

第十二章 优化算法

第 34 讲 无约束优化算法：二阶优化算法

黄定江

DaSE @ ECNU

djhuang@dase.ecnu.edu.cn

1 34.1 牛顿法

2 34.2 拟牛顿法

1 34.1 牛顿法

2 34.2 拟牛顿法

- 从上一讲中可以看出，梯度法仅使用了目标函数的一阶信息。
- 如果函数足够光滑，那么就可以使用更多的信息，例如二阶信息。
- 直观上，可以期望得到更好的优化算法。这就是本讲我们将探究的牛顿类方法。

34.1.1 牛顿法

函数 f 在 \mathbf{x} 处的二阶 Taylor 近似 (或模型) \hat{f} 为

$$\hat{f}(\mathbf{x} + \mathbf{v}) \approx f(\mathbf{x}) + \nabla f(\mathbf{x})^T \mathbf{v} + \frac{1}{2} \mathbf{v}^T \nabla^2 f(\mathbf{x}) \mathbf{v} \quad (1)$$

这是 \mathbf{v} 的二次凸函数, 当 Hessian 矩阵为半正定时, 显然在 $\mathbf{v} = \Delta \mathbf{x}_{nt}$ 处达到最小值。如下图所示:

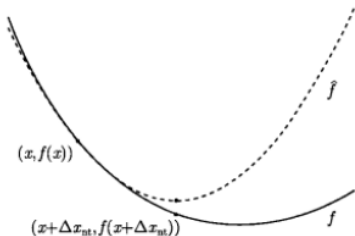


图 1

Newton 步径

若 Hessian 矩阵 $\nabla^2 f(\mathbf{x})$ 正定, 对二阶近似求极小, 可得牛顿方程 $\nabla^2 f(\mathbf{x}) \Delta \mathbf{x}_{nt} = -\nabla f(\mathbf{x})$, 并解得:

$$\Delta \mathbf{x}_{nt} = -\nabla^2 f(\mathbf{x})^{-1} \nabla f(\mathbf{x}).$$

它被称之为 (f 在 \mathbf{x} 处的) **Newton 步径**。

由正定性可知, 除非 $\nabla f(\mathbf{x}) = 0$, 否则就有

$$\nabla f(\mathbf{x})^T \Delta \mathbf{x}_{nt} = -\nabla f(\mathbf{x})^T \nabla^2 f(\mathbf{x})^{-1} \nabla f(\mathbf{x}) < 0$$

因此, Newton 步径是下降方向 (除非 \mathbf{x} 是最优点)。将 \mathbf{x} 加上 Newton 步径 $\Delta \mathbf{x}_{nt}$ 能够极小化 f 在 \mathbf{x} 处的二阶近似。

注: 当 $f(\mathbf{x})$ 就是正定二次函数时, 只需要一步便能求出最优解。

Newton 减量估计

通过简单计算，可以得到牛顿法的单步下降量约为 $f(\mathbf{x}) - \hat{f}(\mathbf{x} + \Delta \mathbf{x}_{nt})$ ，即：

$$f(\mathbf{x}) - \hat{f}(\mathbf{x} + \Delta \mathbf{x}_{nt}) \triangleq \frac{1}{2} \delta(\mathbf{x})^2,$$

其中 \hat{f} 仍是 f 在 \mathbf{x} 处的二阶近似。

注：这一量可作为牛顿法的终止判定准则。

牛顿法步骤

因此，可以得到牛顿法具体步骤如下：

Algorithm 1 牛顿法

- 1: 给定初始点 $\mathbf{x} \in \text{dom } f$, 误差阈值 $\epsilon > 0$
- 2: 计算 Newton 步径和减量。

$$\Delta \mathbf{x}_{nt} := -\nabla^2 f(\mathbf{x})^{-1} \nabla f(\mathbf{x}); \quad \delta^2 := \nabla f(\mathbf{x})^T \nabla^2 f(\mathbf{x})^{-1} \nabla f(\mathbf{x})$$

- 3: 停止准则。如果 $\delta^2/2 \leq \epsilon$, 退出。
 - 4: 直线搜索。通过回溯直线搜索确定步长 λ .
 - 5: 改进。 $\mathbf{x} := \mathbf{x} + \lambda \Delta \mathbf{x}_{nt}$
 - 6: 重复上述步骤，直至退出。
-

经典牛顿法有很好的局部收敛性质. 实际上我们有如下定理:

定理 1

(经典牛顿法的收敛性) 假设目标函数 f 是二阶连续可微的函数, 且海瑟矩阵在最优点 \mathbf{x}^* 的一个邻域 $N_\delta(\mathbf{x}^*)$ 内是利普希茨连续的, 即存在常数 $L > 0$ 使得

$$\|\nabla^2 f(\mathbf{x}) - \nabla^2 f(\mathbf{y})\| \leq L \|\mathbf{x} - \mathbf{y}\|, \quad \forall \mathbf{x}, \mathbf{y} \in N_\delta(\mathbf{x}^*)$$

如果函数 $f(\mathbf{x})$ 在点 \mathbf{x}^* 处满足 $\nabla f(\mathbf{x}^*) = 0, \nabla^2 f(\mathbf{x}^*) \succ 0$, 则对于步长为 1 时的迭代有如下结论:

- (1) 如果初始点离 \mathbf{x}^* 足够近, 则牛顿法产生的迭代点列 $\{\mathbf{x}^k\}$ 收敛到 \mathbf{x}^* ;
- (2) $\{\mathbf{x}^k\}$ 收敛到 \mathbf{x}^* 的速度是 Q-二次的;
- (3) $\{\|\nabla f(\mathbf{x}^k)\|\}$ Q-二次收敛到 0.

从上述定理可以看出牛顿法收敛速度快，但牛顿法也有如下缺陷：

- 函数的 Hessian 矩阵本身计算代价大，难以存储。
- Hessian 矩阵还可能面临着不正定的问题，应该如何修正？
- 在高维问题中，求解 Hessian 矩阵的逆（或者是解大规模线性方程组）的计算量更大。
- 能否以较小的代价找到 Hessian 矩阵的一个较好的近似？

这就是接下来要介绍的修正牛顿法和拟牛顿法（变尺度法）。

34.1.2 修正牛顿法

前面也已提及经典牛顿法有如下缺陷：

- 每一步迭代需要求解一个 n 维线性方程组, 这导致在高维问题中计算量很大. 海瑟矩阵 $\nabla^2 f(\mathbf{x})$ 既不容易计算又不容易储存.
- 当 $\nabla^2 f(\mathbf{x})$ 不正定时, 由牛顿方程给出的牛顿步径 $\Delta \mathbf{x}_{nt}$ 的性质通常比较差. 例如可以验证当海瑟矩阵正定时, $\Delta \mathbf{x}_{nt}$ 是一个下降方向, 而在其他情况下 $\Delta \mathbf{x}_{nt}$ 不一定为下降方向.
- 当迭代点距最优值较远时, 直接选取步长 $\alpha = 1$ 会使得迭代极其不稳定, 在有些情况下迭代点列会发散.

为了克服这些缺陷, 这里介绍带线搜索的修正牛顿法, 其基本思想是对牛顿方程中的海瑟矩阵 $\nabla^2 f(\mathbf{x})$ 进行修正, 使其变成正定矩阵. 它的一般框架如下所示:

Algorithm 2 修正牛顿法

- 1: 给定初始点 \mathbf{x}^0 .
 - 2: **for** $k = 0, 1, 2, \dots$ **do**
 - 3: 确定矩阵 \mathbf{E}^k 使得矩阵 $\mathbf{B}^k \stackrel{\text{def}}{=} \nabla^2 f(\mathbf{x}^k) + \mathbf{E}^k$ 正定且条件数较小.
 - 4: 求解修正的牛顿方程 $\mathbf{B}^k \mathbf{d}^k = -\nabla f(\mathbf{x}^k)$ 得方向 \mathbf{d}^k .
 - 5: 使用任意一种线搜索准则确定步长 α_k .
 - 6: 更新 $\mathbf{x}^{k+1} = \mathbf{x}^k + \alpha_k \mathbf{d}^k$.
 - 7: **end for**
-

该算法的关键在于修正矩阵 \mathbf{E}^k 如何选取.

显式修正

一个最直接的取法是取

$$\mathbf{E}^k = \tau_k \mathbf{I},$$

即取 \mathbf{E}^k 为单位矩阵的常数倍.

- 根据矩阵理论可以知道, 当 τ_k 充分大时, 总可以保证 \mathbf{B}^k 是正定矩阵.
- τ_k 不宜取得过大, 这是因为当 τ_k 趋于无穷时, \mathbf{d}^k 的方向会接近负梯度方向.
- 比较合适的取法是先估计 $\nabla^2 f(\mathbf{x}^k)$ 的最小特征值, 再适当选择 τ_k .

隐式修正

另一种 E^k 的选取是隐式的:

- 它是通过修正 Cholesky 分解的方式来求解牛顿方程.
- 我们知道当海瑟矩阵正定时, 方程组可以用 Cholesky 分解快速求解.
- 当海瑟矩阵不定或条件数较大时, Cholesky 分解会失败. 而修正 Cholesky 分解算法对基本 Cholesky 分解算法进行修正, 且修正后的分解和原矩阵相差不大.
 - 根据 Cholesky 分解的形式, 如果 A 正定且条件数较小, 矩阵 D 的对角线元素不应该太小.
 - 如果计算过程中发现 d_j 过小就应该及时修正. 同时我们需要保证该修正是有界的,

我们首先回顾 Cholesky 分解的定义. 对任意对称正定矩阵 $\mathbf{A} = (a_{ij})$, 它的 Cholesky 分解可写作

$$\mathbf{A} = \mathbf{L}\mathbf{D}\mathbf{L}^T$$

其中 $\mathbf{L} = (l_{ij})$ 是对角线元素均为 1 的下三角矩阵, $\mathbf{D} = \text{Diag}(d_1, d_2, \dots, d_n)$ 是对角矩阵且对角线元素均为正.

因此, 对修正后的矩阵元素也需要有上界约束. 具体来说, 我们选取两个正参数 δ, β 使得

$$d_j \geq \delta, \quad l_{ij} \sqrt{d_j} \leq \beta, \quad i = j+1, j+2, \dots, n.$$

因此, 我们只需要修改 Cholesky 分解时 d_j 的更新即可保证上述条件成立. 具体更新方式为

$$d_j = \max \left\{ |c_{jj}|, \left(\frac{\theta_j}{\beta} \right)^2, \delta \right\}, \quad \theta_j = \max_{i>j} |c_{ij}|.$$

可以证明, 修正的 Cholesky 分解算法实际上是计算修正矩阵 $\nabla^2 f(\mathbf{x}^k) + \mathbf{E}^k$ 的 Cholesky 分解, 其中 \mathbf{E}^k 是对角矩阵且对角线元素非负. 当 $\nabla^2 f(\mathbf{x}^k)$ 正定且条件数足够小时有 $\mathbf{E}^k = 0$.

34.1.3 应用：牛顿法求解 Logistic 回归

在前面我们已经介绍了二分类的 Logistic 回归模型：

$$\min_{\mathbf{x}} L(\mathbf{x}) \stackrel{\text{def}}{=} \frac{1}{m} \sum_{i=1}^m \ln(1 + \exp(-b_i \mathbf{a}_i^T \mathbf{x})) + \lambda \|\mathbf{x}\|_2^2.$$

接下来推导求解该问题的牛顿法，这转化为计算目标函数 $L(\mathbf{x})$ 的梯度和 Hessian 矩阵的问题。根据第六章介绍的向量值函数求导法，容易算出梯度为

$$\begin{aligned} \nabla L(\mathbf{x}) &= \frac{1}{m} \sum_{i=1}^m \frac{1}{1 + \exp(-b_i \mathbf{a}_i^T \mathbf{x})} \cdot \exp(-b_i \mathbf{a}_i^T \mathbf{x}) \cdot (-b_i \mathbf{a}_i) + 2\lambda \mathbf{x} \\ &= -\frac{1}{m} \sum_{i=1}^m (1 - p_i(\mathbf{x})) b_i \mathbf{a}_i + 2\lambda \mathbf{x} \end{aligned}$$

其中 $p_i(\mathbf{x}) = \frac{1}{1 + \exp(-b_i \mathbf{a}_i^T \mathbf{x})}$ 。

引入矩阵 $\mathbf{A} = [\mathbf{a}_1, \mathbf{a}_2, \dots, \mathbf{a}_m]^T \in \mathbb{R}^{m \times n}$, 向量 $\mathbf{b} = (b_1, b_2, \dots, b_m)^T$, 以及

$$\mathbf{p}(\mathbf{x}) = (p_1(\mathbf{x}), p_2(\mathbf{x}), \dots, p_m(\mathbf{x}))^T.$$

此时梯度可简写为:

$$\nabla L(\mathbf{x}) = -\frac{1}{m} \mathbf{A}^T (\mathbf{b} - \mathbf{b} \odot \mathbf{p}(\mathbf{x})) + 2\lambda \mathbf{x}.$$

再对梯度求导, 并写成更为紧凑的矩阵形式, 可得到 Hessian 矩阵

$$\nabla^2 L(\mathbf{x}) = \frac{1}{m} \mathbf{A}^T \mathbf{P}(\mathbf{x}) \mathbf{A} + 2\lambda \mathbf{I}$$

其中 $\mathbf{P}(\mathbf{x})$ 为由 $\{p_i(\mathbf{x})(1 - p_i(\mathbf{x}))\}_{i=1}^m$ 生成的对角矩阵。因此, 牛顿法可以写作

$$\mathbf{x}^{k+1} = \mathbf{x}^k + \left(\frac{1}{m} \left(\mathbf{A}^T \mathbf{P}(\mathbf{x}^k) \mathbf{A} + 2\lambda \mathbf{I} \right)^{-1} \left(\frac{1}{m} \mathbf{A}^T (\mathbf{b} - \mathbf{b} \odot \mathbf{p}(\mathbf{x}^k)) - 2\lambda \mathbf{x}^k \right) \right).$$

在实际中, λ 经常取为 $\frac{1}{100m}$ 。另外, 当变量规模不是很大时, 可以利用正定矩阵的 Cholesky 分解来求解牛顿方程; 当变量规模较大时, 可以使用共轭梯度法进行不精确求解。这里采用 LIBSVM 网站的数据集, 包括: a9a、ijcnn1 和 CINA 数据集。然后使用牛顿法进行求解, 其求解结果参见图2。从中可以看出, 在精确解附近梯度范数具有 Q- 超线性收敛性。

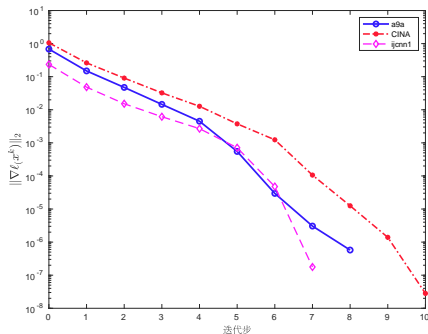


图 2: 牛顿法求解 Logistic 回归模型

1 34.1 牛顿法

2 34.2 拟牛顿法

- 拟牛顿法（变尺度法）是近 40 多年来发展起来的，它是求解无约束优化问题的一种有效方法。
- 由于它既避免了计算二阶导数矩阵及其求逆过程，又比梯度法的收敛速度快，特别是对高维问题具有显著的优越性，因而使拟牛顿法获得了很高的声誉，至今仍被公认为求解无约束优化问题最有效的算法之一。
- 下面我们就来简要地介绍拟牛顿法的基本原理及其计算过程。

34.2.1 割线方程

若想找到 Hessian 矩阵或它的逆的一个近似，就应当寻找到它需要满足的条件。以便构造近似矩阵，并令其也满足同样的条件即可。

显然，这样的一个近似应当满足 Taylor 展开。根据 Taylor 展开，梯度函数 $\nabla f(\mathbf{x})$ 在点 $\mathbf{x}^{(k+1)}$ 处的近似为

$$\nabla f(\mathbf{x}) \approx \nabla f(\mathbf{x}^{(k+1)}) + \nabla^2 f(\mathbf{x}^{(k+1)})(\mathbf{x} - \mathbf{x}^{(k+1)})$$

$\mathbf{x} = \mathbf{x}^{(k)}$ 即有

$$\nabla f(\mathbf{x}^{(k+1)}) - \nabla f(\mathbf{x}^{(k)}) = \nabla^2 f(\mathbf{x}^{(k+1)})(\mathbf{x}^{(k+1)} - \mathbf{x}^{(k)}) \quad (2)$$

或

$$\mathbf{x}^{(k+1)} - \mathbf{x}^{(k)} = \nabla^2 f(\mathbf{x}^{(k+1)})^{-1} [\nabla f(\mathbf{x}^{(k+1)}) - \nabla f(\mathbf{x}^{(k)})] \quad (3)$$

这两个式子并称为割线方程，即拟牛顿条件。

割线方程

若令

$$\begin{cases} \Delta \mathbf{g}^{(k)} = \nabla f(\mathbf{x}^{(k+1)}) - \nabla f(\mathbf{x}^{(k)}) \\ \Delta \mathbf{x}^{(k)} = \mathbf{x}^{(k+1)} - \mathbf{x}^{(k)} \end{cases} \quad (4)$$

则式(2) 变为

$$\Delta \mathbf{g}^{(k)} = \nabla^2 f(\mathbf{x}^{(k+1)}) \Delta \mathbf{x}^{(k)}$$

式(3) 变为

$$\Delta \mathbf{x}^{(k)} = \nabla^2 f(\mathbf{x}^{(k+1)})^{-1} \Delta \mathbf{g}^{(k)}$$

34.2.2 拟牛顿算法框架

如果得到满足割线方程的 **Hessian** 矩阵的近似（下面均用 \mathbf{H} 表示），或者 **Hessian** 矩阵的逆的近似（下面均用 $\overline{\mathbf{H}}$ 表示），则可以得到拟牛顿方法的一般求解框架。

Algorithm 3 拟牛顿法计算框架

- 1: 给定 $\mathbf{x}^0 \in \mathbb{R}^n$, 初始矩阵 $\mathbf{H}^0 \in \mathbb{R}^{n \times n}$ (或 $\overline{\mathbf{H}}^0$), 令 $k = 0$.
 - 2: **while** $k = 0, 1, 2, \dots$ **do**
 - 3: 计算方向 $\mathbf{d}^k = -(\mathbf{H}^k)^{-1} \nabla f(\mathbf{x}^k)$ 或 $\mathbf{d}^k = -\overline{\mathbf{H}}^k \nabla f(\mathbf{x}^k)$.
 - 4: 通过线搜索找到合适的步长 $\lambda_k > 0$, 令 $\mathbf{x}^{k+1} = \mathbf{x}^k + \lambda_k \mathbf{d}^k$.
 - 5: 更新海瑟矩阵的近似矩阵 \mathbf{H}^{k+1} 或其逆矩阵的近似 $\overline{\mathbf{H}}^{k+1}$.
 - 6: $k \leftarrow k + 1$.
 - 7: **end while**
-

下面，我们将讨论如何借助割线方程（拟牛顿条件）具体的构造 Hessian 矩阵或其逆的近似。这里补充几点说明：

- 基于(2) 式得到 Hessian 矩阵的近似，具有较好的理论性质，迭代序列较为稳定，但仍然可能在大规模问题上是非常耗时的。
- 基于(3) 式得到 Hessian 矩阵的逆的近似，更加实用。
- 上述两种方式之间具有很好的形式对称性，下面仅基于 Hessian 矩阵逆的近似进行探究。

34.2.3 秩一更新

现设 $\overline{\mathbf{H}}^{(k)}$ 是第 k 步 Hessian 矩阵的逆的近似，现需构造出 $\overline{\mathbf{H}}^{(k+1)}$ ，则直观的想法是对它进行秩一修正。考虑到对称性，则可设

$$\overline{\mathbf{H}}^{(k+1)} = \overline{\mathbf{H}}^{(k)} + a\mathbf{u}\mathbf{u}^T \quad (5)$$

其中 $\mathbf{u} \in \mathbb{R}^n$, $a \in \mathbb{R}$ 待定。

秩一更新

利用割线方程，有

$$\begin{aligned}\Delta \mathbf{x}^{(k)} &= \overline{\mathbf{H}}^{(k+1)} \Delta \mathbf{g}^{(k)} \\ &= (\overline{\mathbf{H}}^{(k)} + a \mathbf{u} \mathbf{u}^T) \Delta \mathbf{g}^{(k)}\end{aligned}\tag{6}$$

整理得

$$a(\mathbf{u}^T \Delta \mathbf{g}^{(k)}) \mathbf{u} = \Delta \mathbf{x}^{(k)} - \overline{\mathbf{H}}^{(k)} \Delta \mathbf{g}^{(k)}.$$

秩一更新

因此 \mathbf{u} 与 $\Delta \mathbf{x}^{(k)} - \overline{\mathbf{H}}^{(k)} \Delta \mathbf{g}^{(k)}$ 共线。不妨令 $\mathbf{u} = \Delta \mathbf{x}^{(k)} - \overline{\mathbf{H}}^{(k)} \Delta \mathbf{g}^{(k)}$ ，则代入可得

$$a = \frac{1}{(\Delta \mathbf{x}^{(k)} - \overline{\mathbf{H}}^{(k)} \Delta \mathbf{g}^{(k)})^T \Delta \mathbf{g}^{(k)}}.$$

从而，得到秩一更新公式：

$$\overline{\mathbf{H}}^{(k+1)} = \overline{\mathbf{H}}^{(k)} + \frac{(\Delta \mathbf{x}^{(k)} - \overline{\mathbf{H}}^{(k)} \Delta \mathbf{g}^{(k)})(\Delta \mathbf{x}^{(k)} - \overline{\mathbf{H}}^{(k)} \Delta \mathbf{g}^{(k)})^T}{(\Delta \mathbf{x}^{(k)} - \overline{\mathbf{H}}^{(k)} \Delta \mathbf{g}^{(k)})^T \Delta \mathbf{g}^{(k)}}.$$

上述矩阵也称之为尺度矩阵。

秩一更新

如果考虑的是 Hessian 矩阵本身的近似, 则同理可得对应的秩一更新公式如下:

$$\mathbf{H}^{(k+1)} = \mathbf{H}^{(k)} + \frac{(\Delta \mathbf{g}^{(k)} - \mathbf{H}^{(k)} \Delta \mathbf{x}^{(k)})(\Delta \mathbf{g}^{(k)} - \mathbf{H}^{(k)} \Delta \mathbf{x}^{(k)})^T}{(\Delta \mathbf{g}^{(k)} - \mathbf{H}^{(k)} \Delta \mathbf{x}^{(k)})^T \Delta \mathbf{x}^{(k)}}.$$

通过对比发现, 实际上二者之间只是做了如下形式的替换:

$$\bar{\mathbf{H}} \rightarrow \mathbf{H}, \quad \Delta \mathbf{x}^{(k)} \leftrightarrow \Delta \mathbf{g}^{(k)}.$$

注: 这样的形式对称性对我们了解拟牛顿法大有裨益。

秩一更新

秩一更新存在的重大缺陷：

- 秩一更新虽然结构简单，易计算；
- 但不能保证在迭代过程中保持正定。
- 需要寻求更好的近似。

34.2.4 DFP

为克服秩一更新的缺陷，直观的改进是对它进行秩二修正。同样地考虑对称性，则可设

$$\overline{\mathbf{H}}^{(k+1)} = \overline{\mathbf{H}}^{(k)} + a\mathbf{u}\mathbf{u}^T + b\mathbf{v}\mathbf{v}^T \quad (7)$$

其中 $\mathbf{u}, \mathbf{v} \in \mathbb{R}^n$, $a, b \in \mathbb{R}$ 待定。

利用割线方程，有

$$\begin{aligned}\Delta \mathbf{x}^{(k)} &= \overline{\mathbf{H}}^{(k+1)} \Delta \mathbf{g}^{(k)} \\ &= (\overline{\mathbf{H}}^{(k)} + a \mathbf{u} \mathbf{u}^T + b \mathbf{v} \mathbf{v}^T) \Delta \mathbf{g}^{(k)}\end{aligned}\tag{8}$$

整理得

$$a(\mathbf{u}^T \Delta \mathbf{g}^{(k)}) \mathbf{u} + b(\mathbf{v}^T \Delta \mathbf{g}^{(k)}) \mathbf{v} = \Delta \mathbf{x}^{(k)} - \overline{\mathbf{H}}^{(k)} \Delta \mathbf{g}^{(k)}.$$

因此 \mathbf{u}, \mathbf{v} 的线性组合等于 $\Delta \mathbf{x}^{(k)} - \overline{\mathbf{H}}^{(k)} \Delta \mathbf{g}^{(k)}$ 。同样地，不妨令 $\mathbf{u} = \Delta \mathbf{x}^{(k)}$, $\mathbf{v} = \overline{\mathbf{H}}^{(k)} \Delta \mathbf{g}^{(k)}$ ，则代入可得

DFP

$$a = \frac{1}{(\Delta \mathbf{x}^{(k)})^T \Delta \mathbf{g}^{(k)}},$$
$$b = -\frac{1}{(\Delta \mathbf{g}^{(k)})^T \overline{\mathbf{H}}^{(k)} \Delta \mathbf{g}^{(k)}}.$$

从而，得到更新公式：

$$\overline{\mathbf{H}}^{(k+1)} = \overline{\mathbf{H}}^{(k)} + \frac{\Delta \mathbf{x}^{(k)} \Delta \mathbf{x}^{(k)T}}{(\Delta \mathbf{x}^{(k)})^T \Delta \mathbf{g}^{(k)}} - \frac{\overline{\mathbf{H}}^{(k)} \Delta \mathbf{g}^{(k)} (\overline{\mathbf{H}}^{(k)} \Delta \mathbf{g}^{(k)})^T}{(\Delta \mathbf{g}^{(k)})^T \overline{\mathbf{H}}^{(k)} \Delta \mathbf{g}^{(k)}}. \quad (9)$$

这种迭代公式由 Davidon 发现，并由 Fletcher 以及 Powell 进一步发展。因此被称为 **DFP** 公式。

34.2.5 BFGS

利用前面提及的形式对称性，可得如下更新公式：

$$\mathbf{H}^{(k+1)} = \mathbf{H}^{(k)} + \frac{\Delta \mathbf{g}^{(k)} \Delta \mathbf{g}^{(k)T}}{(\Delta \mathbf{g}^{(k)})^T \Delta \mathbf{x}^{(k)}} - \frac{\mathbf{H}^{(k)} \Delta \mathbf{x}^{(k)} (\mathbf{H}^{(k)} \Delta \mathbf{x}^{(k)})^T}{(\Delta \mathbf{x}^{(k)})^T \mathbf{H}^{(k)} \Delta \mathbf{x}^{(k)}}.$$

这种迭代格式就是著名的 **BFGS** 公式，它是由 Broyden, Fletcher, Goldfarb 和 Shanno 四人的名字首字母组成。

尽管 DFP 格式和 BFGS 格式存在这种对偶关系，但实际上，BFGS 格式效果更好些。因此，在实际中 BFGS 格式被使用的更多。

34.2.6 拟牛顿法计算步骤

综上，可将拟牛顿法（以 DFP 为例）的计算步骤总结如下：

① 给定初始点 $\mathbf{x}^{(0)}$ 及梯度允许误差 $\varepsilon > 0$;

② 若

$$\|\nabla f(\mathbf{x}^{(0)})\|^2 \neq \varepsilon$$

则 $\mathbf{x}^{(0)}$ 即为近似极小点，停止迭代。否则，转向下一步；

拟牛顿法计算步骤

③ 令

$$\overline{\mathbf{H}}^{(0)} = \mathbf{I}(\text{单位阵})$$

$$\mathbf{p}^{(0)} = -\overline{\mathbf{H}}^{(0)} \nabla f(\mathbf{x}^{(0)})$$

在 $\mathbf{p}^{(0)}$ 方向进行一维搜索，确定最佳步长 λ_0

$$\min_{\lambda} f(\mathbf{x}^{(0)} + \lambda \mathbf{p}^{(0)}) = f(\mathbf{x}^{(0)} + \lambda_0 \mathbf{p}^{(0)})$$

如此可得下一个近似点

$$\mathbf{x}^{(1)} = \mathbf{x}^{(0)} + \lambda_0 \mathbf{p}^{(0)}$$

拟牛顿法计算步骤

- ④ 一般地, 设已得到近似点 $\mathbf{x}^{(k)}$, 算出 $\nabla f(\mathbf{x}^{(k)})$, 若

$$\|\nabla f(\mathbf{x}^{(k)})\|^2 \leq \varepsilon$$

则 $\mathbf{x}^{(k)}$ 即为所求的近似解, 停止迭代; 否则, 按式(9) 计算 $\overline{\mathbf{H}}^{(k)}$, 并令

$$\mathbf{p}^{(k)} = -\overline{\mathbf{H}}^{(k)} \nabla f(\mathbf{x}^{(k)})$$

在 $\mathbf{p}^{(k)}$ 方向进行一维搜索, 确定最佳最长 λ_k

$$\min_{\lambda} f(\mathbf{x}^{(k)} + \lambda \mathbf{p}^{(k)}) = f(\mathbf{x}^{(k)} + \lambda_k \mathbf{p}^{(k)})$$

其下一个近似点为

$$\mathbf{x}^{(k+1)} = \mathbf{x}^{(k)} + \lambda_k \mathbf{p}^{(k)}$$

拟牛顿法计算步骤

- ⑤ 若 $\mathbf{x}^{(k+1)}$ 点满足精度要求, 则 $\mathbf{x}^{(k+1)}$ 即为所求的近似解。否则, 转回第 (4) 步, 直到求出某点满足精度要求为止。

注: 与共轭梯度法相类似, 如果迭代 n 次仍不收敛, 则以 $\mathbf{x}^{(n)}$ 为新的 $\mathbf{x}^{(0)}$, 以这时的 $\mathbf{x}^{(0)}$ 为起点重新开始一轮新的迭代。

对于拟牛顿法, 我们也可以得到其基本的收敛性以及收敛速度.

定理 2

(BFGS 全局收敛性) 假设初始矩阵 \mathbf{H}^0 是对称正定矩阵, 目标函数 $f(\mathbf{x})$ 是二阶连续可微函数, 且下水平集

$$\mathcal{L} = \{\mathbf{x} \in \mathbb{R}^n \mid f(\mathbf{x}) \leq f(\mathbf{x}^0)\}$$

是凸的, 并且存在正数 m 以及 M 使得对于任意的 $\mathbf{z} \in \mathbb{R}^n$ 以及任意的 $\mathbf{x} \in \mathcal{L}$ 有

$$m\|\mathbf{z}\|^2 \leq \mathbf{z}^T \nabla^2 f(\mathbf{x}) \mathbf{z} \leq M\|\mathbf{z}\|^2,$$

则采用 BFGS 格式并结合 Wolfe 线搜索的拟牛顿算法全局收敛到 $f(\mathbf{x})$ 的极小值点 \mathbf{x}^* .

上述定理叙述了 BFGS 格式的全局收敛性, 但没有说明以什么速度收敛. 下面这个定理介绍了在一定条件下 BFGS 格式会达到 Q -超线性收敛速度. 这里仍然只给出定理结果, 感兴趣的读者可以查阅相关文献, 了解详细的证明过程.

定理 3

(BFGS 收敛速度) 设 $f(\mathbf{x})$ 二阶连续可微, 在最优点 \mathbf{x}^* 的一个邻域内海瑟矩阵利普希茨连续, 且使用 BFGS 迭代格式收敛到 f 的最优值点 \mathbf{x}^* . 若迭代点列 $\{\mathbf{x}^k\}$ 满足

$$\sum_{k=1}^{\infty} \|\mathbf{x}^k - \mathbf{x}^*\| < +\infty,$$

则 $\{\mathbf{x}^k\}$ 以 Q -超线性收敛到 \mathbf{x}^* .

例 1

试用 DFP 法重新计算下述二次函数的极小点

$$f(\mathbf{x}) = \frac{3}{2}x_1^2 + \frac{1}{2}x_2^2 - x_1x_2 - 2x_1$$

解

我们从 $\mathbf{x}^{(0)} = (-2, 4)^T$ 开始, 并取

$$\overline{\mathbf{H}}^{(0)} = \begin{pmatrix} 1 & 0 \\ 0 & 1 \end{pmatrix}$$

$$\nabla f(\mathbf{x}) = [(3x_1 - x_2 - 2), (x_2 - x_1)]^T$$

$$\nabla f(\mathbf{x}^{(0)}) = (-12, 6)^T$$

$$\mathbf{p}^{(0)} = -\overline{\mathbf{H}}^{(0)} \nabla f(\mathbf{x}^{(0)}) = - \begin{pmatrix} 1 & 0 \\ 0 & 1 \end{pmatrix} \begin{pmatrix} -12 \\ 6 \end{pmatrix} = \begin{pmatrix} 12 \\ -6 \end{pmatrix}$$

利用一维搜索，即 $\min_{\lambda} f(\mathbf{x}^{(0)} + \lambda \mathbf{p}^{(0)})$ ，可算得

$$\lambda_0 = \frac{5}{17}$$

$$\mathbf{x}^{(1)} = \mathbf{x}^{(0)} + \lambda_0 \mathbf{p}^{(0)} = \begin{pmatrix} -2 \\ 4 \end{pmatrix} + \frac{5}{17} \begin{pmatrix} 12 \\ -6 \end{pmatrix} = \left(\frac{26}{17}, \frac{38}{17} \right)^T$$

$$\nabla f(\mathbf{x}^{(1)}) = \left(\frac{6}{17}, \frac{12}{17} \right)^T$$

$$\Delta \mathbf{x}^{(0)} = \mathbf{x}^{(1)} - \mathbf{x}^{(0)} = \left(\frac{26}{17}, \frac{38}{17} \right)^T - (-2, 4)^T = \left(\frac{60}{17}, -\frac{30}{17} \right)^T$$

$$\Delta \mathbf{g}^{(0)} = \nabla f(\mathbf{x}^{(1)}) - \nabla f(\mathbf{x}^{(0)}) = \left(\frac{6}{17}, \frac{12}{17} \right)^T - (-12, 6)^T = \left(\frac{210}{17}, -\frac{90}{17} \right)^T$$

$$\begin{aligned} \overline{\mathbf{H}}^{(1)} &= \overline{\mathbf{H}}^{(0)} + \frac{\Delta \mathbf{x}^{(0)} (\Delta \mathbf{x}^{(0)})^T}{(\Delta \mathbf{g}^{(0)})^T \Delta \mathbf{x}^{(0)}} - \frac{\overline{\mathbf{H}}^{(0)} \Delta \mathbf{g}^{(0)} (\Delta \mathbf{g}^{(0)})^T \overline{\mathbf{H}}^{(0)}}{(\Delta \mathbf{g}^{(0)})^T \overline{\mathbf{H}}^{(0)} \Delta \mathbf{g}^{(0)}} \\ &= \begin{pmatrix} 1 & 0 \\ 0 & 1 \end{pmatrix} + \frac{\left(\frac{60}{17}, -\frac{30}{17} \right)^T \left(\frac{60}{17}, -\frac{30}{17} \right)}{\left(\frac{210}{17}, -\frac{90}{17} \right) \left(\frac{60}{17}, -\frac{30}{17} \right)^T} - \frac{\begin{pmatrix} 1 & 0 \\ 0 & 1 \end{pmatrix} \left(\frac{210}{17}, -\frac{90}{17} \right)^T \left(\frac{210}{17}, -\frac{90}{17} \right) \begin{pmatrix} 1 & 0 \\ 0 & 1 \end{pmatrix}}{\left(\frac{210}{17}, -\frac{90}{17} \right) \begin{pmatrix} 1 & 0 \\ 0 & 1 \end{pmatrix} \left(\frac{210}{17}, -\frac{90}{17} \right)^T} \\ &= \frac{1}{986} \begin{pmatrix} 385 & 241 \\ 241 & 891 \end{pmatrix} \end{aligned}$$

$$\mathbf{p}^{(1)} = -\overline{\mathbf{H}}^{(1)} \nabla f(\mathbf{x}^{(1)}) = -\frac{1}{986} \begin{pmatrix} 385 & 241 \\ 241 & 891 \end{pmatrix} \begin{pmatrix} \frac{6}{17} \\ \frac{12}{17} \end{pmatrix} = -\begin{pmatrix} \frac{9}{29} \\ \frac{21}{29} \end{pmatrix}$$

再由一维搜索 $\min_{\lambda} f(\mathbf{x}^{(1)} + \lambda \mathbf{p}^{(1)})$, 得

$$\lambda_1 = \frac{29}{17}$$

从而

$$\mathbf{x}^{(2)} = \mathbf{x}^{(1)} + \lambda_1 \mathbf{p}^{(1)} = \begin{pmatrix} \frac{26}{17} \\ \frac{38}{17} \end{pmatrix} + \frac{29}{17} \begin{pmatrix} -\frac{9}{29} \\ -\frac{21}{29} \end{pmatrix} = \begin{pmatrix} 1 \\ 1 \end{pmatrix}$$

$$\nabla f(\mathbf{x}^{(2)}) = (0, 0)^T$$

可知 $\mathbf{x}^{(2)} = (1, 1)^T$ 为极小点。

- 在以上讨论中, 若取第一个尺度矩阵 $\overline{\mathbf{H}}^{(0)}$ 为对称正定阵, 则可以证明, 由式(9) 迭代形成的尺度矩阵均为对称正定阵。
- 由此可知其搜索方向 $\mathbf{p}^{(k)} = -\overline{\mathbf{H}}^{(k)} \nabla f(\mathbf{x}^{(k)})$ 为下降方向, 这就可以保证每次迭代均能使目标函数值有所改善。
- 当把 DFP 拟牛顿法用于正定二次函数时, 产生的搜索方向为共轭方向, 因而也具有有限步收敛的性质。若将初始尺度矩阵也取为单位矩阵, 对这种函数来说, DFP 法就与共轭梯度法一样了。

34.2.7 应用：拟牛顿法求解压缩感知问题

考虑压缩感知问题:

$$\min_{\mathbf{x} \in \mathbb{R}^n} \|\mathbf{x}\|_1, \quad \text{s.t.} \quad \mathbf{A}\mathbf{x} = \mathbf{b}, \quad (10)$$

其中 $\mathbf{A} \in \mathbb{R}^{m \times n}$, $\mathbf{b} \in \mathbb{R}^m$ 为给定的矩阵和向量. 这是一个约束优化问题, 如何将其转化为一个无约束优化问题呢? 自然地, 我们可以考虑其对偶问题. 由于问题 (10) 的对偶问题的无约束优化形式不是可微的, 即无法计算梯度 (读者可以自行验证)。

我们考虑如下正则化问题:

$$\min_{\mathbf{x} \in \mathbb{R}^n} \|\mathbf{x}\|_1 + \frac{1}{2\alpha} \|\mathbf{x}\|_2^2, \quad \text{s.t.} \quad \mathbf{A}\mathbf{x} = \mathbf{b}, \quad (11)$$

这里 $\alpha > 0$ 为正则化参数. 显然, 当 α 趋于无穷大时, 问题 (11) 的解会逼近 (10) 的解. 由于问题 (11) 的目标函数是强凸的, 其对偶问题的无约束优化形式的目标函数是可微的. 具体地, 问题 (11) 的对偶问题为

$$\min_{\mathbf{y} \in \mathbb{R}^m} f(\mathbf{y}) = -\mathbf{b}^T \mathbf{y} + \frac{\alpha}{2} \left\| \mathbf{A}^T \mathbf{y} - \mathcal{P}_{[-1,1]^n}(\mathbf{A}^T \mathbf{y}) \right\|_2^2, \quad (12)$$

其中 $\mathcal{P}_{[-1,1]^n}(\mathbf{x})$ 为 \mathbf{x} 到集合 $[-1, 1]^n$ 的投影.

通过简单计算, 可知

$$\nabla f(\mathbf{y}) = -\mathbf{b} + \alpha \mathbf{A} \left(\mathbf{A}^T \mathbf{y} - \mathcal{P}_{[-1,1]^n} \left(\mathbf{A}^T \mathbf{y} \right) \right)$$

那么, 我们可以利用 L-BFGS 方法来求解问题 (12). 在得到该问题的解 \mathbf{y}^* 之后, 问题 (11) 的解 \mathbf{x}^* 可通过下式近似得到:

$$\mathbf{x}^* \approx \alpha \left(\mathbf{A}^T \mathbf{y}^* - \mathcal{P}_{[-1,1]^n} \left(\mathbf{A}^T \mathbf{y}^* \right) \right)$$

进一步地, 当 α 充分大时, 问题 (11) 的解就是原问题 (10) 的解. 因此, 我们可以通过选取合适的 α , 通过求解问题 (12) 来得到问题 (10) 的解.

我们用 LASSO 问题中的 \mathbf{A} 和 \mathbf{b} , 分别选取 $\alpha = 5, 10$, 调用 BFGS 方法¹求解问题 (12), 其中内存长度取为 5. 迭代收敛过程见图 3. 从图中我们可以看到, 当靠近最优解时, BFGS 方法的迭代点列呈 Q-线性收敛.

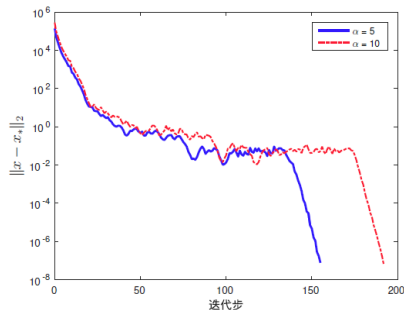


图 3: 压缩感知问题

¹实际中, 使用的是有限内存 BFGS, 感兴趣的同学可以查阅相关文献。

本讲小结

牛顿法

- 二阶近似
- 牛顿步径
- 牛顿法、修正的牛顿法

拟牛顿法（变尺度法）

- 割线方程
- 秩一更新
- DFP、BFGS

注意，牛顿法中还有非精确牛顿法以及拟牛顿法中还有有限内存 BFGS 方法等改进方法，本课程没有详细介绍。另外，本节我们主要讨论了无约束优化问题的二阶求解算法。那么我们是如何处理约束优化问题的呢？