

ADMM 算法

王璐

Alternating direction method of multipliers (ADMM) 算法建立在在凸优化算法的基础上, 如 dual ascent, augmented Lagrangian method 等, 它在统计和机器学习问题中有广泛应用, 比如 lasso, group lasso, 稀疏协方差矩阵的估计等 (Boyd et al., 2011).

1 Dual Ascent

考虑如下带等式限制条件的凸优化问题:

$$\begin{aligned} \min_{\mathbf{x}} \quad & f(\mathbf{x}) \\ \text{s.t.} \quad & A\mathbf{x} = \mathbf{b} \end{aligned} \tag{1}$$

其中 $\mathbf{x} \in \mathbb{R}^n$, $A \in \mathbb{R}^{m \times n}$, $f: \mathbb{R}^n \rightarrow \mathbb{R}$ 是一个凸函数。

(1)的 Lagrangian 为

$$L(\mathbf{x}, \boldsymbol{\lambda}) = f(\mathbf{x}) + \boldsymbol{\lambda}^\top (A\mathbf{x} - \mathbf{b}). \tag{2}$$

(1)的 dual objective 为

$$\Theta_D(\boldsymbol{\lambda}) = \min_{\mathbf{x}} L(\mathbf{x}, \boldsymbol{\lambda}). \tag{3}$$

当 strong duality 成立时, primal optimization 和 dual optimization 的函数最优值相等:

$$\min_{\mathbf{x}} \left[\max_{\boldsymbol{\lambda}} L(\mathbf{x}, \boldsymbol{\lambda}) \right] = \max_{\boldsymbol{\lambda}} \Theta_D(\boldsymbol{\lambda}). \tag{4}$$

Dual ascent 方法是使用梯度上升法求解 dual optimization (4). 令

$$\mathbf{x}^* = \operatorname{argmin}_{\mathbf{x}} L(\mathbf{x}, \boldsymbol{\lambda})$$

则

$$\nabla \Theta_D(\boldsymbol{\lambda}) = A\mathbf{x}^* - \mathbf{b} \tag{5}$$

即 dual objective (3)的梯度为等式限制条件的残差 (residual). 因此 dual ascent 方法可总结为按如下迭代不断更新 \mathbf{x} 和 $\boldsymbol{\lambda}$:

$$\begin{aligned}\mathbf{x}^{(t+1)} &= \underset{\mathbf{x}}{\operatorname{argmin}} L(\mathbf{x}, \boldsymbol{\lambda}^{(t)}) \\ \boldsymbol{\lambda}^{(t+1)} &= \boldsymbol{\lambda}^{(t)} + \alpha_t (A\mathbf{x}^{(t+1)} - \mathbf{b})\end{aligned}\tag{6}$$

其中 $\alpha_t > 0$ 是第 t 步在梯度方向上移动的步长。如果每步选择合适的步长 α_t , dual objective (3)随迭代进行会不断增大, 即 $\Theta_D(\boldsymbol{\lambda}^{(t+1)}) > \Theta_D(\boldsymbol{\lambda}^{(t)})$. 当算法(6)收敛时, $(A\mathbf{x}^{(t+1)} - \mathbf{b})$ 会收敛到 0, 保证得到的解 \mathbf{x}^* 是 primal feasible. 在一些假设条件成立的情况下, 比如 f 是有界的严格凸函数, dual ascent 算法(6)最终会收敛到 $(\mathbf{x}, \boldsymbol{\lambda})$ 的最优解, 但是在很多应用中, 这些假设条件并不满足, 导致 dual ascent 失效。

2 Augmented Lagrangian and the Method of Multipliers

Augmented Lagrangian 方法可以增强 dual ascent 的稳定性 (robustness), 它可以放松 dual ascent 的一些假设条件, 比如严格凸或有界。

优化问题(1)的 augmented Lagrangian 定义为:

$$L_\rho(\mathbf{x}, \boldsymbol{\lambda}) = f(\mathbf{x}) + \boldsymbol{\lambda}^\top (A\mathbf{x} - \mathbf{b}) + \frac{\rho}{2} \|A\mathbf{x} - \mathbf{b}\|_2^2\tag{7}$$

其中 $\rho > 0$ 是惩罚系数。当 $\rho = 0$ 时, (7)退化为标准的 Lagrangian (2)。

Augmented Lagrangian (7)可以看作以下优化问题的 Lagrangian 函数:

$$\begin{aligned}\min_{\mathbf{x}} \quad & f(\mathbf{x}) + \frac{\rho}{2} \|A\mathbf{x} - \mathbf{b}\|_2^2 \\ \text{s.t.} \quad & A\mathbf{x} = \mathbf{b}\end{aligned}\tag{8}$$

注意到该优化问题(8)与(1)是等价的, 因为在满足等式限制条件的 \mathbf{x} 中, (8)中的 $\rho/2 \|A\mathbf{x} - \mathbf{b}\|_2^2$ 对目标函数的贡献为 0. 将(8)的 dual objective 记为:

$$\Theta_{D,\rho}(\boldsymbol{\lambda}) = \min_{\mathbf{x}} L_\rho(\mathbf{x}, \boldsymbol{\lambda}).\tag{9}$$

加入 $\rho/2 \|A\mathbf{x} - \mathbf{b}\|_2^2$ 的目的是使(8)中的目标函数变为严格凸 (strictly convex) 函数, 以避免在(6)中更新 \mathbf{x} 时出现 \mathbf{x} 的某些分量为 $\pm\infty$ 的情况 (此时可能导致无法更新 $\boldsymbol{\lambda}$)。

令

$$\mathbf{x}^* = \underset{\mathbf{x}}{\operatorname{argmin}} L_\rho(\mathbf{x}, \boldsymbol{\lambda})$$

则 augmented dual objective (9)的梯度也为

$$\nabla \Theta_{D,\rho}(\boldsymbol{\lambda}) = A\mathbf{x}^* - \mathbf{b}.$$

所以对优化问题(8)使用 dual ascent 算法可总结为以下迭代:

$$\begin{aligned}\mathbf{x}^{(t+1)} &= \underset{\mathbf{x}}{\operatorname{argmin}} L_\rho(\mathbf{x}, \boldsymbol{\lambda}^{(t)}) \\ \boldsymbol{\lambda}^{(t+1)} &= \boldsymbol{\lambda}^{(t)} + \rho(A\mathbf{x}^{(t+1)} - \mathbf{b})\end{aligned}\tag{10}$$

这里我们将每步步长 α_t 取为 ρ , 原因如下: 当 $\mathbf{x}^{(t+1)}$ 最小化 $L_\rho(\mathbf{x}, \boldsymbol{\lambda}^{(t)})$ 时,

$$\begin{aligned}0 &= \nabla_{\mathbf{x}} L_\rho(\mathbf{x}^{(t+1)}, \boldsymbol{\lambda}^{(t)}) \\ &= \nabla_{\mathbf{x}} f(\mathbf{x}^{(t+1)}) + A^\top [\boldsymbol{\lambda}^{(t)} + \rho(A\mathbf{x}^{(t+1)} - \mathbf{b})] \\ &= \nabla_{\mathbf{x}} f(\mathbf{x}^{(t+1)}) + A^\top \boldsymbol{\lambda}^{(t+1)}.\end{aligned}$$

因此在算法(10)中使用步长 ρ 可以保证每步更新的 $(\mathbf{x}^{(t+1)}, \boldsymbol{\lambda}^{(t+1)})$ 满足优化问题(1)的 Lagrangian stationary 条件:

$$\nabla_{\mathbf{x}} L(\mathbf{x}^{(t+1)}, \boldsymbol{\lambda}^{(t+1)}) = \nabla_{\mathbf{x}} f(\mathbf{x}^{(t+1)}) + A^\top \boldsymbol{\lambda}^{(t+1)} = 0.\tag{11}$$

算法(10)被称为 method of multipliers. 虽然该方法使 dual ascent 算法更稳定, 它也有一个缺点: 当目标函数 f separable 时, 即

$$f(\mathbf{x}) = \sum_{i=1}^n f_i(x_i)\tag{12}$$

对应的 augmented Lagrangian L_ρ (7) 并不 separable. 因此在(10)中不能单独更新 \mathbf{x} 的各个分量 x_i 进行求解. ADMM 算法通过轮流更新 \mathbf{x} 的各个分量将直接求解高维优化问题转化为求解一系列较简单的低维优化问题, 提高每步迭代的计算效率。

3 ADMM 算法

考虑具有如下形式的优化问题:

$$\begin{aligned}\min_{\mathbf{x}, \mathbf{z}} \quad & f(\mathbf{x}) + g(\mathbf{z}) \\ \text{s.t.} \quad & A\mathbf{x} + B\mathbf{z} = \mathbf{c}\end{aligned}\tag{13}$$

其中 f 和 g 都是凸函数, $\mathbf{x} \in \mathbb{R}^n$, $\mathbf{z} \in \mathbb{R}^m$, $A \in \mathbb{R}^{p \times n}$, $B \in \mathbb{R}^{p \times m}$, $\mathbf{c} \in \mathbb{R}^p$. 注意此时(13)中的目标函数关于 \mathbf{x} 和 \mathbf{z} separable.

写出(13)的 augmented Lagrangian:

$$L_\rho(\mathbf{x}, \mathbf{z}, \boldsymbol{\lambda}) = f(\mathbf{x}) + g(\mathbf{z}) + \boldsymbol{\lambda}^\top (A\mathbf{x} + B\mathbf{z} - \mathbf{c}) + \frac{\rho}{2} \|A\mathbf{x} + B\mathbf{z} - \mathbf{c}\|_2^2.\tag{14}$$

ADMM 算法求解(13)的过程可总结为以下迭代：

$$\begin{aligned} \mathbf{x}^{(t+1)} &= \underset{\mathbf{x}}{\operatorname{argmin}} L_{\rho}(\mathbf{x}, \mathbf{z}^{(t)}, \boldsymbol{\lambda}^{(t)}) \\ \mathbf{z}^{(t+1)} &= \underset{\mathbf{z}}{\operatorname{argmin}} L_{\rho}(\mathbf{x}^{(t+1)}, \mathbf{z}, \boldsymbol{\lambda}^{(t)}) \\ \boldsymbol{\lambda}^{(t+1)} &= \boldsymbol{\lambda}^{(t)} + \rho(A\mathbf{x}^{(t+1)} + B\mathbf{z}^{(t+1)} - \mathbf{c}) \end{aligned} \quad (15)$$

其中 $\rho > 0$. ADMM 算法(15)与 method of multipliers (10)十分相似, 不同之处是 ADMM 每步轮流更新变量 \mathbf{x} 和 \mathbf{z} , 而算法(10)联合更新 (\mathbf{x}, \mathbf{z}) :

$$\begin{aligned} (\mathbf{x}^{(t+1)}, \mathbf{z}^{(t+1)}) &= \underset{\mathbf{x}, \mathbf{z}}{\operatorname{argmin}} L_{\rho}(\mathbf{x}, \mathbf{z}, \boldsymbol{\lambda}^{(t)}) \\ \boldsymbol{\lambda}^{(t+1)} &= \boldsymbol{\lambda}^{(t)} + \rho(A\mathbf{x}^{(t+1)} + B\mathbf{z}^{(t+1)} - \mathbf{c}). \end{aligned}$$

这也正是 ADMM 中 alternating direction 得名的原因。

3.1 Scaled Form

如果对(14)中的 $\boldsymbol{\lambda}$ 做一些放缩, ADMM 算法(15)可写为更容易求解的形式。

将残差记为 $\mathbf{r} = A\mathbf{x} + B\mathbf{z} - \mathbf{c}$, 则有

$$\begin{aligned} \boldsymbol{\lambda}^{\top} \mathbf{r} + \frac{\rho}{2} \|\mathbf{r}\|_2^2 &= \frac{\rho}{2} \left(\mathbf{r} + \frac{1}{\rho} \boldsymbol{\lambda} \right)^{\top} \left(\mathbf{r} + \frac{1}{\rho} \boldsymbol{\lambda} \right) - \frac{1}{2\rho} \boldsymbol{\lambda}^{\top} \boldsymbol{\lambda} \\ &= \frac{\rho}{2} \|\mathbf{r} + \mathbf{u}\|_2^2 - \frac{\rho}{2} \|\mathbf{u}\|_2^2 \end{aligned} \quad (16)$$

其中新变量 $\mathbf{u} \triangleq \boldsymbol{\lambda}/\rho$, 被称为 scaled dual variables.

将(16)代入 augmented Lagrangian (14), ADMM 算法(15)可写为如下更方便计算的 scaled form:

$$\begin{aligned} \mathbf{x}^{(t+1)} &= \underset{\mathbf{x}}{\operatorname{argmin}} f(\mathbf{x}) + \frac{\rho}{2} \left\| A\mathbf{x} + B\mathbf{z}^{(t)} - \mathbf{c} + \mathbf{u}^{(t)} \right\|_2^2 \\ \mathbf{z}^{(t+1)} &= \underset{\mathbf{z}}{\operatorname{argmin}} g(\mathbf{z}) + \frac{\rho}{2} \left\| A\mathbf{x}^{(t+1)} + B\mathbf{z} - \mathbf{c} + \mathbf{u}^{(t)} \right\|_2^2 \\ \mathbf{u}^{(t+1)} &= \mathbf{u}^{(t)} + A\mathbf{x}^{(t+1)} + B\mathbf{z}^{(t+1)} - \mathbf{c}. \end{aligned} \quad (17)$$

将每步的残差记为 $\mathbf{r}^{(t)} = A\mathbf{x}^{(t)} + B\mathbf{z}^{(t)} - \mathbf{c}$, 由(17)可得

$$\mathbf{u}^{(T)} = \mathbf{u}^{(0)} + \sum_{t=1}^T \mathbf{r}^{(t)}$$

即 $\mathbf{u}^{(T)}$ 是前 T 步残差的累加。

3.2 ADMM 收敛性

Boyd et al. (2011) 证明了如下有关 ADMM 收敛性的定理。

Theorem 1. 当优化问题(13)满足以下假设条件时:

- **假设 1.** 函数 $f: \mathbb{R}^n \rightarrow \mathbb{R}$ 和 $g: \mathbb{R}^m \rightarrow \mathbb{R}$ 是闭凸函数 (*closed convex functions*).
- **假设 2.** (13)的 *unaugmented Lagrangian* L_0 至少有一个驻点。

ADMM 算法(15)可以保证:

- 残差收敛: $t \rightarrow \infty$ 时, $\mathbf{r}^{(t)} \rightarrow \mathbf{0}$. 即迭代可以保证 $(\mathbf{x}^{(t)}, \mathbf{z}^{(t)})$ 趋于 *primal feasible*.
- 目标函数收敛: $t \rightarrow \infty$ 时, $f(\mathbf{x}^{(t)}) + g(\mathbf{z}^{(t)}) \rightarrow p^*$, 其中 $p^* = \inf \{f(\mathbf{x}) + g(\mathbf{z}) : A\mathbf{x} + B\mathbf{z} = \mathbf{c}\}$.
- *Dual variable convergence*: $t \rightarrow \infty$ 时, $\boldsymbol{\lambda}^{(t)} \rightarrow \boldsymbol{\lambda}^*$, 其中 $\boldsymbol{\lambda}^*$ 是(13)的一个 *dual optimal point*, 即 $\boldsymbol{\lambda}^* \in \underset{\boldsymbol{\lambda}}{\operatorname{argmax}} \Theta_D(\boldsymbol{\lambda})$.

Remarks.

1. 称函数 $f: \mathbb{R}^n \rightarrow \mathbb{R}$ 为**闭凸函数**当且仅当集合

$$\{(\mathbf{x}, t) \in \mathbb{R}^n \times \mathbb{R} : f(\mathbf{x}) \leq t\}$$

是一个非空的闭凸集 (closed convex set)。

2. 假设 1 保证了 ADMM 每步迭代(15)中 \mathbf{x} 和 \mathbf{z} 都是可解的。
3. 假设 1 并不要求函数 f 或 g 有界, 比如 f 可以是如下的 indicator function:

$$f = \begin{cases} 0, & \mathbf{x} \in \mathcal{C} \\ +\infty, & \mathbf{x} \notin \mathcal{C} \end{cases}$$

其中 \mathcal{C} 是一个非空的闭凸集。此时(17)中对 \mathbf{x} 的更新即是寻找一个二次函数在 \mathcal{C} 中的最优解。

4. 假设 2 表明存在 $(\mathbf{x}^*, \mathbf{z}^*, \boldsymbol{\lambda}^*)$ 使得

$$\nabla_{\mathbf{x}} L_0(\mathbf{x}^*, \mathbf{z}^*, \boldsymbol{\lambda}^*) = \nabla f(\mathbf{x}^*) + A^\top \boldsymbol{\lambda}^* = \mathbf{0} \quad (18)$$

$$\nabla_{\mathbf{z}} L_0(\mathbf{x}^*, \mathbf{z}^*, \boldsymbol{\lambda}^*) = \nabla g(\mathbf{z}^*) + B^\top \boldsymbol{\lambda}^* = \mathbf{0} \quad (19)$$

$$\nabla_{\boldsymbol{\lambda}} L_0(\mathbf{x}^*, \mathbf{z}^*, \boldsymbol{\lambda}^*) = A\mathbf{x}^* + B\mathbf{z}^* - \mathbf{c} = \mathbf{0}. \quad (20)$$

(20)表明此时 $(\mathbf{x}^*, \mathbf{z}^*)$ 是 primal feasible. (18) - (19)表明 $(\mathbf{x}^*, \mathbf{z}^*, \boldsymbol{\lambda}^*)$ 满足 Lagrangian stationary 条件, 由 KKT 条件可得 $(\mathbf{x}^*, \mathbf{z}^*)$ 是(13)的 primal optimal point, $\boldsymbol{\lambda}^*$ 是(13)的 dual optimal point.

3.3 ADMM 算法的终止条件 (Stopping Criterion)

根据 KKT 条件, $(\mathbf{x}^*, \mathbf{z}^*, \boldsymbol{\lambda}^*)$ 是优化问题(13)最优解的充分条件是(18) - (20). 接下来我们检查 ADMM 算法每步更新的 $(\mathbf{x}^{(t+1)}, \mathbf{z}^{(t+1)}, \boldsymbol{\lambda}^{(t+1)})$ 是否满足这些条件。

在 ADMM 算法(15)中, 由于 $\mathbf{z}^{(t+1)}$ 最小化 $L_\rho(\mathbf{x}^{(t+1)}, \mathbf{z}, \boldsymbol{\lambda}^{(t)})$, 所以

$$\begin{aligned} 0 &= \nabla_{\mathbf{z}} L_\rho(\mathbf{x}^{(t+1)}, \mathbf{z}^{(t+1)}, \boldsymbol{\lambda}^{(t)}) \\ &= \nabla g(\mathbf{z}^{(t+1)}) + B^\top \boldsymbol{\lambda}^{(t)} + \rho B^\top (A\mathbf{x}^{(t+1)} + B\mathbf{z}^{(t+1)} - \mathbf{c}) \\ &= \nabla g(\mathbf{z}^{(t+1)}) + B^\top (\boldsymbol{\lambda}^{(t)} + \rho \mathbf{r}^{(t+1)}) \\ &= \nabla g(\mathbf{z}^{(t+1)}) + B^\top \boldsymbol{\lambda}^{(t+1)}. \end{aligned} \quad (21)$$

(21)表明 $\mathbf{z}^{(t+1)}$ 和 $\boldsymbol{\lambda}^{(t+1)}$ 总满足条件(19). 这与 method of multipliers 的解总是满足 Lagrangian stationary 条件(11)的证明类似。

在(15)中, $\mathbf{x}^{(t+1)}$ 最小化 $L_\rho(\mathbf{x}, \mathbf{z}^{(t)}, \boldsymbol{\lambda}^{(t)})$, 则有

$$\begin{aligned} 0 &= \nabla_{\mathbf{x}} L_\rho(\mathbf{x}^{(t+1)}, \mathbf{z}^{(t)}, \boldsymbol{\lambda}^{(t)}) \\ &= \nabla f(\mathbf{x}^{(t+1)}) + A^\top \boldsymbol{\lambda}^{(t)} + \rho A^\top (A\mathbf{x}^{(t+1)} + B\mathbf{z}^{(t)} - \mathbf{c}) \\ &= \nabla f(\mathbf{x}^{(t+1)}) + A^\top [\boldsymbol{\lambda}^{(t)} + \rho \mathbf{r}^{(t+1)} + \rho B(\mathbf{z}^{(t)} - \mathbf{z}^{(t+1)})] \\ &= \nabla f(\mathbf{x}^{(t+1)}) + A^\top \boldsymbol{\lambda}^{(t+1)} + \rho A^\top B(\mathbf{z}^{(t)} - \mathbf{z}^{(t+1)}) \end{aligned}$$

即

$$\rho A^\top B(\mathbf{z}^{(t+1)} - \mathbf{z}^{(t)}) = \nabla f(\mathbf{x}^{(t+1)}) + A^\top \boldsymbol{\lambda}^{(t+1)}. \quad (22)$$

将(22)中等式左边的项记为

$$\mathbf{s}^{(t+1)} = \rho A^\top B(\mathbf{z}^{(t+1)} - \mathbf{z}^{(t)}). \quad (23)$$

称 $\mathbf{s}^{(t+1)}$ 为 *dual residual*, 称此时的 $\mathbf{r}^{(t+1)} = A\mathbf{x}^{(t+1)} + B\mathbf{z}^{(t+1)} - \mathbf{c}$ 为 *primal residual*.

总结一下, ADMM 的解是优化问题(13)的最优解的充分条件是(18) - (20), 第二个条件(19)对每步得到的 $(\mathbf{x}^{(t+1)}, \mathbf{z}^{(t+1)}, \boldsymbol{\lambda}^{(t+1)})$ 总成立, 其它两个条件与 dual residual $\mathbf{s}^{(t+1)}$ 和 primal residual $\mathbf{r}^{(t+1)}$ 有关. 定理1表明随着 ADMM 迭代的进行, 残差项 $\mathbf{s}^{(t+1)}$ 和 $\mathbf{r}^{(t+1)}$ 都会收敛到 $\mathbf{0}$ (Boyd et al. (2011) 在 Appendix A 中证明了 $t \rightarrow \infty$ 时, $B(\mathbf{z}^{(t+1)} - \mathbf{z}^{(t)}) \rightarrow \mathbf{0}$). 因此 ADMM 算法一个合理的终止条件是保证 primal and dual residuals 都很小, 即

$$\|\mathbf{r}^{(t+1)}\|_2 < \epsilon_p \text{ 且 } \|\mathbf{s}^{(t+1)}\|_2 < \epsilon_d$$

其中 $\epsilon_p > 0$ 和 $\epsilon_d > 0$ 是一些很小的临界值。

3.4 ADMM 应用举例

在很多机器学习的问题中，最小化的目标函数通常由一个可导的损失函数和一个惩罚项组成。ADMM 算法的优势在于它可以将原目标函数拆分为两个独立的函数 f 和 g 分别优化，这为处理不可导的惩罚项，比如 L_1 惩罚项，提供了一种新思路。

3.4.1 加入 L_1 惩罚的统计模型的估计

令 $l(\beta)$ 表示一个统计模型的对数似然函数，比如 logistic regression, Poisson regression 或任意广义线性模型，它是一个可导的凹函数。为保证模型的稀疏性，我们在参数估计中加入对系数向量 $\beta \in \mathbb{R}^p$ 的 L_1 范数的惩罚，则 β 的最优估计为以下凸优化问题的解：

$$\min_{\beta} -l(\beta) + \lambda \|\beta\|_1. \quad (24)$$

(24)可以写为如下等价形式，然后使用 ADMM 求解：

$$\begin{aligned} \min_{\beta, z} \quad & -l(\beta) + g(z) \\ \text{s.t.} \quad & \beta = z \end{aligned} \quad (25)$$

其中 $g(z) = \lambda \|z\|_1$.

ADMM 求解(25)的迭代格式 (scaled form) 为：

$$\begin{aligned} \beta^{(t+1)} &= \operatorname{argmin}_{\beta} -l(\beta) + \frac{\rho}{2} \left\| \beta - z^{(t)} + u^{(t)} \right\|_2^2 \\ z^{(t+1)} &= \operatorname{argmin}_z \lambda \|z\|_1 + \frac{\rho}{2} \left\| \beta^{(t+1)} - z + u^{(t)} \right\|_2^2 \\ u^{(t+1)} &= u^{(t)} + \beta^{(t+1)} - z^{(t+1)}. \end{aligned} \quad (26)$$

其中 $z^{(t+1)}$ 的每个分量都有解析解：

$$z_j^{(t+1)} = \operatorname{sign} \left(\beta_j^{(t+1)} + u_j^{(t)} \right) \left(\left| \beta_j^{(t+1)} + u_j^{(t)} \right| - \frac{\lambda}{\rho} \right)_+, \quad j = 1, \dots, p.$$

在(26)中，对 β 的更新涉及最小化一个可导的凸函数，因此总可以用一些经典方法求解，比如 Newton 迭代法。

3.4.2 逆协方差矩阵的稀疏估计

假设样本 $x_i \in \mathbb{R}^p$ 服从期望为 $\mathbf{0}$ 的多元正态分布：

$$x_i \stackrel{iid}{\sim} N_p(\mathbf{0}, \Sigma), \quad i = 1, \dots, n. \quad (27)$$

在逆协方差矩阵 (precision matrix) Σ^{-1} 中, 第 (j, k) 元素 $(\Sigma^{-1})_{jk} = 0$ 表明: 给定随机向量 \mathbf{x} 其它分量的取值, \mathbf{x} 的第 j 和第 k 个分量是条件独立的。很多 graphical models 假设 Σ^{-1} 是一个稀疏矩阵 (Friedman et al., 2008) 以便对条件相关的变量做一个筛选。筛选后, 把随机向量 \mathbf{x} 的 p 个分量看作 p 个节点, 将条件相关的分量用边连接就得到了一个无向图, 这是 graphical models 得名的原因。此时 Σ^{-1} 越稀疏得到的无向图就越稀疏。当样本 n 较小时 ($n < p$), 对 Σ^{-1} 做稀疏性假设可以减少参数个数, 使估计结果更稳定。

对模型(27)重新参数化, 令

$$\Theta = \Sigma^{-1}. \quad (28)$$

则 n 个样本的似然函数为

$$\begin{aligned} l(\Theta) &\propto [\det(\Theta)]^{n/2} \exp \left\{ -\frac{1}{2} \sum_{i=1}^n \mathbf{x}_i^\top \Theta \mathbf{x}_i \right\} \\ &\propto [\det(\Theta)]^{n/2} \exp \left\{ -\frac{1}{2} \sum_{i=1}^n \text{tr}(\mathbf{x}_i^\top \Theta \mathbf{x}_i) \right\} \\ &\propto [\det(\Theta)]^{n/2} \exp \left\{ -\frac{1}{2} \text{tr} \left(\sum_{i=1}^n \mathbf{x}_i \mathbf{x}_i^\top \Theta \right) \right\}. \end{aligned} \quad (29)$$

将样本协方差矩阵记为

$$S = \frac{1}{n} \sum_{i=1}^n \mathbf{x}_i \mathbf{x}_i^\top. \quad (30)$$

为得到 Θ 的稀疏估计, 考虑惩罚矩阵 Θ 中非对角线 (off-diagonal) 元素的绝对值, 则 Θ 的估计值即为以下优化问题的解:

$$\begin{aligned} \min_{\Theta \in \mathcal{S}_+} \quad & -\frac{2}{n} \log l(\Theta) + \lambda \|\Theta\|_1 \\ = \min_{\Theta \in \mathcal{S}_+} \quad & -\log [\det(\Theta)] + \text{tr}(S\Theta) + \lambda \|\Theta\|_1 \end{aligned} \quad (31)$$

其中 \mathcal{S}_+ 表示所有 $p \times p$ 对称正定矩阵的集合, $\|\Theta\|_1$ 表示 Θ 非对角线元素绝对值的和。(31)是一个凸优化问题, 因为函数 $f(\Theta) = -\log [\det(\Theta)] + \text{tr}(S\Theta)$ 是凸函数。文献中有很多算法可以求解(31), 比如 neighborhood selection with lasso (Meinshausen et al., 2006), graphical lasso (Friedman et al., 2008), interior point algorithm (Yuan and Lin, 2007), projected subgradient method (Duchi et al., 2012), smoothing method (Lu, 2009) 等。Scheinberg et al. (2010) 使用 ADMM 算法求解(31), 并展示它的效率超越后两种算法。

将优化问题(31)写为以下等价形式:

$$\begin{aligned} \min_{\Theta \in \mathcal{S}_+} \quad & -\log [\det(\Theta)] + \text{tr}(S\Theta) + \lambda \|\Theta\|_1 \\ \text{s.t.} \quad & \Theta = Z. \end{aligned} \quad (32)$$

使用 ADMM 求解(32)的迭代格式 (scaled form) 为:

$$\begin{aligned}\Theta^{(t+1)} &= \underset{\Theta \in \mathcal{S}_+}{\operatorname{argmin}} -\log[\det(\Theta)] + \operatorname{tr}(S\Theta) + \frac{\rho}{2} \left\| \Theta - Z^{(t)} + U^{(t)} \right\|_F^2 \\ Z^{(t+1)} &= \underset{Z}{\operatorname{argmin}} \lambda \|Z\|_1 + \frac{\rho}{2} \left\| \Theta^{(t+1)} - Z + U^{(t)} \right\|_F^2 \\ U^{(t+1)} &= U^{(t)} + \Theta^{(t+1)} - Z^{(t+1)}\end{aligned}\tag{33}$$

其中 $\|\cdot\|_F$ 为矩阵的 Frobenius norm, 即矩阵中所有元素的平方和再开根号。

(33)中对矩阵 Z 每个元素的更新存在解析解:

$$Z_{jk}^{(t+1)} = \begin{cases} \operatorname{sign}(\Theta_{jk}^{(t+1)} + U_{jk}^{(t)}) \left(\left| \Theta_{jk}^{(t+1)} + U_{jk}^{(t)} \right| - \frac{\lambda}{\rho} \right)_+, & j \neq k \\ \Theta_{jk}^{(t+1)} + U_{jk}^{(t)}, & j = k. \end{cases}\tag{34}$$

(33)中对矩阵 Θ 的更新也存在解析解。令 Θ 的一阶导数为零得

$$-\Theta^{-1} + S + \rho(\Theta - Z^{(t)} + U^{(t)}) = \mathbf{0}.\tag{35}$$

由于 Θ 是对称正定矩阵, 我们需要找到满足

$$\rho\Theta - \Theta^{-1} = \rho(Z^{(t)} - U^{(t)}) - S\tag{36}$$

的对称正定矩阵 Θ . 由于(36)等式两端都是对称矩阵, 所以存在特征值分解, 令

$$\rho(Z^{(t)} - U^{(t)}) - S = Q\Lambda Q^\top$$

其中 Q 是单位正交矩阵, $Q^\top Q = QQ^\top = I$, Λ 是对角阵 $\Lambda = \operatorname{diag}(\lambda_1, \dots, \lambda_p)$. 对等式(36)两端左乘 Q^\top 右乘 Q 得

$$\rho Q^\top \Theta Q - Q^\top \Theta^{-1} Q = \Lambda.\tag{37}$$

令 $\tilde{\Theta} = Q^\top \Theta Q$, 则 $\tilde{\Theta}^{-1} = Q^\top \Theta^{-1} Q$. 由(37)得

$$\rho \tilde{\Theta} - \tilde{\Theta}^{-1} = \Lambda.\tag{38}$$

我们可以找到满足(38)的一个对角阵 $\tilde{\Theta}$. 假设 $\tilde{\Theta} = \operatorname{diag}(\theta_1, \dots, \theta_p)$, 则对角线元素满足

$$\rho\theta_j - \frac{1}{\theta_j} = \lambda_j, \quad j = 1, \dots, p.\tag{39}$$

所以 $\theta_j > 0$ 的解为

$$\theta_j = \frac{\lambda_j + \sqrt{\lambda_j^2 + 4\rho}}{2\rho}, \quad j = 1, \dots, p.$$

由此得到的 $\Theta = Q\tilde{\Theta}Q^\top$ 一定是对称正定矩阵且满足一阶条件(35), 记为 $\Theta^{(t+1)}$.

References

- Boyd, S., Parikh, N., Chu, E., Peleato, B., Eckstein, J., et al. (2011). Distributed optimization and statistical learning via the alternating direction method of multipliers. *Foundations and Trends® in Machine learning*, 3(1):1–122.
- Duchi, J., Gould, S., and Koller, D. (2012). Projected subgradient methods for learning sparse gaussians. *arXiv preprint arXiv:1206.3249*.
- Friedman, J., Hastie, T., and Tibshirani, R. (2008). Sparse inverse covariance estimation with the graphical lasso. *Biostatistics*, 9(3):432–441.
- Lu, Z. (2009). Smooth optimization approach for sparse covariance selection. *SIAM Journal on Optimization*, 19(4):1807–1827.
- Meinshausen, N., Bühlmann, P., et al. (2006). High-dimensional graphs and variable selection with the lasso. *The annals of statistics*, 34(3):1436–1462.
- Scheinberg, K., Ma, S., and Goldfarb, D. (2010). Sparse inverse covariance selection via alternating linearization methods. In *Advances in neural information processing systems*, pages 2101–2109.
- Yuan, M. and Lin, Y. (2007). Model selection and estimation in the gaussian graphical model. *Biometrika*, 94(1):19–35.