

最优化算法作业 3

陈文轩

更新: April 14, 2025

1. • 若 $f(x)$ 是二阶连续可微, 证明 $\nabla f(x)$ 是 L-Lipschitz 连续等价于 $LI \succeq \nabla^2 f(x) \succeq -LI$ 。

解

$$\Rightarrow: \forall x, v \in \mathbb{R}^n, \|\nabla^2 f(x)v\| = \lim_{t \rightarrow 0} \frac{\|\nabla f(x+tv) - \nabla f(x)\|}{t} \leq \lim_{t \rightarrow 0} \frac{L\|tv\|}{t} = L\|v\|,$$

故 $\|\nabla^2 f(x)\| \leq L$, 又 $\nabla^2 f(x)$ 是对称矩阵, 故 $\rho(\nabla^2 f(x)) = \|\nabla^2 f(x)\|_2 \leq L$,

即 $\nabla^2 f(x)$ 的特征值模长均不大于 L , 即 $LI \succeq \nabla^2 f(x) \succeq -LI$ 。

$$\Leftarrow: \forall x, y \in \mathbb{R}^n, \nabla f(y) - \nabla f(x) = \int_0^1 \nabla^2 f(x+t(y-x))(y-x) dt, \|\nabla^2 f(x)\| \leq L,$$

$$\|\nabla f(y) - \nabla f(x)\| \leq \int_0^1 \|\nabla^2 f(x+t(y-x))\| \cdot \|y-x\| dt \leq \int_0^1 L\|y-x\| dt = L\|y-x\|$$

即 $\nabla f(x)$ 是 L-Lipschitz 连续的。

- 估计逻辑回归函数的梯度的 Lipschitz 常数:

$$\min_x l(x) := \frac{1}{m} \sum_{i=1}^m \log(1 + \exp(-b_i a_i^\top x))$$

其中 $b_i \in \{-1, 1\}, a_i \in \mathbb{R}^n, i = 1, \dots, m$ 是给定的数据。

解

$$\text{记 } f_i(x) = \log(1 + e^{-b_i a_i^\top x}), \sigma(x) = \frac{1}{1 + e^{-x}}, \nabla f_i(x) = \frac{-b_i a_i}{1 + e^{-b_i a_i^\top x}} = -b_i a_i \sigma(-b_i a_i^\top x)。$$

记 $\sigma_i = \sigma(-b_i a_i^\top x) \in [0, 1]$, 则 $\nabla^2 f_i(x) = -b_i a_i (-b_i a_i^\top) \sigma_i (1 - \sigma_i) = a_i a_i^\top \sigma_i (1 - \sigma_i)$,

$$\|\nabla^2 l(x)\| = \left\| \frac{1}{m} \sum_{i=1}^m \nabla^2 f_i(x) \right\| = \frac{1}{m} \sum_{i=1}^m \|a_i a_i^\top \sigma_i (1 - \sigma_i)\| \leq \frac{1}{4m} \sum_{i=1}^m \|a_i\|^2$$

故梯度的一个 Lipschitz 常数为 $\frac{1}{4m} \sum_{i=1}^m \|a_i\|^2$ 。

2. 对于次梯度算法, 请构造一个非光滑函数例子, 说明常数步长不收敛。

解

$$\text{取 } f(x) = \|x\|, \text{ 则 } \partial f(x) = \begin{cases} \{\text{sgn}(x)\}, & x \neq 0 \\ [-1, 1], & x = 0 \end{cases}。 \text{ 对任意步长 } \alpha, \text{ 取初值 } x_0 = \frac{\alpha}{2},$$

则迭代序列为 $x_n = (-1)^n \frac{\alpha}{2}$ 不收敛。

3. 计算下面函数的邻近点映射，即 $\text{prox}_h(x) = \arg \min_y \left(h(y) + \frac{1}{2} \|y - x\|^2 \right)$

- $h(x) = \|x\|_\infty$ (需要求解一个一维子问题)。

解

对 $h(x) = \|x\|_\infty$, $h^*(z) = \sup_x (z^\top x - h(x)) = I_{\{z \mid \|z\|_1 \leq 1\}}(z)$ 。由 $x = \text{prox}_h(x) + \text{prox}_{h^*}(x)$

考虑 $\text{prox}_{h^*}(x) = \arg \min_z \left(I_{\{z \mid \|z\|_1 \leq 1\}}(z) + \frac{1}{2} \|z - x\|^2 \right)$ 是 x 关于 L^1 范数球的投影。

$$\text{而对 } C = \{z \mid \|z\|_1 \leq 1\}, P_C(x)_k = \begin{cases} x_k - \lambda, & x_k > \lambda \\ 0, & -\lambda \leq x_k \leq \lambda \\ x_k + \lambda, & x_k < -\lambda \end{cases}$$

其中对 $\|x\| \leq 1, \lambda = 0$, 否则 λ 是 $\sum_{k=1}^n \max\{|x_k| - \lambda, 0\} = 1$ 的解。

综上所述, $\text{prox}_h(x) = x - P_C(x)$, $P_C(x)$ 如上定义。

- $h(x) = \max\{0, \|x\|_2 - 1\}$ 。

解

对 $\|y\|_2 \leq 1, h(y) = 0$, 此时原问题化为 x 关于 L^1 范数球的投影问题。此时有

$$\text{prox}_h(x) = \arg \min_{\|y\|_2 \leq 1} \frac{1}{2} \|y - x\|^2 = \begin{cases} x, & \|x\|_2 \leq 1 \\ \frac{x}{\|x\|_2}, & \|x\|_2 > 1 \end{cases}.$$

对 $\|y\|_2 > 1$, 问题化为 $\arg \min_{\|y\|_2 > 1} \left(\|y\| - 1 + \frac{1}{2} \|y - x\|^2 \right)$, 由对称性, 取 x, y 共线。

$$\begin{aligned} \arg \min_{\|y\|_2 > 1} \left(\|y\| - 1 + \frac{1}{2} \|y - x\|^2 \right) & \stackrel{u = \frac{x}{\|x\|_2}}{=} \arg \min_{t > 1} \left(t - 1 + \frac{1}{2} (t - \|x\|_2)^2 \right) \\ & = \frac{x}{\|x\|_2} \arg \min_{t > 1} \left(\frac{1}{2} t^2 + (1 - \|x\|_2)t + \frac{1}{4} \|x\|_2^2 - 1 \right) = \begin{cases} \frac{x}{\|x\|_2}, & \|x\|_2 < 2 \\ \frac{x}{\|x\|_2} (\|x\|_2 - 1), & \|x\|_2 \geq 2 \end{cases} \end{aligned}$$

比较两种情况下 $f(y) = \max\{0, \|y\|_2 - 1\} + \frac{1}{2} \|y - x\|^2$ 的函数值:

对 $\|x\|_2 < 1, f(y_1) = 0, f(y_2) = \frac{1}{2} (\|x\|_2 - 1)^2$, 故使用第一种情况的解 $y = x$ 。

对 $1 < \|x\|_2 < 2$, 两种情况下解都为 $y = \frac{x}{\|x\|_2}$, 即为最终答案。

对 $\|x\|_2 > 2, f(y_1) = \frac{1}{2} (\|x\|_2 - 1)^2, f(y_2) = \|x\|_2 - \frac{3}{2}, f(y_1) - f(y_2) = \frac{1}{2} (\|x\|_2 - 2)^2 \geq 0$, 故使用第二种情况的解 $y = \frac{x}{\|x\|_2} (\|x\|_2 - 1)$ 。

$$\text{综上所述, } \text{prox}_h(x) = \begin{cases} x, & \|x\|_2 \leq 1 \\ \frac{x}{\|x\|_2}, & 1 < \|x\|_2 \leq 2 \\ \frac{x}{\|x\|_2} (\|x\|_2 - 1), & \|x\|_2 > 2 \end{cases}$$

4. 考虑 D-最优实验设计 (D-optimal experimental design), 其目标是最大化估计量的信息内容, 通过差分香农熵测量, 具体到最大化 $\det V(m_1, \dots, m_n)$, 具体背景参考《convex optimization: 7.5 节》。

该问题需要求解下述约束问题:

$$\min_{x \in \Delta_n} -\log \det V(x)$$

其中 $V(x) = \sum_{i=1}^n x_i a_i a_i^\top$, $a_i \in \mathbb{R}^d$, $i = 1, \dots, d$ 是给定的数据, $\Delta_n = \left\{ x : \sum_{i=1}^n x_i = 1, x \geq 0 \right\}$

请使用条件梯度法求解该问题, 写出迭代公式, 并且给出子问题的解。

解
 $\frac{\partial}{\partial x_i} -\log \det V(x) = -\text{tr} \left(\nabla_V \log \det V(x) \frac{\partial V(x)}{\partial x_i} \right) = -\text{tr}(V(x)^{-1} a_i a_i^\top) = -a_i^\top V(x)^{-1} a_i$
 $\Rightarrow \nabla(-\log \det V(x)) = \left(-a_1^\top V(x)^{-1} a_1, \dots, -a_n^\top V(x)^{-1} a_n \right)^\top$, 记 $f(x) = -\log \det V(x)$,
 需要求解子问题 $x_k = \arg \min_{x \in \Delta_n} \langle \nabla f(y_{k-1}), x \rangle$, 这是一个线性问题, 最优解在顶点 e_j 上取。

其中 $j = \arg \min_{x \in \Delta_n} \langle \nabla f(y_{k-1}), e_i \rangle = \arg \min_{i \in \{1, \dots, n\}} (-a_i^\top V(y_{k-1}) a_i) = \arg \max_{i \in \{1, \dots, n\}} a_i^\top V(y_{k-1}) a_i$ 。

$y_k = (1 - \alpha_k) y_{k-1} + \alpha_k e_j$, α_k 取消失步长或通过精确线搜索得到。

5. 求解问题

$$\min f(x) \quad \text{s.t. } x \in \Delta$$

其中, $\Delta_n = \left\{ x \in \mathbb{R}^n : \sum_{i=1}^n x_i = 1, x_n \geq 0 \right\}$ 。使用镜像梯度法, 迭代公式为

$$x^{k+1} = \arg \min_{x \in \Delta} \left(\nabla f(x^k)^\top (x - x^k) + \frac{1}{\alpha_k} \sum_{i=1}^n x_i \log \frac{x_i}{x_i^k} \right)$$

证明:

$$x_i^{k+1} = \frac{x_i^k \exp(-\alpha_k \nabla f(x^k)_i)}{\sum_{j=1}^n x_j^k \exp(-\alpha_k \nabla f(x^k)_j)}$$

解

令 $L(x, \lambda) = \nabla f(x^k)^\top (x - x^k) + \frac{1}{\alpha_k} \sum_{i=1}^n x_i \log \frac{x_i}{x_i^k} + \lambda \left(\sum_{i=1}^n x_i - 1 \right)$, 对 x 求导,

$$\frac{\partial L(x, \lambda)}{\partial x_i} = \nabla f(x^k)_i + \frac{1}{\alpha_k} \left(\log \frac{x_i}{x_i^k} + 1 \right) + \lambda, \quad \text{求解 } \frac{\partial L(x, \lambda)}{\partial x_i} = 0, \text{ 得到}$$

$$x_i = x_i^k \exp \left(-\alpha_k \nabla f(x^k)_i - \alpha_k \lambda - 1 \right) \stackrel{C = \exp(-\alpha_k \lambda - 1)}{=} C x_i^k \exp \left(-\alpha_k \nabla f(x^k)_i \right) \geq 0.$$

需要调整 λ 使得 C 满足 $\sum_{i=1}^n x_i = 1 \Rightarrow C \sum_{j=1}^n x_j^k \exp \left(-\alpha_k \nabla f(x^k)_j \right) = 1$

故 $C = \frac{1}{\sum_{j=1}^n x_j^k \exp(-\alpha_k \nabla f(x^k)_j)}$, $x_i = \frac{x_i^k \exp(-\alpha_k \nabla f(x^k)_i)}{\sum_{j=1}^n x_j^k \exp(-\alpha_k \nabla f(x^k)_j)}$ 即为所求。