

BCD & ADMM 算法收敛理论

Lecture 5: Block Coordinate Descent (BCD) 算法

罗自炎

北京交通大学数统学院

E-mail: zyluo@bjtu.edu.cn

参考资料

- 教材与参考文献：
 - 最优化：建模、算法与理论
 - Bolte et al., MP 2014
 - Fazel et al., SIMAX 2013
 - Rockafellar & Wets, Variational Analysis
- 致谢： 北京大学文再文教授; 清华大学张立平教授

Outline of BCD

- 问题描述
- 分块坐标下降法
- 应用举例

BCD算法求解的优化模型

► 典型优化问题形式:

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

- 可行域: \mathcal{X}
- 决策变量: $x = (x_1, x_2, \dots, x_s)$, $x_i \in \mathbb{R}^{n_i}$, $\sum_{i=1}^s n_i = n$
- 目标函数: f 是关于 x 的可微函数(coupled term), $r_i(x_i)$ 关于 x_i 是适当的闭凸函数(separable terms), 但不一定可微.
- 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

- **LASSO 模型** 待估参数 $x = (x_1, x_2, \dots, x_p)^T \in \mathbb{R}^p$ 可以分成 p 组

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

- **组LASSO 模型** 待估参数 $x = (\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_G)^T \in \mathbb{R}^p$ 可以分成 G 组, 且 $\{\mathbf{x}_i\}_{i=1}^G$ 中只有少数的非零向量.

$$\min_x \frac{1}{2n} \|b - Ax\|_2^2 + \lambda \sum_{i=1}^G \|\mathbf{x}_i\|_2.$$

- ***K*-means 聚类问题** 给定 p 维空间中 n 个数据点 a_1, a_2, \dots, a_n , 聚类问题就是要寻找 k 个不相交的非空集合 S_1, S_2, \dots, S_k , 使得

$$\{a_1, a_2, \dots, a_n\} = S_1 \cup S_2 \cup \dots \cup S_k,$$

并且使得组内距离平方和最小, 即

$$\begin{aligned} \min_{S_1, S_2, \dots, S_k} \quad & \sum_{i=1}^k \sum_{a \in S_i} \|a - c_i\|^2, \\ \text{s.t.} \quad & S_1 \cup S_2 \cup \dots \cup S_k = \{a_1, a_2, \dots, a_n\}, \\ & S_i \cap S_j = \emptyset, \quad \forall i \neq j \end{aligned}$$

- 等价转化:

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

$$\text{其中: } A = \begin{pmatrix} a_1^T \\ \vdots \\ a_n^T \end{pmatrix} \in \mathbb{R}^{n \times p}, H = \begin{pmatrix} h_1^T \\ \vdots \\ h_k^T \end{pmatrix} \in \mathbb{R}^{k \times p},$$

$$\Phi = (\phi_{ij}) \in \mathbb{R}^{n \times k}, \phi_{ij} = \begin{cases} 1, & a_i \in S_j; \\ 0, & a_i \notin S_j. \end{cases}$$

典型问题举例

- 非负矩阵分解

$$\min_{X, Y \geq 0} \frac{1}{2} \|M - XY\|_F^2$$

其中 $M \in \mathbb{R}^{m \times n}$ 是已知矩阵, $X \in \mathbb{R}^{m \times k}$, $Y \in \mathbb{R}^{k \times n}$ 为待求非负矩阵.

- 非负张量分解

$$\min_{A_1, A_2, \dots, A_N \geq 0} \frac{1}{2} \|\mathcal{M} - \llbracket A_1, A_2, \dots, A_N \rrbracket\|_F^2$$

其中 $\mathcal{M} \in \mathbb{R}^{I_1 \times \dots \times I_N}$ 是已知张量,

$\llbracket A_1, A_2, \dots, A_N \rrbracket = \sum_{i=1}^k a_i^{(1)} \circ \dots \circ a_i^{(N)}$, “ \circ ” 表示张量的外积运算,

$A_j = [a_1^{(j)} \ \dots \ a_k^{(j)}] \in \mathbb{R}^{I_j \times k}$, $j = 1, \dots, N$.

典型问题举例

- **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 , 要从 A 中学习出一个字典 $D \in \mathbb{R}^{m \times k}$ 和系数矩阵 $X \in \mathbb{R}^{k \times n}$:

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

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

- **交替极小**: 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

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

- 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 $(k-1)$.

三种变量更新方式

(i) 固定其他分量然后对单一变量求极小:

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

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

$$x_i^k = \arg \min_{x_i \in \mathcal{X}_i^k} \left\{ f_i^k(x_i) + \frac{L_i^{k-1}}{2} \|x_i - x_i^{k-1}\|_2^2 + r_i(x_i) \right\}. \quad (3)$$

(iii) 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:

$$x_i^k = \arg \min_{x_i \in \mathcal{X}_i^k} \left\{ \langle \hat{g}_i^k, x_i - \hat{x}_i^{k-1} \rangle + \frac{L_i^{k-1}}{2} \|x_i - \hat{x}_i^{k-1}\|_2^2 + r_i(x_i) \right\}, \quad (4)$$

其中 \hat{x}_i^{k-1} 采用**外推**定义:

$$\hat{x}_i^{k-1} = x_i^{k-1} + \omega_i^{k-1}(x_i^{k-1} - x_i^{k-2}), \quad (5)$$

$\omega_i^k \geq 0$ 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端梯度法(Proximal Gradient)的更新.

BCD 算法框架

► 分块坐标下降法(Block Coordinate Descent, BCD)

- 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(2) 或(3) 或(4) 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(2)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(2)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(3) (4) 则是对格式(2)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(3)可使得算法收敛性在函数 F 为非严格凸时有所改善.
- 格式(4)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(4)的计算量很小，比较容易实现.

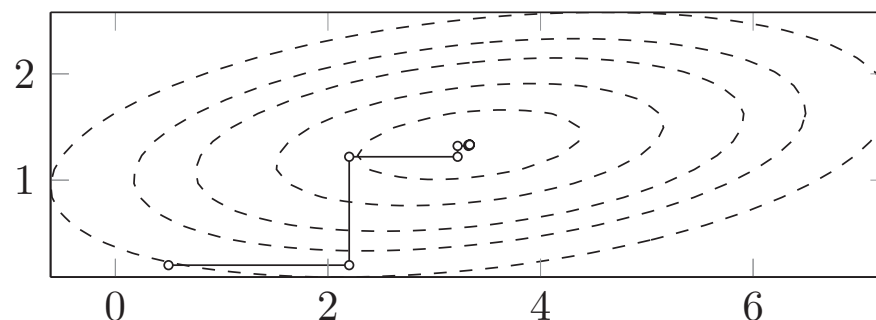
► 【例】 $\min f(x, y) = x^2 - 2xy + 10y^2 - 4x - 20y$

- 采用格式(2)的分块坐标下降法:

$$x^{k+1} = \arg \min_{x \in \mathbb{R}} \{x^2 - 2xy^k - 4x\} = 2 + y^k,$$

$$y^{k+1} = \arg \min_{y \in \mathbb{R}} \{-2x^{k+1}y + 10y^2 - 20y\} = 1 + \frac{x^{k+1}}{10}.$$

- 初始点: $(x, y) = (0.5, 0.2)$ 时迭代点轨迹如图, 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.



非凸 $f(x)$, BCD可能失效!

► 不收敛反例 (Powell, 1973): (使用格式(2))

- 目标函数:

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

$$f(x_1, x_2, x_3) = -x_1x_2 - x_2x_3 - x_3x_1 \text{ 非凸}$$

- 梯度分量: $\nabla_{x_1} F(x_1, x_2, x_3) = -x_2 - x_3 + 2(x_1 - 1)_+ - 2(-x_1 - 1)_+$
- 初始点: $x^0 = (-1 - \varepsilon, 1 + \frac{\varepsilon}{2}, -1 - \frac{\varepsilon}{4})$, 其中 $\varepsilon > 0$,
容易验证迭代序列满足

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

两个聚点: $u = (-1, 1, -1)^T, v = (1, -1, 1)^T$,

$$\nabla F(u) = (0, 2, 0)^T = -\nabla F(v)$$

因此 u, v 均不是 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(2)更新第 i 块, 令 $c_i = b - \bar{A}_i \bar{x}_i$, 则求解

$$\min_{x_i} f_i(x_i) \stackrel{\text{def}}{=} \mu |x_i| + \frac{1}{2} \|a_i\|^2 x_i^2 - a_i^\top c_i x_i. \quad (6)$$

$$= \|a_i\|^2 \min_{x_i} \left\{ \frac{\mu}{\|a_i\|^2} |x_i| + \frac{1}{2} (x_i - a_i^\top c_i / \|a_i\|^2)^2 + \text{const} \right\}$$

- (6)的最小值点

$$x_i^k = \arg \min_{x_i} f_i(x_i) = \text{Soft}_{\lambda_i}(\hat{u}_i) := \begin{cases} \hat{u}_i - \lambda_i, & \text{if } \hat{u}_i > \lambda_i, \\ \hat{u}_i + \lambda_i, & \text{if } \hat{u}_i < -\lambda_i, \\ 0, & \text{otherwise.} \end{cases}$$

$$\text{其中 } \lambda_i = \frac{\mu}{\|a_i\|^2}, \hat{u}_i = \frac{a_i^T c_i}{\|a_i\|^2}.$$

► K -均值聚类问题

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

- 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(4), 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(4). 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(4)的BCD:

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

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(2). 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(2)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:

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