

FIML Estimation for Multivariate Normal Data with Newton's Algorithm

This document provides the mathematical building blocks for applying Newton's algorithm to obtain maximum likelihood estimates of a mean vector and covariance matrix, μ and Σ . The accompanying R program shows the computations for a bivariate normal estimation problem.

Derivative Equations

The Matrix Cookbook is good resource that gives details behind some of the expressions and matrix derivatives in this document. The Cookbook is available at www.math.uwaterloo.ca/~hwolkowi/matrixcookbook.pdf. To begin, the first derivatives of the log-likelihood with respect to μ and Σ are as follows.

$$\frac{\partial LL}{\partial \mu} = -N\Sigma^{-1}\mu + \Sigma^{-1} \sum_{i=1}^N Y_i \quad (1)$$

$$\frac{\partial LL}{\partial \Sigma} = -\frac{1}{2} \sum_{i=1}^N (\Sigma^{-1} - \Sigma^{-1}(Y_i - \mu)(Y_i - \mu)'\Sigma^{-1}) \quad (2)$$

Note that Equation 1 is a vector but Equation 2 is itself a matrix (a 2 by 2 matrix in the bivariate illustration).

Second derivatives quantify the curvature or steepness of the log-likelihood function near its peak (i.e., the rate at which the first-order slopes change across the range of parameter values). Second derivatives are obtained by applying matrix calculus rules to the previous equations, and the Hessian collects these equations in a symmetric matrix with P rows and columns, where P is the number of unique parameters in μ and Σ . The Hessian consists of three unique blocks: a block each for μ and Σ and a block for the cross-product derivatives involving an element from μ and an element from Σ .

$$H_O(\theta) = \begin{pmatrix} \frac{\partial^2 LL}{\partial \mu^2} & \frac{\partial^2 LL}{\partial \mu \partial \Sigma} \\ \frac{\partial^2 LL}{\partial \Sigma \partial \mu} & \frac{\partial^2 LL}{\partial \Sigma^2} \end{pmatrix} \quad (3)$$

The second derivative equations below are the building blocks for the observed information matrix, the inverse of which is the variance–covariance matrix of the estimates with squared standard errors on the diagonal. The blocks of the Hessian matrix are as follows.

$$\frac{\partial^2 LL}{\partial \mu^2} = -N\Sigma^{-1} \quad (4)$$

$$\frac{\partial^2 LL}{\partial \Sigma^2} = - \sum_{i=1}^N \mathbf{D}'_V (\Sigma^{-1} \otimes \{ \Sigma^{-1} (Y_i - \mu)(Y_i - \mu)' \Sigma^{-1} - .5\Sigma^{-1} \}) \mathbf{D}_V \quad (5)$$

$$\frac{\partial^2 LL}{\partial \mu \partial \Sigma} = - \sum_{i=1}^N (\Sigma^{-1} \otimes (Y_i - \mu)' \Sigma^{-1}) \mathbf{D}_V \quad (6)$$

The \otimes symbol is a Kronecker product that multiplies one matrix by each element of another matrix, and \mathbf{D}_V is the duplication matrix (Magnus & Neudecker, 1999). Each covariance term appears only once in the Hessian (and similarly, only once in the variance–covariance matrix of the estimates). The duplication matrix combines redundant off-diagonal terms into a single value. Substituting the maximum likelihood estimates into the derivative expressions, multiplying $H_O(\hat{\theta})$ by -1 , then taking its inverse gives the variance–covariance matrix of the estimates.

Newton's Algorithm

With Newton's algorithm, the jump from the current to the updated parameter value was as follows.

$$\theta^{(t+1)} = \theta^{(t)} - \left(\frac{\partial^2 LL}{\partial \theta^2} \right)^{-1} \left(\frac{\partial LL}{\partial \theta} \right) \text{current estimate} + \text{step size} \quad (7)$$

The step size, computed as the ratio of the first and second derivatives at the current parameter value θ^t , corresponds to the horizontal distance between the current estimate and the peak of the projected quadratic curve (the point at which the first derivative of the quadratic approximation equals 0). The updating step extends to more complex models with multiple parameters. In this case, the multivariate updating equation is

$$\theta^{(t+1)} = \theta^{(t)} - H_O^{-1}(\theta^{(t)}) \nabla LL^{(t)} \quad (8)$$

where θ is a vector of parameter values, t indexes the iterations, $H_O^{-1}(\theta^{(t)})$ is the inverse of the Hessian (the matrix of second derivatives) evaluated at the current parameter estimates at iteration t , and $\nabla LL^{(t)}$ is a vector of first derivatives (the gradient vector) evaluated at the current estimates.

$$\nabla LL^{(t)} = \begin{pmatrix} \frac{\partial LL}{\partial \mu_X^{(t)}} \\ \frac{\partial LL}{\partial \mu_Y^{(t)}} \\ \vdots \\ \frac{\partial LL}{\partial \sigma_Y^{2(t)}} \end{pmatrix} \quad (9)$$

Recall that Equation 1 (the first derivatives with respect to μ) returned a vector, and these quantities appear as the first two elements of the gradient vector. However, the first derivative of the log-likelihood with respect to Σ in Equation 2 returns matrix (a 2 by 2 matrix in the bivariate example). The vec function in the equation below stacks all matrix elements into a vector (e.g., a 4-element vector if Σ is a 2 by 2 matrix), and the duplication matrix combines the two identical off-diagonal elements into a single quantity.

$$\nabla LL_{\Sigma} = -\frac{1}{2} \sum_{i=1}^N D'_V \text{vec}(\Sigma^{-1} - \Sigma^{-1}(Y_i - \mu)(Y_i - \mu)' \Sigma^{-1}) \quad (10)$$

For example, vector if Σ is a 2 by 2 matrix, Equation 10 returns a 3-element vector, one element for each unique estimate in Σ .