

Newton-Raphson 算法

王璐

对很多常用的统计模型, 比如 exponential family, 数据的 log-likelihood $l(\boldsymbol{\beta})$ 是参数 $\boldsymbol{\beta}$ 的凹函数, 这种情况下计算 $\boldsymbol{\beta}$ 的 MLE 只需令一阶导数 (score function) 等于 $\mathbf{0}$:

$$\frac{\partial l(\boldsymbol{\beta})}{\partial \boldsymbol{\beta}} = \mathbf{0}. \quad (1)$$

由于 $l(\boldsymbol{\beta})$ 的一阶导数一般是非线性函数, (1)经常需要用到求解非线性方程组的算法 — Newton-Raphson 算法.

1 Newton-Raphson Algorithm

对于一般的非线性方程组

$$\mathbf{h}(\mathbf{x}) = \mathbf{0}$$

其中 $\mathbf{h}: \mathbb{R}^p \rightarrow \mathbb{R}^p$, 求解该方程组的 Newton-Raphson 算法如下:

1. 选取一个初始值 $\mathbf{x}^{(0)}$.
2. 在 $\mathbf{x}^{(0)}$ 的邻域内用线性函数近似 $\mathbf{h}(\cdot)$

$$\mathbf{h}(\mathbf{x}) \approx \mathbf{h}(\mathbf{x}^{(0)}) + \nabla \mathbf{h}(\mathbf{x}^{(0)})(\mathbf{x} - \mathbf{x}^{(0)}).$$

3. 令上式右边等于 $\mathbf{0}$, 得到 $\mathbf{h}(\mathbf{x}) = \mathbf{0}$ 的一个近似解:

$$\mathbf{x}^{(1)} = \mathbf{x}^{(0)} - \left[\nabla \mathbf{h}(\mathbf{x}^{(0)}) \right]^{-1} \mathbf{h}(\mathbf{x}^{(0)})$$

4. 重复此过程, 第 $(t+1)$ 次迭代得到

$$\mathbf{x}^{(t+1)} = \mathbf{x}^{(t)} - \left[\nabla \mathbf{h}(\mathbf{x}^{(t)}) \right]^{-1} \mathbf{h}(\mathbf{x}^{(t)})$$

5. 如果 $\|\mathbf{h}(\mathbf{x}^{(t+1)})\| < \varepsilon$, 算法终止, 输出 $\mathbf{x}^{(t+1)}$.

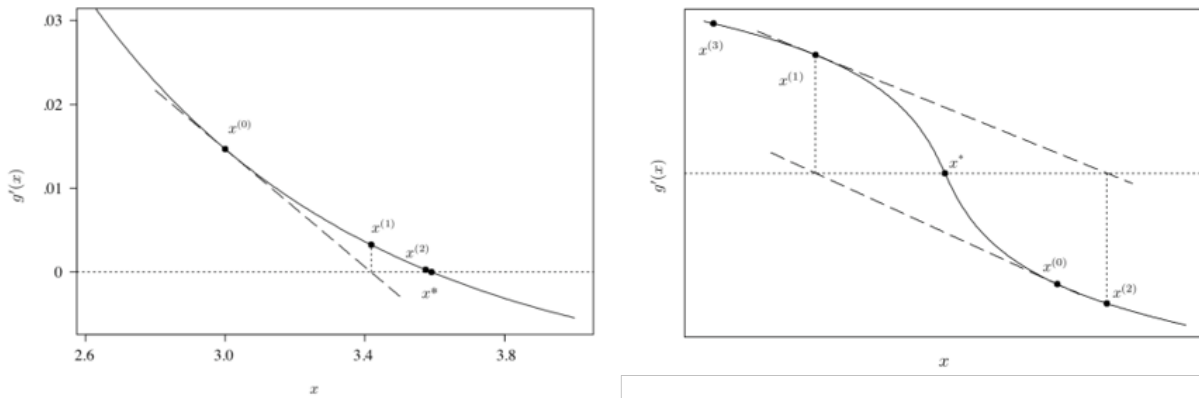


Figure 1: Newton-Raphson 算法收敛和不收敛的情形举例。

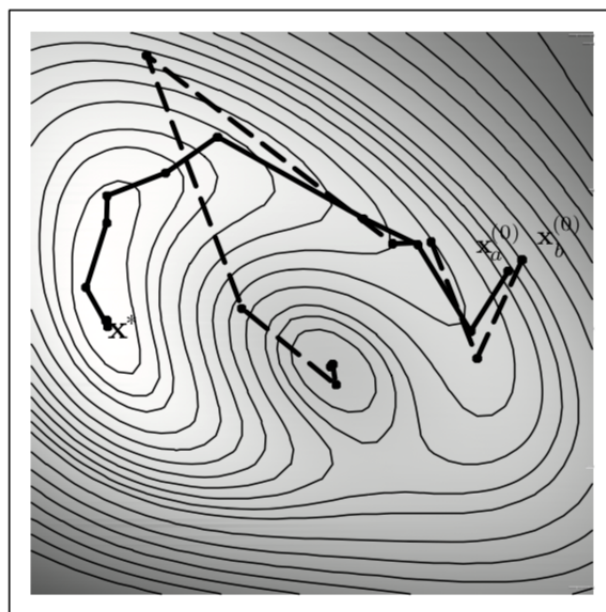


Figure 2: 使用 Newton 算法最大化一个二元函数，图中越亮的区域代表函数值越大。在 Newton 算法中使用两个不同初始值 $\mathbf{x}_a^{(0)}$ 和 $\mathbf{x}_b^{(0)}$ ，最终一个收敛到极大值点，一个收敛到极小值点。Picture source: *Computational Statistics* (2nd edition) by Geof H.Givens (2012)

Remarks

1. Newton 算法是否收敛取决于 $\mathbf{h}(\cdot)$ 的形状和选取的初始值 $\mathbf{x}^{(0)}$ ，如图1 - 2所示。
2. Memory cost $O(p^2)$: 每一步 Newton 迭代需要储存一个 $p \times p$ 矩阵 $\nabla \mathbf{h}(\mathbf{x}^{(t)})$.
Computation cost $O(p^3)$: 每一步 Newton 迭代需要计算一个 $p \times p$ (稠密) 矩阵的逆或解一组线性方程组。

2 收敛性分析

以一元函数为例, 假设 x^* 是 $h(x) = 0$ 的一个根, $h(x)$ 二阶连续可导且 $h'(x^*) \neq 0$. 由于 h' 连续且 $h'(x^*) \neq 0$, 则存在 x^* 的一个邻域, 使得这个邻域内的所有点 x 都满足 $h'(x) \neq 0$. 假设 $x^{(t)}$ 在这个邻域内, 将 $h(x^*)$ 在 $x^{(t)}$ 处 Taylor 展开

$$0 = h(x^*) = h(x^{(t)}) + h'(x^{(t)})(x^* - x^{(t)}) + \frac{1}{2}h''(\tilde{x})(x^* - x^{(t)})^2 \quad (2)$$

其中 \tilde{x} 介于 x^* 和 $x^{(t)}$ 之间。整理(2)得

$$\underbrace{x^{(t)} - \frac{h(x^{(t)})}{h'(x^{(t)})}}_{x^{(t+1)}} - x^* = (x^* - x^{(t)})^2 \frac{h''(\tilde{x})}{2h'(x^{(t)})}. \quad (3)$$

令 $\epsilon_t = x^{(t)} - x^*$, 由(3)得

$$\epsilon_{t+1} = \epsilon_t^2 \frac{h''(\tilde{x})}{2h'(x^{(t)})}. \quad (4)$$

考虑 x^* 的一个 δ -邻域, $\mathcal{N}_\delta(x^*) = (x^* - \delta, x^* + \delta)$, 和下面这个跟 δ 有关的函数:

$$c(\delta) = \max_{x_1, x_2 \in \mathcal{N}_\delta(x^*)} \left| \frac{h''(x_1)}{2h'(x_2)} \right|.$$

由(4)得

$$\begin{aligned} |\epsilon_{t+1}| &\leq \epsilon_t^2 c(\delta) \\ |c(\delta)\epsilon_{t+1}| &\leq (c(\delta)\epsilon_t)^2. \end{aligned} \quad (5)$$

如果 $|\epsilon_0| = |x^{(0)} - x^*| < \delta$, 则由(5)得

$$|c(\delta)\epsilon_t| \leq (c(\delta)\epsilon_{t-1})^2 \leq (c(\delta)\epsilon_{t-2})^{2^2} \leq \cdots \leq (c(\delta)\epsilon_0)^{2^t} = (\delta c(\delta))^{2^t}. \quad (6)$$

进一步研究函数 $c(\delta)$ 的性质。注意到当 $\delta \rightarrow 0$,

$$c(\delta) \rightarrow \left| \frac{h''(x^*)}{2h'(x^*)} \right|$$

即 $c(\delta)$ 收敛到一个有限值。因此 $\delta \rightarrow 0$, $\delta c(\delta) \rightarrow 0$. 此时可以找到一个 δ_1 满足 $\delta_1 c(\delta_1) < 1$. 假设初始值 $x^{(0)}$ 满足 $|\epsilon_0| = |x^{(0)} - x^*| < \delta_1$, 根据(6), 当 $t \rightarrow \infty$,

$$|\epsilon_t| \leq \frac{(\delta_1 c(\delta_1))^{2^t}}{c(\delta_1)} \rightarrow 0$$

则 $x^{(t)} \rightarrow x^*$.

以上我们证明了一个定理: 如果函数 h 二阶连续可导, x^* 是 $h(x) = 0$ 的一个根且 $h'(x^*) \neq 0$, 那么存在 x^* 的一个邻域, 在这个邻域内选取任意初始值 $x^{(0)}$, Newton 算法都会收敛。

事实上, 当函数 h 二阶连续可导且是凸函数或凹函数, 并且有一个根, 则 Newton 算法对任意初始值都收敛。

2.1 收敛阶 (convergence order)

算法的收敛速度可以用算法的**收敛阶**衡量。

Definition 2.1 (收敛阶). 如果算法第 t 步误差 ϵ_t 满足

$$\lim_{t \rightarrow \infty} \epsilon_t = 0 \text{ 且 } \lim_{t \rightarrow \infty} \frac{|\epsilon_{t+1}|}{|\epsilon_t|^\alpha} = c$$

其中常数 $c \neq 0, \alpha > 0$, 称算法的收敛阶为 α .

算法的收敛阶越高意味着算法可以越快逼近真实解。如果 Newton 算法收敛, 由(4)得

$$\lim_{t \rightarrow \infty} \frac{|\epsilon_{t+1}|}{|\epsilon_t|^2} = \left| \frac{h''(x^*)}{2h'(x^*)} \right| = c.$$

因此 Newton 算法是二阶收敛 ($\alpha = 2$, quadratic convergence).

3 估计 Logistic 回归的 MLE

3.1 Logistic 回归模型

被解释变量 (response)

$$Y_i | \mathbf{x}_i \stackrel{ind}{\sim} \text{Bernoulli}(\pi_i), \quad i = 1, \dots, n.$$

解释变量 (covariates) $\mathbf{x}_i \in \mathbb{R}^p$ 与概率 $\pi_i = P(Y_i = 1 | \mathbf{x}_i)$ 有以下关系:

$$\text{logit}(\pi_i) \triangleq \log \left(\frac{\pi_i}{1 - \pi_i} \right) = \mathbf{x}_i^\top \boldsymbol{\beta}. \quad (7)$$

Remarks

1. 模型 (7)为什么合理? 与线性回归类似, 我们希望将 $E(Y_i | \mathbf{x}_i) = P(Y_i = 1 | \mathbf{x}_i) = \pi_i$ 与被解释变量 \mathbf{x}_i 的线性组合建立关系。但是 $\mathbf{x}_i^\top \boldsymbol{\beta}$ 可以是任何实数, 因此需要将 $\mathbf{x}_i^\top \boldsymbol{\beta}$ 映射到 $(0, 1)$ 区间。

- odds ratio $\frac{\pi_i}{1 - \pi_i}$ 将概率 π_i 转化为一个正实数, 再取 \log 将它转化为任意实数。

2. 将 π_i 直接表示为 $\mathbf{x}_i^\top \boldsymbol{\beta}$ 的函数:

$$\begin{aligned} \frac{\pi_i}{1 - \pi_i} &= e^{\mathbf{x}_i^\top \boldsymbol{\beta}} \\ \pi_i &= \frac{1}{1 + e^{-\mathbf{x}_i^\top \boldsymbol{\beta}}}. \end{aligned}$$

3. logistic function

$$\pi(x) = \frac{1}{1 + e^{-x}}$$

- $\pi : \mathbb{R} \rightarrow (0, 1)$
- $\pi'(x) = \pi(x)(1 - \pi(x))$

3.2 MLE of β

在模型 (7) 下, 数据的似然函数 (likelihood) 为:

$$\begin{aligned} L(\beta) &= P(Y_1 = y_1, \dots, Y_n = y_n \mid \mathbf{x}_1, \dots, \mathbf{x}_n, \beta) \\ &= \prod_{i=1}^n P(Y_i = y_i \mid \mathbf{x}_i, \beta) \\ &= \prod_{i=1}^n \pi_i^{y_i} (1 - \pi_i)^{1-y_i} \end{aligned}$$

根据定义, β 的 MLE (maximum likelihood estimation) 为

$$\hat{\beta} = \operatorname{argmax} L(\beta) = \operatorname{argmax} \log L(\beta)$$

为了计算方便, 一般选择最大化对数似然函数 (log-likelihood). 它是似然函数的单调变化, 不改变最大值点 (argmax) 的位置。

$$\begin{aligned} \log L(\beta) &= \sum_{i=1}^n y_i \log(\pi_i) + (1 - y_i) \log(1 - \pi_i) \\ &= \sum_{i=1}^n y_i \log\left(\frac{\pi_i}{1 - \pi_i}\right) + \log(1 - \pi_i) \\ &= \sum_{i=1}^n y_i \mathbf{x}_i^T \beta + \log(1 - \pi_i) \end{aligned}$$

- 一阶导数 (score function)

$$\begin{aligned} \frac{\partial \log L(\beta)}{\partial \beta} &= \sum_{i=1}^n y_i \mathbf{x}_i - \frac{1}{1 - \pi_i} \frac{\partial \pi_i}{\partial \beta} \\ &= \sum_{i=1}^n y_i \mathbf{x}_i - \frac{1}{1 - \pi_i} \pi_i (1 - \pi_i) \mathbf{x}_i \\ &= \sum_{i=1}^n (y_i - \pi_i) \mathbf{x}_i \end{aligned}$$

- 二阶导数 (negative of the observed Fisher's information)

$$\frac{\partial^2 \log L(\beta)}{\partial \beta \partial \beta^\top} = \sum_{i=1}^n \underbrace{-\pi_i(1 - \pi_i) \mathbf{x}_i \mathbf{x}_i^\top}_{\substack{\geq 0 \\ p.s.d}}$$

注意到 $\frac{\partial^2 \log L(\boldsymbol{\beta})}{\partial \boldsymbol{\beta} \partial \boldsymbol{\beta}^\top}$ 与 $y_{1:n}$ 无关, 且是半负定矩阵, 因此 $\log L(\boldsymbol{\beta})$ 是 $\boldsymbol{\beta}$ 的凹函数, 则 $\boldsymbol{\beta}_{MLE}$ 是以下方程的根:

$$\frac{\partial \log L(\boldsymbol{\beta})}{\partial \boldsymbol{\beta}} = \mathbf{0}. \quad (8)$$

使用 Section 1 的 Newton-Raphson 算法求解。令

$$\mathbf{h}(\boldsymbol{\beta}) = \frac{\partial \log L(\boldsymbol{\beta})}{\partial \boldsymbol{\beta}} = X^\top \mathbf{z}.$$

其中 $X = (\mathbf{x}_1, \dots, \mathbf{x}_n)^\top$, $\mathbf{z} = (z_1, \dots, z_n)^\top$, $z_i = y_i - \pi_i$, $i = 1, \dots, n$.

迭代中还需计算:

$$\begin{aligned} \nabla \mathbf{h}(\boldsymbol{\beta}) &= \frac{\partial^2 \log L(\boldsymbol{\beta})}{\partial \boldsymbol{\beta} \partial \boldsymbol{\beta}^\top} = -X^\top W X \\ [\nabla \mathbf{h}(\boldsymbol{\beta})]^{-1} &= -(X^\top W X)^{-1} \end{aligned}$$

其中 $W = \text{diag}(w_1, \dots, w_n)$, $w_i = \pi_i(1 - \pi_i)$, $i = 1, \dots, n$. 此时 Newton 迭代的公式为

$$\begin{aligned} \boldsymbol{\beta}^{(t+1)} &= \boldsymbol{\beta}^{(t)} - [\nabla \mathbf{h}(\boldsymbol{\beta}^{(t)})]^{-1} \mathbf{h}(\boldsymbol{\beta}^{(t)}) \\ &= \boldsymbol{\beta}^{(t)} + (X^\top W^{(t)} X)^{-1} X^\top \mathbf{z}^{(t)} \end{aligned} \quad (9)$$

其中 $W^{(t)}$, $\mathbf{z}^{(t)}$ 为 W , \mathbf{z} 在第 t 步的取值, 因为 π_i 在第 t 步的估计值依赖 $\boldsymbol{\beta}^{(t)}$.

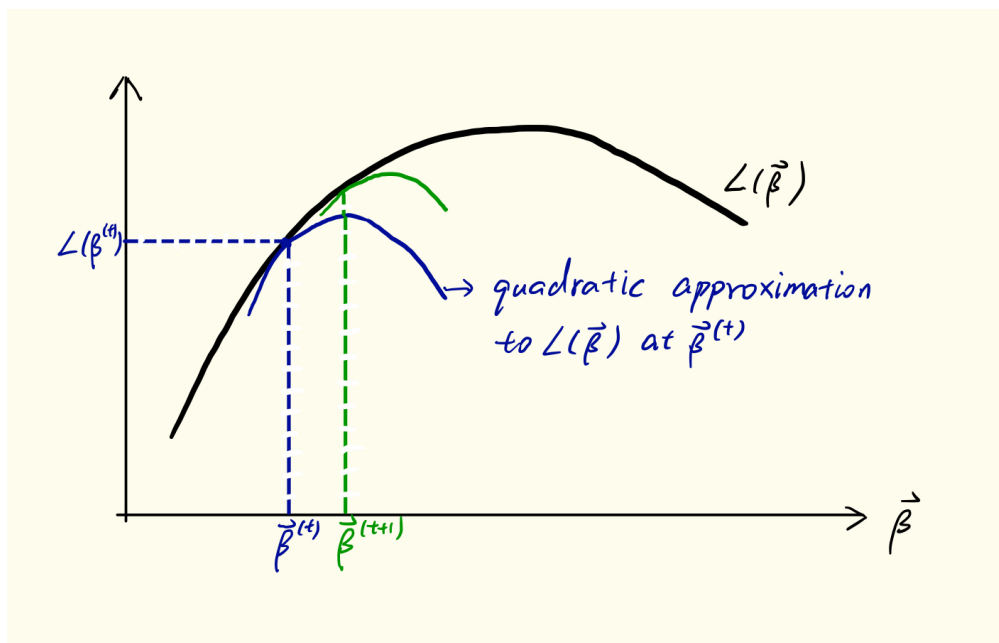
Remarks

1. 以上使用 Newton-Raphson 求解 MLE 的过程需要用到 score function 和 observed Fisher's information, 因此这种算法在统计中也被称为 Fisher-scoring method.
2. 迭代(9)也常写为以下形式:

$$\boldsymbol{\beta}^{(t+1)} = \boldsymbol{\beta}^{(t)} + (X^\top W^{(t)} X)^{-1} X^\top W^{(t)} \tilde{\mathbf{z}}^{(t)} \quad (10)$$

其中 $\tilde{\mathbf{z}}^{(t)} = [W^{(t)}]^{-1} \mathbf{z}^{(t)}$. 此时等号右边第二项跟 weighted least square (WLS) 估计量的形式相同。这样从初始值 $\boldsymbol{\beta}^{(0)}$ 开始, 每一步迭代只需先计算出 $W^{(t)}$, $\mathbf{z}^{(t)}$ 和 $\tilde{\mathbf{z}}^{(t)}$, 再做一个 WLS 就可以得到 $\boldsymbol{\beta}^{(t+1)}$. 因此迭代(10)也被称为 iteratively reweighted least square (IRLS).

3. 对于 Logistic 回归, Newton-Raphson 通常从任意初始值开始都收敛得很快。
4. 从几何角度考虑, 使用 Newton-Raphson 求解方程(8)时, 每一步迭代使用一阶 Taylor 展开 (线性函数) 近似 score function, 这其实相当于用一个二次函数近似 log-likelihood. 更新的 $\boldsymbol{\beta}$ 是这个二次函数的最大值点, 对应更大的 log-likelihood, 最终会收敛到 MLE.



习题:数据 `facerecognition.RData` (<https://github.com/wangronglu/statcomp/blob/gh-pages/facerecognition.RData>) 记录了一个人脸识别算法的测试结果。该实验将每个受试者的一张图片输入人脸识别算法, 如果算法能从剩余的 2143 张图片中正确挑选出受试者的另一张图片, 称为一次成功的匹配 (match). 如果算法成功匹配受试者 i 的两张图片, 则响应变量 $y_i = 1$, 否则 $y_i = 0$. 解释变量 `eyediff` 记录了受试者两张图片眼部区域像素强度之差的绝对值。

1. 对数据 `facerecognition.RData` 做 Logistic 回归, 自己编写 Newton-Raphson 算法估计回归系数. 初始值取为 $\beta^{(0)} = (\beta_0^{(0)}, \beta_1^{(0)})^\top = (0.96, 0)^\top$, 迭代停止容限 $\varepsilon = 10^{-5}$, 作图展示梯度模长 $\left\| \frac{\partial \log L(\beta)}{\partial \beta} \right\|_2$ 随迭代变化情况。
2. 作图展示向量 $\beta^{(t)}$ 随迭代的移动轨迹。将最终结果与 R 函数 `glm(..., family="binomial")` 的系数估计结果进行比较。
3. 在 2 的轨迹图中加入 $\log L(\beta)$ 的 contour plot.

References

Geof H.Givens, J. A. H. (2012). *Computational Statistics*. John Wiley & Sons, Inc, second edition.