

交替方向乘子法(ADMM)

宋晓良

大连理工大学数学科学学院

教材《最优化：建模、算法与理论》配套电子教案

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

1 交替方向乘子法

2 常见变形和技巧

3 应用举例

典型问题形式

考虑如下凸问题：

$$\begin{aligned} \min_{x_1, x_2} \quad & f_1(x_1) + f_2(x_2), \\ \text{s.t.} \quad & A_1 x_1 + A_2 x_2 = b, \end{aligned} \tag{1}$$

- f_1, f_2 是适当的闭凸函数，但不要求是光滑的， $x_1 \in \mathbb{R}^n, x_2 \in \mathbb{R}^m$, $A_1 \in \mathbb{R}^{p \times n}, A_2 \in \mathbb{R}^{p \times m}, b \in \mathbb{R}^p$.
- 问题特点：目标函数可以分成彼此分离的两块，但是变量被线性约束结合在一起。常见的一些无约束和带约束的优化问题都可以表示成这一形式。

问题形式举例

- 可以分成两块 of 无约束优化问题

$$\min_x f_1(x) + f_2(x).$$

引入一个新的变量 z 并令 $x = z$, 将问题转化为

$$\begin{aligned} \min_{x,z} \quad & f_1(x) + f_2(z), \\ \text{s.t.} \quad & x - z = 0. \end{aligned}$$

- 带线性变换的无约束优化问题

$$\min_x f_1(x) + f_2(Ax).$$

可以引入一个新的变量 z , 令 $z = Ax$, 则问题变为

$$\begin{aligned} \min_{x,z} \quad & f_1(x) + f_2(z), \\ \text{s.t.} \quad & Ax - z = 0. \end{aligned}$$

问题形式举例

- 凸集 $C \subset \mathbb{R}^n$ 上的约束优化问题

$$\begin{aligned} \min_x \quad & f(x), \\ \text{s.t.} \quad & Ax \in C, \end{aligned}$$

$I_C(z)$ 是集合 C 的示性函数，引入约束 $z = Ax$ ，那么问题转化为

$$\begin{aligned} \min_{x,z} \quad & f(x) + I_C(z), \\ \text{s.t.} \quad & Ax - z = 0. \end{aligned}$$

- 全局一致性问题

$$\min_x \sum_{i=1}^N \phi_i(x).$$

令 $x = z$ ，并将 x 复制 N 份，分别为 x_i ，那么问题转化为

$$\begin{aligned} \min_{x_i, z} \quad & \sum_{i=1}^N \phi_i(x_i), \\ \text{s.t.} \quad & x_i - z = 0, \quad i = 1, 2, \dots, N. \end{aligned}$$

增广拉格朗日函数法

- 首先写出问题(1)的增广拉格朗日函数

$$\begin{aligned} L_{\rho}(x_1, x_2, y) = & f_1(x_1) + f_2(x_2) + y^T(A_1x_1 + A_2x_2 - b) \\ & + \frac{\rho}{2} \|A_1x_1 + A_2x_2 - b\|_2^2, \end{aligned} \quad (2)$$

其中 $\rho > 0$ 是二次罚项的系数.

- 常见的求解带约束问题的增广拉格朗日函数法为如下更新:

$$(x_1^{k+1}, x_2^{k+1}) = \underset{x_1, x_2}{\operatorname{argmin}} L_{\rho}(x_1, x_2, y^k), \quad (3)$$

$$y^{k+1} = y^k + \tau \rho (A_1x_1^{k+1} + A_2x_2^{k+1} - b), \quad (4)$$

其中 τ 为步长.

交替方向乘子法

Alternating direction method of multipliers, ADMM

- 交替方向乘子法的基本思路: 第一步迭代(3)同时对 x_1 和 x_2 进行优化有时候比较困难, 而固定一个变量求解关于另一个变量的极小问题可能比较简单, 因此我们可以考虑对 x_1 和 x_2 交替求极小
- 其迭代格式可以总结如下:

$$x_1^{k+1} = \underset{x_1}{\operatorname{argmin}} L_\rho(x_1, x_2^k, y^k), \quad (5)$$

$$x_2^{k+1} = \underset{x_2}{\operatorname{argmin}} L_\rho(x_1^{k+1}, x_2, y^k), \quad (6)$$

$$y^{k+1} = y^k + \tau \rho (A_1 x_1^{k+1} + A_2 x_2^{k+1} - b), \quad (7)$$

其中 τ 为步长, 通常取值于 $\left(0, \frac{1+\sqrt{5}}{2}\right]$

原问题最优性条件

- 因为 f_1, f_2 均为闭凸函数，约束为线性约束，所以当Slater条件成立时，可以使用凸优化问题的KKT条件来作为交替方向乘子法的收敛准则。问题(1)的拉格朗日函数为

$$L(x_1, x_2, y) = f_1(x_1) + f_2(x_2) + y^T(A_1x_1 + A_2x_2 - b).$$

- 根据最优性条件定理，若 x_1^*, x_2^* 为问题(1)的最优解， y^* 为对应的拉格朗日乘子，则以下条件满足：

$$0 \in \partial_{x_1} L(x_1^*, x_2^*, y^*) = \partial f_1(x_1^*) + A_1^T y^*, \quad (8a)$$

$$0 \in \partial_{x_2} L(x_1^*, x_2^*, y^*) = \partial f_2(x_2^*) + A_2^T y^*, \quad (8b)$$

$$A_1 x_1^* + A_2 x_2^* = b. \quad (8c)$$

在这里条件(8c)又称为原始可行性条件，条件(8a)和条件(8b)又称为对偶可行性条件。

ADMM单步迭代最优性条件

- 由 x_2 的更新步骤

$$x_2^k = \operatorname{argmin}_x \left\{ f_2(x) + \frac{\rho}{2} \left\| A_1 x_1^k + A_2 x - b + \frac{y^{k-1}}{\rho} \right\|^2 \right\},$$

根据最优性条件不难推出

$$0 \in \partial f_2(x_2^k) + A_2^T [y^{k-1} + \rho(A_1 x_1^k + A_2 x_2^k - b)]. \quad (9)$$

当 $\tau = 1$ 时, 根据(7)可知上式方括号中的表达式就是 y^k , 最终有

$$0 \in \partial f_2(x_2^k) + A_2^T y^k,$$

- 由 x_1 的更新公式

$$x_1^k = \operatorname{argmin}_x \left\{ f_1(x) + \frac{\rho}{2} \|A_1 x + A_2 x_2^{k-1} - b + \frac{y^{k-1}}{\rho}\|^2 \right\},$$

假设子问题能精确求解, 根据最优性条件

$$0 \in \partial f_1(x_1^k) + A_1^T [\rho(A_1 x_1^k + A_2 x_2^{k-1} - b) + y^{k-1}].$$

ADMM单步迭代最优性条件

- 根据ADMM的第三式(7)取 $\tau = 1$ 有

$$0 \in \partial f_1(x_1^k) + A_1^T(y^k + \rho A_2(x_2^{k-1} - x_2^k)). \quad (10)$$

对比条件(8a)可知多出来的项为 $A_1^T A_2(x_2^{k-1} - x_2^k)$ 。因此要检测对偶可行性只需要检测残差

$$s^k = A_1^T A_2(x_2^{k-1} - x_2^k)$$

- 综上当 x_2 更新取到精确解且 $\tau = 1$ 时，判断ADMM是否收敛只需要检测前述两个残差 r^k, s^k 是否充分小：

$$\begin{aligned} 0 &\approx \|r^k\| = \|A_1 x_1^k + A_2 x_2^k - b\| && \text{(原始可行性)}, \\ 0 &\approx \|s^k\| = \|A_1^T A_2(x_2^{k-1} - x_2^k)\| && \text{(对偶可行性)}. \end{aligned} \quad (11)$$

提纲

1 交替方向乘子法

2 常见变形和技巧

3 应用举例

线性化

- 线性化技巧使用近似点项对子问题目标函数进行二次近似.
- 不失一般性, 我们考虑第一个子问题, 即

$$\min_{x_1} f_1(x_1) + \frac{\rho}{2} \|A_1 x_1 - v^k\|^2, \quad (12)$$

其中 $v^k = b - A_2 x_2^k - \frac{1}{\rho} y^k$.

- 当子问题目标函数可微时, 线性化将问题(12)变为

$$x_1^{k+1} = \operatorname{argmin}_{x_1} \left\{ (\nabla f_1(x_1^k) + \rho A_1^T (A_1 x_1^k - v^k))^T x_1 + \frac{1}{2\eta_k} \|x_1 - x^k\|_2^2 \right\},$$

其中 η_k 是步长参数, 这等价于做一步梯度下降.

- 当目标函数不可微时, 可以考虑只将二次项线性化, 即

$$x_1^{k+1} = \operatorname{argmin}_{x_1} \left\{ f_1(x_1) + \rho (A_1^T (A_1 x_1^k - v^k))^T x_1 + \frac{1}{2\eta_k} \|x_1 - x^k\|_2^2 \right\},$$

这等价于做一步近似点梯度步.

缓存分解

- 如果目标函数中含二次函数，例如 $f_1(x_1) = \frac{1}{2} \|C\mathbf{x}_1 - d\|_2^2$ ，那么针对 x_1 的更新(5)等价于求解线性方程组

$$(C^T C + \rho A_1^T A_1) x_1 = C^T d + \rho A_1^T v^k.$$

- 虽然子问题有显式解，但是每步求解的复杂度仍然比较高，这时候可以考虑用**缓存分解**的方法。首先对 $C^T C + \rho A_1^T A_1$ 进行**Cholesky** 分解并缓存分解的结果，在每步迭代中只需要求解简单的三角形方程组
- 当 ρ 发生更新时，就要重新进行分解。特别地，当 $C^T C + \rho A_1^T A_1$ 一部分容易求逆，另一部分是低秩的情形时，可以用**SMW**公式来求逆。

优化转移

- 有时候为了方便求解子问题，可以用一个性质好的矩阵 D 近似二次项 $A_1^T A_1$ ，此时子问题(12)替换为

$$\begin{aligned} x_1^{k+1} = \operatorname{argmin}_{x_1} & \left\{ f_1(x_1) + \frac{\rho}{2} \|A_1 x_1 - v^k\|_2^2 \right. \\ & \left. + \frac{\rho}{2} (x_1 - x^k)^T (D - A_1^T A_1) (x_1 - x^k) \right\}. \end{aligned}$$

这种方法也称为优化转移.

- 通过选取合适的 D ，当计算 $\operatorname{argmin}_{x_1} \left\{ f_1(x_1) + \frac{\rho}{2} x_1^T D x_1 \right\}$ 明显比计算 $\operatorname{argmin}_{x_1} \left\{ f_1(x_1) + \frac{\rho}{2} x_1^T A_1^T A_1 x_1 \right\}$ 要容易时，优化转移可以极大地简化子问题的计算. 特别地，当 $D = \frac{\eta_k}{\rho} I$ 时，优化转移等价于做单步的近似点梯度步.

二次罚项系数的动态调节

- 原始可行性和对偶可行性分别用 $\|r^k\|$ 和 $\|s^k\|$ 度量.
- 求解过程中二次罚项系数 ρ 太大会导致原始可行性 $\|r^k\|$ 下降很快, 但是对偶可行性 $\|s^k\|$ 下降很慢; 二次罚项系数太小, 则会有相反的效果. 这样都会导致收敛比较慢或得到的解的可行性很差.
- 一个自然的想法是在每次迭代时动态调节惩罚系数 ρ 的大小, 从而使得原始可行性和对偶可行性能够以比较一致的速度下降到零. 一个简单有效的方式是令

$$\rho^{k+1} = \begin{cases} \gamma_p \rho^k, & \|r^k\| > \mu \|s^k\|, \\ \rho^k / \gamma_d & \|s^k\| > \mu \|r^k\|, \\ \rho^k, & \text{其他,} \end{cases}$$

其中 $\mu > 1, \gamma_p > 1, \gamma_d > 1$ 是参数, 常见的选择

为 $\mu = 10, \gamma_p = \gamma_d = 2$. 在迭代过程中将原始可行性 $\|r^k\|$ 和对偶可行性 $\|s^k\|$ 保持在彼此的 μ 倍内. 如果发现 $\|r^k\|$ 或 $\|s^k\|$ 下降过慢就应该相应增大或减小二次罚项系数 ρ^k .

- 在(6)式与(7)式中, $A_1x_1^{k+1}$ 可以被替换为

$$\alpha_k A_1 x_1^{k+1} + (1 - \alpha_k)(A_2 x_2^k - b),$$

其中 $\alpha_k \in (0, 2)$ 是一个松弛参数.

- 当 $\alpha_k > 1$ 时, 这种技巧称为超松弛; 当 $\alpha_k < 1$ 时, 这种技巧称为欠松弛. 实验表明 $\alpha_k \in [1.5, 1.8]$ 的超松弛可以提高收敛速度.

多块问题的ADMM

- 考虑有多块变量的情形

$$\begin{aligned} \min_{x_1, x_2, \dots, x_N} \quad & f_1(x_1) + f_2(x_2) + \dots + f_N(x_N), \\ \text{s.t.} \quad & A_1 x_1 + A_2 x_2 + \dots + A_N x_N = b. \end{aligned} \tag{13}$$

这里 $f_i(x_i)$ 是闭凸函数, $x_i \in \mathbb{R}^{n_i}$, $A_i \in \mathbb{R}^{m \times n_i}$.

- 同样写出增广拉格朗日函数 $L_\rho(x_1, x_2, \dots, x_N, y)$, 相应的多块ADMM迭代格式为

$$\begin{aligned} x_1^{k+1} &= \operatorname{argmin}_x L_\rho(x, x_2^k, \dots, x_N^k, y^k), \\ x_2^{k+1} &= \operatorname{argmin}_x L_\rho(x_1^{k+1}, x, \dots, x_N^k, y^k), \\ &\dots\dots\dots \\ x_N^{k+1} &= \operatorname{argmin}_x L_\rho(x_1^{k+1}, x_2^{k+1}, \dots, x, y^k), \\ y^{k+1} &= y^k + \tau \rho (A_1 x_1^{k+1} + A_2 x_2^{k+1} + \dots + A_N x_N^{k+1} - b), \end{aligned}$$

其中 $\tau \in (0, (\sqrt{5} + 1)/2)$ 为步长参数.

提纲

1 交替方向乘子法

2 常见变形和技巧

3 应用举例

LASSO 问题的Primal 形式

- LASSO 问题

$$\min \quad \mu \|x\|_1 + \frac{1}{2} \|Ax - b\|^2.$$

转换为标准问题形式：

$$\begin{aligned} \min_{x,z} \quad & \frac{1}{2} \|Ax - b\|^2 + \mu \|z\|_1, \\ \text{s.t.} \quad & x = z. \end{aligned}$$

- 交替方向乘子法迭代格式为

$$\begin{aligned} x^{k+1} &= \operatorname{argmin}_x \left\{ \frac{1}{2} \|Ax - b\|^2 + \frac{\rho}{2} \|x - z^k + y^k/\rho\|^2 \right\}, \\ &= (A^T A + \rho I)^{-1} (A^T b + \rho z^k - y^k), \\ z^{k+1} &= \operatorname{argmin}_z \left\{ \mu \|z\|_1 + \frac{\rho}{2} \|x^{k+1} - z + y^k/\rho\|^2 \right\}, \\ &= \operatorname{prox}_{(\mu/\rho)\|\cdot\|_1} (x^{k+1} + y^k/\rho), \\ y^{k+1} &= y^k + \tau \rho (x^{k+1} - z^{k+1}). \end{aligned}$$

LASSO 问题的Primal 形式

- 注意，因为 $\rho > 0$ ，所以 $A^T A + \rho I$ 总是可逆的。 x 迭代本质上是计算一个岭回归问题（ ℓ_2 范数平方正则化的最小二乘问题）；而对 z 的更新为 ℓ_1 范数的邻近算子，同样有显式解。在求解 x 迭代时，若使用固定的罚因子 ρ ，我们可以缓存矩阵 $A^T A + \rho I$ 的初始分解，从而减小后续迭代中的计算量。
- 需要注意的是，在 LASSO 问题中，矩阵 $A \in \mathbb{R}^{m \times n}$ 通常有较多的列（即 $m \ll n$ ），因此 $A^T A \in \mathbb{R}^{n \times n}$ 是一个低秩矩阵，二次罚项的作用就是将 $A^T A$ 增加了一个正定项。该 ADMM 主要运算量来自更新 x 变量时求解线性方程组，复杂度为 $O(n^3)$ （若使用缓存分解技术或 SMW 公式则可进一步降低每次迭代的运算量）

LASSO 问题的对偶形式

- 考虑LASSO 问题的对偶问题

$$\begin{aligned} \min \quad & b^T y + \frac{1}{2} \|y\|^2, \\ \text{s.t.} \quad & \|A^T y\|_\infty \leq \mu. \end{aligned} \quad (14)$$

- 引入约束 $A^T y + z = 0$, 可以得到如下等价问题:

$$\begin{aligned} \min \quad & \underbrace{b^T y + \frac{1}{2} \|y\|^2}_{f(y)} + \underbrace{I_{\|z\|_\infty \leq \mu}(z)}_{h(z)}, \\ \text{s.t.} \quad & A^T y + z = 0. \end{aligned} \quad (15)$$

- 对约束 $A^T y + z = 0$ 引入乘子 x , 对偶问题的增广拉格朗日函数为

$$L_\rho(y, z, x) = b^T y + \frac{1}{2} \|y\|^2 + I_{\|z\|_\infty \leq \mu}(z) - x^T (A^T y + z) + \frac{\rho}{2} \|A^T y + z\|^2.$$

LASSO 问题的对偶形式

- 当固定 y, x 时, 对 z 的更新即向无穷范数球 $\{z \mid \|z\|_\infty \leq \mu\}$ 做欧几里得投影, 即将每个分量截断在区间 $[-\mu, \mu]$ 中; 当固定 z, x 时, 对 y 的更新即求解线性方程组

$$(I + \rho AA^T)y = A(x^k - \rho z^{k+1}) - b.$$

- 因此得到ADMM 迭代格式为

$$\begin{aligned} z^{k+1} &= \mathcal{P}_{\|z\|_\infty \leq \mu} (x^k / \rho - A^T y^k), \\ y^{k+1} &= (I + \rho AA^T)^{-1} (A(x^k - \rho z^{k+1}) - b), \\ x^{k+1} &= x^k - \tau \rho (A^T y^{k+1} + z^{k+1}). \end{aligned}$$

- 虽然ADMM 应用于对偶问题也需要求解一个线性方程组, 但由于LASSO 问题的特殊性 ($m \ll n$), 求解 y 更新的线性方程组需要的计算量是 $O(m^3)$, 使用缓存分解技巧后可进一步降低至 $O(m^2)$, 这大大小于针对原始问题的ADMM.

References

Douglas-Rachford method, ADMM, Spingarn's method

- J. E. Spingarn, *Applications of the method of partial inverses to convex programming: decomposition*, Mathematical Programming (1985)
- J. Eckstein and D. Bertsekas, *On the Douglas-Rachford splitting method and the proximal algorithm for maximal monotone operators*, Mathematical Programming (1992)
- P.L. Combettes and J.-C. Pesquet, *A Douglas-Rachford splitting approach to nonsmooth convex variational signal recovery*, IEEE Journal of Selected Topics in Signal Processing (2007)
- S. Boyd, N. Parikh, E. Chu, B. Peleato, J. Eckstein, *Distributed optimization and statistical learning via the alternating direction method of multipliers* (2010)
- N. Parikh, S. Boyd, *Block splitting for distributed optimization* (2013)

image deblurring: the example is taken from

D. O'Connor and L. Vandenberghe, *Primal-dual decomposition by operator splitting and applications to image deblurring* (2014)