# Applied Panel Data Analysis – Lecture 8

Christopher F. Parmeter

AGRODEP
September 9-13[th], 2013
Dakar, Senegal

U miami

- We will discuss endogeneity in the unobserved effects model
- Traditional endogeneity where some of the covariates are correlated with the idiosyncratic shocks in the model
- Endogeneity induced when some of the covariates are correlated with the unobserved effect while others are not

- In most applied economic settings endogeneity is considered the key hurdle for providing credible results
- Controlling for endogeneity allows one to progress from make statements about correlations to causation
- Endogeneity can arise from omitted variables, measurement error, sample selectivity, or self selection

- Consider our structural unobserved unobserved effects model

$$y_{it} = Y'_{1,it}\gamma + x'_{1,it}\beta + c_i + \varepsilon_{it}$$
$$= Z_{it}\delta + c_i + \varepsilon_{it} \tag{1}$$

  where $Y_{1,it}$ are $g_1$ endogenous variables and $x_{1,it}$ are $k_1$ exogenous variables; $Z = [Y_1, X_1]$

- We also have $k_2 > g_1$ additional instrumental variables, $x_{2,it}$

- Let $x_{it} = [x_{1,it}, x_{2,it}]$ (so that $X = [X_1, X_2]$) denote the collection of all exogenous variables

- Endogeneity enters in the model when
  $E[\varepsilon_{it}|Z_{it}, c_i] \neq E[\varepsilon_{it}] = 0$
- Our instruments will satisfy the standard exogeneity condition
  $E[\varepsilon_{it}|X_{it}, c_i] \neq E[\varepsilon_{it}] = 0$
- Notice that our focus at the moment is with correlation
  between $Y_1$ and $\varepsilon$, the unobserved effect will be controlled
  through either the fixed or random effects framework

- Suppose that we assume the fixed effects framework for our structural unobserved effects model
- Using the within transformation on (1) we have

$$Qy_{it} = QZ_{it}\delta + Q\varepsilon_{it} \tag{2}$$

- Using the instruments $QX_{it}$, two stage least squares estimation produces

$$\hat{\delta}_{W2SLS} = \left(\tilde{Z}'P_{\tilde{X}}\tilde{Z}\right)^{-1}\tilde{Z}'P_{\tilde{X}}\tilde{y} \tag{3}$$

where $P_{\tilde{X}} = \tilde{X}(\tilde{X}'\tilde{X})^{-1}\tilde{X}'$

- This is nothing more than 2SLS except we use as instruments $\tilde{X}$ instead of $X$

- The reason for the transformation is that the correlation between the fixed effects and the regressors needs to be removed prior to controlling for the correlation between $Y_1$ and $X_2$
- The variance-covariance matrix for $\hat{\delta}_{W2SLS}$ is

$$Var(\hat{\delta}_{W2SLS}) = \sigma_\varepsilon^2 \left( \tilde{Z}' P_{\tilde{X}} \tilde{Z} \right)^{-1} \tag{4}$$

- An alternative set of instruments is $PX$, which constitutes between estimation
- Using the between transformation on (1) we have

$$Py_{it} = PZ_{it}\delta + P\varepsilon_{it} \tag{5}$$

and two stage least squares estimation produces

$$\hat{\delta}_{B2SLS} = \left(\bar{Z}'P_{\bar{X}}\bar{Z}\right)^{-1}\bar{Z}'P_{\bar{X}}\bar{y} \tag{6}$$

where $P_{\bar{X}} = \bar{X}(\bar{X}'\bar{X})^{-1}\bar{X}'$

- The variance-covariance matrix for $\hat{\delta}_{B2SLS}$ is

$$Var(\hat{\delta}_{B2SLS}) = \sigma_1^2\left(\bar{Z}'P_{\bar{X}}\bar{Z}\right)^{-1} \tag{7}$$

where $\sigma_1^2 = T\sigma_c^2 + \sigma_\varepsilon^2$

- If we focus on the one-way error component model, in the single equation setting our variance-covariance setting is identical to the fully exogenous setting we discussed in Lecture 4

- Recall that GLS estimation of

$$\begin{pmatrix} Qy \\ Py \end{pmatrix} = \begin{pmatrix} QX \\ PX \end{pmatrix} \beta + \begin{pmatrix} Qu \\ Pu \end{pmatrix}$$

  produced the one-way random effects estimator

- Now consider the following system of $2NT$ observations

$$\begin{pmatrix} X'Qy \\ X'Py \end{pmatrix} = \begin{pmatrix} X'QZ \\ X'PZ \end{pmatrix} \delta + \begin{pmatrix} X'Qu \\ X'Pu \end{pmatrix} \qquad (8)$$

- Given the validity of the instruments we have

$$E \left( \begin{array}{c} X'Qu \\ X'Pu \end{array} \right) = 0$$

and

$$Var \left( \begin{array}{c} X'Qu \\ X'Pu \end{array} \right) = \left[ \begin{array}{cc} \sigma_\varepsilon^2 X'QX & 0 \\ 0 & \sigma_1^2 X'PX \end{array} \right] \quad (9)$$

- Thus, GLS estimation will produce an unbiased and consistent estimator of $\delta$ from (8)

- Baltagi (1981) derived the error component two-stage least squares estimator based on GLS estimation of (8):

$$\hat{\delta}_{EC2SLS} = \left[ \frac{\tilde{Z}'P_{\tilde{X}}\tilde{Z}}{\sigma_\varepsilon^2} + \frac{\bar{Z}'P_{\bar{X}}\bar{Z}}{\sigma_1^2} \right]^{-1} \left[ \frac{\tilde{Z}'P_{\tilde{X}}\tilde{y}}{\sigma_\varepsilon^2} + \frac{\bar{Z}'P_{\bar{X}}\bar{y}}{\sigma_1^2} \right] \tag{10}$$

- As in the fully exogenous case, the EC2SLS estimator can be succinctly written as a matrix weighted average of the B2SLS and FE2SLS estimators

$$\hat{\delta}_{EC2SLS} = W_1\hat{\delta}_{FE2SLS} + W_2\hat{\delta}_{B2SLS} \tag{11}$$

- Baltagi (1981) suggests consistent estimators for $\sigma_\varepsilon^2$ and $\sigma_1^2$ using the residual sum of squares from fixed effects two-stage least squares estimation and between two-stage least squares

$$\widehat{\sigma}_\varepsilon^2 = \frac{\hat{\varepsilon}'_{FE2SLS} Q \hat{\varepsilon}_{FE2SLS}}{N(T-1)} \tag{12}$$

$$\widehat{\sigma}_1^2 = \frac{\hat{\varepsilon}'_{B2SLS} P \hat{\varepsilon}_{B2SLS}}{N} \tag{13}$$

where $\hat{\varepsilon}_{FE2SLS} = y - Z\hat{\delta}_{FE2SLS}$ and $\hat{\varepsilon}_{B2SLS} = y - Z\hat{\delta}_{B2SLS}$

- An alternative two-stage least squares estimator for the one-way error component model is from Balestra and Varadharajan-Krisnakumar (1987)
- They suggest direct GLS estimation of (1) using $\Omega^{-1/2}$ with instruments $\Omega^{-1/2}X = \frac{\tilde{X}}{\sigma_\varepsilon} + \frac{\bar{X}}{\sigma_1}$
- Their generalized two-stage least squares estimator is

$$\hat{\delta}_{G2SLS} = \left(Z^{*\prime}P_{X^*}Z^*\right)^{-1}Z^{*\prime}P_{X^*}y^* \qquad (14)$$

- Note that Baltagi's (1981) EC2SLS estimator uses as instruments $[\tilde{X}, \bar{X}]$ while Balestra and Varadharajan-Krisnakumar's (1987) G2SLS estimator uses as instruments $\frac{\tilde{X}}{\sigma_\varepsilon} + \frac{\bar{X}}{\sigma_1}$

- How do these instrument sets differ?

- $[\tilde{X}, \bar{X}]$ spans a linear space of dimension $2(k_1 + k_2)$ while $\frac{\tilde{X}}{\sigma_\varepsilon} + \frac{\bar{X}}{\sigma_1}$ spans a linear space of dimension $k_1 + k_2$, i.e. the instrument set of Balestra and Varadharajan-Krisnakumar's (1987) is a subset of that of Baltagi (1981)

- Baltagi and Li (1992) show that in the single equation setting, these extra instruments do not yield reductions in the variance covariance matrix of $\hat{\delta}_{EC2SLS}$

- Moreover, $\hat{\delta}_{EC2SLS}$ and $\hat{\delta}_{G2SLS}$ have the same asymptotic variance-covariance matrix

- However, the use of $\hat{\delta}_{EC2SLS}$ is still common because while the variance-covariance matrix is asymptotically the same as $\hat{\delta}_{G2SLS}$ in the single equation setting, in the full system setting, Baltagi's approach yields gains in efficiency
- While not common, one should check estimates across the two instrument sets to see if there are any perceptible differences (there should not be except in perverse settings)

- Recall that the distinction between the random effects framework and the fixed effects framework was an all or nothing proposition
- Either all of the regressors were independent from the unobserved effect (random effects framework) or all regressors were allowed to be correlated with the unobserved effect (fixed effects framework)
- There was no middle ground for estimation between these two frameworks
- Hausman and Taylor (1981) proposed an estimation strategy that accomplished just this

- To begin, endogeneity will now exist in the unobserved effects model through correlation amongst a subset of the regressors and the unobserved effects

- An example is a wage regression where work experience and years of education are correlated with ability (lets assume it is time constant), which is part of the unobserved effect

- If we assume the fixed effects framework (which is feasible), but only these two variables are correlated with ability, then our structure is too strict

- The unobserved effects model of Hausman and Taylor (1981) is

$$y_i t = x'_{it}\beta + z'_i\gamma + c_i + \varepsilon_{it} \tag{15}$$

- We further partition $x_{it}$ and $z_i$ into two pieces: $x_{it} = [x_{1,it}, x_{2,it}]$ and $z_i = [z_{1,i}, z_{2,i}]$
- $x_{1,it}$ is $k_1 \times 1$, $x_{2,it}$ is $k_2 \times 1$, $z_{1,i}$ is $g_1 \times 1$ and $z_{2,i}$ is $g_2 \times 1$
- We assume that $x_1$ and $z_1$ are exogenous with respect to both $c$ and $\varepsilon$ while $x_2$ and $z_2$ are exogenous with respect to $\varepsilon$ but are endogenous with respect to $c$

- Notice that the within transformation would eliminate the endogeneity of $x_2$, but it also removes $z_1$ and $z_2$ from the model
- Hausman and Taylor's (1981) approach is to control for endogeneity without eliminating the time constant covariates from the model
- How do they do this?

- They suggest the standard random effects framework transformation $\Omega^{-1/2}$ to (15) and then application of two-stage least squares using as instruments

$$A = \left[ \tilde{X}, PX_1, Z_1 \right] \qquad (16)$$

- Note that $Z_1$ instruments itself (since it is exogenous), while $X_1$ and $X_2$ are instrumented by $\tilde{X}$

- $Z_2$ is instrumented by $PX_1$; given the panel structure $X_1$ can be used in two different dimensions as an instrument

- As it stands the Hausman and Taylor (1981) procedure is infeasible given that the elements of $\Omega$ are unknown
- To construct a feasible estimator Hausman and Taylor propose the following approach
- First, estimate the model in (15) using the within transformation; this will naturally eliminate $Z$ from the model so $\gamma$ is not identified

- Second, average the residuals from within estimation of (15) across time

$$\hat{u}_{i\cdot} = \bar{y}_{i\cdot} - \bar{X}_{i\cdot}\tilde{\beta} \qquad (17)$$

- Third, perform two-stage least squares using instrument matrix $A = [X_1, Z_1]$ on the model

$$\hat{u}_{i\cdot} = Z_i\gamma + \omega_i \qquad (18)$$

- The estimator from this regression is

$$\hat{\gamma}_{2SLS} = (Z'P_A Z)^{-1} Z' P_A \hat{u} \qquad (19)$$

- These steps provide consistent estimates of $\beta$ and $\gamma$, which can be used to construct consistent estimates of $\sigma_c^2$ and $\sigma_\varepsilon^2$
- Consistent estimators of the variance components are

$$\hat{\sigma}_\varepsilon^2 = \frac{y'Q\left(I - P_{QX}\right)Qy}{N(T-1)} \tag{20}$$

and

$$\hat{\sigma}_1^2 = \frac{\left(y - X\tilde{\beta} - Z\hat{\gamma}_{2SLS}\right)' P \left(y - X\tilde{\beta} - Z\hat{\gamma}_{2SLS}\right)}{N} \tag{21}$$

- Using the variance component estimates the original unobserved effects model in (15) is transformed with $\hat{\Omega}^{-1/2}$ and two-stage least squares is performed using instrument matrix $A$ from (16)
- If $k_1 < g_2$ then the model is under-identified, $\hat{\beta}_{HT} = \tilde{\beta}$ and $\hat{\gamma}_{HT}$ does not exist
- If $k_1 = g_2$ then the model is exactly-identified, $\hat{\beta}_{HT} = \tilde{\beta}$ and $\hat{\gamma}_{HT} = \hat{\gamma}_{2SLS}$
- If $k_1 > g_2$ then the model is over-identified, and $\hat{\beta}_{HT}$ is more efficient than $\tilde{\beta}$

- An over-identification test follows along the lines of the Hausman test of the random effects framework

$$\hat{m} = \left(\hat{\beta}_{HT} - \tilde{\beta}\right)' \left(var(\tilde{\beta}) - var(\hat{\beta}_{HT})\right)^{-1} \left(\hat{\beta}_{HT} - \tilde{\beta}\right) \quad (22)$$

- This statistic has limiting distribution $\chi^2_\ell$ where $\ell = \min[k_1 - g_2, NT - (k_1 + k_2)]$
- This test allows one to discern if endogeneity in $X_2$ is severe
- Note that one only tests using $\beta$ as the within transformation cannot identify $\gamma$

- Both Amemiya and MaCurdy (1986) and Breusch, Mizon and Schmidt (1989) proposed instrument sets that produce an estimator more efficient than the original Hausman and Taylor (1981) estimator

- These instrument sets are more likely to ensure identification of $\gamma$, however, they come at the expense of rapidly increasing the instrument set based on the time dimension of the panel

- Too many instruments can also be viewed negatively, even though there are efficiency gains to be had

- Further, there are additional exogeneity conditions that must be satisfied with these expanded instrument sets; whereas Hausman and Taylor (1981) only require that the time averaged $X_1$s are uncorrelated with $c$, the Amemiya and MaCurdy (1986) instrument set requires conditional strict exogeneity, a much stronger condition

- Discussed accounting for endogeneity in the unobserved effects model
- Estimation covered both the fixed and random effects framework
- Beyond endogeneity with the idiosyncratic error term, also discussed compromise between fixed and random effects framework that can allow for time constant variables
- Hausman-Taylor estimator allows endogeneity between covariates and unobserved effect; can identify time constant effects