

Non-stochastic Optimization

Background: likelihood inference

- let x_1, \dots, x_n be an iid sample from $f(x|\theta^*)$, where true parameter value θ^* is unknown
- the likelihood function is $L(\theta) = \prod_{i=1}^n f(x_i|\theta)$
- the maximum likelihood estimator (MLE) of θ is the maximizer of $L(\theta)$
- usually it is easier to work with the loglikelihood $l(\theta) = \log L(\theta)$
- typically maximization of $l(\theta)$ is done by solving $l'(\theta) = 0$
 - $l'(\theta)$ is called the score function
- for any θ , $E_{\theta}\{l'(\theta)\} = 0$

$$E_{\theta}\{l'(\theta) l'(\theta)^T\} = -E_{\theta}\{l''(\theta)\}$$

where E_{θ} is expectation wrt $f(x|\theta)$

- Fisher information: $I(\theta) = E_{\theta}\{l'(\theta) l'(\theta)^T\}$
- observed Fisher information: $-l''(\theta)$
 - if $\dim(\theta) = 1$, $I(\theta)$ is a nonnegative number
 - if $\dim(\theta) > 1$, $I(\theta)$ is a nonnegative definite matrix

- Importance of $I(\theta)$: it sets the limit on how accurate an unbiased estimate of θ can be

- as $n \rightarrow \infty$, $\sqrt{n}(\hat{\theta}_{MLE} - \theta^*) \Rightarrow N_p(0, I(\theta^*)^{-1})$

Working with Derivatives

- suppose $g(x)$ is a differentiable function, where $x = (x_1, \dots, x_n)$

- to find its maximum/minimum, one method is to solve the equation $g'(x) = 0$, where $g'(x) = \left(\frac{\partial g}{\partial x_1}, \dots, \frac{\partial g}{\partial x_n} \right)^T$

- that is, maximization/minimization is equivalent to solving $f(x) = 0$ where $f = g'$

Univariate Case

Newton's method

- a fast approach to solve $f(x) = 0$

- steps: (i) start with an initial estimate x_0
(ii) for $t = 0, 1, \dots$, compute

$$x_{t+1} = x_t + h_t \quad \text{with } h_t = -\frac{f(x_t)}{f'(x_t)}$$

(iii) continue until convergence

- also known as Newton-Raphson
- need to specify x_0
- if $f(x)=0$ has multiple roots, end result will depend on x_0
- iteration cannot continue if $f'(x_e)=0$

Why it works?

- let x^0 be true solution, \bar{x} be an approximation of x^0
- Taylor expansion: $f(x) = f(\bar{x}) + (x - \bar{x}) f'(\bar{x}) + \frac{(x - \bar{x})^2}{2} f''(\tilde{x})$
where \tilde{x} lies between x and \bar{x}
- since $f(x^0)=0$, we have $0 = f(\bar{x}) + (x^0 - \bar{x}) f'(\bar{x}) + \frac{(x^0 - \bar{x})^2}{2} f''(\tilde{x})$
- since x^0 and \bar{x} are close, the last term can be ignored:
$$0 \approx f(\bar{x}) + (x^0 - \bar{x}) f'(\bar{x}) \Rightarrow x^0 \approx \bar{x} - \frac{f(\bar{x})}{f'(\bar{x})}$$

Optimization with Newton

- can be applied to optimize g by applying to $f=g'$
- both g' (gradient) and g'' (Hessian) are needed
- many variants of Newton's method avoid the computation of g'' , which can be difficult, especially for multivariate functions

Example: To maximize $g(x) = \frac{\log x}{1+x}$

- first find $f(x) = g'(x) = \frac{1 + \frac{1}{x} - \log x}{(1+x)^2}$

$$f'(x) = g''(x) = \frac{-(3 + 4/x + 1/x^2 - 2 \log x)}{(1+x)^3}$$

- therefore $h_t = \frac{(x_t + 1)(1 + \frac{1}{x_t} - \log x_t)}{3 + \frac{4}{x_t} + \frac{1}{x_t^2} - 2 \log x_t}$

- a simpler formula: note that solving $f(x) = 0$ is the same as solving $1 + \frac{1}{x} - \log x$.

- treat $1 + \frac{1}{x} - \log x$ as a new f function

- then $h_t = x_t - \frac{x_t^2 \log x_t}{1 + x_t} \Rightarrow x_{t+1} = 2x_t - \frac{x_t^2 \log x_t}{1 + x_t} \neq$

Example:

to maximize loglikelihood $l(\theta)$, $\theta_{t+1} = \theta_t - \frac{l'(\theta_t)}{l''(\theta_t)}$

- consider the model with shift $p(x|\theta) = p(x-\theta)$.

- given observations $x_1, \dots, x_n \text{ iid } \sim p(x|\theta)$,

$$l(\theta) = \sum_{i=1}^n \log p(x_i - \theta), \quad l'(\theta) = - \sum_{i=1}^n \frac{p'(x_i - \theta)}{p(x_i - \theta)}$$

$$l''(\theta) = \sum_{i=1}^n \frac{p''(x_i - \theta)}{p(x_i - \theta)} - \sum_{i=1}^n \left\{ \frac{p'(x_i - \theta)}{p(x_i - \theta)} \right\}^2$$

- note that we update θ , not x_1, \dots, x_n

- In R, to minimize a function, one can use

$$Z = \text{nlmminb}(x_0, g, \text{gr.g}, \text{hess.g})$$

x_0 : initial value

g : function being minimized

$\left. \begin{array}{l} \text{gr.g: gradient of } g \\ \text{hess.g: Hessian of } g \end{array} \right\} \text{ have to be analytically calculated}$

- one can also use

$$Z = \text{nlmminb}(x_0, g, \text{gr.g}) \text{ or } Z = \text{nlmminb}(x_0, g),$$

where $\text{gr.g} / \text{hess.g}$ will be numerically approximated

Secant Method

- approximating $f'(x_t)$ by $\frac{f(x_t) - f(x_{t-1})}{x_t - x_{t-1}}$, the Newton method becomes the secant method:

$$x_{t+1} = x_t - \frac{f(x_t)(x_t - x_{t-1})}{f(x_t) - f(x_{t-1})}$$

- need to specify x_0 and x_1

Fisher Scoring

- another variant of Newton's method
- specific for MLE
- replace the Hessian $l''(\theta)$ by its expectation; i.e., Fisher information $I(\theta)$

$$\theta_{t+1} = \theta_t + \frac{l'(\theta_t)}{I(\theta_t)}$$

- in practice, use Fisher scoring in the beginning to make rapid improvements, then Newton's method for refinement near the end

Example continue with the previous example on $p(x|\theta) = p(x-\theta)$

- to use Fisher scoring, need to compute $I(\theta) = -E_{\theta}(l''(\theta))$

$$\begin{aligned}\Rightarrow I(\theta) &= -n E_{\theta} \left[\frac{p''(x-\theta)}{p(x-\theta)} - \left\{ \frac{p'(x-\theta)}{p(x-\theta)} \right\}^2 \right] \\&= -n \int \left[\frac{p''(x-\theta)}{p(x-\theta)} - \left\{ \frac{p'(x-\theta)}{p(x-\theta)} \right\}^2 \right] p(x-\theta) dx \\&= -n \int p''(x-\theta) dx + n \int \frac{[p'(x-\theta)]^2}{p(x-\theta)} dx \\&= -n \frac{d^2}{d\theta^2} \int p(x-\theta) dx + n \int \frac{[p'(x)]^2}{p(x)} dx \\&= -n \frac{d^2}{d\theta^2} 1 + n \int \frac{p'(x)^2}{p(x)} dx = n \int \frac{p'(x)^2}{p(x)} dx \quad \# \end{aligned}$$

Multivariate Case now g is a function in $\underline{x} = (x_1, \dots, x_p)^T$.

Newton's Method

- generalization is straight forward
- to maximize/minimize $g(\underline{x})$, use

$$\underline{x}_{t+1} = \underline{x}_t - [g''(\underline{x}_t)]^{-1} g'(\underline{x}_t)$$

- $g''(\underline{x}_t)$: $p \times p$ matrix with (i,j) th element as $\frac{\partial^2 g(\underline{x})}{\partial x_i \partial x_j}$

- $g'(\underline{x}) = \left[\frac{\partial g(\underline{x})}{\partial x_1}, \dots, \frac{\partial g(\underline{x})}{\partial x_p} \right]^T$, a $p \times 1$ vector

- note: need to compute the inverse of $g''(\underline{x}_t)$

Fisher Scoring

- use $\underline{\theta}_{t+1} = \underline{\theta}_t + I(\underline{\theta}_t)^{-1} l'(\underline{\theta}_t)$

Other Newton-like Methods

- computing $g''(\underline{x})$ or $[g''(\underline{x})]^{-1}$ could be hard.
- the idea is to replace $g''(\underline{x})$ by some easily-computable matrix, say $M(\underline{x})$
- $\underline{x}_{t+1} = \underline{x}_t - M_t^{-1} g'(\underline{x}_t)$

Steepest Ascent Method

- set $M_t = -\alpha_t^{-1} I_p$ (I_p : identity matrix)
- $\underline{x}_{t+1} = \underline{x}_t + \alpha_t g'(\underline{x}_t)$
- $\alpha_t > 0$: step size at t which can shrink to ensure ascent
- if at step t , the original step turns out to be downhill, the updating can be back track by halving α_t
- also known as steepest Descent (for minimization)

Gauss-Newton Method

- want to maximize $g(\underline{\theta}) = -\sum_{i=1}^n \{y_i - f_i(\underline{\theta})\}^2$
where each $f_i(\underline{\theta})$ is differentiable
- first consider linear regression $y_i = \underline{x}_i^T \underline{\theta} + \varepsilon, i=1, \dots, n$
- the least-squares estimator of $\underline{\theta}$ maximizes $g(\underline{\theta})$ with $f_i(\hat{\underline{\theta}}) = \underline{x}_i^T \hat{\underline{\theta}}$
- $\hat{\underline{\theta}} = (\underline{X}^T \underline{X})^{-1} \underline{X}^T \underline{y}$ where $\underline{X} = \begin{pmatrix} \underline{x}_1^T \\ \underline{x}_2^T \\ \vdots \\ \underline{x}_n^T \end{pmatrix}, \underline{y} = \begin{pmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{pmatrix}$
- Gauss Newton uses a similar idea for nonlinear $f_i(\underline{\theta})$

- let $\underline{\theta}^*$ be the unknown maximizer of $g(\underline{\theta})$

- consider $h(\underline{u}) = - \sum_{i=1}^n \{y_i - f_i(\underline{\theta} + \underline{u})\}^2$

- $h(\underline{u})$ is maximized by $\underline{u}^* = \underline{\theta}^* - \underline{\theta}$ (\underline{u}^* unknown)

- if $\underline{\theta}$ is near $\underline{\theta}^*$, $\underline{u}^* \approx 0$ and by Taylor expansion of $h(\underline{u})$, \underline{u}^* should be close to the maximizer of

$$- \sum_{i=1}^n \{y_i - f_i(\underline{\theta}) - f'_i(\underline{\theta})^T \underline{u}\}^2$$

- treat $y_i - f_i(\underline{\theta})$ as y_i as in linear regression

- $f'_i(\underline{\theta})$ as x_i

- we have $\underline{u}^* = \underline{\theta}^* - \underline{\theta} \approx (\underline{A}^T \underline{A})^{-1} \underline{A}^T \underline{z}$

$$\text{where } \underline{A} = \underline{A}(\underline{\theta}) = \begin{pmatrix} f'_1(\underline{\theta})^T \\ \vdots \\ f'_n(\underline{\theta})^T \end{pmatrix}, \quad \underline{z} = \underline{z}(\underline{\theta}) = \begin{pmatrix} y_1 - f_1(\underline{\theta}) \\ \vdots \\ y_n - f_n(\underline{\theta}) \end{pmatrix}$$

- the updating formula is

$$\underline{\theta}_{t+1} = \underline{\theta}_t + (\underline{A}_t^T \underline{A}_t)^{-1} \underline{A}_t^T \underline{z}_t$$

$$\text{where } \underline{A}_t = \underline{A}(\underline{\theta}_t), \quad \underline{z}_t = \underline{z}(\underline{\theta}_t)$$

some tricks for Newton / Fisher Scoring:

- calculate $g''(\underline{x}_t) / I(\underline{\theta}_t)$ every, say, 3 iteration

- use $M_t = \alpha I + g''(\underline{x}_t)$

$$\approx M_t = \alpha I + I(\underline{\theta}_t)$$

if $g''(\underline{x}_t) / I(\underline{\theta}_t)$ is near singular, where $\alpha > 0$