

机器学习中的优化算法

Lecture09: 复合优化算法-分块坐标下降法

张立平

清华大学数学科学系

办公室：理科楼#A302, Tel: 62798531

E-mail: lipingzhang@tsinghua.edu.cn

Contents and Acknowledgement

- 教材：最优化：建模、算法与理论

<http://bicmr.pku.edu.cn/wenzw/bigdata2021.html>

- 致谢：北京大学文再文教授

Outline of BCD

- 问题描述
- 分块坐标下降法
- 应用举例
- 收敛性分析

BCD针对的问题描述

► 典型优化问题形式:

$$\min_{x \in \mathcal{X}} F(x_1, x_2, \dots, x_s) = f(x_1, x_2, \dots, x_s) + \sum_{i=1}^s r_i(x_i),$$

- \mathcal{X} 是函数的可行域, 自变量 x 拆分成 s 个变量块 x_1, x_2, \dots, x_s , 每个变量块 $x_i \in \mathbb{R}^{n_i}$.
- 函数 f 是关于 x 的可微函数, 每个 $r_i(x_i)$ 关于 x_i 是适当的闭凸函数, 但不一定可微.
- 目标函数 F 的性质体现在 f , 每个 r_i 以及自变量的分块上. 通常情况下, f 对于所有变量块 x_i 不可分, 但单独考虑每一块自变量时, f 有简单结构; r_i 只和第 i 个自变量块有关, 因此 r_i 在目标函数中是一个可分项.
- 求解该问题的难点在于如何利用分块结构处理不可分的函数 f .

典型问题举例

- **分组LASSO 模型** 参数 $x = (x_1, x_2, \dots, x_G) \in \mathbb{R}^p$ 可以分成 G 组, 且 $\{x_i\}_{i=1}^G$ 中只有少数的非零向量.

$$\min_x \frac{1}{2n} \|b - Ax\|_2^2 + \lambda \sum_{i=1}^G \sqrt{p_i} \|x_i\|_2.$$

- **K -均值聚类问题**

$$\begin{aligned} \min_{\Phi, H} \quad & \|A - \Phi H\|_F^2, \\ \text{s.t.} \quad & \Phi \in \mathbb{R}^{n \times k}, \text{ 每一行只有一个元素为1, 其余为0,} \\ & H \in \mathbb{R}^{k \times p}. \end{aligned}$$

- **低秩矩阵恢复** 设 $b \in \mathbb{R}^m$ 是已知的观测向量, \mathcal{A} 是线性映射.

$$\min_{X, Y} \frac{1}{2} \|\mathcal{A}(XY) - b\|_2^2 + \alpha \|X\|_F^2 + \beta \|Y\|_F^2, \quad (\alpha, \beta > 0 \text{ 为正则化参数}).$$

- 非负矩阵分解

$$\min_{A_1, A_2, \dots, A_N \geq 0} \frac{1}{2} \|\mathcal{M} - A_1 \circ A_2 \circ \dots \circ A_N\|_F^2 + \sum_{i=1}^N \lambda_i r_i(A_i),$$

其中 \mathcal{M} 是已知张量, “ \circ ”表示张量的外积运算.

- 字典学习 设 $A \in \mathbb{R}^{m \times n}$ 为 n 个观测, 每个观测的信号维数是 m , 现在我们要从 A 中学习出一个字典 $D \in \mathbb{R}^{m \times k}$ 和系数矩阵 $X \in \mathbb{R}^{k \times n}$:

$$\begin{aligned} \min_{D, X} \quad & \frac{1}{2n} \|DX - A\|_F^2 + \lambda \|X\|_1, \\ \text{s.t.} \quad & \|D\|_F \leq 1. \end{aligned}$$

挑战和难点

- 函数 f 关于变量全体一般是非凸的，这使得问题求解具有挑战性
- 应用在非凸问题上的算法收敛性不易分析，很多针对凸问题设计的算法通常会失效
- 目标函数的整体结构十分复杂，变量的更新需要很大计算量
- **目标**: 发展一种更新方式简单且有全局收敛性（收敛到稳定点）的有效算法

变量划分

- **分块坐标下降法更新方式**: 按照 x_1, x_2, \dots, x_s 的次序依次固定其他 $(s - 1)$ 块变量极小化 F , 完成一块变量的极小化后, 它的值便立即被更新到变量空间中, 更新下一块变量时将使用每个变量最新的值.
- 变量划分

$$\mathcal{X}_i^k = \{x \in \mathbb{R}^{n_i} \mid (x_1^k, \dots, x_{i-1}^k, x, x_{i+1}^{k-1}, \dots, x_s^{k-1}) \in \mathcal{X}\}.$$

- 辅助函数

$$f_i^k(x_i) = f(x_1^k, \dots, x_{i-1}^k, x_i, x_{i+1}^{k-1}, \dots, x_s^{k-1}),$$

其中 x_j^k 表示在第 k 次迭代中第 j 块自变量的值, 函数 f_i^k 表示在第 k 次迭代更新第 i 块变量时所需要考虑的目标函数的光滑部分. 考虑第 i 块变量时前 $(i - 1)$ 块变量已经完成更新, 因此上标为 k ; 而后面下标从 $(i + 1)$ 起的变量仍为旧的值, 因此上标为 $(k - 1)$.

变量更新方式

- 固定其他分量然后对单一变量求极小:

$$x_i^k = \arg \min_{x_i \in \mathcal{X}_i^k} \left\{ f_i^k(x_i) + r_i(x_i) \right\}. \quad (1)$$

- 增加了一个近似点项 $\frac{L_i^{k-1}}{2} \|x_i - x_i^{k-1}\|_2^2$ 来限制下一步迭代不应该与当前位置相距过远, 增加近似点项的作用是使得算法能够收敛($L_i^k > 0$ 为常数):

$$x_i^k = \arg \min_{x_i \in \mathcal{X}_i^k} \left\{ f_i^k(x_i) + \frac{L_i^{k-1}}{2} \|x_i - x_i^{k-1}\|_2^2 + r_i(x_i) \right\}. \quad (2)$$

- 对 $f_i^k(x)$ 进行线性化以简化子问题的求解, 并引入了 Nesterov 加速算法的技巧加快收敛:

$$x_i^k = \arg \min_{x_i \in \mathcal{X}_i^k} \left\{ \langle \hat{g}_i^k, x_i - \hat{x}_i^{k-1} \rangle + \frac{L_i^{k-1}}{2} \|x_i - \hat{x}_i^{k-1}\|_2^2 + r_i(x_i) \right\}, \quad (3)$$

其中 \hat{x}_i^{k-1} 采用外推定义:

$$\hat{x}_i^{k-1} = x_i^{k-1} + \omega_i^{k-1}(x_i^{k-1} - x_i^{k-2}), \quad (4)$$

$\omega_i^k \geq 0$ 为外推的权重, $\hat{g}_i^k \stackrel{\text{def}}{=} \nabla f_i^k(\hat{x}_i^{k-1})$ 为外推点处的梯度. 取权重 $\omega_i^k = 0$ 即可得到不带外推的更新格式, 此时等价于进行一次近似点梯度法的更新.

►分块坐标下降法

- 1: 初始化: 选择两组初始点 $(x_1^{-1}, x_2^{-1}, \dots, x_s^{-1}) = (x_1^0, x_2^0, \dots, x_s^0)$.
- 2: **for** $k = 1, 2, \dots$ **do**
- 3: **for** $i = 1, 2, \dots$ **do**
- 4: 使用格式(1) 或(2) 或(3) 更新 x_i^k .
- 5: **end for**
- 6: **if** 满足停机条件 **then**
- 7: 返回 $(x_1^k, x_2^k, \dots, x_s^k)$, 算法终止.
- 8: **end if**
- 9: **end for**

BCD算法格式解释

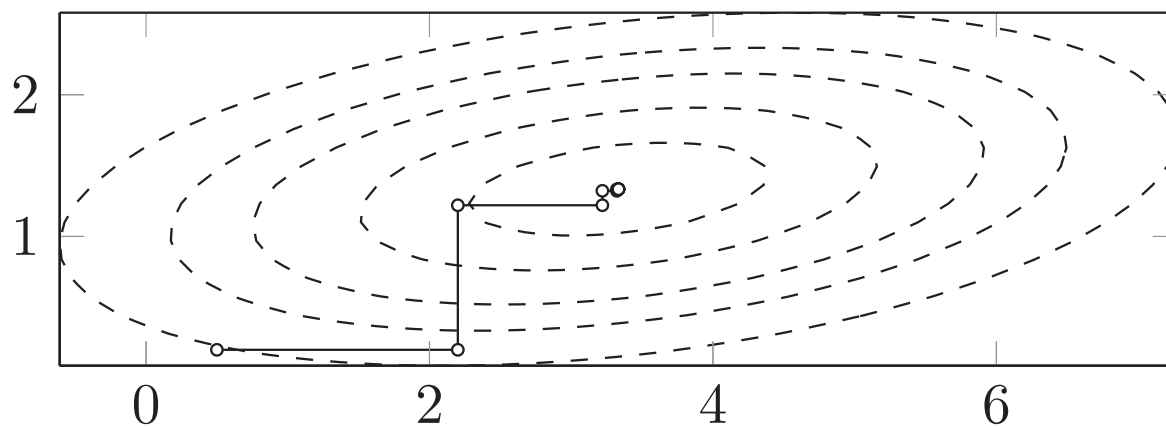
- BCD算法的子问题可采用三种不同的更新格式，这三种格式可能会产生不同的迭代序列，可能会收敛到不同的解，坐标下降算法的数值表现也不相同.
- 格式(1)是最直接的更新方式，它严格保证了整个迭代过程的目标函数值是下降的. 然而由于 f 的形式复杂，子问题求解难度较大. 在收敛性方面，格式(1)在强凸问题上可保证目标函数收敛到极小值，但在非凸问题上不一定收敛.
- 格式(2) (3) 则是对格式(1)的修正，不保证迭代过程目标函数的单调性，但可以改善收敛性结果. 使用格式(2)可使得算法收敛性在函数 F 为非严格凸时有所改善.
- 格式(3)实质上为目标函数的一阶泰勒展开近似，在一些测试问题上有更好的表现，可能的原因是使用一阶近似可以避开一些局部极小值点. 此外，格式(3)的计算量很小，比较容易实现.

► 【例】 $\min f(x, y) = x^2 - 2xy + 10y^2 - 4x - 20y$

- 采用格式(1)的分块坐标下降法:

$$x^{k+1} = 2 + y^k, \quad y^{k+1} = 1 + \frac{x^{k+1}}{10}.$$

- 下图描绘了当初始点为 $(x, y) = (0.5, 0.2)$ 时的迭代点轨迹, 可以看到在进行了7次迭代后迭代点与最优解已充分接近.
- 对于比较病态的问题, 由于分块坐标下降法是对逐个分量处理, 它能较好地捕捉目标函数的各向异性, 而梯度法则会受到很大影响.



► **不收敛反例** 对于非凸函数 $f(x)$, 分块坐标下降法可能失效. Powell 在1973年就给出了一个使用格式(1)但不收敛的例子:

$$F(x_1, x_2, x_3) = -x_1x_2 - x_2x_3 - x_3x_1 + \sum_{i=1}^3 [(x_i - 1)_+^2 + (-x_i - 1)_+^2],$$

设 $\varepsilon > 0$, 初始点取为:

$$x^0 = \left(-1 - \varepsilon, 1 + \frac{\varepsilon}{2}, -1 - \frac{\varepsilon}{4}\right),$$

容易验证迭代序列满足

$$x^k = (-1)^k \cdot (-1, 1, -1) + \left(-\frac{1}{8}\right)^k \cdot \left(-\varepsilon, \frac{\varepsilon}{2}, -\frac{\varepsilon}{4}\right),$$

这个迭代序列有两个聚点 $(-1, 1, -1)$ 与 $(1, -1, 1)$, 但这两个点都不是 F 的稳定点.

BCD应用举例

► **LASSO问题** $\min_x \mu \|x\|_1 + \frac{1}{2} \|Ax - b\|^2$

- 将自变量 x 记为 $x = [x_i, \bar{x}_i^\top]^\top$, 其中 \bar{x}_i 为 x 去掉第 i 个分量而形成的列向量. 相应地, 矩阵 A 在第 i 块的更新记为 $A = [a_i \ \bar{A}_i]$, 其中 \bar{A}_i 为矩阵 A 去掉第 i 列而形成的矩阵. LASSO问题可以写为

$$\min_{x_i} \mu |x_i| + \mu \|\bar{x}_i\|_1 + \frac{1}{2} \|a_i x_i - (b - \bar{A}_i \bar{x}_i)\|^2.$$

- 利用格式(1)更新第 i 块, 令 $c_i = b - \bar{A}_i \bar{x}_i$, 则求解

$$\min_{x_i} f_i(x_i) \stackrel{\text{def}}{=} \mu |x_i| + \frac{1}{2} \|a_i\|^2 x_i^2 - a_i^\top c_i x_i. \quad (5)$$

- (5)的最小值点

$$x_i^k = \arg \min_{x_i} f_i(x_i) = \begin{cases} \frac{a_i^T c_i - \mu}{\|a_i\|^2}, & a_i^T c_i > \mu, \\ \frac{a_i^T c_i + \mu}{\|a_i\|^2}, & a_i^T c_i < -\mu, \\ 0, & \text{otherwise.} \end{cases}$$

► K -均值聚类问题

$$\begin{aligned} \min_{\Phi, H} \quad & \|A - \Phi H\|_F^2, \\ \text{s.t.} \quad & \Phi \in \mathbb{R}^{n \times k}, \text{ 每一行只有一个元素为1, 其余为0,} \\ & H \in \mathbb{R}^{k \times p}. \end{aligned}$$

- 当固定 H 时, 设 Φ 的每一行为 ϕ_i^T , 则

$$A - \Phi H = \begin{pmatrix} a_1^T \\ a_2^T \\ \vdots \\ a_n^T \end{pmatrix} - \begin{pmatrix} \phi_1^T \\ \phi_2^T \\ \vdots \\ \phi_n^T \end{pmatrix} H = \begin{pmatrix} a_1^T - \phi_1^T H \\ a_2^T - \phi_2^T H \\ \vdots \\ a_n^T - \phi_n^T H \end{pmatrix}.$$

注意到 ϕ_i 只有一个分量为1, 其余分量为0, 不妨设其第 j 个分量为1, 此时 $\phi_i^T H$ 相当于将 H 的第 j 行取出, 因此 $\|a_i^T - \phi_i^T H\|$ 为 a_i^T 与 H 的第 j 个行向量的距离. 我们的最终目的是极小化 $\|A - \Phi H\|_F^2$, 所以 j 应该选矩阵 H 中距离 a_i^T 最近的那一行:

$$\Phi_{ij} = \begin{cases} 1, & j = \arg \min_l \|a_i - h_l\|, \\ 0, & \text{otherwise.} \end{cases}$$

其中 h_l^T 表示矩阵 H 的第 l 行.

- 当固定 Φ 时, 考虑 H 的每一行 h_j^T , 则有

$$\|A - \Phi H\|_F^2 = \sum_{j=1}^k \sum_{a \in S_j} \|a - h_j\|^2,$$

因此只需对每个 h_j 求最小. 设 \bar{a}_j 是目前第 j 类所有点的均值,

则 $\sum_{a \in S_j} \langle a - \bar{a}_j, \bar{a}_j - h_j \rangle = 0$ 且

$$\begin{aligned} \sum_{a \in S_j} \|a - h_j\|^2 &= \sum_{a \in S_j} \|a - \bar{a}_j + \bar{a}_j - h_j\|^2 \\ &= \sum_{a \in S_j} (\|a - \bar{a}_j\|^2 + \|\bar{a}_j - h_j\|^2 + 2 \langle a - \bar{a}_j, \bar{a}_j - h_j \rangle) \\ &= \sum_{a \in S_j} (\|a - \bar{a}_j\|^2 + \|\bar{a}_j - h_j\|^2). \end{aligned}$$

故 $h_j = \bar{a}_j$ 可达到最小值.

► 非负矩阵分解问题 $\min_{X, Y \geq 0} f(X, Y) = \frac{1}{2} \|XY - M\|_F^2$

- $f(X, Y)$ 的梯度

$$\frac{\partial f}{\partial X} = (XY - M)Y^T, \quad \frac{\partial f}{\partial Y} = X^T(XY - M).$$

- 利用格式(3), 注意到当 $r_i(X)$ 为凸集示性函数时即是求解到该集合的投影, 因此得到分块坐标下降法:

$$X^{k+1} = \max\{X^k - t_k^x (X^k Y^k - M)(Y^k)^T, 0\},$$

$$Y^{k+1} = \max\{Y^k - t_k^y (X^k)^T (X^k Y^k - M), 0\},$$

其中 t_k^x, t_k^y 是步长.

►字典学习 $\min_{D, X} \frac{1}{2n} \|DX - A\|_F^2 + \lambda \|X\|_1 + \frac{\mu}{2} \|D\|_F^2$

- 当固定变量 D 时, 考虑函数

$$f_D(X) = \frac{1}{2n} \|DX - A\|_F^2 + \lambda \|X\|_1.$$

使用格式(3). 计算 $f_D(X)$ 中光滑部分的梯度:

$$G = \frac{1}{n} D^T (DX - A),$$

因此格式(3)的BCD:

$$X^{k+1} = \text{prox}_{t_k \lambda \|\cdot\|_1} \left(X^k - \frac{t_k}{n} (D^k)^T (D^k X^k - A) \right),$$

其中 t_k 为步长.

- 当固定变量 X 时, 考虑函数

$$f_X(D) = \frac{1}{2n} \|DX - A\|_F^2 + \frac{\mu}{2} \|D\|_F^2.$$

使用格式(1). 计算关于 D^T 的梯度:

$$\nabla_{D^T} f_X(D) = \frac{1}{n} X(X^T D^T - A^T) + \mu D^T.$$

于是可得

$$D = AX^T(XX^T + n\mu I)^{-1}.$$

因为 $X \in \mathbb{R}^{k \times n}$, $k \ll n$, 故可方便地求出 XX^T 的逆. 故格式(1)的BCD:

$$D^{k+1} = A(X^{k+1})^T(X^{k+1}(X^{k+1})^T + n\mu I)^{-1}.$$

◆ 若先更新 X 再更新 D , 则最终可以得到如下的分块坐标下降法:

$$\begin{aligned} X^{k+1} &= \text{prox}_{t_k \lambda \|\cdot\|_1} \left(X^k - \frac{t_k}{n} (D^k)^T (D^k X^k - A) \right), \\ D^{k+1} &= A(X^{k+1})^T(X^{k+1}(X^{k+1})^T + n\mu I)^{-1}. \end{aligned}$$

BCD收敛性分析: 近似点交替线性化方法

考虑优化问题

$$\min \Psi(x, y) \stackrel{\text{def}}{=} f(x) + g(y) + H(x, y), \quad (x, y) \in \mathbb{R}^n \times \mathbb{R}^m$$

其中 f 和 g 为适当(非凸)闭函数, H 为连续可微函数.

- 基于BCD格式(3), 其分块坐标下降法迭代如下:

$$\begin{aligned} x^{k+1} &\in \text{prox}_{c_k f} \left(x^k - c_k \nabla_x H \left(x^k, y^k \right) \right), \\ y^{k+1} &\in \text{prox}_{d_k g} \left(y^k - d_k \nabla_y H \left(x^{k+1}, y^k \right) \right), \end{aligned} \tag{6}$$

其中 c_k, d_k 为步长参数. 由于自变量只有两块, 对光滑部分 H 采用线性化处理, 因此该BCD又称为近似点交替线性化方法.

► **非凸函数的邻近算子** 设 h 是适当闭函数(可以非凸), 且具有有限的下界, $\inf_{x \in \text{dom } h} h(x) > -\infty$, 定义 h 的邻近算子为

$$\text{prox}_h(x) = \arg \min_{u \in \text{dom } h} \left\{ h(u) + \frac{1}{2} \|u - x\|^2 \right\}.$$

► **非凸函数的次微分** 设 $f: \mathbb{R}^n \rightarrow (-\infty, +\infty]$ 是适当下半连续函数.

- ① 对给定的 $x \in \text{dom } f$, 满足如下条件的所有向量 $u \in \mathbb{R}^n$ 的集合定义为 f 在点 x 处的 *Fréchet* 次微分:

$$\liminf_{y \rightarrow x, y \neq x} \frac{f(y) - f(x) - \langle u, y - x \rangle}{\|y - x\|} \geq 0,$$

记为 $\hat{\partial}f(x)$. 当 $x \notin \text{dom } f$ 时, 将 $\hat{\partial}f(x)$ 定义为空集 \emptyset .

- ② f 在点 $x \in \mathbb{R}^n$ 处的 **极限次微分**(或简称为次微分)定义为

$$\partial f(x) = \{u \in \mathbb{R}^n : \exists x^k \rightarrow x, f(x^k) \rightarrow f(x), u^k \in \hat{\partial}f(x^k) \rightarrow u\}.$$

极限次微分通过对 x 附近的点处的 *Fréchet* 次微分取极限得到.

- 设 h 是适当闭函数且 $\inf_{x \in \text{dom } h} h(x) > -\infty$, 则 $\forall x \in \text{dom } h$, $\text{prox}_h(x)$ 是 \mathbb{R}^n 上的非空紧集.
- $\hat{\partial}f(x) \subseteq \partial f(x)$, 前者是闭凸集, 后者是闭集. 并非在所有的 $x \in \text{dom } f$ 处都存在 *Fréchet* 次微分.
- 设 h 是适当闭函数(可非凸)且有下界, $u \in \text{prox}_h(x)$, 则 $x - u \in \partial h(u)$

► 近似点交替线性化方法(6)收敛的假设条件: 假设A

- (1) $f: \mathbb{R}^n \rightarrow (-\infty, +\infty]$, $g: \mathbb{R}^m \rightarrow (-\infty, +\infty]$ 均为适当下半连续函数,
 $\inf_{\mathbb{R}^n \times \mathbb{R}^m} \Psi > -\infty$, $\inf_{\mathbb{R}^n} f > -\infty$, 以及 $\inf_{\mathbb{R}^m} g > -\infty$
- (2) $H: \mathbb{R}^n \times \mathbb{R}^m \rightarrow \mathbb{R}$ 是连续可微函数, 且 ∇H 在有界集上是联合利普希茨连续的: 对于任意的 $B_1 \times B_2 \subset \mathbb{R}^n \times \mathbb{R}^m$, 存在 $L > 0$ 使得对于任意的 $(x_1, y_1), (x_2, y_2) \in B_1 \times B_2$ 有

$$\begin{aligned} & \|(\nabla_x H(x_1, y_1) - \nabla_x H(x_2, y_2), \nabla_y H(x_1, y_1) - \nabla_y H(x_2, y_2))\| \\ & \leq L \|(x_1 - x_2, y_1 - y_2)\|. \end{aligned}$$

由假设A可得:

- 在有界集上 H 关于每个分量都是梯度 L -利普希茨连续的

$$\begin{aligned} \|\nabla_x H(x_1, y) - \nabla_x H(x_2, y)\| & \leq L \|x_1 - x_2\|, \\ \|\nabla_y H(x, y_1) - \nabla_y H(x, y_2)\| & \leq L \|y_1 - y_2\|. \end{aligned} \tag{7}$$

- $\Psi(x, y)$ 的次微分:

$$\partial \Psi(x, y) = (\nabla_x H(x, y) + \partial f(x), \nabla_y H(x, y) + \partial g(y)). \tag{8}$$

Kurdyka-Łojasiewicz (KL) 性质

非凸的情况下进行收敛性分析的主要工具是Kurdyka-Łojasiewicz (KL) 性质.

- 设 $\sigma : \mathbb{R}^d \rightarrow (-\infty, +\infty]$ 是适当下半连续函数且 $\inf_{\mathbb{R}^d} \sigma > -\infty$. 给定实数 $\alpha \leq \beta$, 定义 $[\alpha, \beta]$ 关于函数 σ 的原像为

$$[\alpha \leq \sigma \leq \beta] = \{x \in \mathbb{R}^d \mid \alpha \leq \sigma(x) \leq \beta\}.$$

- **Φ_η 函数类**: 定义 Φ_η 是凹连续函数 $\varphi : [0, \eta) \rightarrow \mathbb{R}_+$ 的集合且满足:
 - (i) $\varphi(0) = 0$; (ii) φ 在 $(0, \eta)$ 内连续可微, 在点0处连续;
 - (iii) 对任意的 $s \in (0, \eta)$, 都有 $\varphi'(s) > 0$.
- **KL性质** 设 $\sigma : \mathbb{R}^d \rightarrow (-\infty, +\infty]$ 是适当下半连续函数.
 - (1) 称函数 σ 在给定点 $\bar{u} \in \mathbf{dom} \partial\sigma \stackrel{\text{def}}{=} \{u \mid \partial\sigma(u) \neq \emptyset\}$ 处具有**KL 性质**, 若存在 $\eta \in (0, +\infty]$ 和 \bar{u} 的一个邻域 U 以及函数 $\varphi \in \Phi_\eta$, 使得

$$\varphi'(\sigma(u) - \sigma(\bar{u})) \cdot \text{dist}(0, \partial\sigma(u)) \geq 1, \quad \forall u \in U \cap [\sigma(\bar{u}) < \sigma < \sigma(\bar{u}) + \eta],$$

其中 $\text{dist}(x, S)$ 表示点 x 到集合 S 的距离.

(2) 若 σ 在 $\text{dom } \partial\sigma$ 上处处满足KL 性质, 则称 σ 是一个KL 函数.

►KL性质的解释

- 一大类函数都具有KL 性质, 该性质刻画了函数本身在给定点 \bar{u} 处的某种行为.
- 如果点 \bar{u} 不是函数 σ 的临界点, 那么KL 性质在点 \bar{u} 处自然成立. 因此KL 性质成立的不平凡情形是 \bar{u} 是 σ 的临界点, 即 $0 \in \partial\sigma(\bar{u})$.
- 这种情况下KL 性质保证了“函数 σ 可被锐化”. 直观上来说, 令

$$\tilde{\varphi}(u) = \varphi(\sigma(u) - \sigma(\bar{u})),$$

KL 性质在某种条件下可以改写成

$$\text{dist}(0, \partial\tilde{\varphi}(u)) \geq 1,$$

其中 u 的取法需要保证 $\sigma(u) > \sigma(\bar{u})$.

- 以上性质表明, 无论 u 多么接近临界点 \bar{u} , $\tilde{\varphi}(u)$ 的次梯度的模长均大于1. 所以KL 性质也被称为是函数 σ 在重参数化子 φ 下的一个锐化, 这种几何性质在分析一阶算法的收敛性时起到关键作用.

半代数与KL函数

半代数, 对数指数函数是KL函数

- \mathbb{R}^d 的子集 S 是一个半代数集, 如果存在有限个实多项式函数 $g_{ij}, h_{ij} : \mathbb{R}^d \rightarrow \mathbb{R}$ 使得

$$S = \bigcup_{j=1}^p \bigcap_{i=1}^q \{u \in \mathbb{R}^d : g_{ij}(u) = 0, h_{ij}(u) < 0\}.$$

- 函数 $h : \mathbb{R}^d \rightarrow (-\infty, +\infty]$ 称为半代数的, 如果它的图 $\{(u, t) \in \mathbb{R}^{d+1} : h(u) = t\}$ 是 \mathbb{R}^{d+1} 的半代数子集.
- 设 $\sigma(u) : \mathbb{R}^d \rightarrow (-\infty, +\infty)$ 是适当的下半连续函数. 若 σ 是半代数的, 则它在 $\text{dom } \sigma$ 中任一点处满足KL性质.

► 半代数函数举例:

- 实多项式函数.
- 半代数集的指示函数.
- 半代数函数的有限和与有限乘积; 半代数函数的复合.

- 上极限/下极限类函数. 例如, 当 g 是半代数函数并且 C 是半代数集时, $\sup\{g(u, v) : v \in C\}$ 是半代数的.
- 半正定矩阵锥, Stiefel流形以及恒秩矩阵都是半代数集.
- S 是 \mathbb{R}^d 中的非空半代数子集, 则函数 $\text{dist}(x, S)^2$ 是半代数的.
- $\|\cdot\|_0, \|\cdot\|_p$ 是半代数函数, 其中 p 是有理数.

Lemma 1 (一致KL性质). 设 Ω 是紧集, $\sigma : \mathbb{R}^d \rightarrow (-\infty, +\infty]$ 是适当下半连续函数, 在 Ω 上为常数且在 Ω 的每个点处都满足KL性质, 则存在 $\varepsilon > 0, \eta > 0, \varphi \in \Phi_\eta$ 使得对任意 $\bar{u} \in \Omega$ 和所有满足以下条件的 u :

$$\{u \in \mathbb{R}^d : \text{dist}(u, \Omega) < \varepsilon\} \cap [\sigma(\bar{u}) < \sigma < \sigma(\bar{u}) + \eta],$$

有

$$\varphi'(\sigma(u) - \sigma(\bar{u}))\text{dist}(0, \partial\sigma(u)) \geq 1.$$

Proof. 因为 \mathbb{R}^d 上的紧集可以由有限多个开集覆盖, 因此该问题可在有限个点上进行讨论.

- 设 μ 是 σ 在 Ω 上的取值. 由于 Ω 是紧集, 根据有限覆盖定理, 存在有限多个开球 $B(u_i, \varepsilon_i)$ (其中 $u_i \in \Omega, i = 1, 2, \dots, p$) 使得 $\Omega \subset \bigcup_{i=1}^p B(u_i, \varepsilon_i)$.

- 现在考虑这些点 u_i . 在点 u_i 上 KL 性质成立, 设 $\varphi_i : [0, \eta_i) \rightarrow \mathbb{R}_+$ 是对应的重参数化子, 则对任意 $u \in B(u_i, \varepsilon_i) \cap [\mu < \sigma < \mu + \eta_i]$, 有(逐点的) KL 性质:

$$\varphi'_i(\sigma(u) - \mu) \text{dist}(0, \partial\sigma(u)) \geq 1.$$

取充分小的 $\varepsilon > 0$ 使得

$$U_\varepsilon \stackrel{\text{def}}{=} \{u \in \mathbb{R}^d \mid \text{dist}(u, \Omega) \leq \varepsilon\} \subset \bigcup_{i=1}^p B(u_i, \varepsilon_i).$$

- 取 $\eta = \min_i \eta_i$, 以及

$$\varphi(s) = \int_0^s \max_i \varphi'_i(t) dt, \quad s \in [0, \eta).$$

容易验证 $\varphi \in \Phi_\eta$.

- 对任意的 $u \in U_\varepsilon \cap [\mu < \sigma < \mu + \eta]$, u 必定落在某个球 $B(u_{i_0}, \varepsilon_{i_0})$ 中, 于是我们有

$$\begin{aligned} \varphi'(\sigma(u) - \mu) \text{dist}(0, \partial\sigma(u)) &= \max_i \varphi'_i(\sigma(u) - \mu) \text{dist}(0, \partial\sigma(u)) \\ &\geq \varphi'_{i_0}(\sigma(u) - \mu) \text{dist}(0, \partial\sigma(u)) \geq 1. \end{aligned}$$

即一致 KL 性质成立.

□

近似点交替线性化方法的收敛性

分析近似点交替线性化方法(6)的收敛性, 主要分为三个步骤:

① **充分下降**: 找到一个正常数 ρ_1 使得

$$\rho_1 \|z^{k+1} - z^k\|^2 \leq \Psi(z^k) - \Psi(z^{k+1}).$$

② **次梯度上界**: 假设算法产生的迭代序列有界, 找到一个常数 ρ_2 , 使得次梯度有一个上界估计:

$$\|w^{k+1}\| \leq \rho_2 \|z^{k+1} - z^k\|, \quad w^k \in \partial\Psi(z^k).$$

③ **利用KL 性质证明全序列收敛**: 假设 Ψ 是一个KL函数, 证明迭代序列 $\{z^k\}$ 是一个柯西列.

★ **注**: 前两个步骤是证明多数算法的基本步骤, 当这两个性质成立时, 对任意的算法产生的迭代序列的聚点集合都为非空连通紧集, 且这些聚点都是 Ψ 的临界点.

►充分下降

Lemma 2. 设 $h : \mathbb{R}^d \rightarrow \mathbb{R}$ 是连续可微函数, 梯度 ∇h 是 L_h -利普希茨连续的, $\sigma : \mathbb{R}^d \rightarrow (-\infty, +\infty]$ 是适当下半连续函数且 $\inf_{\mathbb{R}^d} \sigma > -\infty$. 固定 $t < \frac{1}{L_h}$, 则对任意的 $u \in \text{dom } \sigma$ 和 $\tilde{u} \in \text{prox}_{t\sigma}(u - t\nabla h(u))$, 有

$$h(\tilde{u}) + \sigma(\tilde{u}) \leq h(u) + \sigma(u) - \frac{1}{2} \left(\frac{1}{t} - L_h \right) \|\tilde{u} - u\|^2.$$

Proof. 根据 σ 的假设, \tilde{u} 是良定义的. 根据 \tilde{u} 的最优性, 有

$$\langle \tilde{u} - u, \nabla h(u) \rangle + \frac{1}{2t} \|\tilde{u} - u\|^2 + \sigma(\tilde{u}) \leq \sigma(u).$$

再结合二次上界, 有

$$\begin{aligned} h(\tilde{u}) + \sigma(\tilde{u}) &\leq h(u) + \langle \tilde{u} - u, \nabla h(u) \rangle + \frac{L_h}{2} \|\tilde{u} - u\|^2 + \sigma(\tilde{u}) \\ &\leq h(u) + \frac{L_h}{2} \|\tilde{u} - u\|^2 + \sigma(u) - \frac{1}{2t} \|\tilde{u} - u\|^2 \\ &= h(u) + \sigma(u) - \frac{1}{2} \left(\frac{1}{t} - L_h \right) \|\tilde{u} - u\|^2. \end{aligned}$$

□

Theorem 1 (充分下降). 设假设A成立, $\{z^k\} = \{(x^k, y^k)\}$ 为BCD(6) 产生的迭代序列, 且假设 $\{z^k\}$ 有界. 取步长 $c_k = d_k = \frac{1}{\gamma L}$, 其中 $\gamma > 1$ 是常数, L 为 ∇H 的利普希茨系数, 则以下结论成立:

(1) 函数值序列 $\{\Psi(z^k)\}$ 是单调下降的, 且

$$\frac{\rho_1}{2} \|z^{k+1} - z^k\|^2 \leq \Psi(z^k) - \Psi(z^{k+1}), \quad \forall k \geq 0,$$

其中 $\rho_1 = (\gamma - 1)L$.

(2) 序列 $\{\|z^{k+1} - z^k\|\}_{k=1}^{\infty}$ 平方可和:

$$\sum_{k=1}^{\infty} \left(\|x^{k+1} - x^k\|^2 + \|y^{k+1} - y^k\|^2 \right) = \sum_{k=1}^{\infty} \|z^{k+1} - z^k\|^2 < +\infty,$$

且 $\lim_{k \rightarrow \infty} \|z^{k+1} - z^k\| = 0$.

Proof. (1) 根据假设A (2), (7) 和引理2可知, 每一步关于 x^k 和 y^k 的下降量估计:

$$\begin{aligned} H(x^{k+1}, y^k) + f(x^{k+1}) &\leq H(x^k, y^k) + f(x^k) - \frac{1}{2} \left(\frac{1}{c_k} - L \right) \|x^{k+1} - x^k\|^2 \\ &= H(x^k, y^k) + f(x^k) - \frac{1}{2} (\gamma - 1) L \|x^{k+1} - x^k\|^2, \end{aligned}$$

$$\begin{aligned} H(x^{k+1}, y^{k+1}) + g(y^{k+1}) &\leq H(x^{k+1}, y^k) + g(y^k) - \frac{1}{2} \left(\frac{1}{d_k} - L \right) \|y^{k+1} - y^k\|^2 \\ &= H(x^{k+1}, y^k) + g(y^k) - \frac{1}{2} (\gamma - 1) L \|y^{k+1} - y^k\|^2. \end{aligned}$$

$$\Rightarrow \Psi(z^k) - \Psi(z^{k+1}) \geq \frac{1}{2} (\gamma - 1) L \left(\|x^{k+1} - x^k\|^2 + \|y^{k+1} - y^k\|^2 \right).$$

由此可得

$$\frac{\rho_1}{2} \|z^{k+1} - z^k\|^2 \leq \Psi(z^k) - \Psi(z^{k+1}). \quad (9)$$

因此, 根据假设 $\inf \Psi > -\infty$ 可知, 函数值序列 $\{\Psi(z^k)\}$ 单调下降收敛到一个有限的数 Ψ^* .

(2) 设 N 为任意的整数, 在(9) 中对 k 求和, 得

$$\sum_{k=0}^{N-1} \|z^{k+1} - z^k\|^2 \leq \frac{2}{\rho_1} (\Psi(z^0) - \Psi(z^N)) \leq \frac{2}{\rho_1} (\Psi(z^0) - \Psi^*).$$

令 $N \rightarrow \infty$, 可得 $\sum_{k=0}^{\infty} \|z^{k+1} - z^k\|^2 < +\infty$, 从而 $\lim_{k \rightarrow \infty} \|z^{k+1} - z^k\| = 0$. \square

► **次梯度上界** 讨论序列 $\{z^k\}$ 是否会趋于某个临界点.

Lemma 3 (次梯度上界). 设定理1中的假设条件成立, 且 $\{z^k\}$ 是BCD(6)产生的有界序列. 对任意的整数 k , 定义

$$A_x^k = \frac{1}{c_{k-1}} (x^{k-1} - x^k) + \nabla_x H(x^k, y^k) - \nabla_x H(x^{k-1}, y^{k-1}),$$

$$A_y^k = \frac{1}{d_{k-1}} (y^{k-1} - y^k) + \nabla_y H(x^k, y^k) - \nabla_y H(x^k, y^{k-1}).$$

则 $(A_x^k, A_y^k) \in \partial \Psi(x^k, y^k)$ 且

$$\|(A_x^k, A_y^k)\| \leq \|A_x^k\| + \|A_y^k\| \leq \rho_2 \|z^k - z^{k-1}\|,$$

其中 $\rho_2 = (2\gamma + 3)L$.

Proof. 由(6)中更新 x^k 的一阶最优性条件可知

$$\nabla_x H(x^{k-1}, y^{k-1}) + \frac{1}{c_{k-1}}(x^k - x^{k-1}) + u^k = 0,$$

其中 $u^k \in \partial f(x^k)$ 为 f 的一个次梯度. 因此,

$$u^k = \frac{1}{c_{k-1}}(x^{k-1} - x^k) - \nabla_x H(x^{k-1}, y^{k-1}).$$

同理, 由(6)中 y^k 的更新可知

$$v^k = \frac{1}{d_{k-1}}(y^{k-1} - y^k) - \nabla_y H(x^k, y^{k-1}),$$

其中 $v^k \in \partial g(y^k)$ 为 g 的一个次梯度. 由 A_x^k, A_y^k 的定义, 及(8)可得

$$A_x^k = \nabla_x H(x^k, y^k) + u^k \in \partial_x \Psi(x^k, y^k),$$

$$A_y^k = \nabla_y H(x^k, y^k) + v^k \in \partial_y \Psi(x^k, y^k).$$

故有 $(A_x^k, A_y^k) \in \partial \Psi(x^k, y^k)$.

下面估计 A_x^k 和 A_y^k 的模长. 由假设A (2)知, 对 $\|A_x^k\|$ 有

$$\begin{aligned}\|A_x^k\| &\leq \frac{1}{c_{k-1}} \|x^{k-1} - x^k\| + \|\nabla_x H(x^k, y^k) - \nabla_x H(x^{k-1}, y^{k-1})\| \\ &\leq \frac{1}{c_{k-1}} \|x^{k-1} - x^k\| + L(\|x^{k-1} - x^k\| + \|y^{k-1} - y^k\|) \\ &= \left(L + \frac{1}{c_{k-1}}\right) \|x^{k-1} - x^k\| + L\|y^{k-1} - y^k\| \\ &= (\gamma + 1)L\|x^{k-1} - x^k\| + L\|y^{k-1} - y^k\| \\ &\leq (\gamma + 2)L\|z^{k-1} - z^k\|.\end{aligned}$$

由(7)知, 对 $\|A_y^k\|$ 有

$$\begin{aligned}
 \|A_y^k\| &\leq \frac{1}{d_{k-1}} \|y^k - y^{k-1}\| + \|\nabla_y H(x^k, y^k) - \nabla_y H(x^k, y^{k-1})\| \\
 &\leq \frac{1}{d_{k-1}} \|y^k - y^{k-1}\| + L \|y^k - y^{k-1}\| \\
 &= \left(\frac{1}{d_{k-1}} + L \right) \|y^k - y^{k-1}\| \\
 &\leq (\gamma + 1)L \|z^k - z^{k-1}\|.
 \end{aligned}$$

由此可得

$$\|(A_x^k, A_y^k)\| \leq \|A_x^k\| + \|A_y^k\| \leq (2\gamma + 3)L \|z^k - z^{k-1}\| = \rho_2 \|z^k - z^{k-1}\|. \quad \square$$

►子列收敛性

- 由引理3知, $\partial\Psi(z^k)$ 将会包含一个模长不断趋于0的向量, 这暗示着某种收敛性. 由于有界序列 $\{z^k\}$ 一定有收敛的子列, 故可猜想 $\{z^k\}$ 的极限点应和 Ψ 的临界点有一定的关系.

- **极限点集的性质:** 定义 $\omega(z^0)$ 为近似点交替线性化方法(6)从点 z^0 出发产生迭代序列的所有极限点集, 且 $\{z^k\}$ 是有界序列, 则以下结论成立:

(1) $\emptyset \neq \omega(z^0) \subset \text{crit } \Psi$, 其中 $\text{crit } \Psi$ 定义为 Ψ 所有的临界点.

(2) z^k 与集合 $\omega(z^0)$ 的距离趋于0,

$$\lim_{k \rightarrow \infty} \text{dist}(z^k, \omega(z^0)) = 0.$$

(3) $\omega(z^0)$ 是非空的连通紧集.

(4) Ψ 在 $\omega(z^0)$ 上是一个有限的常数.

Ref: BOLTE J, SABACH S, TEBOULLE M. *Proximal alternating linearized minimization for nonconvex and nonsmooth problems*. Mathematical Programming, 2014, 146(1): 459-494.

►利用KL 性质证明全序列收敛

Theorem 2 (有限长度性质). 设 Ψ 是KL 函数, 假设 \mathbf{A} 满足, $\{z^k\}$ 是有界序列, 则以下结论成立:

(1) 序列 $\{z^k\}$ 的长度有限,

$$\sum_{k=1}^{\infty} \|z^{k+1} - z^k\| < +\infty. \quad (10)$$

(2) 序列 $\{z^k\}$ 收敛到 Ψ 的一个临界点 $z^* = (x^*, y^*)$.

Proof. 由于 $\{z^k\}$ 是有界序列, 故存在收敛子列 $\{z^{k_q}\} \rightarrow \bar{z} \ (q \rightarrow \infty)$. 由充分下降定理1(1)知, 函数值序列 $\{\Psi(z^k)\}$ 总是收敛的, 且

$$\lim_{k \rightarrow \infty} \Psi(z^k) = \Psi(\bar{z}). \quad (11)$$

不妨设 $\Psi(\bar{z}) < \Psi(z^{\bar{k}})$. 这是因为若存在 \bar{k} 使得 $\Psi(z^{\bar{k}}) = \Psi(\bar{z})$, 由充分下降性可知 $z^{\bar{k}+1} = z^{\bar{k}}$, 进而有 $z^k = z^{\bar{k}}, \forall k > \bar{k}$. 结论自然成立.

由极限(11) 和极限点集 $\omega(z^0)$ 的性质(2): $\lim_{k \rightarrow \infty} \text{dist}(z^k, \omega(z^0)) = 0$ 可知, 对

任意的 $\varepsilon, \eta > 0$, 存在充分大的正整数 l , 使得对任意的 $k > l$,

$$\Psi(z^k) < \Psi(\bar{z}) + \eta, \quad \text{dist}(z^k, \omega(z^0)) < \varepsilon.$$

因此, 由假设 Ψ 是KL函数, 和极限点集 $\omega(z^0)$ 的性质(4): Ψ 在 $\omega(z^0)$ 上是一个有限的常数 可知, 一致KL性质成立.

(1) 根据极限点集 $\omega(z^0)$ 的性质, $\omega(z^0)$ 是非空紧集, 且 Ψ 在 $\omega(z^0)$ 上是常数. 故由一致KL性质引理1知, 对任意的 $k > l$ 有

$$\varphi'(\Psi(z^k) - \Psi(\bar{z})) \text{dist}(0, \partial\Psi(z^k)) \geq 1. \quad (12)$$

根据次梯度上界引理3可知

$$\text{dist}(0, \partial\Psi(z^k)) \leq \|(A_x^k, A_y^k)\| \leq \rho_2 \|z^k - z^{k-1}\|.$$

又由(12)得

$$\varphi'(\Psi(z^k) - \Psi(\bar{z})) \geq \frac{1}{\rho_2} \|z^k - z^{k-1}\|^{-1}. \quad (13)$$

由 φ 的凹性知,

$$\varphi(\Psi(z^k) - \Psi(\bar{z})) - \varphi(\Psi(z^{k+1}) - \Psi(\bar{z})) \geq \varphi'(\Psi(z^k) - \Psi(\bar{z}))(\Psi(z^k) - \Psi(z^{k+1})). \quad (14)$$

为了表示方便, 定义

$$\Delta_{p,q} = \varphi(\Psi(z^p) - \Psi(\bar{z})) - \varphi(\Psi(z^q) - \Psi(\bar{z})),$$

其中 p, q 为任意正整数. 显然 $\Delta_{p,q} + \Delta_{q,r} = \Delta_{p,r}$. 定义常数

$$C = \frac{2\rho_2}{\rho_1} > 0.$$

根据不等式(14), (13) 和充分下降定理1, 可知

$$\begin{aligned} \Delta_{k,k+1} &\geq \varphi'(\Psi(z^k) - \Psi(\bar{z}))(\Psi(z^k) - \Psi(z^{k+1})) \\ &\geq \frac{1}{\rho_2} \|z^k - z^{k-1}\|^{-1} \cdot \frac{\rho_1}{2} \|z^{k+1} - z^k\|^2 \\ &= \frac{\|z^{k+1} - z^k\|^2}{C\|z^k - z^{k-1}\|} \\ \Rightarrow \quad &\|z^{k+1} - z^k\| \leq \sqrt{C\Delta_{k,k+1}\|z^k - z^{k-1}\|}. \end{aligned}$$

根据基本不等式 $2\sqrt{ab} \leq a + b, \forall a, b > 0$, 取 $a = \|z^k - z^{k-1}\|$, $b = C\Delta_{k,k+1}$, 则

$$2\|z^{k+1} - z^k\| \leq \|z^k - z^{k-1}\| + C\Delta_{k,k+1}. \quad (15)$$

对任意的 $k > l$, 在(15)中把 k 替换成 i 并对 $i = l + 1, l + 2, \dots, k$ 求和, 得

$$\begin{aligned} 2 \sum_{i=l+1}^k \|z^{i+1} - z^i\| &\leq \sum_{i=l+1}^k \|z^i - z^{i-1}\| + C \sum_{i=l+1}^k \Delta_{i,i+1} \\ &\leq \sum_{i=l+1}^k \|z^{i+1} - z^i\| + \|z^{l+1} - z^l\| + C\Delta_{l+1,k+1}. \end{aligned}$$

故有

$$\begin{aligned} \sum_{i=l+1}^k \|z^{i+1} - z^i\| &\leq \|z^{l+1} - z^l\| + C \left(\varphi(\Psi(z^{l+1}) - \Psi(\bar{z})) - \varphi(\Psi(z^{k+1}) - \Psi(\bar{z})) \right) \\ &\leq \|z^{l+1} - z^l\| + C\varphi(\Psi(z^{l+1}) - \Psi(\bar{z})) \\ \Rightarrow \sum_{k=1}^{\infty} \|z^{k+1} - z^k\| &< +\infty. \end{aligned}$$

(2) 由(10)知, 要证 $\{z^k\}$ 收敛, 只需证明 $\{z^k\}$ 是柯西列.

对任意的 $q > p > l$,

$$z^q - z^p = \sum_{k=p}^{q-1} (z^{k+1} - z^k),$$

根据三角不等式,

$$\|z^q - z^p\| = \left\| \sum_{k=p}^{q-1} (z^{k+1} - z^k) \right\| \leq \sum_{k=p}^{q-1} \|z^{k+1} - z^k\|.$$

因此, 由(10)知, $\{z^k\}$ 是一个柯西列, 故收敛.

□