lambda_i^* f_i(\mathbf{x}) + \sum_{j=1}^p \nu_j^* h_j(\mathbf{x})$$

- 如果此问题的解是原问题可行解，那么它就是原问题最优解
- 如果不是，那么原问题不存在最优点，即原问题的最优解无法达到

熵的最大化

例 18

考虑问题

$$\begin{aligned} \min \quad & f_0(\mathbf{x}) = \sum_{i=1}^n x_i \log x_i \\ \text{s.t.} \quad & \mathbf{A}\mathbf{x} \leq \mathbf{b} \\ & \mathbf{1}^T \mathbf{x} = 1 \end{aligned}$$

其中定义域为 \mathbb{R}_{++}^n ，求其最优解。

熵的最大化

解

对偶问题：

$$\begin{aligned} \max \quad & -\mathbf{b}^T \boldsymbol{\lambda} - \nu - e^{-\nu-1} \sum_{i=1}^n e^{-\mathbf{a}_i^T \boldsymbol{\lambda}} \\ \text{s.t.} \quad & \boldsymbol{\lambda} \geq \mathbf{0} \end{aligned}$$

- 假设 Slater 条件的弱化形式成立，即存在 $\mathbf{x} > \mathbf{0}$ 使得 $\mathbf{A}\mathbf{x} \leq \mathbf{b}$ 以及 $\mathbf{1}^T \mathbf{x} = 1$ ，因此强对偶性成立，存在一个对偶最优解 $(\boldsymbol{\lambda}^*, \nu^*)$

熵的最大化

解

(续) 设对偶问题已经解出。 $(\boldsymbol{\lambda}^*, \nu^*)$ 处的拉格朗日函数为

$$L(\mathbf{x}, \boldsymbol{\lambda}^*, \nu^*) = \sum_{i=1}^n x_i \log x_i + \boldsymbol{\lambda}^{*T} (\mathbf{A}\mathbf{x} - \mathbf{b}) + \nu^* (\mathbf{1}^T \mathbf{x} - 1)$$

它在 \mathcal{D} 上严格凸且有下界, 因此有一个唯一解 \mathbf{x}^* ,

$$x_i^* = 1 / \exp(\mathbf{a}_i^T \boldsymbol{\lambda}^* + \nu^* + 1), \quad i = 1, \dots, n$$

其中 \mathbf{a}_i 是矩阵 \mathbf{A} 的列向量。如果 \mathbf{x}^* 是原问题可行解, 则其必是原问题的最优解。如果 \mathbf{x}^* 不是原问题可行解, 则原问题的最优解不能达到。

在等式约束下极小化可分函数

例 19

考虑问题

$$\begin{aligned} \min \quad & f_0(\mathbf{x}) = \sum_{i=1}^n f_i(x_i) \\ \text{s.t.} \quad & \mathbf{a}^T \mathbf{x} = b \end{aligned}$$

其中 $\mathbf{a} \in \mathbb{R}^n$, $b \in \mathbb{R}$, 函数 $f_i: \mathbb{R} \rightarrow \mathbb{R}$ 是可微函数, 严格凸, 求其最优解。

在等式约束下极小化可分函数

解

拉格朗日函数:

$$L(\mathbf{x}, \nu) = \sum_{i=1}^n f_i(x_i) + \nu(\mathbf{a}^T \mathbf{x} - b) = -b\nu + \sum_{i=1}^n (f_i(x_i) + \nu a_i x_i)$$

对偶函数:

$$\begin{aligned} g(\nu) &= -b\nu + \inf_{\mathbf{x}} \left(\sum_{i=1}^n (f_i(x_i) + \nu a_i x_i) \right) \\ &= -b\nu + \sum_{i=1}^n \inf_{x_i} (f_i(x_i) + \nu a_i x_i) \\ &= -b\nu - \sum_{i=1}^n f_i^*(-\nu a_i) \end{aligned}$$

在等式约束下极小化可分函数

解

(续) 对偶问题:

$$\max \quad -b\nu - \sum_i^n f_i^*(-\nu a_i)$$

- 假设找到一个对偶最优解 ν^* 。 $L(\mathbf{x}, \nu^*)$ 具有唯一的最小点 $\tilde{\mathbf{x}}$, 有 $\tilde{\mathbf{x}} = \mathbf{x}^*$ 。 这可以通过求解 $\nabla_{\mathbf{x}} L(\mathbf{x}, \nu^*) = 0$ 得到 \mathbf{x}^* , 即求解方程组

$$f_i(x_i^*) = -\nu^* a_i, \quad i = 1, \dots, n.$$

- 1 32.1 无约束优化最优性条件
- 2 32.2 拉格朗日对偶函数
- 3 32.3 拉格朗日对偶问题
- 4 32.4 Slater 约束准则与强对偶原理
- 5 32.5 KKT 最优性条件
- 6 32.6 对偶问题在数据科学中优化问题的应用**

32.6.1 线性可分支支持向量机

分类学习最基本的想法就是基于训练集 D 在样本空间中找到一个划分超平面 $\mathbf{w}^T \mathbf{x} + b = 0$ ，将不同类别的样本分开。但对于线性可分数据集，能将训练样本分开的划分超平面可能有很多，我们应该去找到哪一个呢？

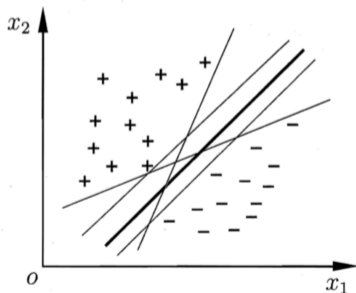


图 1: 存在多个划分超平面将两类训练样本分开

线性可分支持向量机

因此，期望找到具有“最大间隔” (maximum margin) 的划分超平面，将两类样本尽可能地分开。如下所示：

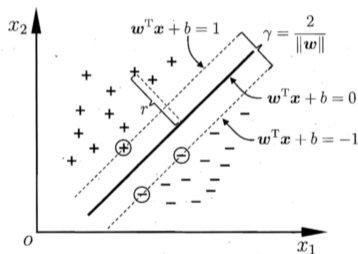


图 2: 支持向量与最大间隔

此时最大间隔为 $\gamma = \frac{2}{\|w\|}$ 。

线性可分支持向量机

因此，该问题为要找到能分类正确，使得 γ 最大的参数 \mathbf{w} 和 b ，即解如下优化问题：

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

这等价于

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

这便是线性可分支持向量机（Support vector Machine，简称 SVM）的模型。下面讨论其对偶问题。

线性可分支持向量机对偶问题

- 设拉格朗日乘子 $\alpha \geq 0$ ，则该问题的拉格朗日函数可写为

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

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

$$\mathbf{w} = \sum_{i=1}^N \alpha_i y_i x_i, \quad 0 = \sum_{i=1}^N \alpha_i y_i$$

- 回代可得对偶问题

$$\begin{aligned} \max_{\alpha} \quad & \sum_{i=1}^N \alpha_i - \frac{1}{2} \sum_{i=1}^N \sum_{j=1}^N \alpha_i \alpha_j y_i y_j x_i^T x_j \\ \text{s.t.} \quad & \sum_{i=1}^N \alpha_i y_i = 0, \quad \alpha_i \geq 0, \quad i = 1, 2, \dots, N \end{aligned}$$

线性可分支持向量机对偶问题

假设通过对偶问题求得 α ，则最终模型：

$$f(\mathbf{x}) = \mathbf{sign}(\mathbf{w}^T \mathbf{x} + b) = \mathbf{sign} \left(\sum_{i=1}^N \alpha_i y_i x_i^T \mathbf{x} + b \right)$$

KKT 条件

$$\begin{cases} \alpha_i \geq 0 \\ y_i f(x_i) - 1 \geq 0 \\ \alpha_i (y_i f(x_i) - 1) = 0 \end{cases}$$

总有 $\alpha_i = 0$ 或 $y_i f(x_i) = 1$

解的稀疏性：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，支持向量机因此而得名。

32.6.2 线性支持向量机

现实中的训练样本集，通常是线性不可分的，或者是近似线性可分。对于线性近似可分数据集，在第六章已经介绍，通过引入松弛变量的方式使其“可分”。因此得到如下优化问题：

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

线性支持向量机对偶问题

- 设拉格朗日乘子 $\alpha \geq 0$, $\mu \geq 0$, 则该问题的拉格朗日函数可写为

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

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

$$\mathbf{w} = \sum_{i=1}^N \alpha_i y_i x_i, \quad 0 = \sum_{i=1}^N \alpha_i y_i, \quad \boldsymbol{\alpha} + \boldsymbol{\mu} = C \mathbf{1}$$

线性支持向量机对偶问题

• 回代可得对偶问题

$$\begin{aligned} \max_{\alpha} \quad & \sum_{i=1}^N \alpha_i - \frac{1}{2} \sum_{i=1}^N \sum_{j=1}^N \alpha_i \alpha_j y_i y_j x_i^T x_j \\ \text{s.t.} \quad & \sum_{i=1}^N \alpha_i y_i = 0 \\ & 0 \leq \alpha_i \leq C, \quad i = 1, 2, \dots, N \end{aligned}$$

同样地，求解对偶问题，然后获得最终模型。

本讲小结

无约束优化的最优性条件

- 无约束可微优化问题：一阶条件、二阶条件
- 无约束不可微优化问题：凸优化的一阶条件
- 无约束不可微优化问题：复合优化的一阶条件

约束优化问题的最优性条件

- 拉格朗日对偶函数与拉格朗日对偶问题
- 弱对偶性、Slater 约束准则与强对偶原理
- KKT 最优性条件：原始约束、对偶约束、稳定性条件、互补松弛

这些最优性条件将用于指导我们设计优化问题的求解算法。