

where  $\mathbf{f}_i(\cdot)$  is an  $n$ -vector of nonlinear functions,  $\mathbf{u}_i$  is an  $n$ -vector of error terms, and  $\boldsymbol{\theta}$  is a  $p$ -vector of parameters to be estimated. In general, subject to whatever restrictions need to be imposed for the system to be identified, all the endogenous and exogenous variables and all the parameters may appear in any equation.

The first step in any sort of IV procedure is to choose the instruments to be used. If the model is nonlinear only in the parameters, the matrix of optimal instruments is  $\mathbf{X}$ . As we have seen, however, there is no simple way to choose the instruments for models that are nonlinear in one or more of the endogenous variables. The theory of Section 17.4 can be applied, of course, but the result that it yields is not very practical. Under the usual assumptions about the error terms, namely, that they are homoskedastic and independent across observations but correlated across equations for each observation, one finds that a matrix of instruments  $\mathbf{W}$  will be optimal if  $\mathcal{S}(\mathbf{W})$  is equal to the subspace spanned by the union of the columns of the  $E(\partial \mathbf{f}_i / \partial \boldsymbol{\theta})$ . This result was originally derived by Amemiya (1977). It makes sense but is generally not very useful in practice. For now, we simply assume that *some* valid  $n \times m$  matrix of instruments  $\mathbf{W}$  is available, with  $m \geq p$ .

A nonlinear IV procedure for full-system estimation, similar in spirit to the single-equation NL2SLS procedure based on minimizing (18.78), was first proposed by Jorgenson and Laffont (1974) and called **nonlinear three-stage least squares**, or **NL3SLS**. The name is somewhat misleading, for the same reason that the name “NL2SLS” is misleading. By analogy with (18.60), the criterion function we would really like to minimize is

$$\sum_{i=1}^g \sum_{j=1}^g \sigma^{ij} \mathbf{f}_i^\top(\mathbf{Y}, \mathbf{X}, \boldsymbol{\theta}) \mathbf{P}_W \mathbf{f}_j(\mathbf{Y}, \mathbf{X}, \boldsymbol{\theta}). \quad (18.80)$$

In practice, however, the elements  $\sigma^{ij}$  of the inverse of the contemporaneous covariance matrix  $\boldsymbol{\Sigma}$  will not be known and will have to be estimated. This may be done in several ways. One possibility is to use NL2SLS for each equation separately. This will generally be easy, but it may not be possible if some parameters are identified only by cross-equation restrictions. Another approach which will work in that case is to minimize the criterion function

$$\sum_{i=1}^g \mathbf{f}_i^\top(\mathbf{Y}, \mathbf{X}, \boldsymbol{\theta}) \mathbf{P}_W \mathbf{f}_i(\mathbf{Y}, \mathbf{X}, \boldsymbol{\theta}), \quad (18.81)$$

in which the unknown covariance matrix  $\boldsymbol{\Sigma}$  is replaced by the identity matrix. The estimator obtained by minimizing (18.81) will evidently be a valid GMM estimator and thus will be consistent even though it is inefficient. Whichever inefficient estimator is used initially, it will yield  $g$  vectors of residuals  $\hat{\mathbf{u}}_i$  from which the matrix  $\boldsymbol{\Sigma}$  may be estimated consistently in exactly the same way as for linear models; see (18.62). Replacing the unknown  $\sigma^{ij}$ 's in (18.80) by

the elements  $\hat{\sigma}^{ij}$  of the inverse of the estimate of  $\Sigma$  then yields the criterion function

$$\sum_{i=1}^g \sum_{j=1}^g \hat{\sigma}^{ij} \mathbf{f}_i^\top(\mathbf{Y}, \mathbf{X}, \boldsymbol{\theta}) \mathbf{P}_W \mathbf{f}_j(\mathbf{Y}, \mathbf{X}, \boldsymbol{\theta}), \quad (18.82)$$

which can actually be minimized in practice.

As usual, the minimized value of the criterion function (18.82) provides a test statistic for overidentifying restrictions; see Sections 7.8 and 17.6. If the model and instruments are correctly specified, this test statistic will be asymptotically distributed as  $\chi^2(m - p)$ ; recall that  $m$  is the number of instruments and  $p$  is the number of free parameters. Moreover, if the model is estimated unrestrictedly and subject to  $r$  distinct restrictions, the difference between the two values of the criterion function will be asymptotically distributed as  $\chi^2(r)$ . If the latter test statistic is to be employed, it is important that the same estimate of  $\Sigma$  be used for both estimations, since otherwise the test statistic may not even be positive in finite samples.

When the sample size is large, it may be less computationally demanding to obtain one-step efficient estimates rather than actually to minimize (18.82). Suppose the initial consistent estimates, which may be either NL2SLS estimates or systems estimates based on (18.81), are denoted  $\hat{\boldsymbol{\theta}}$ . Then a first-order Taylor-series approximation to  $\mathbf{f}_i(\boldsymbol{\theta}) \equiv \mathbf{f}_i(\mathbf{Y}, \mathbf{X}, \boldsymbol{\theta})$  around  $\hat{\boldsymbol{\theta}}$  is

$$\mathbf{f}_i(\hat{\boldsymbol{\theta}}) + \mathbf{F}_i(\hat{\boldsymbol{\theta}})(\boldsymbol{\theta} - \hat{\boldsymbol{\theta}}),$$

where  $\mathbf{F}_i$  is an  $n \times p$  matrix of the derivatives of  $\mathbf{f}_i(\boldsymbol{\theta})$  with respect to the  $p$  elements of  $\boldsymbol{\theta}$ . If certain parameters do not appear in the  $i^{\text{th}}$  equation, the corresponding columns of  $\mathbf{F}_i$  will be identically zero. The one-step estimates, which will be asymptotically equivalent to NL3SLS estimates, are simply  $\hat{\boldsymbol{\theta}} = \hat{\boldsymbol{\theta}} - \hat{\boldsymbol{\epsilon}}$ , where  $\hat{\boldsymbol{\epsilon}}$  denotes the vector of linear 3SLS estimates

$$\hat{\boldsymbol{\epsilon}} = \left[ \sum_{i=1}^g \sum_{j=1}^g \hat{\sigma}^{ij} \hat{\mathbf{F}}_i^\top \mathbf{P}_W \hat{\mathbf{F}}_j \right]^{-1} \left[ \sum_{i=1}^g \sum_{j=1}^g \hat{\sigma}^{ij} \hat{\mathbf{F}}_i^\top \mathbf{P}_W \hat{\mathbf{f}}_j \right]. \quad (18.83)$$

Compare expression (18.64), for the case with no cross-equation restrictions.

It is clear that NL3SLS can be generalized to handle heteroskedasticity of unknown form, serial correlation of unknown form, or both. For example, to handle heteroskedasticity one would simply replace the matrix  $\mathbf{P}_W$  in (18.82) and (18.83) by the matrix

$$\mathbf{W}(\mathbf{W}^\top \hat{\boldsymbol{\Omega}}_{ij} \mathbf{W})^{-1} \mathbf{W}^\top,$$

where, by analogy with (18.76),  $\hat{\boldsymbol{\Omega}}_{ij} = \text{diag}(\hat{u}_{ti} \hat{u}_{tj})$  for  $i, j = 1, \dots, g$ . The initial estimates  $\hat{\boldsymbol{\theta}}$  need not take account of heteroskedasticity. For a more detailed discussion of this sort of procedure, and of NL3SLS in general, see Gallant (1987, Chapter 6).