

# 最优化理论

陈鸿峥

2019.04\*

## 目录

|          |                   |           |
|----------|-------------------|-----------|
| <b>1</b> | <b>简介</b>         | <b>1</b>  |
| 1.1      | 优化概述 . . . . .    | 1         |
| 1.2      | 分类 . . . . .      | 2         |
| 1.3      | 历史 . . . . .      | 2         |
| <b>2</b> | <b>凸集</b>         | <b>3</b>  |
| <b>3</b> | <b>凸函数</b>        | <b>6</b>  |
| <b>4</b> | <b>凸优化问题</b>      | <b>11</b> |
| 4.1      | 标准型 . . . . .     | 11        |
| 4.2      | 线性规划 . . . . .    | 13        |
| <b>5</b> | <b>对偶理论</b>       | <b>17</b> |
| <b>6</b> | <b>优化算法</b>       | <b>27</b> |
| 6.1      | 简介 . . . . .      | 27        |
| 6.2      | 梯度下降法 . . . . .   | 28        |
| 6.3      | 非光滑优化问题 . . . . . | 33        |

## 1 简介

### 1.1 优化概述

优化(optimization): 从一个可行解的集合中寻找出最好的元素

---

\*Build 20190423

### 例 1. • 最小二乘线性拟合（凸问题）

- 深度神经网络（非凸，见下）

$$\mathbf{x}_1^{(i)} = f_1(\mathbf{x}_0^{(i)}, \mathbf{w}_1)$$

... ..

$$\mathbf{x}_n^{(i)} = f_n(\mathbf{x}_{n-1}^{(i)}, \mathbf{w}_n)$$

$$\min \sum_{i=1}^m (\mathbf{y}^{(i)} - \mathbf{x}_n^{(i)})^2$$

- 图像处理，自然图像通常都是分块光滑的，原图 $\Phi_0$ ，有噪声的新图 $\Phi$   
全变参( $TV$ , *Total Variation*)范数，计算图像每个像素点左侧和下侧的差异

$$\|\Phi\|_{TV} = \sum_y \sum_x \sqrt{(\Phi(x, y) - \Phi(x, y-1))^2 + (\Phi(x, y) - \Phi(x-1, y))^2}$$

可得优化目标：近似自然图像，而且跟原图不能差太远

$$\min(\|\Phi\|_{TV} + \lambda \|\Phi - \Phi_0\|_F^2)$$

- 推荐系统：*Netflix*问题

矩阵横向为用户，纵向为电影，值为评分值(1~5)，问题是把矩阵补全，这样就可以做推荐了→低秩矩阵补全

电影很多，但类型不多，关联关系有限→近似低秩<sup>1</sup>

低秩本来需要最小化 $\mathbf{z}$ 的非零奇异值数目 $\|\mathbf{z}\|_0$ ，但是非凸的；转化为最小化和范数<sup>2</sup> $\|\mathbf{z}\|_*$

$$\min \quad \|\mathbf{z}\|_* := \|\mathbf{z}\|_1$$

$$s.t. \quad \mathbf{z}_{ij} = \mathbf{M}_{ij}, (i, j) \in \Omega$$

## 1.2 分类

- 线性规划/非线性规划
- 凸规划/非凸规划（更好的分类）

目标函数凸函数，可行解集为凸集则是凸优化，一般容易求解

## 1.3 历史

- Newton-Raphson算法：求零点，等价于求 $\min f^2(x)$
- Gauss-Seidel算法：求解线性方程组 $A\mathbf{x} = \mathbf{b}$ ，等价于求 $\min \|A\mathbf{x} - \mathbf{b}\|_2^2$
- Lagrange

---

<sup>1</sup> $A$ 的秩等于非零奇异值 $\sqrt{\text{eig}(A^T A)}$ 数目

<sup>2</sup>矩阵所有奇异值之和

- Kantorov: 苏联, 线性规划, 诺贝尔经济学奖
- Dantzig: 美国, 优化决策, 线性规划单纯形
- Von Neumann: 线性规划问题对偶理论
- Karmarkar: 80年代, 线性规划内点法
- Nesterov: 后80年代, 非线性凸优化内点法
- 现代: 并行、随机算法

## 2 凸集

定义 1. 一些集合概念如下

- 仿射集 (*affine set*)

$C$  为仿射集  $\iff$  过  $C$  内任意两点的直线都在  $C$  内

$$\iff \forall x_1, x_2 \in C, \theta \in \mathbb{R}, \theta x_1 + (1 - \theta)x_2 \in C$$

例 2. 用定义易证线性方程组的解集  $C = \{x \mid Ax = b\}$  是仿射集; 反过来, 每一个仿射集都可以用线性方程组的解集表示

- 仿射组合

$$\forall x_1, x_2, \dots, x_k \in C, \theta_1, \dots, \theta_k \in \mathbb{R}, \theta_1 + \dots + \theta_k = 1 : \theta_1 x_1 + \dots + \theta_k x_k \in C$$

- 仿射包 (*hull*): 所有仿射组合的集合

$$\text{aff } C := \{\theta_1 x_1 + \dots + \theta_k x_k \mid \forall x_1, \dots, x_k \in C, \theta_1 + \dots + \theta_k = 1\}$$

- 凸集 (*convex set*)

$C$  为凸集  $\iff$  过  $C$  内任意两点的线段都在  $C$  内

$$\iff \forall x_1, x_2 \in C, \theta \in [0, 1], \theta x_1 + (1 - \theta)x_2 \in C$$

- 凸组合

$$\forall x_1, x_2, \dots, x_k \in C, \theta_1, \dots, \theta_k \in [0, 1], \theta_1 + \dots + \theta_k = 1 : \theta_1 x_1 + \dots + \theta_k x_k \in C$$

- 凸包: 最小的凸集

$$\text{conv } C := \{\theta_1 x_1 + \dots + \theta_k x_k \mid \forall x_1, \dots, x_k \in C, \theta_1, \dots, \theta_k \in [0, 1], \theta_1 + \dots + \theta_k = 1\}$$

- 凸锥(*convex cone*)

$$\mathcal{C} \text{ 为凸锥} \iff \forall x_1, x_2 \in \mathcal{C}, \theta_1, \theta_2 \geq 0, \theta_1 x_1 + \theta_2 x_2 \in \mathcal{C}$$

除了空集的凸锥都得包含原点 (取  $\theta_1 = \theta_2 = 0$ )

- 凸锥组合/非负线性组合:

$$\forall x_1, x_2, \dots, x_k \in \mathcal{C}, \theta_1, \dots, \theta_k \geq 0: \theta_1 x_1 + \dots + \theta_k x_k \in \mathcal{C}$$

- 凸锥包: 类似前面定义

由上面的定义易知, 仿射组合/凸锥组合 (强条件) 一定是凸组合。

**定义 2** (超平面(hyperplane)与半空间(halfspace)). 超平面都是比原空间低一维

$$\{\mathbf{x} \mid \mathbf{a}^T \mathbf{x} = b, \mathbf{x}, \mathbf{a} \in \mathbb{R}^n, b \in \mathbb{R}, \mathbf{a} \neq 0\}$$

超平面将空间划分为两个部分, 即半空间

$$\{\mathbf{x} \mid \mathbf{a}^T \mathbf{x} \leq b, \mathbf{a} \neq 0\}$$

若方程特解为  $\mathbf{x}_0$ , 则  $\mathbf{a} \perp (\mathbf{x} - \mathbf{x}_0)$

**定义 3** (欧式球(Euclidean ball)).

$$B(x_c, r) = \{x \mid \|x - x_c\|_2 \leq r\}$$

范数(*norm*)球可类似定义

**定义 4** (椭球(ellipsoid)).

$$\varepsilon(x_c, P) = \{x \mid (x - x_c)^T P^{-1} (x - x_c) \leq 1\}, P \succ 0$$

其中  $P \succ 0$  代表  $P$  对称且正定 ( $P = P^T$ )

分析. 定义内积  $\langle x^T P^{-1} y \rangle$  (需证满足内积条件), 进而  $P$ -范数  $\|x\|_P := \sqrt{x^T P x}$  是范数, 而椭球不过是  $P$ -范数意义下的球, 由定理得椭球是凸的

**定义 5** (多面体(polyhedron)).

$$P = \{\mathbf{x} \mid \mathbf{a}_i^T \mathbf{x} \leq b_i, \mathbf{c}_j^T \mathbf{x} = d_j, i = 1, \dots, m, j = 1, \dots, p\}$$

**例 3.** • 空集、点、 $\mathbb{R}^n$  空间均为仿射

- 任意直线为仿射; 若过原点则为凸锥

- $\mathbb{R}^n$ 空间的子空间<sup>3</sup>为仿射和凸锥
- 超平面为仿射
- 半空间、欧式球、椭球、多面体为凸集

定义 6 (仿射函数).

$$f: \mathbb{R}^n \mapsto \mathbb{R}^m \quad f(\mathbf{x}) = A\mathbf{x} + \mathbf{b}, A \in \mathbb{R}^{m \times n}, \mathbf{b} \in \mathbb{R}^m$$

性质如下:

- $S \subset \mathbb{R}^n$ 为凸  $\implies f(S) = \{f(\mathbf{x}) \mid \mathbf{x} \in S\}$ 为凸
- $C \subset \mathbb{R}^m$ 为凸  $\implies f^{-1}(C) = \{\mathbf{x} \in \mathbb{R}^n \mid f(\mathbf{x}) \in C\}$ 为凸

例 4. 两个集合的和  $S_1 + S_2 = \{x + y \mid x \in S_1, y \in S_2\}$ 保凸

分析. 直积  $S_1 \times S_2 = \{(x, y) \mid x \in S_1, y \in S_2\}$ 显然可以保凸 (相当于在两个集合同时画线)

令  $A = \begin{bmatrix} I & I \end{bmatrix}, \mathbf{x} = \begin{bmatrix} x & y \end{bmatrix}^T, \mathbf{b} = 0$ , 由仿射函数性质知

定义 7 (透视(perspective)函数<sup>4</sup>). 透视函数  $P: \mathbb{R}^{n+1} \mapsto \mathbb{R}^n, \text{dom } P = \mathbb{R}^n \times \mathbb{R}_{++}$  定义如下

$$P(z, t) = \frac{z}{t}, z \in \mathbb{R}^n, t \in \mathbb{R}_{++}$$

反透视函数

$$P^{-1}(c) := \{(x, t) \in \mathbb{R}^{n+1} \mid \frac{x}{t} \in c, t > 0\}$$

若  $c \in \text{dom } P$  为凸, 则  $P(c) := \{P(x), x \in c\}$  为凸; 反透视函数仍保持  $c$  的凸性。

考虑  $\mathbb{R}^{n+1}$  内的线段,  $x = (\tilde{x} \in \mathbb{R}^n, x_{n+1} \in \mathbb{R}_{++}), y = (\tilde{y}, y_{n+1})$  则经过透视函数仍是线段

分析.

$$P(\theta x + (1 - \theta)y) = \frac{\theta \tilde{x} + (1 - \theta)\tilde{y}}{\theta x_{n+1} + (1 - \theta)y_{n+1}} = \frac{\theta x_{n+1}}{\theta x_{n+1} + (1 - \theta)y_{n+1}} \frac{\tilde{x}}{x_{n+1}} + \frac{(1 - \theta)y_{n+1}}{\theta x_{n+1} + (1 - \theta)y_{n+1}} \frac{\tilde{y}}{y_{n+1}}$$

定义 8 (线性分数函数). 仿射函数

$$g(x) = \begin{bmatrix} A \\ C^T \end{bmatrix} x + \begin{bmatrix} b \\ d \end{bmatrix}, A \in \mathbb{R}^{m \times n}, b \in \mathbb{R}^m, c \in \mathbb{R}^n, d \in \mathbb{R}$$

线性分数函数  $f: \mathbb{R}^n \mapsto \mathbb{R}^m = p \circ g$

$$f(x) = \frac{Ax + b}{c^T x + d}, \text{dom } f = \{x \mid c^T x + d > 0\}$$

保凸性

- 凸集的交

---

<sup>3</sup>零元、加法封闭、数乘封闭

<sup>4</sup>++代表  $\geq 0$ , +++代表  $> 0$

- 仿射、逆仿射
- 透视函数
- 线性分数函数

### 3 凸函数

定义 9 (凸函数). 1.  $f: \mathbb{R}^n \rightarrow \mathbb{R}$  为凸  $\iff \text{dom } f$  为凸且  $\forall x, y \in \text{dom } f, \theta \in [0, 1]$

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

- 严格凸:  $\theta \in (0, 1)$ , 不等式不能取等
- 凹函数: 若  $-f$  为凸

2. 高维定义:  $f: \mathbb{R}^n \mapsto \mathbb{R}$  为凸  $\iff \text{dom } f$  为凸

$$\forall x \in \text{dom } f, v \in \mathbb{R}^n : g(t) := f(x + tv) \text{ 为凸, } \text{dom } g = \{t \mid x + tv \in \text{dom } f\}$$

相当于每一个剖面上的低维函数都是凸的

3. 一阶条件 (first-order condition)<sup>5</sup>

$$f(y) \geq f(x) + \nabla^T f(x)(y - x)$$

4. 二阶条件:  $f: \mathbb{R}^n \mapsto \mathbb{R}$  为凸  $\iff \text{dom } f$  为凸

$$\forall x \in \text{dom } f : \nabla^2 f(x) \succeq 0$$

- 凹函数:  $\nabla^2 f(x) \preceq 0$
- 严格凸:  $\iff \nabla^2 f(x) \succ 0$ , 反例  $f(x) = x^4$  (在一个点斜率不变并不要紧)

例 5.  $f(x) = a^T x + b$

分析. 有  $\nabla f(x) = a$ , 进而

$$f(y) = a^T y + b \geq a^T x + b + a^T (y - x) = a^T y + b$$

定义 10 (凸函数的扩展(extended-value)). 尽管凸函数的定义域为凸, 但往往不好处理, 那就将其扩展到全空间。  $x \in \text{dom } f \subset \mathbb{R}^n, \text{dom } \tilde{f} = \mathbb{R}^n$ , 会有

$$\tilde{f}(x) = \begin{cases} f(x) & x \in \text{dom } f \\ +\infty & x \notin \text{dom } f \end{cases}$$

---

<sup>5</sup> $\nabla^T f(x) = [\nabla f(x)]^T$

指示/示信(indicator)函数不一定是凸的

$$f(x) = \begin{cases} 0 & x \in C \\ +\infty & x \notin C \end{cases}$$

定理 1. 若  $f$  为凸, 可微, 则  $\exists x \in \text{dom } f, \nabla f(x) = 0$

例 6. 二次函数  $f(x) = \frac{1}{2}x^T Px + q^T x + r$ ,  $P \in S^n$  (对称矩阵),  $q^T \in \mathbb{R}^n$ ,  $r \in \mathbb{R}$

分析.  $\nabla^2 f(x) = P$

$P \in S_+^n$  凸,  $P \in S_{++}^n$  严格凸

例 7.  $f(x) = \frac{1}{x^2}$ ,  $\text{dom } f = \{x \in \mathbb{R}, x \neq 0\}$

分析. 注意  $\text{dom } f$  不是凸集

- 指数函数  $f(x) = e^{ax}$
- 幂函数  $f(x) = x^a$
- 绝对值的幂函数  $f(x) = |x|^p, x \in \mathbb{R}, p > 0$ :  $p \in [1, +\infty)$  凸,  $p \in (0, 1)$  既不凸又不凹

分析.

$$f''(x) = \begin{cases} p(p-1)x^{p-2} & x > 0 \\ p(p-1)(-x)^{p-2} & x < 0 \end{cases}$$

- 对数函数  $f(x) = \log x$
- 熵  $f(x) = -x \log x$
- 极大值函数  $f(x) = \max\{x_1, \dots, x_n\}, x \in \mathbb{R}^n$

定义 11 (解析近似). 无穷阶可微

极大值函数的解析近似是  $f(x) = \log(e^{x_1} + \dots + e^{x_n})$

$$\max\{x_1, \dots, x_n\} \leq f(x) \leq \max\{x_1, \dots, x_n\} + \log n$$

分析.

$$\begin{aligned} \frac{\partial f}{\partial x_i} &= \frac{e^{x_i}}{e^{x_1} + \dots + e^{x_n}} \\ \frac{\partial^2 f}{\partial x_i \partial x_j} &= \begin{cases} \frac{-e^{x_i} e^{x_i}}{(e^{x_1} + \dots + e^{x_n})^2} = -\frac{e^{2x_i}}{(e^{x_1} + \dots + e^{x_n})^2} & i = j \\ \frac{-e^{x_i} e^{x_j}}{(e^{x_1} + \dots + e^{x_n})^2} & i \neq j \end{cases} \\ z &:= [e^{x_1} \quad \dots \quad e^{x_n}]^T \end{aligned}$$

求 Hessian 矩阵

$$H = \frac{1}{(\mathbb{1}^T z)^2} (-z \cdot z^T + (\mathbb{1}^T z) \text{diag}(z))$$

将前面常量丢弃<sup>6</sup>

$$\begin{aligned}
 a_i &:= v_i \sqrt{z_i} = \begin{bmatrix} a_1 & \cdots & a_n \end{bmatrix}^T, b_i = \sqrt{z_i} \\
 v^T H v &= (\mathbb{1}^T z) v^T \text{diag}(z) v - v^T z z^T v \\
 &= \left( \sum_i z_i \right) \left( \sum_i v_i^2 z_i \right) - \left( \sum_i v_i z_i \right)^2 \\
 &= (b^T b)(a^T a) - (a^T b)^2 \quad \text{Cauchy} \\
 &\geq 0
 \end{aligned}$$

**定义 12** (范数).  $p(x)$  为范数

1.  $p(ax) = |a|p(x)$
2.  $p(x+y) \leq p(x) + p(y)$
3.  $p(x) = 0 \iff x = 0$

零范数  $\|x\|_0$ : 非零元素数目, 是伪范数 (不符合第一个定义)

$\mathbb{R}^n$  中的范数都是凸函数, 正则化!

分析.

$$\forall x, y, \theta \in [0, 1] p(\theta x + (1 - \theta)y) \leq \theta p(x) + (1 - \theta)p(y)$$

行列式的对数  $f(x) = \log \det(x)$ ,  $\text{dom } f = S_{++}^n$   $n=1$  凹函数证  $n > 1$  也为凹, 用高维定义

$$\begin{aligned}
 g(t) &= f(z + tv) \\
 &= \log \det(z + tv) \\
 &= \log \det(z^{1/2}(I + tz^{1/2}vz^{-1/2})z^{1/2}), \quad z^{1/2} \in S_{++}^n, z^{1/2}z^{1/2} = z \\
 &= \log \det(z) + \log \det(I + tz^{1/2}vz^{-1/2}) \\
 &= \log \det(z) + \sum_{i=1}^n \log(1 + t\lambda_i), \quad \lambda_i = z^{-1/2}vz^{1/2} \text{ 的特征值}
 \end{aligned}$$

$$\begin{aligned}
 g'(t) &= \sum_{i=1}^n \frac{\lambda_i}{1 + t\lambda_i} \\
 g''(t) &= \sum_{i=1}^n -\frac{\lambda_i^2}{(1 + t\lambda_i)^2}
 \end{aligned}$$

补充证明: 对对称阵特征值分解  $tz^{1/2}vz^{1/2} = tQ\Lambda Q^T$ , 对角阵  $\Lambda$  即为  $QQ^T = I$ ,  $Q$  为酉矩阵

$$I + tz^{-1/2}vz^{-1/2} = QQ^T + tQ\Lambda Q^T = Q(I + t\Lambda)Q^T$$

$$\log \det(I + tz^{-1/2}vz^{-1/2}) = \log \det(Q) + \log \det(I + t\Lambda) + \log \det(Q^T)$$

保持函数凸性

---

<sup>6</sup>  $H$  半正定, 则  $\forall v \in \mathbb{R}^n : v^T H v \geq 0$



- 非负加权和  $f_1, \dots, f_m$  为凸, 定义域  $\mathbb{R}^n$

$$f := \sum_{i=1}^m w_i f_i, w_i \geq 0$$

- 非负积分  $f(x, y)$  对  $y \in A$  均为凸 ( $A$  不一定为凸),  $w(y) \geq 0$

$$g(x) := \int_{y \in A} w(y) f(x, y) dy$$

- 仿射映射  $f: \mathbb{R}^n \mapsto \mathbb{R}$  为凸,  $A \in \mathbb{R}^{m \times n}, b \in \mathbb{R}^m$ ,  $\text{dom } g = \{x \mid Ax + b \in \text{dom } f\}$

$$g(x) := f(Ax + b)$$

分析. —  $\text{dom } f$  为凸, 则  $\text{dom } g$  为凸

$$- \forall x, y \in \text{dom } g, \forall \theta \in [0, 1]$$

$$\begin{aligned} g(\theta x + (1 - \theta)y) &= f(A(\theta x + (1 - \theta)y) + b) \\ &= f(\theta(Ax + b) + (1 - \theta)(Ay + b)) \\ &\leq \theta f(Ax + b) + (1 - \theta)f(Ay + b) \\ &= \theta g(x) + (1 - \theta)g(y) \end{aligned}$$

— 其实只是在定义域上改变, 而不是改变值域, 因而函数凸性不会改变

- 两个函数的极大值函数,  $f_1, f_2$  为凸

$$f(x) := \max\{f_1(x), f_2(x)\}, \text{dom } f = \text{dom } f_1 \cap \text{dom } f_2$$

- 任意个凸函数极大值函数为凸

$$f(x) = \max\{a_1^T x + b_1, \dots, a_m^T x + b_m\}$$

- 无限个凸函数,  $y \in A$ ,  $f(x, y)$  对于  $x$  为凸, 则

$$g(x) := \sup_{y \in A} f(x, y)$$

例 8. 点  $x$  到集合  $C$  的最远距离

$$f(x) = \sup_{y \in A} \|x - y\|$$

位移对于范数凸性不会有影响

例 9.  $x \in \mathbb{R}^n$ ,  $x[i]$  为第  $i$  大元素,  $x[1] \geq x[2] \geq \dots \geq x[r] \geq \dots \geq x[n]$

$$f(x) := \sum_{i=1}^r x[i]$$

–  $r = 1$ :  $f(x) = x[1] = \max\{x_1, \dots, x_n\}$ , 每一项都是  $\mathbf{e}_i^T x_i$

–  $r > 1$ :  $f(x) = \max\{x_{i_1} + \dots + x_{i_r} \mid 1 \leq i_1 < i_2 < \dots < i_r \leq n\}$

- 函数的组合:  $h: \mathbb{R}^k \mapsto \mathbb{R}, g: \mathbb{R}^n \mapsto \mathbb{R}^k$

$$f := h \circ g: \mathbb{R}^n \mapsto \mathbb{R}$$

先考虑  $n = k = 1, \text{dom } g = \mathbb{R}^n, \text{dom } h = \mathbb{R}^k, \text{dom } f = \mathbb{R}$ ,  $h, g$  二阶可微

$$f'(x) = h'(g(x)) \cdot g'(x)$$

$$f''(x) = h''(g(x))(g'(x))^2 + h'(g(x))g''(x) > 0$$

即当  $g$  为凸,  $h$  为凸且不降;  $g$  为凹,  $h$  为凸且不增时,  $f(x)$  为凸

(若定义域非全空间) 当  $g$  为凸,  $h$  为凸, 扩展值函数  $\tilde{h}$  不降;  $g$  为凹,  $h$  为凸,  $\tilde{h}$  不增时,  $f(x)$  为凸

例 10.  $g$  为凸,  $\exp g(x)$  为凸;  $g$  为凹,  $g > 0$ ,  $\log g(x)$  为凹;  $g$  为凸,  $g > 0$ ,  $1/g(x)$  为凸

例 11.  $g(x) = x^2, \text{dom } g = \mathbb{R}, h(y) = 0, \text{dom } h = [1, 2], f = h \circ g$ , 注意  $\tilde{h}$  并非不降!

- 函数透视:  $P: \mathbb{R}^{n+1} \mapsto \mathbb{R}^n, \text{dom } P \in \mathbb{R}^n \times \mathbb{R}_{++}, P(z, t) = \frac{z}{t}$

$$f: \mathbb{R}^n \mapsto \mathbb{R}, g(x, t) = tf\left(\frac{x}{t}\right), \text{dom } g = \{(x, t) \mid \frac{x}{t} \in \text{dom } f\}, g: \mathbb{R}^n \times \mathbb{R}_{++} \mapsto \mathbb{R}$$

若  $f(x)$  为凸, 则  $g(x, t)$  相对于  $(x, t)$  联合凸

例 12. –  $f(x) = x^T x, g(x, t) = x^T x/t$

–  $f(x) = -\log x, g(x, t) = t \log(t/x)$

–  $u, v \in \mathbb{R}_{++}^n, g(u, v) = \sum_{i=1}^n u_i \log(u_i/v_i)$ , 信息论常用, 衡量相似性, KL 散度

$$D_{KL} := \sum_{i=1}^n \left( u_i \log \frac{u_i}{v_i} - u_i + v_i \right)$$

定义 13 ( $\alpha$  次水平集( $\alpha$ -sub level set)).  $f: \mathbb{R}^n \mapsto \mathbb{R}, C_\alpha = \{x \in \text{dom } f \mid f(x) \leq \alpha\}$

定义 14 (拟凸函数(quasi-convex)).  $\alpha$  次水平集为凸集  $\iff f$  为拟凸函数

拟凸函数有很好的性质  $\rightarrow$  单模态/单峰函数

凸函数与凸集联系

- 凸函数定义域为凸集
- 凸函数的  $\alpha$  次水平集为凸集

## 4 凸优化问题

### 4.1 标准型

广义定义：极小化凸函数，约束为凸集

$$\begin{aligned} & \text{minimize} && f_0(\mathbf{x}) \\ & \text{subject to} && f_i(\mathbf{x}) \leq 0 \quad i = 1, \dots, m \\ & && h_j(\mathbf{x}) = 0 \quad j = 1, \dots, p \end{aligned}$$

- 优化变量  $\mathbf{x} \in \mathbb{R}^n$
- 目标/损失函数  $f_0 : \mathbb{R}^n \mapsto \mathbb{R}$
- 不等式约束函数  $f_i : \mathbb{R}^n \mapsto \mathbb{R}$
- 等式约束函数  $h_j : \mathbb{R}^n \mapsto \mathbb{R}$
- 域  $\mathcal{D} = \bigcap_{i=1}^m \text{dom } f_i \cap \bigcap_{j=1}^p \text{dom } h_j$
- 可行解  $\mathcal{X} = \{\mathbf{z} \mid f_i(\mathbf{z}) \leq 0, h_j(\mathbf{z}) = 0, i = 1, \dots, m, j = 1, \dots, p\}$
- 最优值  $P^* = \inf\{f_0(\mathbf{x}) \mid \mathbf{x} \in \mathcal{X}\}$
- 最优解  $\mathbf{x}^* \iff \forall \mathbf{z} \in \mathbb{R}^n, \mathbf{z} \in \mathcal{X} : f_0(\mathbf{z}) \geq f_0(\mathbf{x}^*)$
- 最优解集  $X^* = \{\mathbf{x}^* \mid f_0(\mathbf{x}^*) = P^*, \mathbf{x}^* \in \mathcal{X}\}$
- $\varepsilon$ -次最优解集  $X_\varepsilon = \{\mathbf{x} \mid f_0(\mathbf{x}) \leq P^* + \varepsilon, \mathbf{x} \in \mathcal{X}\}$
- 局部最优  $\exists R > 0, f_0(x) = \inf\{f_0(\mathbf{z}) \mid \mathbf{x} \in \mathcal{X}, \mathbf{z} \in \mathcal{X}, \|\mathbf{x} - \mathbf{z}\| \leq R\}$
- 局部最优解集  $x_{local} = \{\mathbf{x} \mid \mathbf{x} \text{ 为局部最优}\}$

狭义定义：  $f_i(x), i = 0, 1, \dots$  为凸函数，  $h_i(x)$  为仿射函数

例 13.

$$\begin{aligned} \min & f_0(x) = x_1^2 + x_2^2 \\ \text{s.t.} & f_1(x) = \frac{x_1}{1+x_2^2} \leq 0 \implies x_1 \leq 0 \\ & h_1(x) = (x_1 + x_2)^2 = 0 \implies x_1 + x_2 = 0 \end{aligned}$$

定理 2. 凸问题局部最优等价于全局最优

分析. 若  $x$  为局部最优

$$\exists R > 0 : f_0(x) = \inf\{f_0(z) \mid z \in \mathcal{X}, x \in \mathcal{X}, \|x - z\|_2 \leq R\}$$

反证法，设  $x$  不是全局最优，  $y$  为全局最优，  $f_0(x) > f_0(y)$

$$z = \theta x + (1 - \theta)y, \theta = \frac{R}{2\|y - x\|_2}$$

$$\|z - x\|_2 = \frac{R\|y - x\|_2}{2\|y - x\|_2} = \frac{R}{2}$$

由  $\|z - x\|_2 \leq R \implies f_0(x) \leq f_0(z)$ , 有

$$f_0(z) \leq \theta f_0(x) + (1 - \theta)f_0(y) < \theta f_0(z) + (1 - \theta)f_0(z) = f_0(z)$$

矛盾

可微凸目标函数

无约束  $\min f_0(x), \nabla f_0^*(x) = 0$

$$\forall x, y : f_0(y) \geq f_0(x) + \langle \nabla f_0(x), y - x \rangle$$

$$f_0(y) \geq f_0(x^*) + \langle \nabla f_0(x^*), y - x^* \rangle = f_0(x^*)$$

有约束  $\min f_0(x), s.t. x \in \mathcal{X}$

$$x^* \in \mathcal{X}, \langle \nabla f_0(x^*), y - x^* \rangle \geq 0, \forall y \in \mathcal{X}$$

例 14. 等式约束  $\min f_0(x), \text{dom } f_0 \subset \mathbb{R}^n, f_0$  可微, 使得  $Ax = b$

分析.  $x^*$  最优,  $Ax^* = b, \forall y \in \mathcal{X}, Ay = b$

$$\langle \nabla f_0(x^*), y - x^* \rangle \geq 0$$

$$\begin{cases} y = x^* + v \\ Av = 0 \end{cases}, v \in \text{Nul } A$$

$$\forall v \in \text{Nul } A, \langle \nabla f_0(x^*), v \rangle \geq 0$$

1.  $\text{Nul } A = \{0\}$

2.  $A$  不可逆,  $\nabla f_0(x^*) \perp \text{Nul } A$

例 15. 正约束  $\min f_0(x), s.t. x \geq 0$

分析. 若  $x^*$  最优,  $\iff x^* \geq 0, \forall y \geq 0, \langle \nabla f_0(x^*), y - x^* \rangle \geq 0$

$$\iff \langle \nabla f_0(x^*), y \rangle \geq \langle \nabla f_0(x^*), x^* \rangle$$

1. 若  $\nabla f_0(x^*) \not\geq 0$  有矛盾 (负数行乘上正无穷), 故  $\nabla f_0(x^*) \geq 0$

2. 令  $y = 0$ , 有  $0 \geq \langle \nabla f_0(x^*), x^* \rangle \implies \sum_{i=1}^n (\nabla f_0(x^*))_i x_i^* \leq 0$   
前面  $\geq 0$ , 进而互补松弛条件

3.  $x^* \geq 0$

## 4.2 线性规划

$$\begin{aligned} \min \quad & c^T \mathbf{x} + \mathbf{d} \\ \text{s.t.} \quad & G\mathbf{x} \leq \mathbf{h} \\ & A\mathbf{x} = \mathbf{b} \end{aligned}$$

$$\begin{aligned} \min \quad & c^T \mathbf{x} + \mathbf{d} \\ \text{s.t.} \quad & G\mathbf{x} + \mathbf{s} = \mathbf{h} \\ & \mathbf{s} \geq 0 \end{aligned}$$

$\mathbf{s}$ 为松弛变量(slack variable)

用 $\mathbf{x}^+$ 和 $\mathbf{x}^-$ 拆分, 得到 $\mathbf{x} = \mathbf{x}^+ - \mathbf{x}^-$ ,  $\mathbf{x}^+ \geq 0, \mathbf{x}^- \geq 0, \mathbf{s} \geq 0$

**例 16** (食谱问题).  $m$ 种营养元素不小于 $b_1, \dots, b_m$ ,  $n$ 种食物, 单位含量 $a_{1j}, \dots, a_{mj}$ , 食物量 $x_1, \dots, x_n$ , 价格 $c_1, \dots, c_n$

$$\begin{aligned} \min \quad & \sum_{j=1}^n c_j x_j \\ \text{s.t.} \quad & \sum_{j=1}^n a_{ij} x_j \geq b_i \\ & x_j \geq 0 \end{aligned}$$

其中 $i = 1, \dots, m, j = 1, \dots, n$

线性分数规划

$$\begin{aligned} \min \quad & f_0(x) = \frac{c^T x + d}{e^T x + f}, \text{dom } f = \{x \mid e^T x + f > 0\} \\ \text{s.t.} \quad & Gx \leq h \\ & Ax = b \end{aligned}$$

等价于

$$\begin{aligned}
 \min \quad & c^T y + dz \\
 \text{s.t.} \quad & Gy - hz \leq 0 \\
 & Ay - bz = 0 \\
 & e^T y + fz = 1 \\
 & z \geq 0
 \end{aligned}$$

分析. 证明两个问题等价,  $P_0$ 与 $P_1$

若 $x$ 在 $P_0$ 内可行

$$y = \frac{x}{e^T x + f}, z = \frac{1}{e^T x + f}$$

若 $(y, z)$ 在 $P_1$ 中可行

$$x = \frac{y}{z} (z \neq 0)$$

若 $z = 0$ ,  $x_0$ 为 $P_0$ 的可行解

$$\begin{aligned}
 x &= x_0 + ty, t \geq 0 \\
 \lim_{t \rightarrow \infty} \frac{c^T(x_0 + ty) + d}{e^T(x_0 + ty) + f} &= c^T y
 \end{aligned}$$

代入看所有条件结论都相同

二次规划(Quadratic Programming)

$$\begin{aligned}
 \min \quad & \frac{1}{2} x^T p x + q^T x + r, \quad p \succ 0 \\
 \text{s.t.} \quad & Gx \leq h \\
 & Ax = b
 \end{aligned}$$

二次约束二次规划(QCQP)

$$\begin{aligned}
 \min \quad & \frac{1}{2} p_0 x + q_0^T x + r_0, \quad p \succ 0 \\
 \text{s.t.} \quad & \frac{1}{2} x^T p_i x + q_i^T x + r_i \leq 0, \quad i = 1, \dots, m, p_i \succ 0 \\
 & Ax = b
 \end{aligned}$$

最小二乘问题

$$\begin{aligned}
 \min_x \quad & \frac{1}{2} \|Ax - b\|_2^2 \\
 \text{s.t.} \quad & Ax + e = b
 \end{aligned}$$

$$\frac{1}{2} (x^T A^T A x - 2b^T A x + b^T b)$$

一范数规范化最小二乘

$$\min \frac{1}{2} \|Ax - b\|_2^2 + \lambda_1 \|x\|_1$$

本来用零范数，但用一范数拟合

改写

$$\|x\|_1 = \mathbb{1}^T \mathbf{x}^+ + \mathbb{1}^T \mathbf{x}^-$$

Basic Pursuit

$$\begin{aligned} \min \quad & \frac{1}{2} \|Ax - b\|_2^2 \\ \text{s.t.} \quad & \|x\|_1 \leq \varepsilon_1 \end{aligned}$$

原式很难平衡两者，下式只需考虑 $\|x\|_1$ 的影响

岭回归(Ridge): 所有 $x$ 差距不要太大

$$\min \frac{1}{2} \|Ax - b\|_2^2 + \frac{1}{2} \lambda_2 \|x\|_2^2$$

$$\begin{aligned} \min \quad & \frac{1}{2} \|Ax - b\|_2^2 \\ \text{s.t.} \quad & \|x\|_2^2 \leq \varepsilon_2 \end{aligned}$$

投资组合问题(portfolio optimization): 初始价格 $x_1, \dots, x_n$ , 最终价格 $P_1 x_1, \dots, P_n x_n$

$$\begin{aligned} \max \quad & P_1 x_1 + \dots + P_n x_n \\ \text{s.t.} \quad & x_1 + \dots + x_n = B \\ & x_1, \dots, x_n \geq 0 \end{aligned}$$

$\bar{P} = \mathbb{E}(P)$ 已知,  $\Sigma = \mathbb{D}(P)$

$$\begin{aligned} \min \quad & x^T \Sigma x \\ \text{s.t.} \quad & p^T x \geq r_{\min} \\ & x_1 + \dots + x_n = B \\ & x_1, \dots, x_n \geq 0 \end{aligned}$$

半定规划(semi-definite programming, SDP) (矩阵意义下的线性规划问题):  $X \in \mathbb{R}^{n \times n}, C \in \mathbb{R}^{n \times n}, A_i \in$

$$\mathbb{R}^{n \times n}, b_i \in \mathbb{R}$$

$$\begin{aligned} \min \quad & \text{tr}(CX) \\ \text{s.t.} \quad & \text{tr}(A_i X) = b_i, i = 1, \dots, p \\ & X \succeq 0 \end{aligned}$$

**例 17** (谱范数极小化问题). 矩阵多项式  $A(x) = A_0 + x_1 A_1 + \dots + x_n A_n, A_i \in \mathbb{R}^{p \times q}$

$$\min_x \|A(x)\|_2$$

谱范数代表  $A(x)$  的最大奇异值<sup>7</sup>

$$\begin{aligned} \min_{x, s} \quad & S \\ \text{s.t.} \quad & A^T(x)A(x) \preceq SI \end{aligned}$$

**例 18** (最快分布式线性平均).

$$\begin{aligned} x(t) &= Px(t-1) \\ P &= \begin{bmatrix} P_{11} & \dots & P_{1n} \\ \vdots & & \vdots \\ P_{n1} & \dots & P_{nn} \end{bmatrix}, P\mathbb{1} = \mathbb{1} \end{aligned}$$

其中  $(i, j) \in E$  或  $i = j$ ,  $P_{ij} \neq 0$ ; 否则  $P_{ij} = 0$

$$P = P^T, P_{ij} = P_{ji}, P_{ij} > 0$$

$$P \succeq 0, P_{ij} \geq 0$$

只要图是连通图, 则一定会收敛

收敛速度与第二大/特征值绝对值/有关

$$1 = \lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_n \geq -1$$

即收敛速度由  $\max\{|\lambda_2|, |\lambda_n|\}$  决定

$$\begin{aligned} \min \max\{|\lambda_2|, |\lambda_n|\} \\ \max\{|\lambda_2|, |\lambda_n|\} = \left\| P - \frac{1}{n} \mathbb{1} \mathbb{1}^T \right\| \end{aligned}$$

---

<sup>7</sup>谱范数是诱导范数, F-范数(Frobenias)  $\|A(x)\|_F$  才算是矩阵意义下的2范数



$$\begin{aligned}
\min \quad & t := \left\| P - \frac{1}{n} \mathbb{1} \mathbb{1}^T \right\|_2 \\
\text{s.t.} \quad & P \mathbb{1} = \mathbb{1} \\
& P = P^T \\
& P \succeq 0 \\
& P_{ij} = 0, \quad (i, j) \neq E \wedge i \neq j \\
& -tI \preceq P - \frac{1}{n} \mathbb{1} \mathbb{1}^T \preceq tI
\end{aligned}$$

多目标优化问题：帕累托最优解

若有另一解在某个指标上更好，则必有指标更差

帕累托最优值/帕累托最优面

$\min f_{01}$  与  $\min f_{02}$  的交点为理想点(oracle)

若  $f_{01}(x), \dots, f_{0q}(x)$  为凸， $\mathcal{X}$  为凸

$$\begin{aligned}
\min \quad & \lambda_1 f_{01}(x) + \dots + \lambda_q f_{0q}(x), \quad \lambda_1, \dots, \lambda_q \geq 0 \\
\text{s.t.} \quad & x \in \mathcal{X}
\end{aligned}$$

1. 能找到一个Pareto最优解
2. 遍历  $\lambda_1, \dots, \lambda_q$ ，可找到全部

岭回归的多目标优化表示

$$\begin{cases} \min \frac{1}{2} \|Ax - b\|_2^2 \\ \min \frac{1}{2} \|x\|_2^2 \end{cases}$$

## 5 对偶理论

拉格朗日函数(Lagrangian function)

$$L(x, \lambda, v) = f_0(x) + \sum_{i=1}^m \lambda_i f_i(x) + \sum_{i=1}^p v_i h_i(x), \text{dom } L = D \times \mathbb{R}^m \times \mathbb{R}^p$$

拉格朗日乘子(multiplier)

- 原变量(primal variable):  $\lambda = [\lambda_1 \quad \dots \quad \lambda_m]^T$
- 对偶变量(dual variable):  $v = [v_1 \quad \dots \quad v_p]^T$

拉格朗日对偶函数

$$\begin{aligned} g(\lambda, v) &= \inf_{x \in \mathcal{D}} L(x, \lambda, v) \\ &= \inf_{x \in \mathcal{D}} \left( f_0(x) + \sum_{i=1}^m \lambda_i f_i(x) + \sum_{i=1}^p v_i h_i(x) \right) \end{aligned}$$

注意遍历域是  $\mathcal{D} = \bigcap_{i=1}^m \text{dom } f_i \cap \bigcap_{i=1}^p \text{dom } h_i$ ，而不是可行解集  $\mathcal{X}$

- $g(\lambda, v)$  一定是关于  $\lambda$  和  $v$  的凹函数（关于  $\lambda$  和  $v$  的仿射函数，注意  $x$  为常数）
  - $\forall \lambda \geq 0, \forall v, g(\lambda, v) \leq P^*$
- 对偶(dual)问题

$$\begin{aligned} \max \quad & g(\lambda, v) \\ \text{s.t.} \quad & \lambda \geq 0 \end{aligned}$$

其最优解记为  $D^*$ ，则  $D^* \leq P^*$ ，即给出了原问题的一个最优下界

$x^*$  原问题最优解

$$\begin{aligned} \sum_{i=1}^m \lambda_i f_0(x^*) + \sum_{i=1}^p v_i h_i(x^*) &\leq 0 \\ L(x^*, \lambda, v) = f_0(x^*) + (\dots) &\leq P^* \\ g(\lambda, v) = \inf_{x \in \mathcal{D}} L(x, \lambda, v) &\leq L(x^*, \lambda, v) \leq P^* \end{aligned}$$

例 19.

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

分析.

$$\begin{aligned} L(x, v) &= x^T x + v^T (Ax - b) \\ g(v) &= \inf_{x \in \mathcal{D}} L(x, v) \\ &= \inf_{x \in \mathcal{D}} x^T x + v^T Ax - v^T b \\ &= \left(-\frac{A^T v}{2}\right)^T \left(-\frac{A^T v}{2}\right) + v^T A \left(-\frac{A^T v}{2}\right) - v^T b \\ &= -\frac{1}{4} v^T A A^T v - b^T v \end{aligned}$$

补充求梯度： $2x + A^T v = 0 \implies x = -\frac{A^T v}{2}$

因而得到对偶问题

$$\max_v -\frac{1}{4} v^T A A^T v - b^T v$$

例 20.

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

分析. 注意 $\lambda$ 前面符号, 要化为一般形式

$$L(x, \lambda, v) = c^T x - \lambda^T x + v^T (Ax - b)$$

$$\begin{aligned} g(\lambda, v) &= \inf_x L(x, \lambda, v) \\ &= \inf_x (c - \lambda + A^T v)^T x - v^T b \\ &= \begin{cases} -\infty & c - \lambda + A^T v \neq 0 \\ -v^T b & c - \lambda + A^T v = 0 \end{cases} \end{aligned}$$

对偶问题, 由于要极大, 故不考虑负无穷部分

$$\begin{aligned} \max_{\lambda, v} \quad & -v^T b \\ \text{s.t.} \quad & c - \lambda + A^T v = 0 \\ & \lambda \geq 0 \end{aligned}$$

逆过来求解

$$\begin{aligned} \min \quad & b^T v \\ \text{s.t.} \quad & A^T v + c \geq 0 \end{aligned}$$

$$L(v, \lambda) = b^T v - \lambda^T (A^T v + c)$$

$$\begin{aligned} g(\lambda) &= \inf_v L(v, \lambda) \\ &= \inf_v (b - A\lambda)^T v - \lambda^T c \\ &= \begin{cases} -\lambda^T c & b - A\lambda = 0 \\ -\infty & b - A\lambda \neq 0 \end{cases} \end{aligned}$$

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

对偶的对偶不一定回去，线性规划才满足

例 21.

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

分析.

$$L(x, v) = x^T w x + \sum_{i=1}^n v_i (x_i^2 - 1)$$

$$\begin{aligned} g(v) &= \inf_x L(x, v) \\ &= \inf_x x^T w x + \sum_{i=1}^n v_i x_i^2 - \sum_{i=1}^n v_i \\ &= \inf_x x^T (w + \text{diag}(v)) x - \mathbb{1}^T v = \begin{cases} -\mathbb{1}^T v & w + \text{diag}(v) \succeq 0 \\ -\infty & \text{otherwise} \end{cases} \end{aligned}$$

补充求梯度： $2(w + \text{diag}(v))x = 0$

$$\begin{aligned} \max_v \quad & -\mathbb{1}^T v \\ \text{s.t.} \quad & w + \text{diag}(v) \succeq 0 \end{aligned}$$

定义 15 (函数的共轭).  $f: \mathbb{R}^n \mapsto \mathbb{R}, f^*(y) = \sup_{x \in \text{dom } f} (y^T x - f(x))$ , 几何意义即到不同斜率直线的距离最大值

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

$$\begin{aligned} L(x, \lambda, v) &= f_0(x) + \lambda^T (Ax - b) + v^T (cx - d) \\ &= f_0(x) + (A^T \lambda + c^T v)^T x - \lambda^T b - v^T d \\ g(\lambda, v) &= \inf_x f_0(x) + (A^T \lambda + c^T v)^T x - \lambda^T b - v^T d \\ &= -\sup_x -(A^T \lambda + c^T v)^T x - f_0 \\ &= -f_0^*(-(A^T \lambda + c^T v)) - \lambda^T b - v^T d \end{aligned}$$

对偶间隙(duality gap):  $p^* - d^* \geq 0$

- 弱对偶：严格大于0
- 强对偶：对偶间隙为0

1. 对于非凸问题, 通常  $p^* \neq d^*$
2. 对于凸问题, 若slater条件满足,  $p^* = d^*$

定义 16 (相对内点(relative interior)).

$$relint D = \{x \in D \mid B(x, r) \cap \text{aff } D \subset v, \exists r > 0\}$$

定理 3 (Slater条件).

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

$\exists x \in relint D$  使得  $f_i(x) < 0, i = 1, \dots, m, Ax = b$

例 22. 二次规划(QP)

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

Slater条件  $\{x \mid Ax = b\}$  非空

例 23. 二次约束二次规划(QCQP)

$$\begin{aligned} \min \quad & \frac{1}{2} x^T P_0 x + q_0^T x + r_0 \\ \text{s.t.} \quad & \frac{1}{2} x^T P_i x + q_i^T x + r_i \leq 0, \quad i = 1, \dots, m \end{aligned}$$

$P_0, \dots, P_i$  半正定

凸问题+Slater条件  $\implies p^* = d^*$ , 但有可能不满足Slater条件也依然强对偶

例 24.

$$\begin{aligned} \min \quad & x, x \in \mathbb{R} \\ \text{s.t.} \quad & x \leq 0 \\ & -x \leq 0 \end{aligned}$$

分析.

$$\begin{aligned} L(x, \lambda_1, \lambda_2) &= x + \lambda_1 x - \lambda_2 x = (1 + \lambda_1 - \lambda_2)x \\ g(\lambda_1, \lambda_2) &= \inf_{x \in \mathbb{R}} (1 + \lambda_1 - \lambda_2)x = \begin{cases} 0 & 1 + \lambda_1 - \lambda_2 = 0 \\ -\infty & \text{otherwise} \end{cases} \end{aligned}$$

$$\begin{aligned} \max_{\lambda_1, \lambda_2} \quad & 0 \\ \text{s.t.} \quad & 1 + \lambda_1 - \lambda_2 = 0 \end{aligned}$$

$$\implies p^* = d^* = 0$$

置信域问题

$$\begin{aligned} \min \quad & x^T A x + b^T x \\ \text{s.t.} \quad & x^T x \leq 1 \\ & A \not\equiv 0 \end{aligned}$$

依然可以得到  $p^* = d^*$   
几何解释

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

$$G = \{(f_1(x), f_0(x)) \mid x \in \mathcal{D}\}$$

$$g(\lambda) = \inf \{t + \lambda u \mid (u, t) \in G\}$$

$$L(x, \lambda) = f_0(x) + \lambda f_1(x)$$

$$g(\lambda) = \inf_{x \in \mathcal{D}} \{f_0(x) + \lambda f_1(x)\}$$

$$p^* = \inf \{t \mid (u, t) \in G, u \leq 0\}$$

$$\lambda \geq 0, \max g(\lambda)$$

注意问题必须要有可行解

经济学解释：满足原材料约束下，利润最多价格  $\lambda_i \geq 0$

$$g(\lambda) = \inf_x f(x) + \lambda_1 f_1(x) + \dots + \lambda_m f_m(x) = \inf_x L(x, \lambda)$$

则  $g(\lambda)$  为对偶函数，市场  $p^*$  损失最小 ( $g(\lambda) \leq p^*$ )

$$d^* = \sup_{\lambda \geq 0} g(\lambda)$$

市场平衡点，均衡市场  $p^* = d^*$ ，最优/影子价格  $\lambda^*$

多目标优化解释

$$\begin{cases} \min f_0(x) & 1 \\ \min f_1(x) & \lambda_1 \\ \vdots & \vdots \\ \min f_m(x) & \lambda_m \end{cases}$$

$$\min_x f_0(x) + \lambda_1 f_1(x) + \cdots + \lambda_m f_m(x)$$

鞍点(saddle point)解释

$$f(w, z), w \in S_w, z \in S_z$$

极小极大不等式

$$\sup_{z \in S_z} \inf_{w \in S_w} f(w, z) \leq \inf_{w \in S_w} \sup_{z \in S_z} f(w, z)$$

若有 $(\tilde{w}, \tilde{z})$ 使得

$$(\tilde{w}, \tilde{z}) = \arg \max_{z \in S_z} \min_{w \in S_w} f(w, z)(\tilde{w}, \tilde{z}) = \arg \min_{w \in S_w} \max_{z \in S_z} f(w, z)$$

则 $(\tilde{w}, \tilde{z})$ 为鞍点

有下面不等式成立

$$f((\tilde{w}, z)) \leq f(\tilde{w}, \tilde{z}) \leq f(w, \tilde{z}), \forall z \in S_z, w \in S_w$$

即从一个方向望过去是最小，从另一个方向望过去是最大

$$\begin{aligned} L(x, \lambda) &= f_0(x) + \sum_{i=1}^m \lambda_i f_i(x) \\ \implies \sup_{\lambda \geq 0} L(x, \lambda) &= \sup_{\lambda \geq 0} \{f_0(x) + \sum_{i=1}^m \lambda_i f_i(x)\} \\ &= \begin{cases} f_0(x) & f_i(x) \leq 0, i = 1, \dots, m \\ +\infty & \text{otherwise} \end{cases} \\ \implies p^* &= \inf_x \{f_0(x) \mid f_i(x) \leq 0, i = 1, \dots, m\} = \inf_x \sup_{\lambda \geq 0} L(x, \lambda) \\ d^* &= \sup_{\lambda \geq 0} g(\lambda) = \sup_{\lambda \geq 0} \inf_x L(x, \lambda) \implies p^* \geq d^* \end{aligned}$$

如果 $L(x, \lambda)$ 有鞍点，则必有 $p^* = d^*$

鞍点在无约束优化问题中是很糟糕的点（所有方向上梯度为0），但是有约束优化问题则是非常好的点

若 $(\tilde{x}, \tilde{\lambda})$ 为 $L(x, \lambda)$ 鞍点  $\iff p^* = d^*$  且 $\tilde{x}, \tilde{\lambda}$ 为原对偶问题最优解  $\implies$  若为鞍点， $p^* = d^*$

$$\sup_{\lambda \geq 0} \inf_x L(x, \lambda) = \inf_x \sup_{\lambda \geq 0} L(x, \lambda)$$

已知 $(\tilde{x}, \tilde{\lambda})$ 为左边最优

$$\tilde{\lambda} = \arg \max_{\lambda \geq 0} \inf_x L(x, \lambda)$$

$$\tilde{x} = \arg \inf_x \sup_{\lambda \geq 0} L(x, \lambda)$$

则 $\tilde{\lambda}$ 对偶最优,  $\tilde{x}$ 为原问题最优

$$\min f_0(x)$$

$$\text{s.t. } f_i(x) \leq 0, \quad i = 1, \dots, m$$

**定理 4.**  $(\tilde{x}, \tilde{\lambda})$ 为拉格朗日函数鞍点  $\iff p^* = d^*$ , 且 $(\tilde{x}, \tilde{\lambda})$ 为原对偶的最优解

分析. 右推左,  $(\tilde{x}, \tilde{\lambda})$ 原对偶可行

$$f_i(\tilde{x}) \leq 0, i = 1, \dots, m, \tilde{\lambda} \geq 0$$

因 $p^* = d^*$ , 有

$$\begin{aligned} f_0(\tilde{x}) &= g(\tilde{\lambda}) \\ &= \inf_x \{f_0(x) + \sum_{i=1}^m \tilde{\lambda}_i f_i(x)\} \\ &\leq f_0(\tilde{x}) + \sum_{i=1}^m \tilde{\lambda}_i f_i(\tilde{x}) \\ &\leq f_0(\tilde{x}) \end{aligned}$$

进而不等号都得为等号

$$1. \inf_x L(x, \tilde{\lambda}) = L(\tilde{x}, \tilde{\lambda})$$

$$2. f_0(\tilde{x}) = \sup_{\lambda \geq 0} \{f_0(\tilde{x}) + \sum_{i=1}^m \lambda_i f_i(\tilde{x})\} = \sup_{\lambda \geq 0} L(\tilde{x}, \lambda)$$

$$\implies L(\tilde{x}, \tilde{\lambda}) = \sup_{\lambda \geq 0} L(\tilde{x}, \lambda)$$

$$\implies (\tilde{x}, \tilde{\lambda}) \text{ 是 } L(x, \lambda) \text{ 的鞍点}$$

一般优化问题的对偶理论

$$\min f_0(x)$$

$$\text{s.t. } f_i(x) \leq 0, \quad i = 1, \dots, m$$

$$h_i(x) = 0, \quad i = 1, \dots, p$$

不一定是凸问题, 但 $p^* = d^*$ , 最优解满足什么条件?



## 对偶问题

$$\begin{aligned} \max \quad & g(\lambda, v) \\ \text{s.t.} \quad & \lambda \geq 0 \end{aligned}$$

分析. 设 $(x^*, \lambda^*, v^*)$ 为原对偶最优解, 则 $(x^*, \lambda^*, v^*)$ 为原对偶可行解

$$f_i(x^*) \leq 0, i = 1, \dots, m, \quad h_i(x^*) = 0, i = 1, \dots, p, \quad \lambda^* \geq 0$$

$$\begin{aligned} p^* = d^* &\implies f_0(x^*) = g(\lambda^*, v^*) \\ &= \inf_x \{f_0(x) + \sum_{i=1}^m \lambda_i^* f_i(x) + \sum_{i=1}^p v_i^* h_i(x)\} \\ &\leq f_0(x^*) + \sum_{i=1}^m \lambda_i^* f_i(x^*) + \sum_{i=1}^p v_i^* h_i(x^*) \\ &\leq f_0(x^*) \end{aligned}$$

同上理, 不等号全取等

$$1. \lambda_i^* f_i(x^*) = 0, \forall i = 1, \dots, m$$

$$2. x^* = \arg \min_x L(x, \lambda^*, v^*)$$

若 $f_0, f_i, h_i$ 均可微, 则必要条件为

$$\left. \frac{\partial L(x, \lambda^*, v^*)}{\partial x} \right|_{x=x^*} = 0$$

可微优化问题的KKT(Karush-Kuhn-Tucker)条件

- $f_i(x^*) \leq 0, i = 1, \dots, m$  primal feasibility
- $h_i(x^*) = 0, i = 1, \dots, p$  primal feasibility
- $\lambda^* \geq 0$  dual feasibility
- $\lambda_i^* f_i(x^*) = 0, i = 1, \dots, m$  complementarity slackness(对偶互斥条件)
- $\left. \frac{\partial L(x, \lambda^*, v^*)}{\partial x} \right|_{x=x^*} = 0$  stability

**定理 5.** 若原问题为凸, 则KKT条件为充要条件

分析. 必要性已证, 证明充分性

若 $(\tilde{x}, \tilde{\lambda}, \tilde{v})$ 满足KKT条件  $\implies (\tilde{x}, \tilde{\lambda}, \tilde{v})$ 最优  $\tilde{x}$ 为原问题可行解,  $(\tilde{\lambda}, \tilde{v})$ 为对偶问题可行解

证明思路:  $g(\tilde{\lambda}, \tilde{v}) = f_0(\tilde{x})$

$L(x, \tilde{\lambda}, \tilde{v})$  为  $x$  的凸函数, 则  $\tilde{x}$  使  $L(x, \tilde{\lambda}, \tilde{v})$  最小

$$\begin{aligned}
g(\tilde{\lambda}, \tilde{v}) &= \inf_x L(x, \tilde{\lambda}, \tilde{v}) \\
&= L(\tilde{x}, \tilde{\lambda}, \tilde{v}) \\
&= f_0(\tilde{x}) + \sum_{i=1}^m \tilde{\lambda}_i f_i(\tilde{x}) + \sum_{i=1}^p \tilde{v}_i h_i(\tilde{x}) \\
&= f_0(\tilde{x})
\end{aligned}$$

**例 25** (Waterfilling 算法). 共  $n$  个信道 (*channel*)

$$source \longleftrightarrow destination$$

$$\begin{aligned}
\min \quad & - \sum_{i=1}^n \log(\alpha_i + x_i) \\
\text{s.t.} \quad & x \geq 0 \\
& \mathbb{1}^T = 1
\end{aligned}$$

分析. *KKT* 条件

- $x^* \geq 0$
- $\mathbb{1}^T x^* = 1$
- $\lambda^* \geq 0$
- $x_i^* \lambda_i^* = 0, \forall i$

$$L(x, \lambda, v) = - \sum_{i=1}^n \log(\alpha_i + x_i) - \lambda^T x + v(\mathbb{1}^T x - 1)$$

$$\left( \frac{\partial L(x, \lambda, v)}{\partial x} \right)_i = - \frac{1}{\alpha_i + x_i} - \lambda_i + v$$

$$- \frac{1}{\alpha_i + x_i^*} - \lambda_i^* + v^* = 0, \forall i$$

$$\implies v^* \frac{1}{\alpha_i + x_i^*}, i = 1, \dots, n$$

$$x_i^* \left( v^* - \frac{1}{\alpha_i + x_i^*} \right) = 0, i = 1, \dots, n$$

$$\text{若 } v^* > \frac{1}{\alpha_i} \implies x_i^* = 0$$

$$\text{若 } v^* < \frac{1}{\alpha_i}$$

$$\frac{1}{\alpha_i} > v^* \geq \frac{1}{\alpha_i + x_i^*}$$

进而

$$\begin{aligned} x_i^* &> 0 \\ v^* &= \frac{1}{\alpha_i + x_i^*} \\ x_i^* &= \frac{1}{v^*} - \alpha_i \\ \implies x_i^* &= \max\{0, \frac{1}{v^*} - \alpha_i\} \end{aligned}$$

结合  $\sum_i x_i^* = 1$ , 即注水算法

Motivation: 误差, 调整参数测灵敏度

$$\begin{aligned} \min \quad & f_0(x) \\ \text{s.t.} \quad & f(x) \leq u_i, \quad i = 1, \dots, m \\ & h_i(x) = w_i, \quad i = 1, \dots, p \end{aligned}$$

新问题的最优解记为  $p^*(\mathbf{u}, \mathbf{w})$

性质: 若原始问题为凸, 则  $p^*(\mathbf{u}, \mathbf{w})$  是  $(u, w)$  的凸函数

布尔线性规划问题做松弛(relaxation)

$$x_i \in \{0, 1\} \implies 1 \geq x_i \geq 0$$

## 6 优化算法

### 6.1 简介

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

罚函数法

$$\begin{aligned} \min \quad & f_0(x) + \frac{\lambda}{2} \|\mathbf{Ax} - \mathbf{b}\|_2^2 \\ \tilde{x} &= \arg \min_x F \\ \nabla f_0(\tilde{x}) + \lambda \mathbf{A}^T (\mathbf{A}\tilde{x} - \mathbf{b}) &= 0 \end{aligned}$$

$$\begin{aligned}
L(x, v) &= f_0(x) + v^T(A\mathbf{x} - \mathbf{b}) \\
\Rightarrow g(v) &= \inf_x f_0(x) + v^T(A\mathbf{x} - \mathbf{b}) \\
v &= \lambda(A\tilde{\mathbf{x}} - \mathbf{b}) \\
\Rightarrow g(\lambda(A\tilde{\mathbf{x}} - \mathbf{b})) &= \inf_x f_0(x) + \lambda(A\tilde{\mathbf{x}} - \mathbf{b})^T(A\mathbf{x} - \mathbf{b}) \\
\nabla f_0(x) + \lambda A^T(A\tilde{\mathbf{x}} - \mathbf{b}) &= 0
\end{aligned}$$

$$\begin{aligned}
\min \quad & f_0(x) \\
\text{s.t.} \quad & A\mathbf{x} \geq \mathbf{b}
\end{aligned}$$

log-barrier

$$\min f_0(x) + \sum_{i=1}^m u_i \log(a_i^T \mathbf{x} - b_i)$$

$\min f_0(x)$ 可微, 凸, 无约束

1. 所有算法都是迭代的

$$x^{(k+1)} = x^{(k)} + \alpha^{(k)} d^{(k)}$$

$\alpha \geq 0$ 为步长,  $d$ 为方向, 所有算法本质上都是选择方向与步长的问题

2. 如何选择步长 $\alpha^{(k)}$

$$\left\{ \begin{array}{l} \text{确定步长} \\ \text{搜索步长} \end{array} \right\} \left\{ \begin{array}{l} \text{固定步长} \\ \text{变化步长 (递减步长)} \end{array} \right.$$

最优步长: 线搜索问题

$$\alpha^{(k)} = \arg \min_{\alpha \geq 0} f_0(x^{(k)} + \alpha d^{(k)})$$

3. 关键问题是选方向

黄金分割法(0.618法)/优选法求解线搜索问题: 这样做的采样复杂度很低, 之前算过的点很容易被再用!

不精确线搜索(Armijo Rule): 一阶泰勒展开

实际上没有必要要求最优步长, 在该方向上的差异并没有太大

## 6.2 梯度下降法

$$d^{(k)} = -\nabla f_0(x^{(k)})$$

- 能否收敛
- 收敛到哪里
- 收敛速度

假设

0. 基本假设:  $f$  为可微的凸函数,

$$x^* = \arg \min_x f_0(x)$$

存在且有限,  $f_0(x^*)$  有限

1. Lipschitz 连续梯度

$$\exists L \geq 0, \|\nabla f_0(x) - \nabla f_0(y)\| \leq L \|x - y\|, \forall x, y$$

等价定义:

a. 若  $f_0(x)$  二阶可微

$$\nabla^2 f_0(x) \preceq LI, \forall x$$

b. 下界

$$\langle \nabla f_0(x) - \nabla f_0(y), x - y \rangle \geq \frac{1}{L} \|\nabla f_0(x) - \nabla f_0(y)\|^2$$

c. 上界

$$\langle \nabla f_0(x) - \nabla f_0(y), x - y \rangle \leq L \|x - y\|^2$$

d. 当函数为凸时

$$0 \leq f_0(y) - f_0(x) - \langle \nabla f_0(x), y - x \rangle \leq \frac{L}{2} \|x - y\|^2$$

2. 强凸性(strong convexity)

$$\exists \mu > 0: f_0(y) \geq f_0(x) + \langle \nabla f_0(x), y - x \rangle + \frac{\mu}{2} \|x - y\|_2^2, \forall x, y$$

二阶可微情况下的等价定义

$$\nabla^2 f(x) \succeq \mu I$$

例 26.

$$\begin{array}{lll} f_0(x) = \mathbb{1}^T x & L = 0 & \times \\ f_0(x) = \frac{1}{2} \|x\|_2^2 & L = 1 & \mu = 1 \\ f_0(x) = \frac{1}{4} \|x\|_2^4 & \times & \times \end{array}$$

区别于严格凸(strictly convex), 强凸一定是严格凸

定理 6. 严格凸函数只有一个最小值点

分析. 反证法, 假设  $x, y$  均为最小值点, 且  $x \neq y$

$$f_0(y) > f_0(x) + \langle \nabla f_0(x), x - y \rangle = f_0(x)$$

定理 7. 若  $f_0(x)$  有 *Lipschitz* 连续梯度, 常数  $L$ , 若  $\alpha \in (0, \frac{2}{L})$ , 则有

$$f_0(x^{(k)}) - f_0(x^*) \leq \frac{2(f_0(x^0) - f_0(x^*)) \|x^0 - x^*\|^2}{2\|x^0 - x^*\|^2 + k\alpha(2 - L\alpha)(f_0(x^0) - f_0(x^*))}, \forall x^*$$

即以  $O(\frac{1}{k})$  速度收敛

分析. 1° 点的单调性: 与任意  $x^*$  的距离在缩小

$$\|x^{(k+1)} - x^*\|^2 \leq \|x^{(k)} - x^*\|^2, \forall x^*$$

$$\begin{aligned} LHS &= \|x^{(k)} - x^* - \alpha \nabla f_0(x^k)\|^2 \\ &= \|x^{(k)} - x^*\|^2 - 2\alpha \langle x^k - x^*, \nabla f_0(x^k) \rangle + \alpha^2 \|\nabla f_0(x^k)\|^2 \\ &\leq \|x^{(k)} - x^*\|^2 + \alpha(\alpha - \frac{2}{L}) \|\nabla f_0(x^k)\|^2 \quad \text{注意到 } \nabla f_0(x^*) = 0, \text{ 利用 } Lipschitz \text{ 连续梯度} \\ &\leq \|x^{(k)} - x^*\|^2 \end{aligned}$$

2° 函数值的单调性:  $f_0(x^{(k+1)}) \leq f_0(x^{(k)})$  (注意下降可能非常缓慢, 并不一定收敛)

$$\begin{aligned} f_0(x^{(k+1)}) &\leq f_0(x^{(k)}) + \langle \nabla f_0(x^{(k)}), x^{(k+1)} - x^{(k)} \rangle + \frac{L}{2} \|x^{(k+1)} - x^{(k)}\|^2 \\ &= f_0(x^{(k)}) - \alpha(1 - \frac{L\alpha}{2}) \|\nabla f_0(x^{(k)})\|^2 \\ &\leq f_0(x^{(k)}) \end{aligned}$$

3° 函数值的充分下降 (即证明收敛性)

$$\begin{aligned} f_0(x^{(k+1)}) - f_0(x^*) &\leq f_0(x^{(k)}) - f_0(x^*) - \omega \|\nabla f_0(x^{(k)})\|^2 \\ f_0(x^{(k)}) - f_0(x^*) &\leq \langle f_0(x^{(k)}), x^{(k)} - x^* \rangle \\ &= \langle \nabla f_0(x^{(k)}) - \nabla f_0(x^*), x^{(k)} - x^* \rangle \\ &\leq \|\nabla f_0(x^{(k)}) - \nabla f_0(x^*)\| \|x^{(k)} - x^*\| \\ &\leq \|\nabla f_0(x^{(k)})\| \|x^{(k)} - x^*\| \\ \Delta^{(k+1)} &\leq \Delta^{(k)} - \frac{\omega}{\|x^0 - x^*\|^2} (\Delta^{(k)})^2 \\ \frac{1}{\Delta^{(k+1)}} &\leq \frac{1}{\Delta^{(k)}} - \frac{\omega}{\|x^0 - x^*\|^2} \frac{\Delta^{(k)}}{\Delta^{(k+1)}} \end{aligned}$$

错位相消可得结论  $O(\frac{1}{k})$  收敛速度

定理 8. 若  $f_0$  有 *Lipschitz* 连续梯度, 常数  $L$ , 强凸函数  $n$ , 步长  $\alpha \in (0, \frac{2}{\mu+L}]$ , 则

$$\|x^{(k)} - x^*\|^2 \leq \left(1 - \frac{2\alpha\mu L}{\mu + L}\right)^k \|x^{(0)} - x^*\|^2$$

分析.

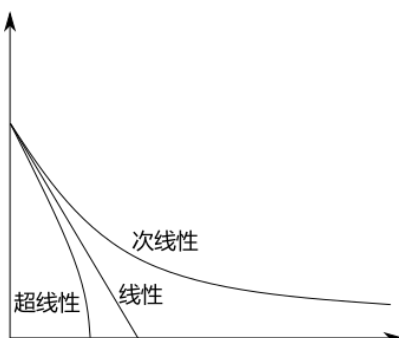
$$\begin{aligned}
 \|x^{(k)} - x^*\|^2 &= \|x^{(k)} - \alpha \nabla f_0(x^{(k)}) - x^*\|^2 \\
 &= \|x^{(k)} - x^*\|^2 - 2\alpha \langle x^{(k)} - x^*, \nabla f_0(x^{(k)}) \rangle + \alpha^2 \|\nabla f_0(x^{(k)})\|^2 \\
 &\leq \|x^{(k)} - x^*\|^2 - \frac{2\alpha}{\mu + L} \|\nabla f_0(x^{(k)})\|^2 + \alpha^2 \|\nabla f_0(x^{(k)})\|^2 \quad \text{内积不等式} \\
 &\leq RHS
 \end{aligned}$$

$$1 - \frac{4\mu L}{(\mu + L)^2} = \frac{(L - \mu)^2}{(L + \mu)^2} = \frac{\left(\frac{L}{\mu} - 1\right)^2}{\left(\frac{L}{\mu} + 1\right)^2}$$

$L$ 为Hessian矩阵的最大特征值,  $\mu$ 为Hessian矩阵的最小特征值, 则 $\frac{L}{\mu}$ 为该矩阵的条件数

不同收敛速度

- 次线性收敛
- 线性收敛
- 超线性收敛



例 27.

$$f_0(x) = \frac{1}{2} \|Ax - b\|_2^2$$

分析.

$$x^{(0)} \rightarrow x^{(1)}$$

$$x^{(1)} = x^{(0)} - \alpha(x^{(0)} - b) = b$$

条件数糟糕的病态矩阵收敛速度是非常糟糕的, 会出现zig-zag的情况

$$A = \begin{bmatrix} 1 & 0 \\ 0 & 10^{-4} \end{bmatrix}$$

可以通过预处理(precondition)来解决条件数糟糕的问题

$f_0$ , Lipschitz连续梯度( $L$ ), 强凸( $\mu$ ), 函数值收敛性

$$\begin{aligned}\tilde{f}_0(\alpha^{(k)}) &= f_0(x^{(k+1)}) = f_0(x^{(k)} - \alpha^{(k)} \nabla f_0(x^{(k)})) \\ &\leq f_0(x^{(k)}) + \langle \nabla f_0(x^{(k)}), -\alpha^{(k)} \nabla f_0(x^{(k)}) \rangle + \frac{L}{2} \left\| -\alpha^{(k)} \nabla f_0(x^{(k)}) \right\|^2 \\ &= f_0(x^{(k)}) + \frac{L(\alpha^{(k)})^2 - 2\alpha^{(k)}}{2} \left\| \nabla f_0(x^{(k)}) \right\|^2\end{aligned}$$

$\alpha^{(k)} = \alpha_{exact}^{(k)}$  精确线搜索

$$\begin{aligned}\tilde{f}_0\left(\frac{1}{L}\right) &= f_0(x^{(k)}) - \frac{1}{2L} \left\| \nabla f_0(x^{(k)}) \right\|^2 \\ \tilde{f}_0(\alpha_{exact}^{(k)}) &\leq \tilde{f}_0\left(\frac{1}{L}\right) \\ \implies \tilde{f}_0(\alpha_{exact}^{(k)}) - f_0(x^{(k)}) - f_0(x^*) &= -\frac{1}{2L} \left\| \nabla f_0(x^{(k)}) \right\|^2 \\ &\leq \left(1 - \frac{\mu}{L}\right)(f_0(x^{(k)}) + f_0(x^*))\end{aligned}$$

$$\begin{aligned}f_0(x^*) &\geq f_0(x^{(k)}) + \langle \nabla f_0(x^{(k)}), x^* - x^{(k)} \rangle + \frac{\mu}{2} \left\| x^{(k)} - x^* \right\|^2 \\ &\geq f_0(x^{(k)}) - \frac{\mu}{2} \left\| x^{(k)} - x^* \right\|^2 - \frac{1}{2\mu} \left\| \nabla f_0(x^{(k)}) \right\|^2 + \frac{\mu}{2} \left\| x^{(k)} - x^* \right\|^2 \quad ab \geq -\frac{\mu}{2}a^2 - \frac{1}{2\mu}b^2 \\ &= f_0(x^{(k)}) - \frac{1}{2\mu} \left\| \nabla f_0(x^{(k)}) \right\|^2\end{aligned}$$

$$f_0(x^{(k)}) - f_0(x^*) \leq \frac{1}{2\mu} \left\| \nabla f_0(x^{(k)}) \right\|^2$$

Armijo Rule

$$\tilde{f}_0(\alpha^{(k)}) = f_0(x^{(k+1)}) \leq f_0(x^{(k)}) + \frac{L(\alpha^{(k)})^2 - 2\alpha^{(k)}}{2} \left\| \nabla f_0(x^{(k)}) \right\|^2$$

$$\tilde{f}_0(\alpha^{(k)}) = f_0(x^{(k+1)}) \leq f_0(x^{(k)}) - \gamma \alpha^{(k)} \left\| \nabla f_0(x^{(k)}) \right\|^2$$

首先说明, 若  $0 \leq \alpha^{(k)} \leq \frac{1}{L}$  时, 则

$$\tilde{f}_0(\alpha^{(k)}) \leq f_0(x^{(k)}) - \gamma \alpha^{(k)} \left\| \nabla f_0(x^{(k)}) \right\|^2$$

当  $\alpha^{(k)} \in [0, \frac{1}{2}]$  时,

$$-\alpha^{(k)} + \frac{L}{2}(\alpha^{(k)})^2 \leq -\frac{\alpha^{(k)}}{2} \iff \frac{L}{2}(\alpha^{(k)})^2 \leq \frac{\alpha^{(k)}}{2} \iff L \cdot \alpha^{(k)} \leq 1$$

$$\begin{aligned}f_0(x^{(k+1)}) &\leq f_0(x^{(k)}) + \frac{L(\alpha^{(k)})^2 - 2\alpha^{(k)}}{2} \left\| \nabla f_0(x^{(k)}) \right\|^2 \\ &\leq f_0(x^{(k)}) - \frac{\alpha^{(k)}}{2} \left\| \nabla f_0(x^{(k)}) \right\|^2 \\ &\leq f_0(x^{(k)}) - \gamma \alpha^{(k)} \left\| \nabla f_0(x^{(k)}) \right\|^2\end{aligned}$$



$$f_0(x^{(k+1)}) \leq f_0(x^{(k)}) - \min\{\gamma\alpha_{\max}, \frac{\gamma\beta}{L}\} \|\nabla f_0(x^{(k)})\|^2$$

$$\implies f_0(x^{(k+1)}) - f_0(x^*) \leq \left(1 - \min\{2\mu\gamma\alpha_{\max}, \frac{2\mu\gamma\beta}{L}\}\right) (f_0(x^{(k)}) - f_0(x^*))$$

梯度下降法的解释1

$$x^{(k+1)} = x^{(k)} - \alpha^{(k)} \nabla f_0(x^{(k)})$$

将 $f_0$ 在 $x^{(k)}$ 处进行一阶Taylor展开

$$f_0(x) \approx f_0(x^{(k)}) + \langle \nabla f_0(x^{(k)}), x - x^{(k)} \rangle + \frac{1}{2\alpha^{(k)}} \|x - x^{(k)}\|^2$$

求梯度

$$\nabla f_0(x^{(k)}) + \frac{1}{\alpha^{(k)}} (x - x^{(k)}) = 0$$

$$\alpha^{(k)} \nabla f_0(x^{(k)}) + x - x^{(k)} = 0$$

$$x = x^{(k)} - \alpha^{(k)} \nabla f_0(x^{(k)})$$

解释2

$$f_0(x^{(k)} + v) \approx f_0(x^{(k)}) + \langle \nabla f_0(x^{(k)}), v \rangle$$

$$d^{(k)} = \arg \min_v \{ \langle \nabla f_0(x^{(k)}), v \rangle \mid \|v\| = 1 \}$$

若采用2-范数，可得标准化的负梯度方向(normalized negative gradient)

$$d^{(k)} = \frac{-\nabla f_0(x^{(k)})}{\|\nabla f_0(x^{(k)})\|_2}$$

通过改变不同的范数，有不同的特性

坐标下降法(coordinate descent/alternating direction)交替极小化

$$d^{(k)} = \mathbf{e}_{\text{mod}(k, n)}$$

注意，这里 $x \in \mathbb{R}^n$ ,  $n \bmod n = 0$

$$\alpha^{(k)} = \arg \min f_0(x^{(k)} + \alpha d^{(k)}), \alpha_{\max} \geq \alpha \geq \alpha_{\min}$$

### 6.3 非光滑优化问题

$$\min f_0(x), \quad f_0 \text{连续, 凸, 不可微}$$

梯度下降法→次梯度(subgradient)法 $g_0(x) \in \partial f_0(x)$  (注意凹函数则对应的是supgradient)

$$f_0(y) \geq f_0(x) + \langle g_0(x), y - x \rangle, \forall y$$

$f(x) = |x|$ 在零点处次梯度为 $[-1, 1]$

$$x^{(k+1)} = x^{(k)} - \alpha^{(k)} g_0(x^{(k)})$$

只要有 $0 \in \partial f_0(x_0)$ 就有最优解 $x = x_0$

如果激活函数为非光滑的（如ReLU），那么出来的函数也是非光滑的，就要用次梯度  
关键在于选择步长

- 固定步长 $\alpha^{(k)} = \alpha$
- 不可加但平方可加

$$\sum_{k=0}^{\infty} (\alpha^{(k)})^2 < \infty \quad \sum_{k=0}^{\infty} \alpha^{(k)} = \infty$$

- 不可加递减

$$\sum_{k=0}^{\infty} \alpha^{(k)} = 0 \quad \lim_{k \rightarrow \infty} \alpha^{(k)} \rightarrow 0$$