

非线性最小二乘问题

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 非线性最小二乘问题
- 2 高斯-牛顿方法
- 3 Levenberg-Marquardt 方法
- 4 大残量问题的拟牛顿算法

最小二乘问题

$$\min_x f(x) = \frac{1}{2} \sum_{j=1}^m r_j^2(x) \quad (1)$$

- 其中 $r_j: \mathbb{R}^n \rightarrow \mathbb{R}$ 是光滑函数, 并且假设 $m \geq n$. 称 r_j 为残差.
- 记 $r: \mathbb{R}^n \rightarrow \mathbb{R}^m$

$$r(x) = (r_1(x), r_2(x), \dots, r_m(x))^T.$$

问题可以表述为 $\min f(x) = \frac{1}{2} \|r(x)\|_2^2$

- 一般情况下不是凸问题
- 问题1是无约束优化问题, 可以直接使用线搜索或拟牛顿法求解

最小二乘问题

- 记 $J(x) \in \mathbb{R}^{m \times n}$ 是向量值函数 $r(x)$ 在点 x 处的雅可比矩阵：

$$J(x) = \begin{bmatrix} \nabla r_1(x)^T \\ \nabla r_2(x)^T \\ \vdots \\ \nabla r_m(x)^T \end{bmatrix}.$$

- $f(x)$ 的梯度和海瑟矩阵：

$$\nabla f(x) = \sum_{j=1}^m r_j(x) \nabla r_j(x) = J(x)^T r(x), \quad (2a)$$

$$\nabla^2 f(x) = \sum_{j=1}^m \nabla r_j(x) \nabla r_j(x) + \sum_{i=1}^m r_i(x) \nabla^2 r_i(x) \quad (2b)$$

$$= J(x)^T J(x) + \sum_{i=1}^m r_i(x) \nabla^2 r_i(x), \quad (2c)$$

最小二乘问题

$$\nabla^2 f(x) = J(x)^T J(x) + \sum_{i=1}^m r_i(x) \nabla^2 r_i(x)$$

- $\nabla^2 f(x)$ 在形式上分为两部分
- 在计算 $f(x)$ 导数时已经求出 $J(x)$ ，第一项可自然得到。第二项需要额外计算。
- 如果在最优解附近，残量值较小或残量函数接近线性函数，第二项可以忽略，可以用 $J(x)^T J(x)$ 近似海瑟矩阵，基于牛顿法，结合线搜索或信赖域方法，可设计出高斯-牛顿方法和 Levenberg-Marquardt 方法
- 如果第二项不可忽略，则需要引入带结构的拟牛顿方法。

最小二乘问题

- 若 $r(x) = Ax - b$ ，则是线性最小二乘问题，
- $\nabla f(x) = A^T(Ax - b)$. 最优解满足正则化方程：

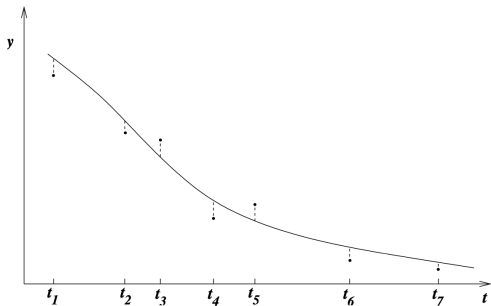
$$A^T A x = A^T b$$

- 线性问题的求解是非线性问题求解的基础. 求解线性最小二乘问题的方法有：
 - 对 $A^T A$ 直接做choloskey分解，简单便捷但受 A 的条件数影响大。
 - 对 A 做QR分解，较稳定，相对误差小。
 - 对 A 做SVD分解，可以获得更精确的敏感性信息。
 - 当问题规模较大时，迭代法更有效，如共轭梯度法。

实例:模型拟合

$$\min_x \frac{1}{2} \sum_{j=1}^n (\phi(t_j; x) - y_j)^2$$

- 模型 $\phi(t; x) = x_1 + tx_2 + t_3^x + x_4 e^{-x_5 t}$ 依赖于参数向量 x 。
- $(t_1, y_1), (t_2, y_2), \dots, (t_n, y_n)$ 是数据点
- 目标为寻找合适的模型参数 x



提纲

1 非线性最小二乘问题

2 高斯-牛顿方法

3 Levenberg-Marquardt 方法

4 大残量问题的拟牛顿算法

高斯-牛顿方法

- 可被看作牛顿法+线搜索
- 在迭代点 x_k ，标准牛顿法为，计算更新方向 d_k ：

$$\nabla^2 f(x_k) d_k = -\nabla f(x_k)$$

然后做一步更新 $x_{k+1} = x_k + d_k$.

- 高斯-牛顿法的迭代方向 d_k^{GN} 满足：

$$J_k^T J_k d_k^{GN} = -J_k^T r_k. \quad (3)$$

其中 J_k 、 r_k 分别是 $J(x_k)$ 、 $r(x_k)$ 的简写

- 使用了近似

$$\nabla^2 f_k \approx J_k^T J_k$$

省略了对 $\nabla^2 r_j$ 的计算，极大的减少了计算量。

高斯-牛顿方法

- 方程(3)与线性最小二乘问题的正则化方程类似，迭代方向是如下问题的解

$$\min_d \frac{1}{2} \|J_k d + r_k\|^2$$

- 求解该问题时，可以直接对 J_k 做QR分解或SVD分解，无需计算出 $J_k^T J_k$ 。
- 若使用共轭梯度法求解，需要计算向量和矩阵 $J_k^T J_k$ 的乘法，可以依次乘 J_k 和 J_k^T ，无需计算 $J_k^T J_k$ 。
- 另一种理解高斯-牛顿方法的方式为，在点 x_k 处，考虑下一步更新 $x_k + d$ ，做近似 $r(x_k + d) \approx r_k + J_k d$ ，原问题近似为：

$$\min_d f(x_k + d) = \frac{1}{2} \|r(x_k + d)\|^2 \approx \frac{1}{2} \|J_k d + r_k\|^2$$

高斯-牛顿方法

算法可总结如下

Algorithm 1 高斯- 牛顿法

- 1: 给定初始值 x_0 , $k \leftarrow 0$.
 - 2: **while** 未达到收敛准则 **do**
 - 3: 计算残差向量 r_k , 雅可比矩阵 J_k .
 - 4: 求解线性最小二乘问题 $\min_d \frac{1}{2} \|J_k d + r_k\|^2$ 确定下降方向 d_k .
 - 5: 使用线搜索准则计算步长 α_k .
 - 6: 更新: $x_{k+1} = x_k + \alpha_k d_k$.
 - 7: $k \leftarrow k + 1$.
 - 8: **end while**
-

全局收敛性分析

- 线搜索条件可选择Armijo或Wolfe.
- 若 J_k 满秩且 ∇f_k 非零, 则 d_k^{GN} 是一个下降方向:

$$\begin{aligned}(d_k)^T \nabla f(x_k) &= d_k^T J_k^T r_k = -d_k^T J_k^T J_k d_k \\ &= -\|J_k d_k\|^2 \leq 0.\end{aligned}$$

- 那么 d^k 是一个合适的线搜索方向, 全局收敛性的证明可以套用线搜索的证明。

全局收敛性分析

- 注意到，雅可比矩阵 J_k 的非奇异性很关键，在这个条件下建立收敛性.
- 具体为：假设雅可比矩阵 $J(x)$ 的奇异值一致地大于0，即存在 $\gamma > 0$ 使得

$$\|J(x)z\| \geq \gamma\|z\|, \quad \forall x \in \mathcal{N}, \quad (4)$$

其中 \mathcal{N} 是下水平集

$$\mathcal{L} = \{x \mid f(x) \leq f(x_0)\} \quad (5)$$

的一个邻域， x_0 是算法的初始点，且假设 \mathcal{L} 是有界的.

全局收敛性

Theorem

全局收敛性 如果每个残差函数 r_j 在有界下水平集(5)的一个邻域 \mathcal{N} 内是利普希茨连续可微的, 并且雅可比矩阵 $J(x)$ 在 \mathcal{N} 内满足一致满秩条件(4), 而步长满足 *Wolfe* 准则, 则对高斯-牛顿法得到的序列 $\{x_k\}$ 有

$$\lim_{k \rightarrow \infty} (J_k)^T r_k = 0.$$

局部收敛性分析

- 在最优值点 x^* 附近, 当 $\sum_{i=1}^m r_i(x) \nabla^2 r_i(x)$ 较小时 ($r(x^*)$ 很小或在 x^* 附近 r 接近仿射函数), $J_k^T J_k$ 占主导位置, 高斯-牛顿方法有类似牛顿法的收敛速度。

Theorem (局部收敛性)

设 $r_i(x)$ 二阶连续可微, x^* 是最小二乘问题(1)的最优解, 海瑟矩阵 $\nabla^2 f(x)$ 和其近似矩阵 $J(x)^T J(x)$ 均在点 x^* 的一个邻域内利普希茨连续, 则当高斯-牛顿算法步长 α_k 恒为1时,

$$\|x_{k+1} - x^*\| \leq C \|((J^*)^T J^*)^{-1} H^*\| \|x_k - x^*\| + \mathcal{O}(\|x_k - x^*\|^2), \quad (6)$$

其中 $H^* = \sum_{i=1}^m r(x^*) \nabla^2 r(x^*)$ 为海瑟矩阵 $\nabla^2 f(x^*)$ 去掉 $J(x^*)^T J(x^*)$ 的部分, $C > 0$ 为常数。

- 1 非线性最小二乘问题
- 2 高斯-牛顿方法
- 3 Levenberg-Marquardt 方法
- 4 大残量问题的拟牛顿算法

Levenberg-Marquardt 方法

- 当 J_k 不满秩时, (3) 有很多个解, 应该怎么更新?
- LM 方法本质为信赖域方法, 更新方向为如下问题的解

$$\min_d \quad \frac{1}{2} \|J^k d + r^k\|^2, \quad \text{s.t.} \quad \|d\| \leq \Delta_k. \quad (7)$$

- LM 方法将如下近似当作信赖域方法中的 m_k :

$$m_k(d) = \frac{1}{2} \|r^k\|^2 + d^T (J^k)^T r^k + \frac{1}{2} d^T (J^k)^T J^k d. \quad (8)$$

- 同样使用 $(J^k)^T J^k$ 来近似海瑟矩阵.

Levenberg-Marquardt 方法

- 类似信赖域方法，引入如下定义来衡量 $m_k(d)$ 近似程度的好坏：

$$\rho_k = \frac{f(x^k) - f(x^k + d^k)}{m_k(0) - m_k(d^k)} \quad (9)$$

为函数值实际下降量与预估下降量（即二阶近似模型下降量）的比值。

- 如果 ρ_k 接近1，说明 $m_k(d)$ 来近似 $f(x)$ 是比较成功的，则应该扩大 Δ_k ；如果 ρ_k 非常小甚至为负，就说明我们过分地相信了二阶模型 $m_k(d)$ ，此时应该缩小 Δ_k 。
- 只有当 ρ_k 足够大，也就是对模型拟合较好时，才进行一步更新，否则不更新。

Levenberg-Marquardt 方法

Algorithm 2 Levenberg-Marquardt 方法

- 1: 给定最大半径 Δ_{\max} , 初始半径 Δ_0 , 初始点 x^0 , $k \leftarrow 0$.
- 2: 给定参数 $0 \leq \eta < \bar{\rho}_1 < \bar{\rho}_2 < 1$, $\gamma_1 < 1 < \gamma_2$.
- 3: **while** 未达到收敛准则 **do**
- 4: 计算子问题(7)得到迭代方向 d^k .
- 5: 根据(9) 计算下降率 ρ_k .
- 6: 更新信赖域半径:

$$\Delta_{k+1} = \begin{cases} \gamma_1 \Delta_k, & \rho_k < \bar{\rho}_1, \\ \min\{\gamma_2 \Delta_k, \Delta_{\max}\}, & \rho_k > \bar{\rho}_2 \text{ 以及 } \|d^k\| = \Delta_k, \\ \Delta_k, & \text{其他.} \end{cases}$$

- 7: 更新自变量:

$$x^{k+1} = \begin{cases} x^k + d^k, & \rho_k > \eta, \\ x^k, & \text{其他.} \end{cases} \quad /* \text{ 只有下降比例足够大才更新} */$$

- 8: $k \leftarrow k + 1$.
- 9: **end while**

子问题求解

Corollary

向量 d^* 是信赖域子问题

$$\min_d \quad \frac{1}{2} \|Jd + r\|^2, \quad \text{s.t.} \quad \|d\| \leq \Delta$$

的解当且仅当 d^* 是可行解并且存在数 $\lambda \geq 0$ 使得

$$(J^T J + \lambda I) d^* = -J^T r, \quad (10)$$

$$\lambda(\Delta - \|d^*\|) = 0. \quad (11)$$

- 问题(10)等价于线性最小二乘问题，具体实现时可利用系数矩阵的结构

$$\min_d \quad \frac{1}{2} \left\| \begin{bmatrix} J \\ \sqrt{\lambda} I \end{bmatrix} p + \begin{bmatrix} r \\ 0 \end{bmatrix} \right\|^2.$$

子问题求解

$$\min_d \frac{1}{2} \left\| \begin{bmatrix} J \\ \sqrt{\lambda} I \end{bmatrix} p + \begin{bmatrix} r \\ 0 \end{bmatrix} \right\|^2.$$

- 在试探 λ 的值时, J 的块不变, 设 $J = QR$, 则

$$\begin{bmatrix} J \\ \sqrt{\lambda} I \end{bmatrix} = \begin{bmatrix} QR \\ \sqrt{\lambda} I \end{bmatrix} = \begin{bmatrix} Q & 0 \\ 0 & I \end{bmatrix} \begin{bmatrix} R \\ \sqrt{\lambda} I \end{bmatrix}.$$

- 矩阵 $\begin{bmatrix} R \\ \sqrt{\lambda} I \end{bmatrix}$ 有较多的零元素, 可以使用Household变换或Givens变换完成QR分解。
- 如果矩阵 J 没有显式形式, 只能提供矩阵乘法, 则仍然可以用截断共轭梯度法。

收敛性分析

Theorem

假设常数 $\eta \in (0, \frac{1}{4})$, 下水平集 \mathcal{L} 是有界的且每个 $r_i(x)$ 在下水平集 \mathcal{L} 的一个邻域 \mathcal{N} 内是利普希茨连续可微的. 假设对于任意的 k , 子问题(7)的近似解 d_k 满足

$$m_k(0) - m_k(d_k) \geq c_1 \|(J_k)^T r_k\| \min \left\{ \Delta_k, \frac{\|(J_k)^T r_k\|}{\|(J_k)^T J_k\|} \right\},$$

其中 $c_1 > 0$ 且 $\|d_k\| \leq \gamma \Delta_k, \gamma \geq 1$, 则

$$\lim_{k \rightarrow \infty} \nabla f(x_k) = \lim_{k \rightarrow \infty} (J_k)^T r_k = 0.$$

- 信赖域型LM方法本质上是固定信赖域半径 Δ , 通过迭代寻找满足条件的乘子 λ , 每一步迭代需要求解线性方程组

$$(J^T J + \lambda I) d = -J^T r$$

该步计算代价较大。

- 注意到在LM方法中, 由于 $J_k^T J_k \succ 0$, 那么有 $-\lambda_1 < 0$, 此时有 $\lambda > -\lambda_1$, 因此若 λ 越大, d 的模长就越小。
- 调整 λ 的大小等价于调整信赖域半径的大小, 这意味着, Δ 被 λ 隐式决定。

- LM的更新基于 Δ ，LMF的更新直接基于 λ ，每一步求解子问题：

$$\min_d \|Jd + r\|_2^2 + \lambda \|d\|_2^2.$$

- 调整 λ 的原则可以参考信赖域半径的调整原则
- 考虑参数 ρ_k

$$\rho_k = \frac{f(x_k) - f(x_k + d_k)}{m_k(0) - m_k(d_k)} \quad (12)$$

较大可以减小下一步的 λ ，较小可以增大下一步的 λ 。

Algorithm 3 LMF 方法

- 1: 给定初始点 x_0 , 初始乘子 λ_0 , $k \leftarrow 0$.
- 2: 给定参数 $0 \leq \eta < \bar{\rho}_1 < \bar{\rho}_2 < 1$, $\gamma_1 < 1 < \gamma_2$.
- 3: **while** 未达到收敛准则 **do**
- 4: 求解LM方程 $((J_k)^T J_k + \lambda I)d = -(J_k)^T r_k$ 得到迭代方向 d_k .
- 5: 根据(12)式计算下降率 ρ_k .
- 6: 更新乘子:

$$\lambda_{k+1} = \begin{cases} \gamma_2 \lambda_k, & \rho_k < \bar{\rho}_1, & /* \text{ 扩大乘子 (缩小信赖域半径) } */ \\ \gamma_1 \lambda_k, & \rho_k > \bar{\rho}_2, & /* \text{ 缩小乘子 (扩大信赖域半径) } */ \\ \lambda_k, & \text{其他.} & /* \text{ 乘子不变 } */ \end{cases}$$

- 7: 更新自变量:

$$x_{k+1} = \begin{cases} x_k + d_k, & \rho_k > \eta, & /* \text{ 只有下降比例足够大才更新 } */ \\ x_k, & \text{其他.} \end{cases}$$

- 8: $k \leftarrow k + 1$.
 - 9: **end while**
-

提纲

- 1 非线性最小二乘问题
- 2 高斯-牛顿方法
- 3 Levenberg-Marquardt 方法
- 4 大残量问题的拟牛顿算法

大残量问题的拟牛顿算法

- 大残量问题中，海瑟矩阵的第二部分不可忽视，此时高斯-牛顿法和LM方法可能只有线性的收敛速度。
- 此时如果直接使用牛顿法，则开销太大；直接使用拟牛顿法，又似乎忽略了问题的特殊结构。
- 重新写出海瑟矩阵：

$$\nabla^2 f(x) = J(x)^T J(x) + \sum_{i=1}^m r_i(x) \nabla^2 r_i(x)$$

第一项是容易求解，可以保留。第二项不易求解但不可忽略，用拟牛顿法进行近似。

大残量问题的拟牛顿算法

- 使用 B_k 来表示 $\nabla^2 f(x_k)$ 的近似矩阵, T_k 表示 $\sum_{j=1}^m r_j(x_k) \nabla^2 r_j(x_k)$ 的近似, 即

$$B_k = (J_k)^T J_k + T_k,$$

- 目标为

$$T_{k+1} \approx \sum_{j=1}^m r_j(x_{k+1}) \nabla^2 r_j(x_{k+1})$$

- 记 $s_k = x_{k+1} - x_k$, T_{k+1} 应该尽量保留原海瑟矩阵的性质

$$\begin{aligned} T_{k+1} s_k &\approx \sum_{j=1}^m r_j(x_{k+1}) (\nabla^2 r_j(x_{k+1})) s_k \\ &\approx \sum_{j=1}^m r_j(x_{k+1}) (\nabla r_j(x_{k+1}) - \nabla r_j(x_k)) \\ &= (J_{k+1})^T r_{k+1} - (J_k)^T r_{k+1}. \end{aligned}$$

大残量问题的拟牛顿算法

- 拟牛顿条件为：

$$T_{k+1}s_k = (J_{k+1})^T r_{k+1} - (J_k)^T r_{k+1}$$

- Dennis, Gay, 和Welsch给出的一种更新格式为：

$$T_{k+1} = T_k + \frac{(y^\# - T_k s_k) y^T + y (y^\# - T_k s_k)^T}{y^T s_k} - \frac{(y^\# - T_k s_k)^T s_k}{(y^T s)^2} y y^T$$

其中

$$s_k = x_{k+1} - x_k$$

$$y = J_{k+1}^T r_{k+1} - J_k^T r_k$$

$$y^\# = J_{k+1}^T r_{k+1} - J_k^T r_{k+1}$$