

On-Line Supp.

How do we proceed to maximize the log-likelihood function of a sample by selecting the optimizing parameters, subject to $\theta \in \Theta$? Appropriate **methods of numerical, constrained optimization** need to be implemented: this is what packages such as Matlab, Gauss, EViews, or Stata are for. For instance (i.e., other, better but more complex methods are feasible), **Newton's method** makes use of the **Hessian**, which is a $K \times K$ matrix $H(\theta) \equiv \partial^2 \ell(\theta) / \partial \theta \partial \theta'$ that collects second partial derivatives of the log-likelihood function with respect to each of the parameters in θ . Similarly the $K \times 1$ **gradient** $\partial \ell(\theta) / \partial \theta$ collects the first partial derivatives of the log-likelihood function with respect to each of the elements in θ . Let $\hat{\theta}_j$ denote the value of the vector of estimates at step j of the algorithm, and let $\partial \ell(\hat{\theta}_j) / \partial \theta$ and $H(\hat{\theta}_j)$ denote, respectively, the gradient and the Hessian evaluated at $\hat{\theta}_j$. Then the fundamental equation to update the estimates according to Newton's algorithm is:

$$\hat{\theta}_{j+1} = \hat{\theta}_j - H^{-1}(\hat{\theta}_j) [\partial \ell(\hat{\theta}_j) / \partial \theta] \quad (0.1)$$

Because the log-likelihood function is to be maximized, the Hessian should be negative definite, at least when $\hat{\theta}_j$ is sufficiently near $\hat{\theta}_r$. This ensures that this step is in an uphill direction. The maximization process therefore proceeds through the following steps:

- Set an initial vector of parameters, $\hat{\theta}_0$, and compute $H^{-1}(\hat{\theta}_0)$ and $\partial \ell(\hat{\theta}_0) / \partial \theta$.
- Compute the new vector of estimated parameters $\hat{\theta}_1 = \hat{\theta}_0 - H^{-1}(\hat{\theta}_0) [\partial \ell(\hat{\theta}_0) / \partial \theta]$ and therefore $H^{-1}(\hat{\theta}_1)$ and $\partial \ell(\hat{\theta}_1) / \partial \theta$; check that the **Euclidean norm** $\|\hat{\theta}_1 - \hat{\theta}_0\|$ (in words, this is the square root of the sum of all squared differences between the elements of $\hat{\theta}_1$ and $\hat{\theta}_0$) is not inferior to some small threshold parameter (typically, 10^{-5}).
- Update the vector of parameter estimates to $\hat{\theta}_2 = \hat{\theta}_1 - H^{-1}(\hat{\theta}_1) [\partial \ell(\hat{\theta}_1) / \partial \theta]$ and check that the norm $\|\hat{\theta}_2 - \hat{\theta}_1\|$ is not inferior to the threshold parameter.
- Continue (unless a maximum number of iteration has been

exceeded, but with fast computers often thousands of iterations are affordable in the space of a few minutes only) until $\hat{\theta}_j = \hat{\theta}_{j-1} - H^{-1}(\hat{\theta}_{j-1})[\partial\ell(\hat{\theta}_{j-1})/\partial\theta]$ is such that $\|\hat{\theta}_j - \hat{\theta}_{j-1}\|$ falls below the fixed convergence threshold, that signals that the optimizing vector has stopped changing.

- Set $\hat{\theta}_T^{ML} = \hat{\theta}_j$.

Numerical optimization is a very sensitive business; a myriad of choices are considered to be crucial to obtain “reliable” results, such as the initial value $\hat{\theta}_0$, the convergence tolerance criterion, and often how much the algorithm is supposed to “travel” in the direction indicated by the inverse Hessian matrix, i.e., the coefficient τ in the iteration in (5.123), generalized to read as $\hat{\theta}_{j+1} = \hat{\theta}_j - \tau H^{-1}(\hat{\theta}_j)[\partial\ell(\hat{\theta}_j)/\partial\theta]$, where $\tau > 0$ (clearly a $\tau < 1$ “dims” the step taken in direction $[\partial\ell(\hat{\theta}_j)/\partial\theta]$, while a $\tau > 1$ acts as a multiplier). Reliability here is often evidence or even taken to offer some guarantee that $\hat{\theta}_T^{ML} = \hat{\theta}_j$ truly represents a **global** (as opposed to local) **maximizer** of the log-likelihood function and as such it is unique, as assumed. For instance, just to get hard evidence on this aspect, it is often advised to start off the maximization algorithm in correspondence of a range of alternative starting values and then retain, for the true and often lengthy iterative Newton-style search, the most promising one(s).

Other numerical optimization methods are of course possible. A few of them are faster than Newton’s method because they replace the Hessian matrix with cheaper to compute negative definite $K \times K$ matrices, for instance $OPG(\theta) \equiv -[\partial\ell(\theta)/\partial\theta][\partial\ell(\theta)/\partial\theta]'$, which is negative definite by construction, unless $\partial\ell(\theta)/\partial\theta = \mathbf{0}$, which would instead show that a stationary point has been reached. The advantage of this expression is that it only requires calculation (often numerically) of first-order derivatives. Moreover, our simplified illustration of Newton’s method ignores the role played by constraints, that may interfere with setting $\hat{\theta}_j = \hat{\theta}_{j-1} - H^{-1}(\hat{\theta}_{j-1})[\partial\ell(\hat{\theta}_{j-1})/\partial\theta]$, when the constraints are violated.