# Lecture 9

## Regression extensions – Multivariate and nonlinear regression

Lectures in SDPE: *Econometrics I* on February 26, 2024

Markus Jäntti
Swedish Institute For Social Research
Stockholm University

- This lecture discusses least squares based regression extensions – e.g., generalized least squares (GLS), esp. Seemingly Unrelated Regressions [SUR, chapter 11 in Hansen (2021)] and non-linear models, estimated using non-linear least squares, [NLS, chapter 23 in Hansen (2021)]
- The treatment is brief, but
  - GLS/SUR and NLS both contribute to understanding an important later topic, generalized method of moments (GMM)
  - NLS is in a way "close" to the linear projection model studied earlier

- We now consider a set of regression equations of the form

$$Y_{ji} = X'_{ji}\beta_j + e_{ji}. \tag{1}$$

  There are $j = 1, \ldots, m$ dependent variables and $i = 1, \ldots, n$ observations with each vector of regressors $X_{ji}$ and associated coefficient vector $\beta_j$ having $k_j$ elements; $e_{ji}$ is a regression error and the total number of coefficients is $\bar{k} = \sum_{j=1}^{m} k_j$.

- Set up in this way, observations $i$ are treated as independent but the variables $j$ as correlated (and not only through $X_{ji}$). A typical example might be the consumption of household $i$ of different goods $j$.

- The dependence across variables is captured by the $m \times m$ covariance matrix $\Sigma_i$:

$$E[e_i e'_i] = \Sigma_i, \tag{2}$$

  where $e_i$ is the vector of $m$ regression errors for observation $i$.

## Regression Systems

- A more compact way of writing the system is in terms of a $m \times 1$ dependent variable and regressions error, a $m \times \overline{k}$ regressor with associated $m \times 1$ coefficients:

$$y_i = \overline{X}_i \beta + e_i \qquad (3)$$

where the $m \times 1$ dependent variable is $y_i = (Y_{1i}, \ldots, Y_{mi})'$ and

$$\overline{X}_i = \begin{bmatrix} X'_{1i} & 0 & \ldots & 0 \\ 0 & X'_{2i} & \ldots & 0 \\ \vdots & \vdots & \vdots & \vdots \\ 0 & 0 & \ldots & X'_{mi} \end{bmatrix}, \qquad (4)$$

or. . .

- . . . by stacking all $n$ observations into $mn \times 1$ and $mn \times \overline{k}$ matrices

$$y = \begin{bmatrix} y_1 \\ \vdots \\ y_m \end{bmatrix}, \qquad e = \begin{bmatrix} e_1 \\ \vdots \\ e_m \end{bmatrix}, \qquad \overline{X} = \begin{bmatrix} \overline{X_1} \\ \vdots \\ \overline{X_m} \end{bmatrix}, \qquad (5)$$

- . . . so we have

$$y = \overline{X}\beta + e \qquad (6)$$

- We might have the same regressors in all equations, so $X_{ji} = X_i$ and $k_j = k$, which can be written in many ways also, e.g., using a $m \times k$ parameter matrix

$$y_i = BX_i + e_i, \quad B = (\beta_1, \ldots, \beta_m) \tag{7}$$

or as

$$Y = XB + E \tag{8}$$

where $Y$ and $E$ are $n \times m$ matrices.

- With the same regressors, we can sometimes use the convenient notation involving the Kronecker product $\otimes$ that

$$\overline{X}_i = \begin{bmatrix} X_i' & \mathbf{0} & \ldots & \mathbf{0} \\ \mathbf{0} & X_i' & \ldots & \mathbf{0} \\ \vdots & \vdots & \vdots & \vdots \\ \mathbf{0} & \mathbf{0} & \ldots & X_i' \end{bmatrix} = I_m \otimes X_i' \tag{9}$$

- One approach to estimation is to apply least squares to each of the $j$ equations in eq 1,

$$\widehat{\beta_j} = \left( \sum_{i=1}^{n} X_{ji} X'_{ji} \right)^{-1} \left( \sum_{i=1}^{n} X_{ji} Y_{ji} \right) \tag{10}$$

and the full set of coefficients is $\widehat{\beta} = (\widehat{\beta'_1}, \ldots, \widehat{\beta'_m})'$.

# Least-Squares Estimator

- To estimate (under homoscedasticity) the error covariance matrix, note that the residuals are

$$\widehat{e_i} = y_i - \overline{X_i}\widehat{\beta}. \tag{11}$$

- The feasible estimator of the $m \times m$ variance matrix is then

$$\widehat{\Sigma} = n^{-1}\sum_{i=1}^{n}\widehat{e_i}\widehat{e_i}'. \tag{12}$$

- In order to determine the conditional mean and variance of $\widehat{\beta}$, we make the strong assumption of conditional mean independence, $E[e_i|X_i] = \mathbf{0}$ (here $X_i$ is the union of all $X_{ji}$). It follows that $E[Y_{ji}|X_i] = X'_{ji}\beta_j$

- To obtain the mean, center the estimator:

$$\widehat{\beta} - \beta = (\overline{X}'\overline{X})^{-1}(\overline{X}'e) = \left(\sum_{i=1}^{n}\overline{X}'_i\overline{X}_i\right)^{-1}\left(\sum_{i=1}^{n}\overline{X}'_ie_i\right). \tag{13}$$

- Now take the conditional expectation:

$$E[\widehat{\beta} - \beta|X] = \beta - \beta = \mathbf{0}. \tag{14}$$

## Mean and Variance of Systems Least-Squares

- To get the variance of the estimator, define (cf. equation 2)

$$\mathrm{E}[e_i e_i' | X_i] = \Sigma_i. \tag{15}$$

- With independence across observations, we have

$$\mathrm{E}[ee'|X] = \mathrm{E}\left( \begin{bmatrix} e_1 e_1' & e_1 e_2' & \dots & e_1 e_n \\ \vdots & \ddots & \dots & \vdots \\ e_n e_1' & e_2 e_1' & \dots & e_n e_n \end{bmatrix} \middle| X \right)$$
$$= \begin{bmatrix} \Sigma_1 & \mathbf{0} & \dots & \mathbf{0} \\ \vdots & \ddots & \dots & \vdots \\ \mathbf{0} & \dots & \mathbf{0} & \Sigma_n \end{bmatrix} \tag{16}$$

- By independence across $i$ we also have

$$\mathrm{Var}\left[ \sum_{i=1}^{n} \overline{X}_i' e_i \middle| X \right] = \sum_{i=1}^{n} \mathrm{Var}[\overline{X}_i' e_i | X_i] = \sum_{i=1}^{n} \overline{X}_i' \Sigma_i \overline{X}_i. \tag{17}$$

The variance of $\widehat{\beta}$ follows as:

$$\mathrm{Var}[\widehat{\beta}|X] = (\overline{X}'\overline{X})^{-1} \left( \sum_{i=1}^{n} \overline{X}_i' \Sigma_i \overline{X}_i \right) (\overline{X}'\overline{X})^{-1}. \tag{18}$$

Mean and Variance of Systems Least-Squares

- With common regressors, matters simplify:

$$\text{Var}[\widehat{\beta}|X] = \left(I_m \otimes (X'X)^{-1}\right)\left(\sum_{i=1}^{n}(\Sigma_i \otimes X_iX_i')\right)\left(I_m \otimes (X'X)^{-1}\right) \quad (19)$$

- With conditionally homoscedastic regressors ($\Sigma_i \equiv \Sigma$),

$$\text{Var}[\widehat{\beta}|X] = \left(\overline{X}'\overline{X}\right)^{-1}\left(\sum_{i=1}^{n}\overline{X}_i'\Sigma\overline{X}_i\right)\left(\overline{X}'\overline{X}\right)^{-1}. \quad (20)$$

- And with both common regressors and homoscedastic errors,

$$\text{Var}[\widehat{\beta}|X] = \Sigma \otimes (X'X)^{-1} \quad (21)$$

## Asymptotic Distribution

- For the asymptotic distribution, we can make do with the equation-by-equation linear projection condition, $E[X_{ji}e_{ji}] = 0$. This makes our $\widehat{\beta}_j$ *consistent* for $\beta_j$ (and all of $\beta$).

- The asymptotic *marginal* distribution of each $\widehat{\beta}_j$ is normal, but we need some additional material to determine the *joint* distribution of $\widehat{\beta}$.

- By our assumptions, the vector

$$\overline{X}_i' e_i = \begin{bmatrix} X_{1i}e_{1i} \\ \vdots \\ X_{mi}e_{mi} \end{bmatrix} \tag{22}$$

is iid and has mean zero, the central limit theorem gives

$$\frac{1}{\sqrt{n}} \sum_{i=1}^{n} \overline{X}_i' e_i \xrightarrow{d} N(0, \boldsymbol{\Omega}), \tag{23}$$

where

$$\boldsymbol{\Omega} = E\left[ \overline{X}_i' e_i e_i' \overline{X}_i \right] = E\left[ \overline{X}_i' \boldsymbol{\Sigma}_i \overline{X}_i \right]. \tag{24}$$

## Asymptotic Distribution

- The rest follows a familiar pattern; denoting $E[\overline{X}_i'\,\overline{X}_i] = Q$, we get for the centered estimator

$$\sqrt{n}(\widehat{\beta} - \beta) \xrightarrow{d} N(0, V_\beta) \qquad (25)$$

  where $V_\beta = Q^{-1}\Omega Q^{-1}$.

- The usefulness (and necessity) of this setup becomes apparent when we consider testing for joint hypotheses of the familiar form $\theta = r(\beta) = r(\beta_1, \ldots, \beta_m)$ with LS estimate $\widehat{\theta} = r(\widehat{\beta}_1, \ldots, \widehat{\beta}_m)$, for which

$$\sqrt{n}(\widehat{\theta} - \theta) \xrightarrow{d} N(0, V_\theta) \qquad (26)$$

  where $V_\theta = R'V_\beta R$ and

$$R = \frac{\partial r(\beta)}{\partial \beta}.$$

  To test for hypotheses across equations thus requires the full system error covariance.

- To estimate the covariance matrix of the estimator, we rely in the general case of conditional heteroscedasticity and non-identical covariates as

$$\widehat{V_{\widehat{\beta}}} = \left(\overline{X}'\overline{X}\right)^{-1}\left(\sum_{i=1}^{n}\overline{X}_i'\widehat{e}_i\widehat{e}_i'\overline{X}_i\right)\left(\overline{X}'\overline{X}\right)^{-1}. \qquad (27)$$

- Under the standard assumptions made for the single-equation case,

$$n\widehat{V_{\widehat{\beta}}}\xrightarrow{p}V_{\beta}. \qquad (28)$$

## Seemingly Unrelated Regression

- A special case of multivariate regression is Seemingly Unrelated Regression (SUR) where "seemingly" unrelated observations are nonetheless related through common shocks.

- Consider the conditionally homoscedastic regression assuming conditional mean independence,

$$\boldsymbol{y}_i = \overline{\boldsymbol{X}_i}\boldsymbol{\beta} + \boldsymbol{e}_i, \quad \text{E}[\boldsymbol{e}_i|\boldsymbol{X}_i] = 0, \quad \text{E}[\boldsymbol{e}_i\boldsymbol{e}_i'|\boldsymbol{X}_i] = \boldsymbol{\Sigma} \tag{29}$$

- The generalized least squares estimator of $\beta$ is

$$\begin{aligned}
\widetilde{\boldsymbol{\beta}} &= \left(\sum_{i=1}^{n}\overline{\boldsymbol{X}_i}'\boldsymbol{\Sigma}^{-1}\overline{\boldsymbol{X}_i}\right)^{-1}\left(\sum_{i=1}^{n}\overline{\boldsymbol{X}_i}'\boldsymbol{\Sigma}^{-1}\boldsymbol{y}_i\right) \\
&= \left(\overline{\boldsymbol{X}}'(\boldsymbol{I}_n \otimes \boldsymbol{\Sigma}^{-1})\overline{\boldsymbol{X}}\right)^{-1}\left(\overline{\boldsymbol{X}}'(\boldsymbol{I}_n \otimes \boldsymbol{\Sigma}^{-1})\boldsymbol{y}\right).
\end{aligned} \tag{30}$$

- Once you replace the unknown $\boldsymbol{\Sigma}$ by its estimator, this is a feasible GLS estimator better known as SUR:

$$\widehat{\boldsymbol{\beta}}_{sur} = \left(\overline{\boldsymbol{X}}'(\boldsymbol{I}_n \otimes \widehat{\boldsymbol{\Sigma}}^{-1})\overline{\boldsymbol{X}}\right)^{-1}\left(\overline{\boldsymbol{X}}'(\boldsymbol{I}_n \otimes \widehat{\boldsymbol{\Sigma}}^{-1})\boldsymbol{y}\right) \tag{31}$$

# An embarrassing example

- In a paper coauthored by yours truly (Jäntti, Pirttilä, and Selin 2015), we estimate among other things labour supply by regressing hours $h$ worked on log net wages $\ln w(1 - \tau)$ and some controls,

$$h = \beta_1 \ln w(1 - \tau) + \cdots + \beta_k + e. \tag{32}$$

- Our interest is not in $\beta_1$ but in the elasticity of hours wrt net wages, which in this semi log specification is

$$\eta = \frac{\beta_1}{\mathrm{E}[h]}.$$

- In the paper, we treat $\widehat{\mathrm{E}[h]}$ as constant, which it is not, since it is an estimator and, moreover, it is an estimator that is correlated with $\widehat{\beta}_1$.

- We should have setup a system as

$$\begin{aligned}
h &= \ln w(1 - \tau)\beta_{11} + \cdots + \beta_{1k} + e_1 \\
h &= \beta_{21} + e_2
\end{aligned} \tag{33}$$

and estimated the elasticity by

$$\widehat{\eta} = \frac{\widehat{\beta}_{11}}{\widehat{\beta}_{21}}$$

and worked out the standard error of this using the delta method (but did not!).

- Consider the CEF (with a single, i.e., scalar, $X$):

$$\mathrm{E}[Y|X] = m(X) = \theta_1 + \theta_2 \exp^{\theta_3 X} \qquad (34)$$

- This is not linear in the parameters and the coefficients can not be estimated by OLS.

- We can formulate this as a non-linear regression:

$$Y = \theta_1 + \theta_2 \exp^{\theta_3 X} + e \qquad (35)$$

- The sum of squared deviations for sample data is

$$S_n(\theta) = \frac{1}{2} \sum_{i=1}^{n} e_i^2 = \frac{1}{2} \sum_{i=1}^{n} [Y_i - m(X_i, \theta)]^2. \qquad (36)$$

## NonLinear Least Squares

- The NLS estimator is the solution to the first-order condition:

$$\frac{\partial S(\theta)}{\partial \theta} = -\sum_{i=1}^{n} [Y_i - m(X_i, \theta)] \frac{\partial m(X_i, \theta)}{\partial \theta} = \mathbf{0}. \tag{37}$$

- For our example equation 34 these are:

$$\frac{\partial S(\theta)}{\partial \theta_1} = -\sum_{i=1}^{n} [Y_i - \theta_1 - \theta_2 e^{\theta_3 X_i}] = 0$$

$$\frac{\partial S(\theta)}{\partial \theta_2} = -\sum_{i=1}^{n} [Y_i - \theta_1 - \theta_2 e^{\theta_3 X_i}] e^{\theta_3 X_i} = 0 \tag{38}$$

$$\frac{\partial S(\theta)}{\partial \theta_3} = -\sum_{i=1}^{n} [Y_i - \theta_1 - \theta_2 e^{\theta_3 X_i}] \theta_2 X_i e^{\theta_3 X_i} = 0$$

These do not have a closed form solution, so $\theta$ needs to be estimated using iterative numerical methods.

- The NLS estimator is an example of an "m-estimator" (see Hansen 2021, ch 22)
- This class includes the LS estimators we have considered hitherto
- For the asymptotic distribution of the NLS estimator (and more generally, m-estimators) requires some additional assumptions (see Hansen 2021, p 777):
  - the parameter set is compact
  - the moment function is finite and continuous in the parameters and bound
  - the criterion function converges uniformly in the parameter set

- For differentiable $m()$, with the vector of partial derivatives

$$\boldsymbol{m}_\theta(\boldsymbol{X}, \theta) = \frac{\partial m(\boldsymbol{X}, \theta)}{\partial \theta} \qquad (39)$$

the non-linear least squares (NLS) estimate is asymptotically normal with

$$\sqrt{n}(\widehat{\theta} - \theta) \xrightarrow{d} N(\boldsymbol{0}, \boldsymbol{V}_\theta) \qquad (40)$$

with $\boldsymbol{V}_\theta = \mathrm{E}[\boldsymbol{m}_\theta \boldsymbol{m}'_\theta]^{-1} \mathrm{E}[\boldsymbol{m}_\theta \boldsymbol{m}'_\theta e^2] \mathrm{E}[\boldsymbol{m}_\theta \boldsymbol{m}'_\theta]^{-1}$ which can be estimated using the NLS residuals and the "plug-in" estimate of $\mathrm{E}[\boldsymbol{m}_\theta \boldsymbol{m}'_\theta]$ at sample data points and the NLS estimate of $\theta$.

## Testing for Omitted NonLinearity

- If the worry is that a given linear regression omits non-linear regressors, this can easily be tested using conventional methods.

- Suppose we first estimate the regression

$$Y = X'\beta + e. \tag{41}$$

$Z = h(X)$ is a set of non-linear functions of $X$. Omitted non-linearity can be tested by estimating

$$Y = X'\widetilde{\beta} + Z'\widetilde{\gamma} + \tilde{e} \tag{42}$$

and testing $\gamma = \mathbf{0}$ using a Wald test.

- A variant is the RESET test, which uses fitted values from the "short" regression $\hat{Y}_i = X_i'\widehat{\beta}$ to form $Z_i' = (\hat{Y}_i^2, \ldots, \hat{Y}_i^m)$,

$$Y_i = X_i'\widetilde{\beta} + Z_i'\widetilde{\gamma} + \tilde{e}_i. \tag{43}$$

Again, the Wald statistic of the hypothesis that $\gamma = \mathbf{0}$ (with a $\chi^2_{m-1}$ distribution) is a test for omitted non-linearity.

Hansen, Bruce E (2021). *Econometrics*. Madison, WI: University of Wisconsin.

Jäntti, Markus, Pirttilä, Jukka, and Selin, Håkan (2015). "Estimating labour supply elasticities based on cross-country micro data: A bridge between micro and macro estimates?" *Journal of Public Economics* 127, pp. 87–99. DOI: doi:10.1016/j.jpubeco.2014.12.006. URL: http://www.sciencedirect.com/science/article/pii/S0047272714002527.