# Model building issues in multilevel linear models with endogenous covariates

Leonardo Grilli and Carla Rampichini

Dipartimento di Statistica "Giuseppe Parenti"
Università di Firenze
grilli@ds.unifi.it, carla@ds.unifi.it

**knem⬡**

KNEMO'06 Conference Anacapri, 3-6 September 2006

---

## Motivation

- Endogeneity is a common problem in applied works

- We explore level 2 endogeneity in linear random effects models, i.e. random effects correlated with covariates

- Measurement error stemming from the use of cluster means is overlooked in the literature

KNEMO'06                                                                 2

---

## Outline

- Endogeneity in multilevel models
- The linear random intercept model
- Between-cluster slope and measurement error: correction for consistent estimates
- Simulations
- Conclusions

KNEMO'06                                                                 3

---

## The linear random intercept model

$$Y_{ij} = \alpha + \beta X_{ij} + v_j + e_{ij}$$

- $i$=1,2,…,$n_j$ elementary (level 1) index
- $j$=1,2,…,$J$ cluster (level 2) index $Y_{ij}$ response
- $X_{ij}$ level 1 covariate
- $v_j$ level 2 errors, or random effects
- $e_{ij}$ level 1 errors

- Examples:
  - Panel data (typically: $n_j$ small, $J$ large)
  - Clustered cross-section data (typically: $n_j$ large, $J$ small)

KNEMO'06                                                                 4

---

## Level 2 endogeneity

- Level 2 endogeneity arises when

$$E(v_j \mid X_{ij}) \neq 0$$

➔ standard estimators are inconsistent for $\beta$

Note that   $Cov(v_j, X_{ij}) \neq 0 \Rightarrow E(v_j \mid X_{ij}) \neq 0$

KNEMO'06                                                                 5

---

## Sources of endogeneity

- Omission of relevant regressors at any level

- Measurement error in the covariates

- Self-selection

- Simultaneity

KNEMO'06                                                                 6

---

## The data generating model

- $X_{ij}$ must be treated as random
- The hierarchical framework requires to specify how varies between and within clusters, assume a variance component model

$$X_{ij} = X_j^B + X_{ij}^W$$

- Under the assumptions

X1    $X_j^B$ are iid with mean $\mu_X$ and variance $\tau_X^2 > 0$

X2    $X_{ij}^W$ are iid with zero mean and variance $\sigma_X^2 > 0$

X3    $X_j^B \perp\!\!\!\perp X_{ij}^W$ , $\forall i, j$

## The Overall model

$$Y_{ij} = \alpha + \beta^W X_{ij}^W + \beta^B X_j^B + u_j + e_{ij}$$

- $\beta^N$ within effect, $\beta^B$ between effect
- in general, $\beta^N \neq \beta^B$
- Assume:
  - Independent clusters
  - Two-stage sampling
  - Unbalanced design

## Assumptions on the errors

In the Overall data generating model:

- $X^W$ and $X^B$ are exogenous

- Errors at different levels are independent

- At each level, errors are i.i.d.

## Example: university effectiveness

- $Y_{ij}$ observed income of the *i*-th student of the *j*-th school one year after graduation

- $X_{ij}$ observed grade of such a student

- $X_{ij} = X_j^B + X_{ij}^W$
  - $X_j^B$      school mean grade
  - $X_{ij}^W = X_{ij} - X_j^B$    student deviation from school mean

- The model decompose the total effect of *X* on *Y*:
  - $\beta^B$      school effect on the income
  - $\beta^N$      student effect on the income

- $u_j$ school level residual (effectiveness)
- $e_{ji}$ student level residual

## An example …cont'd

Between and Within effects are conceptually different and sometimes have opposite signs!

In this example (University of Florence data):
- $\beta^N$ >0 ➔ *within a school*, students with higher grade have higher income
- $\beta^B$ <0 ➔ schools giving higher grades show lower average income (e.g. Humanities)

In many settings we expect $\beta^N \neq \beta^B$

## Nature of level 2 endogeneity

- Overall model alternative parametrization

$$Y_{ij} = \alpha + \beta^W X_{ij} + \delta X_j^B + u_j + e_{ij}$$
$$\delta = \beta^B - \beta^W$$

- If $X^B$ is omitted ➔ it is included in level 2 error

$$Y_{ij} = \eta + \beta^W X_{ij} + v_j + e_{ij}$$
$$v_j = \delta(X_j^B - \mu_X) + u_j$$
$$E(v_j) = 0 \quad Var(v_j) = \tau_{Y|X}^2 = \delta^2 \tau_X^2 + \tau_{Y|X^B X^W}^2$$

## Omission of $X^B$

Note that: $Cov(v_j, X_{ij}) = Cov(v_j, X_j^B) = \delta \tau_X^2$

Between variance of $X$ (assumed >0)

→ X exogenous iff $\delta=0$

level 2 endogeneity can be interpreted as:

- a wrong constraint on the slopes, i.e. $\beta^N = \beta^B = \beta$

$Y_{ij} = \alpha + \beta^W X_{ij}^W + \beta^B X_j^B + u_j + e_{ij}$ → $Y_{ij} = \alpha + \beta X_{ij} + v_j + e_{ij}$

- the omission of a relevant regressor, i.e. $X^B$

$Y_{ij} = \alpha + \beta^W X_{ij} + \delta X_j^B + u_j + e_{ij}$ → $Y_{ij} = \alpha + \beta X_{ij} + v_j + e_{ij}$

KNEMO'06                                                                 13

---

## Correct for level 2 endogeneity

Fit the Overall model

$$Y_{ij} = \alpha + \beta^W X_{ij}^W + \beta^B X_j^B + u_j + e_{ij}$$

or

$$Y_{ij} = \alpha + \beta^W X_{ij} + \delta X_j^B + u_j + e_{ij}$$

BUT $X^B$ and $X^W$ are UNOBSERVABLE!

KNEMO'06                                                                 14

---

## What is observed?

Replace

- $X^B$ with the cluster means $\overline{X}_j = \frac{1}{n_j}\sum_{i=1}^{n_j} X_{ij}$
- $X^W$ with the deviations $\widetilde{X}_{ij} = X_{ij} - \overline{X}_j$

Observable split $X_{ij} = \overline{X}_j + \widetilde{X}_{ij}$

Note that $X^B$ and $X^W$ are measured with error:

$$\overline{X}_j = X_j^B + \overline{X}_j^W \qquad \widetilde{X}_{ij} = X_{ij}^W - \overline{X}_j^W$$

KNEMO'06                                                                 15

---

## Working overall model

$$Y_{ij} = \alpha + \beta^W \widetilde{X}_{ij} + \beta^B \overline{X}_j + z_j + e_{ij}$$

$$z_j = u_j - \delta \overline{X}_j^W$$  measurement error cancel out iff $\delta=0$

Note that $\tilde{X}_{ij}$ :

- is exogenous, i.e. $E(z_j \mid \widetilde{X}_{ij}) = -\delta E(\overline{X}_j^W \mid \widetilde{X}_{ij}^W) = 0$
- and uncorrelated with cluster mean: $Cov(\overline{X}_j, \tilde{X}_{ij}) = 0$

$\beta^W$ consistently estimated

KNEMO'06                                                                 16

---

## Estimable between slope

- The cluster mean is endogenous!

$$Cov(z_j, \overline{X}_j) = -\delta \sigma_X^2 / n_j$$

- We consistently estimate $\beta_{cm}^B \neq \beta^B$
- In the balanced case:

$$\beta_{cm}^B = \lambda_X \beta^B + (1-\lambda_X)\beta^W ,$$

$\lambda_X \in (0,1)$ reliability of X

- The model is correctly specified, but the measurement error causes endogeneity!!!

In the balanced case:
recover $\beta^B$ using the estimated $\beta^B_{cm}$ and $\lambda_X$

KNEMO'06                                                                 17

---

## Reliability

$$\lambda_X = \frac{Var(X_j^B)}{Var(\overline{X}_j)} = \frac{\tau_X^2}{\tau_X^2 + \sigma_X^2/n} = \left(1 + \frac{1}{(\tau_X^2/\sigma_X^2)n}\right)^{-1} , \ \lambda_X \in (0,1)$$

$\lambda_X$ is an increasing function of :

- the cluster size $n$ (sample design)
- the variance ratio $\tau_X^2 / \sigma_X^2$ (model parameters)

It is usual to find values far from 1, e.g.

$\lambda_X = 0.67$ if $\begin{cases} n=2 & \tau_X^2 = \sigma_X^2 & \text{(panel)} \\ n=20 & \tau_X^2 = 0.1\sigma_X^2 & \text{(cross-section)} \end{cases}$

Estimate $\tau_X^2$ and $\sigma_X^2$ by standard ANOVA methods

KNEMO'06                                                                 18

## Variance correction

Estimated residual level 2 variance is inflated!

$$Var(z_j) = \tau^2_{Y|X^B X^W} + \lambda_X \delta^2 \frac{\sigma^2_X}{n}$$

*true* level 2 residual variance

Recover the 'true' variance subtracting this factor from the estimated variance

measurement error bias may be more serious on $\tau^2$ than on $\delta$!!!

---

## In summary …

| Working model | $Y_{ij} = \eta + \beta^W X_{ij} + v_j + e_{ij}$ | $Y_{ij} = \alpha + \beta^W \widetilde{X}_{ij} + \beta^B \overline{X}_j + z_j + e_{ij}$ |
|---|---|---|
| Omission of a regressor | yes | no |
| Measurement error | no | yes |
| Level 2 error cov | $Cov(v_j, X_{ij}) = \delta\tau^2_x$ | $Cov(z_j, \tilde{X}_{ij}) = 0$ <br> $Cov(z_j, \overline{X}_j) = -\delta\sigma^2_X/n_j$ |
| Consistent $\beta^W$ | No for n small | yes |
| Consistent $\beta^B$ | no | yes (correction required) |

"true" model

$$Y_{ij} = \begin{cases} \alpha + \beta^W X_{ij}^W + \beta^B X_j^B + u_j + e_{ij} \\ \alpha + \beta^W X_{ij} + \delta X_j^B + u_j + e_{ij} \end{cases}$$

---

## Simulation: MC means on 1000 replications

| $\delta=\beta^B-\beta^W$ | $Y_{ij} = \alpha + \beta X_{ij} + \dots$ | | $Y_{ij} = \alpha + \beta^W X_{ij} + \delta\overline{X}_j + \dots$ | | |
|---|---|---|---|---|---|
| | $\beta$ | $\tau^2_Y$ | $\beta^W$ | $\delta$ | $\tau^2_Y$ |
| -2 | 0.61 | 3.63 | 1.01 | -1.34 | 2.35 |
| -1.5 | 0.62 | 2.27 | 1.00 | -1.00 | 1.75 |
| -1 | 0.71 | 1.52 | 1.01 | -0.68 | 1.34 |
| -0.5 | 0.84 | 1.12 | 1.00 | -0.33 | 1.09 |
| 0 | 1.00 | 1.00 | 1.00 | 0.00 | 1.00 |
| 0.5 | 1.16 | 1.12 | 1.00 | 0.33 | 1.09 |
| 1 | 1.29 | 1.51 | 0.99 | 0.67 | 1.34 |
| 1.5 | 1.38 | 2.28 | 1.00 | 1.01 | 1.75 |
| 2 | 1.40 | 3.59 | 1.00 | 1.34 | 2.36 |

+ endogeneity ↑ / + endogeneity ↓

$\lambda_X$=2/3, $\beta^W$=1, $\tau^2_Y$=1
$n$=2, $J$=100

---

## What happens when cluster size n→∞

| $\delta=\beta^B-\beta^W$ | $Y_{ij} = \alpha + \beta X_{ij} + \dots$ | | $Y_{ij} = \alpha + \beta^W X_{ij} + \delta\overline{X}_j + \dots$ | | |
|---|---|---|---|---|---|
| | $\beta$ | $\tau^2_Y$ | $\beta^W$ | $\delta$ | $\tau^2_Y$ |
| -2 | 1 | 5 | 1 | -2 | 1 |
| -1.5 | 1 | 3.25 | 1 | -1.5 | 1 |
| -1 | 1 | 2 | 1 | -1 | 1 |
| -0.5 | 1 | 1.25 | 1 | -0.5 | 1 |
| 0 | 1 | 1 | 1 | 0 | 1 |
| 0.5 | 1 | 1.25 | 1 | 0.5 | 1 |
| 1 | 1 | 2 | 1 | 1 | 1 |
| 1.5 | 1 | 3.25 | 1 | 1.5 | 1 |
| 2 | 1 | 5 | 1 | 2 | 1 |

+ endogeneity ↑ / + endogeneity ↓

$\lambda_X$=1, $\beta^W$=1, $\tau^2_Y$=1
$n\approx$50, $J$=100

---

## Alternative RE estimators

- Other methods for consistent estimation of $\beta^W$ exist (e.g CIGLS and IV)

$$Y_{ij} = \alpha + \beta^W X_{ij} + v_j + e_{ij}$$

→ level 2 variance is net (adjusted for *X*)

- If the interest is also on $\beta^B$, the overall model must be used

$$Y_{ij} = \alpha + \beta^W X_{ij} + \delta\overline{X}_j + z_j + e_{ij}$$

▶ level 2 variance is net (adjusted for *X*)
▶ BUT attention to measurement error!

---

## Concluding remarks

- In general RE models are better than FE also in presence of level 2 endogeneity

- More simulations are needed to study efficiency issues

- Many extensions need further investigation:
  - model with two or more covariates (problems of model selection)
  - model with random slopes
  - non-linear models.

## References

Ebbes, P., Bockenholt, U. and Wedel, M. (2004). Regressor and random-effects dependencies in multilevel models, *Statistica Neerlandica*, 58, 161–178.

Fielding, A. (2004). The Role of the Hausman Test and whether Higher Level Effects should be treated as Random or Fixed. *Multilevel Modelling Newsletter*, 16(2), 3–9.

Kim, J. S. and Frees, E. W. (2006). Omitted variables in multilevel models. *Psychometrika*, in press.

Rettore, E. and Martini, A. (2001). Constructing league tables of service providers when the performance of the provider is correlated to the characteristics of the clients. In *Processi e metodi statistici di valutazione*, Proceedings of the Conference of the Italian Statistical Society, Roma.

Rice, N., Jones, A. M. and Goldstein H. (2002). Multilevel models where the random effects are correlated with the fixed predictors. Manuscript available at:
http://www.mlwin.com/hgpersonal/ciqls.pdf.

Snijders, T. A. B. and Berkhof, J. (2006). Diagnostic checks for multilevel models. In Jan de Leeuw (Ed.), *Handbook of Multilevel Analysis*. New York: Springer, to appear.

KNEMO'06          25

---

# Thanks for your attention!



---

## Estimable slope

In model :  $Y_{ij} = \alpha + \beta X_{ij} + v_j + e_{ij}$

the OLS estimable slope is actually

$$\beta = \beta^W + \rho_X \delta = \rho_X \beta^B + (1 - \rho_X)\beta^W$$

$$\rho_X = \tau_X^2 / \left(\sigma_X^2 + \tau_X^2\right)$$

i.e. a mixture of the between and within slope, depending on the value of the ICC of X

KNEMO'06          27

---

## Endogeneity test

- The endogeneity Hausman test is equivalent to the Wald test $H_0$: $\delta$=0 in the Working Overall model $Y_{ij} = \alpha + \beta^W X_{ij} + \delta \overline{X}_j + z_j + e_{ij}$

- The estimable parameter is $\delta_{cm} = \lambda_x \, \delta$
- The test $H_0$: $\delta_{cm}$=0 has a lower power

  do the test using $\delta = \delta_{cm}/\lambda_x$

KNEMO'06          28

---

## FE versus RE

- Standard solution for endogeneity in panel literature is the FE estimator (consistent for $\beta^W$)

but    $\tilde{Y}_{ij} = \beta^W \tilde{X}_{ij} + \varepsilon_{ij}$

  - It doesn't allow any cluster level covariate
  - It can be quite inefficient because the number of parameters grow with the number of clusters

- If interest is only on $\beta^W$, a better solution is the Within RE model

  $Y_{ij} = \alpha + \beta^W \tilde{X}_{ij} + s_j + e_{ij}$

  → level 2 variance is gross (unadjusted for *X*)

KNEMO'06          29