

# Linear Panels and Random Coefficients

Manuel Arellano

CEMFI

September 2013

## Introduction

- Panel data models with fixed effects play an important role in applied econometrics.
- In the linear case several estimation methods are available (within groups, IV & GMM, likelihood methods...).
- Applications of these methods are widespread.
- The purpose of these lectures is to provide an overview of the literature on nonlinear panel data methods, including some emphasis on bias-reduction approaches.
- To set the stage, we begin with a review of some basic concepts of linear panels and random coefficients models.
- The focus is on microeconometrics: individuals, households, and firms, but also cross-country growth and development studies.
- Business cycle and financial volatility studies that relate to time series panels and factor models are out of scope here.

## Linear panels and random coefficients

- Basic motivation in microeconometrics: Identifying models that cannot be identified on single outcome data. Two leading situations:
  - Fixed effects endogeneity (e.g. productivity analysis, price effects in demand models, wage effects in labor supply).
  - Error components, variance decomposition (e.g. inequality, mobility studies, quality-adjusted price indices).

## GMM perspective

- The generalized method of moments has proved very useful for linear panel models as an organizing principle.

General idea:

- Start from a set of moment conditions suggested by the model.
- Use sample counterpart to get estimates of common parameters.
- Invoke a central limit theorem to approximate the distribution of standardized estimates by a normal distribution.
- If more moments than parameters are available, form linear combinations.

## Leading example: within-groups

$$y_{it} = x'_{it}\theta_0 + \alpha_i + v_{it} \quad E(v_{it} \mid x_{i1}, \dots, x_{iT}, \alpha_i) = 0.$$

- In this model  $x_{it}$  may be correlated with  $\alpha_i$  but not with  $v_{is}$  for all  $t, s$ . We say that  $x_{it}$  is endogenous wrt the fixed effect but strictly exogenous wrt the time-varying error.
- Letting  $\tilde{x}_{it} = x_{it} - \bar{x}_i$ , the WG model implies the moment conditions

$$E \left[ \sum_{t=1}^T \tilde{x}_{it} (\tilde{y}_{it} - \tilde{x}'_{it}\theta_0) \right] = 0.$$

- The WG estimator  $\hat{\theta}_{WG}$  solves the sample moments

$$\sum_{i=1}^N \sum_{t=1}^T \tilde{x}_{it} (\tilde{y}_{it} - \tilde{x}'_{it}\hat{\theta}_{WG}) = 0.$$

### Leading example: within-groups (continued)

- Inference can be based on the large  $N$ , fixed  $T$  approximation:

$$\widehat{V}^{-1/2} \left( \widehat{\theta}_{WG} - \theta_0 \right) \approx \mathcal{N}(0, I)$$

where

$$\widehat{V} = H^{-1} \left( \sum_{i=1}^N \sum_{t=1}^T \sum_{s=1}^T \widehat{v}_{it} \widehat{v}_{is} \widetilde{x}_{it} \widetilde{x}'_{is} \right) H^{-1},$$

$$\widehat{v}_{it} = \widetilde{y}_{it} - \widetilde{x}'_{it} \widehat{\theta}_{WG}, \text{ and } H = \sum_{i=1}^N \sum_{t=1}^T \widetilde{x}_{it} \widetilde{x}'_{it}.$$

- The resulting "cluster-robust" standard errors are robust to heteroskedasticity and serial correlation but rely on cross-sectional independence.

## Cluster-robust bootstrap standard errors

- A bootstrap approach is as follows. Let  $W_i = (y_{i1}, x'_{i1}, \dots, y_{iT}, x'_{iT})'$  and regard  $W_1, \dots, W_N$  as a multivariate random sample of size  $N$  according to some cdf  $F$ .
- The WG estimator is a function of the data  $\hat{\theta}_{WG} = h(W_1, \dots, W_N)$  whose distribution we want to estimate

$$\Pr(\hat{\theta}_{WG} \leq r) = \Pr_F[h(W_1, \dots, W_N) \leq r].$$

- A simple candidate is the plug-in estimator. It replaces  $F$  by the empirical cdf  $\hat{F}_N$ :

$$\hat{F}_N(s) = \frac{1}{N} \sum_{i=1}^N \mathbf{1}(W_i \leq s),$$

which assigns probability  $1/N$  to each of the observed values  $w_1, \dots, w_N$  of  $W_1, \dots, W_N$

- Letting  $W_1^*, \dots, W_N^*$  denote a random sample from  $\hat{F}_N$ , the resulting estimator is then

$$\Pr_{\hat{F}_N}[h(W_1^*, \dots, W_N^*) \leq r], \quad (1)$$

which is conceptually simple but prohibitive to calculate.

- The bootstrap method evaluates (1) by simulation.  $M$  of samples  $W_1^*, \dots, W_N^*$  (the bootstrap samples) are drawn from  $\hat{F}_N$ , and the frequency with which

$$h(W_1^*, \dots, W_N^*) \leq r$$

provides the desired approximation to the estimator (1).

### *Cluster-robust bootstrap standard errors (continued)*

- As a result of resampling we have available  $M$  estimates from the artificial samples:  
 $\hat{\theta}_{WG}^{(1)}, \dots, \hat{\theta}_{WG}^{(M)}$ .

- A bootstrap standard error is then obtained as

$$\left[ \frac{1}{M-1} \sum_{m=1}^M \left( \hat{\theta}_{WG}^{(m)} - \overline{\hat{\theta}_{WG}} \right)^2 \right]^{1/2}$$

where  $\overline{\hat{\theta}_{WG}} = \sum_{m=1}^M \hat{\theta}_{WG}^{(m)} / M$ .

- The bootstrap method is very flexible and applicable to many different situations such as the bias and variance of an estimator, the calculation of confidence intervals, etc.
- Under general regularity conditions, using the bootstrap standard error to construct test statistics has the same asymptotic justification as conventional asymptotic procedures.
- Sometimes a data producer will provide users with *replicate weights*, which enable the estimation of the sampling distribution of estimators from complex sample designs without disclosing confidential information.



## Generalizations

*Improved GMM under heteroskedasticity and autocorrelation of unknown form*

- Improved GMM based on the larger set of moments  $E [x_i (\tilde{y}_{it} - \tilde{x}'_{it}\theta_0)] = 0$ , ( $t = 1, \dots, T$ ) or

$$E [x_i (\Delta y_{it} - \Delta x'_{it}\theta_0)] = 0, (t = 2, \dots, T)$$

where  $x_i$  stacks  $x_{i1}, \dots, x_{iT}$ .

*Instrumental variable fixed effects models*

- IV versions where the starting assumption is

$$E (v_{it} \mid z_{i1}, \dots, z_{iT}, \alpha_i) = 0$$

for some strictly exogenous instrument  $z$  (e.g. tax component of price variation).

- The moments become

$$E [z_i (\tilde{y}_{it} - \tilde{x}'_{it}\theta_0)] = 0.$$

- In this case  $x$  is treated as a strictly endogenous variable.

## Generalizations (continued)

### *Testing for correlated effects*

- If  $x$  is uncorrelated with  $\alpha$ , valid moments are  $E [x_i (y_{it} - x'_{it}\theta_0)] = 0$ , ( $t = 1, \dots, T$ ), which include  $E [x_i (\Delta y_{it} - \Delta x'_{it}\theta_0)] = 0$ , ( $t = 2, \dots, T$ ) as a subset.
- Thus, an incremental Sargan test can be used for testing the null of fixed-effects exogeneity (Hausman type testing).

### *Models with both time-invariant and time-varying variables*

- A model with a FE-exogenous time-invariant regressor  $w$  satisfies the moments:

$$\begin{aligned} E [x_i (\tilde{y}_{it} - \tilde{x}'_{it}\theta_0)] &= 0 \\ E [w_i (\bar{y}_i - \bar{x}'_i\theta_0 - w_i\delta_0)] &= 0. \end{aligned}$$

- In an IV version the second moment would specify the orthogonality between the average error and an external time-invariant instrument.

## Error in variables

- In a measurement error version of the WG model where  $x$  is measured with an iid error, valid moments are

$$E \left[ \left( x_{i1}, \dots, x_{i(t-2)}, x_{i(t+1)}, \dots, x_{iT} \right) (\Delta y_{it} - \Delta x'_{it} \theta_0) \right] = 0 \quad (t = 2, \dots, T).$$

- Instruments are relevant as long as there is persistence in latent  $x$ 's.
- If ignored first differencing may exacerbate measurement error bias as illustrated next.
- In a linear regression  $y = \beta x^* + u$  with classical measurement error  $x = x^* + \varepsilon$  where  $u, x^*, \varepsilon$  are mutually independent, the OLS parameter satisfies

$$\frac{\text{Cov}(y, x)}{\text{Var}(x)} = \frac{\text{Cov}(y, x^*)}{\text{Var}(x^*) + \text{Var}(\varepsilon)} = \frac{\beta}{1 + \lambda}$$

where  $\lambda = \text{Var}(\varepsilon) / \text{Var}(x^*)$ .

- Similarly, letting  $\lambda_{\Delta} = \text{Var}(\Delta\varepsilon) / \text{Var}(\Delta x^*)$ , the OLS parameter of the regression in differences satisfies

$$\frac{\text{Cov}(\Delta y, \Delta x)}{\text{Var}(\Delta x)} = \frac{\beta}{1 + \lambda_{\Delta}}.$$

- If  $\text{Cov}(\varepsilon_t, \varepsilon_{t-1}) = 0$  but  $\text{Cov}(x_t^*, x_{t-1}^*) > 0$  then  $\lambda_{\Delta} > \lambda$ . Under these conditions, which are relevant in applications, differencing magnifies measurement error bias.

*Illustration: measuring economies of scale in firm money demand*

- Bover and Watson (2005) estimate firm-level money demand equations of the form

$$\log m_{it} = c(t) \log s_{it} + b(t) + \eta_i + v_{it}.$$

where  $m$  is demand for cash and  $s$  denotes output (or sales).

- The economies of scale coefficient  $c(t)$  is specified as a polynomial in  $t$  to allow for changes over the sample period.
- The year dummies  $b(t)$  capture changes in relative interest rates together with other aggregate effects.
- The individual effect is meant to represent permanent differences across firms in the production of transaction services (so that  $\eta$  varies inversely with the firm's financial sophistication), and  $v$  contains measurement errors in cash holdings and sales.
- We would expect  $\text{Cov}(\log s, \eta) \leq 0$  and a downward unobserved heterogeneity bias in economies of scale.
- We also expect measurement error to account for a larger share of variation in sales growth than in the level of sales.

Firm money demand estimates  
Sample period 1986–1996

	OLS Levels	OLS WG	OLS 1st-diff.	GMM 1st-diff.	GMM 1st-diff. m. error	GMM Levels m. error
Log sales	.72 (30.)	.56 (16.)	.45 (12.)	.49 (16.)	.99 (7.5)	.75 (35.)
Log sales ×trend	−.02 (3.2)	−.03 (9.7)	−.03 (4.9)	−.03 (5.3)	−.03 (5.0)	−.03 (4.0)
Log sales ×trend <sup>2</sup>	.001 (1.2)	.002 (6.6)	.001 (1.9)	.001 (2.0)	.001 (2.3)	.001 (1.4)
Sargan ( <i>p</i> -value)				.12	.39	.00

All estimates include year dummies, and those in levels also include industry dummies. *t*-ratios in brackets robust to heteroskedasticity & serial correlation. *N*=5649. Source: Bover and Watson (2005).

All estimates in the table are obtained from an unbalanced panel of 5649 Spanish firms with at least four consecutive annual observations during the period 1986–1996.

- The comparison between OLS-levels and WG (cols 1 & 2) is consistent with a positive fixed-effects bias (counter to expectation), but the smaller OLS-diff sales effect (col 3) suggests that measurement error bias may be important.
- Col 4 shows GMM estimates based on the moments  $E(\log s_{it} \Delta v_{is}) = 0$  for all  $t, s$ . Absent measurement error, we would expect them to be similar to WG and OLS-diff.
- Col 5 shows GMM estimates based on

$$E(\log s_{it} \Delta v_{is}) = 0 \quad (t = 1, \dots, s - 2, s + 1, \dots, T; s = 1, \dots, T),$$

thus allowing for both correlated firm effects and measurement error in sales.

- Interestingly, now the leading sales coefficient is much higher and close to unity, and the Sargan test has a  $p$ -value close to 40 per cent.
- Finally, col 6 shows GMM estimates based on

$$E(\log s_{it} v_{is}) = 0 \quad (t = 1, \dots, s - 1, s + 1, \dots, T; s = 1, \dots, T),$$

which allow for measurement error in sales but not for correlated effects. The leading sales effect in this case is close to OLS in levels, suggesting that in levels the measurement error bias is not as important as in differences.

### *Conclusion*

- What is interesting about this example is that a comparison between estimates in levels and deviations without consideration of measurement error (e.g. restricted to compare cols 1 & 2, or 1 & 3, as in Hausman-type testing), would lead to the conclusion of correlated effects, but with biases going in entirely the wrong direction.

## Predeterminedness and dynamics

### *Time patterns*

- The previous examples include fixed effects but do not allow for time patterns in the dependence between  $x$  and time-varying errors.
- However, the time dimension makes it possible to go beyond the cross-sectional notions of strict exogeneity and strict endogeneity, whereby the time series of a regressor is either fully independent or fully dependent of the time series of errors.
- Thus,  $x$  may depend on past  $v$ 's but not on future  $v$ 's (predeterminedness), or on  $v$ 's that are close in time but not on  $v$ 's from distant periods.
- A linear model with general predetermined variables replaces the strict exogeneity assumption  $E(v_{it} | x_{i1}, \dots, x_{iT}, \alpha_i) = 0$  with the sequential conditioning assumption

$$E(v_{it} | x_{i1}, \dots, x_{it}, \alpha_i) = 0.$$

Letting  $x_i^t = (x_{i1}, \dots, x_{it})$ , such model implies the moments:

$$E \left[ x_i^{t-1} (\Delta y_{it} - \Delta x_{it}' \theta_0) \right] = 0.$$

- This notion can be generalized to external instruments and to alternative patterns of leads or lags.
- An example is the relationship between the presence of small children at home and female labor supply. Treating children as strictly exogenous in this context is a much more restrictive assumption than treating them as predetermined.

### *First-stage and second-stage regressions in panel GMM*

- In Arellano-Bond GMM estimation there is a sequence of period-by-period first-stage regressions and a pooled second-stage regression.
- Letting for simplicity  $T = 3$  and a single predetermined regressor, the period-by-period first-stage fitted values are

$$\begin{aligned}\widehat{\Delta x_{i2}} &= \widehat{\pi}_{21} x_{i1} \\ \widehat{\Delta x_{i3}} &= \widehat{\pi}_{31} x_{i1} + \widehat{\pi}_{32} x_{i2}\end{aligned}$$

where  $\widehat{\pi}_{21}$  is the cross-sectional OLS coefficient of  $\Delta x_{i2}$  on  $x_{i1}$ , etc. (in practice, orthogonal deviations are preferred to first-differences but the idea is the same).

- The second-stage is a pooled IV regression of  $(\Delta y_{i2}, \Delta y_{i3})$  on  $(\Delta x_{i2}, \Delta x_{i3})$  using  $(\widehat{\Delta x_{i2}}, \widehat{\Delta x_{i3}})$  as instruments.
- The latter is very different to the time-series perspective where instruments would come from a pooled first-stage regression:

$$\begin{pmatrix} \widetilde{\Delta x_{i2}} \\ \widetilde{\Delta x_{i3}} \end{pmatrix} = \widetilde{\pi} \begin{pmatrix} x_{i1} \\ x_{i2} \end{pmatrix}$$

where  $\widetilde{\pi}$  is the pooled OLS coefficient of  $(\Delta x_{i2}, \Delta x_{i3})$  on  $(x_{i1}, x_{i2})$ . The 2nd-stage would be pooled IV of  $(\Delta y_{i2}, \Delta y_{i3})$  on  $(\Delta x_{i2}, \Delta x_{i3})$  using  $(\widetilde{\Delta x_{i2}}, \widetilde{\Delta x_{i3}})$  as instruments.

- In a pooled first-stage regression one cannot easily project on different  $x$ 's at different periods as one does using period-by-period first stage regressions.



## *Dynamic models*

- Time patterns of dependence arise naturally in the context of dynamic models. These are models that consider the effects of lagged outcomes and/or lagged and current independent explanatory variables on current outcomes.
- The simplest example is an autoregressive model, which is a special case of the above with  $x_{it} = y_{i(t-1)}$ .
- The basic moments are:

$$E \left[ y_i^{t-2} \left( \Delta y_{it} - \Delta y_{i(t-1)} \theta_0 \right) \right] = 0,$$

- Under mean stationarity, the following moments for the errors in levels are also available:

$$E \left[ \Delta y_{i(t-1)} \left( y_{it} - y_{i(t-1)} \theta_0 \right) \right] = 0.$$

- Autoregressive models are the workhorse in the analysis of individual earnings and household income dynamics.

## Permanent-transitory income models

- Permanent-transitory models are common in the literature that looks at the relationship between household income and consumption from a life-cycle perspective.
- Examples include Hall & Mishkin (1982) (HM), Blundell, Pistaferri & Preston (2008), and Kaplan & Violante (2010).
- HM used food consumption and labour income from a PSID sample of  $N = 2309$  US households over  $T = 7$  years to test the predictions of a permanent income model.
- We use HM as an illustration of permanent-transitory covariance structures.
- HM specified means of income and consumption changes as regressions on age, age<sup>2</sup>, time, and changes in the number of children and adults in the household.
- They implicitly allowed for unobserved intercept heterogeneity in the levels of the variables, but only for observed heterogeneity in their changes.
- Deviations from the individual means of income and consumption, denoted  $\bar{y}_{it}$  and  $\bar{c}_{it}$  respectively, were specified as follows.

### *Income process*

- HM assumed that income errors  $\bar{y}_{it}$  were the result of two different types of shocks, permanent and transitory:

$$\bar{y}_{it} = y_{it}^L + y_{it}^S.$$

- They also assumed that agents were able to distinguish one type of shock from the other and respond to them accordingly.
- The permanent component  $y_{it}^L$  was specified as a random walk

$$y_{it}^L = y_{i(t-1)}^L + \varepsilon_{it},$$

and the transitory component  $y_{it}^S$  as a moving average process

$$y_{it}^S = \eta_{it} + \rho_1 \eta_{i(t-1)} + \rho_2 \eta_{i(t-2)}.$$

- A limitation was lack of measurement error in observed income (a component to which consumption does not respond). This is important since measurement error in PSID income is large, but identification requires cross-validation information.

### *Consumption process*

- Mean deviations in consumption changes were specified to respond one-to-one to permanent income shocks and by a fraction  $\beta$  to transitory shocks.
- The magnitude of  $\beta$  depends on the persistence in transitory shocks ( $\rho_1$  and  $\rho_2$ ) and real interest rates. Dependence on age is ignored for simplicity.
- This model can be formally derived from an optimization problem with quadratic utility, and constant interest rates that are equal to the subjective discount factor.
- Since only food consumption is observed, an adjustment was made by assuming a constant marginal propensity to consume food  $\alpha$ .
- With these assumptions we have

$$\Delta \bar{c}_{it} = \alpha \varepsilon_{it} + \alpha \beta \eta_{it}.$$

- HM also introduced a measurement error in the level of consumption (or transitory consumption that is independent of income shocks) with an MA(2) specification:

$$c_{it}^S = v_{it} + \lambda_1 v_{i(t-1)} + \lambda_2 v_{i(t-2)}.$$

### *Bivariate covariance structure*

- The model that is taken to the data consists of a joint specification for mean deviations in consumption and income changes as follows:

$$\Delta \bar{c}_{it} = \alpha \varepsilon_{it} + \alpha \beta \eta_{it} + v_{it} - (1 - \lambda_1) v_{i(t-1)} - (\lambda_1 - \lambda_2) v_{i(t-2)} - \lambda_2 v_{i(t-3)}$$

$$\Delta \bar{y}_{it} = \varepsilon_{it} + \eta_{it} - (1 - \rho_1) \eta_{i(t-1)} - (\rho_1 - \rho_2) \eta_{i(t-2)} - \rho_2 \eta_{i(t-3)}.$$

- The three innovations are mutually independent with variances  $\sigma_\varepsilon^2$ ,  $\sigma_\eta^2$  and  $\sigma_v^2$ . Thus, the model contains 9 coefficients:

$$\theta = \left( \alpha \quad \beta \quad \lambda_1 \quad \lambda_2 \quad \rho_1 \quad \rho_2 \quad \sigma_\varepsilon^2 \quad \sigma_\eta^2 \quad \sigma_v^2 \right)'$$

- The model specifies a covariance structure for the  $12 \times 1$  vector

$$w_i = \left( \Delta \bar{c}_{i2} \quad \Delta \bar{c}_{i3} \quad \cdots \quad \Delta \bar{c}_{i7} \quad \Delta \bar{y}_{i2} \quad \Delta \bar{y}_{i3} \quad \cdots \quad \Delta \bar{y}_{i7} \right)'$$

$$E(w_i w_i') = \Omega(\theta).$$

*Bivariate covariance structure (continued)*

- Let us look at the form of some elements of  $\Omega(\theta)$ .

$$\text{Var}(\Delta\bar{y}_{it}) = \sigma_{\varepsilon}^2 + 2 \left(1 - \rho_1 - \rho_1\rho_2 + \rho_1^2 + \rho_2^2\right) \sigma_{\eta}^2 \quad (t = 2, \dots, 7)$$

$$\text{Cov}(\Delta\bar{y}_{it}, \Delta\bar{y}_{i(t-1)}) = -[(1 - \rho_1) - (1 - \rho_1 + \rho_2)(\rho_1 - \rho_2)] \sigma_{\eta}^2$$

and also

$$\text{Cov}(\Delta\bar{c}_{it}, \Delta\bar{y}_{it}) = \alpha\sigma_{\varepsilon}^2 + \alpha\beta\sigma_{\eta}^2 \quad (t = 2, \dots, 7) \quad (2)$$

$$\text{Cov}(\Delta\bar{c}_{it}, \Delta\bar{y}_{i(t-1)}) = 0 \quad (3)$$

$$\text{Cov}(\Delta\bar{c}_{i(t-1)}, \Delta\bar{y}_{it}) = -\alpha\beta(1 - \rho_1)\sigma_{\eta}^2. \quad (4)$$

- A fundamental restriction of the model is lack of correlation between current consumption changes and lagged income changes, as captured by (3).
- The model, nevertheless, predicts correlation between current consumption changes and current and future income changes, as seen from (2) and (4).

### *Empirical results*

- HM estimated their model by Gaussian PML. They estimated  $\hat{\beta} = 0.3$ , which given their estimates of  $\rho_1$  and  $\rho_2$  ( $\hat{\rho}_1 = 0.3$ ,  $\hat{\rho}_2 = 0.1$ ) turned out to be consistent with the model only for unrealistic values of real interest rates (above 30 percent).
- Moreover, they estimated the marginal propensity to consume food as  $\hat{\alpha} = 0.1$ , and the moving average parameters for transitory consumption as  $\hat{\lambda}_1 = 0.2$  and  $\hat{\lambda}_2 = 0.1$ .
- The variance of the permanent income shocks was twice as large as that of the transitory shocks:  $\hat{\sigma}_\varepsilon^2 = 3.4$  and  $\hat{\sigma}_\eta^2 = 1.5$ .
- They tested the covariance structure focusing on the fundamental restriction of lack of correlation between current changes in consumption and lagged changes in income. They found a negative covariance which was significantly different from zero.
- As a result of this finding they considered an extended version of the model in which a fraction of consumers spent their current income.

## GMM estimation of covariance structures

- The previous model specifies a structure on a data covariance matrix. Abstracting from mean components, suppose the covariance matrix of a  $p \times 1$  time series  $y_i$  is a function of a  $k \times 1$  parameter vector  $\theta$  given by

$$E(y_i y_i') = \Omega(\theta).$$

- If  $y_i$  is a scalar time series its dimension will be  $T$ , but in the HM context  $p = 2T$ .
- Vectorizing the expression and eliminating redundant elements (due to symmetry) we obtain a vector of moments of order  $r = (p + 1)p/2$ :

$$\text{vech} E [y_i y_i' - \Omega(\theta)] = E [s_i - \omega(\theta)],$$

where the *vech* operator stacks by rows the lower triangle of a square matrix.

- If  $r > k$  and  $H(\theta) = \partial \omega(\theta) / \partial \theta'$  has full column rank, the model is overidentified. In that case a standard optimal GMM estimator solves:

$$\hat{\theta} = \arg \min_c [\bar{s} - \omega(c)]' \hat{V}^{-1} [\bar{s} - \omega(c)]$$

where  $\bar{s}$  is the sample mean vector of  $s_i$ :

$$\bar{s} = \frac{1}{N} \sum_{i=1}^N s_i$$

and  $\hat{V}$  is some consistent estimator of  $V = \text{Var}(s_i)$ . A natural choice is the sample covariance matrix of  $s_i$ :

$$\hat{V} = \frac{1}{N} \sum_{i=1}^N s_i s_i' - \bar{s} \bar{s}'.$$



*GMM estimation of covariance structures (continued)*

- The first-order conditions from the optimization problem are

$$-H(c)' \widehat{V}^{-1} [\bar{s} - \omega(c)] = 0.$$

- The two standard results for large sample inference are, firstly, asymptotic normality of the scaled estimation error

$$\left[ \frac{1}{N} H(\widehat{\theta})' \widehat{V}^{-1} H(\widehat{\theta}) \right]^{-1/2} (\widehat{\theta} - \theta) \xrightarrow{d} \mathcal{N}(0, I)$$

and, secondly, the asymptotic chi-square distribution of the minimized estimation criterion (*test statistic of overidentifying restrictions*)

$$S = N \left[ \bar{s} - \omega(\widehat{\theta}) \right]' \widehat{V}^{-1} \left[ \bar{s} - \omega(\widehat{\theta}) \right] \xrightarrow{d} \chi_{r-k}^2.$$

## Random coefficients

- Fixed effects methods are a standard way of controlling for endogeneity or unobserved heterogeneity in the estimation of common parameters.
- But sometimes we wish to treat a parameter as a heterogeneous quantity and therefore its mean and other characteristics of its distribution become central objects of interest.
- Examples are random trend earnings models, heterogeneous production functions, and heterogeneous treatment effects.
- The  $T$  equations of the random coefficients model in compact form can be written as

$$y_i = Z_i\delta_0 + X_i\gamma_i + v_i \quad E(v_i | Z_i, X_i, \gamma_i) = 0.$$

- The WG model is a special case in which the only random coefficient is the intercept.
- We assume that  $T > \dim(\gamma_i) = q$  and only consider the subpopulation with  $\det(X_i'X_i) \neq 0$ .
- The parameters of interest are  $\delta_0$  and characteristics of the distribution of  $\gamma_i$ , such as  $\gamma_0 = E(\gamma_i)$  and  $\Sigma_0 = \text{Var}(\gamma_i)$ .
- Now instead of considering LS in deviations from means we consider LS of the residuals in individual-specific regressions of  $y$  and  $z$  on  $x$  ( $\tilde{x}_{it}$  is the residual of a regression of the  $i$ -th time series of  $x$  on an intercept).

### Estimating common parameters and average effects

- The generalized WG operator  $Q_i = I - X_i (X_i' X_i)^{-1} X_i$  leads to the transformed equation

$$Q_i y_i = Q_i Z_i \delta_0 + Q_i v_i$$

and the moments

$$E [Z_i' (Q_i y_i - Q_i Z_i \delta_0)] = 0.$$

- The WG estimator is

$$\hat{\delta} = \left( \sum_{i=1}^N Z_i' Q_i Z_i \right)^{-1} \sum_{i=1}^N Z_i' Q_i y_i$$

- Pre-multiplying the model by the LS operator  $H_i = (X_i' X_i)^{-1} X_i'$  we get

$$H_i (y_i - Z_i \delta_0) = \gamma_i + H_i v_i$$

so that  $\gamma_0$  satisfies the moment

$$\gamma_0 = E [H_i (y_i - Z_i \delta_0)]$$

and a large- $N$  consistent estimator is

$$\hat{\gamma} = \frac{1}{N} \sum_{i=1}^N (X_i' X_i)^{-1} X_i' (y_i - Z_i \hat{\delta}) \equiv \frac{1}{N} \sum_{i=1}^N \hat{\gamma}_i.$$

## Is $\hat{\gamma}_i$ informative about $\gamma_i$ ? An illustration

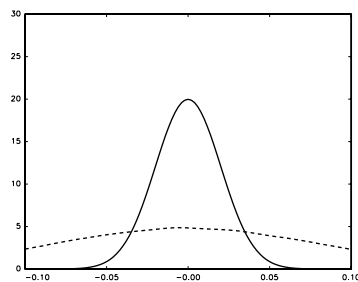
- Consider the random trend model:

$$y_{it} = \alpha_i + \beta_i t + v_{it}$$

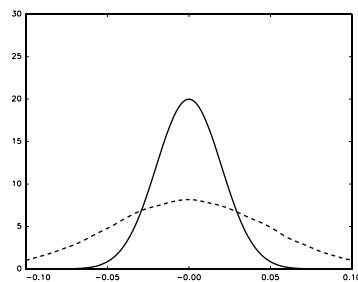
where  $\alpha_i$  and  $\beta_i$  are bivariate normal (or bimodal normal mixture),  $v_{it}$  is normal AR(1) with autoregressive coefficient  $\rho$ .

- Roughly calibrate the parameters to match Guvenen (2008):  $\rho = .8$ ,  $\text{Var}(\alpha_i) = .02$ ,  $\text{Var}(\beta_i) = .0004$  (corr. =  $-.2$ ),  $\sigma_v^2 = .03$ .
- Question: compare the density of  $\hat{\beta}_i$  (resp.  $\hat{\alpha}_i$ ) to that of  $\beta_i$  ( $\alpha_i$ ).

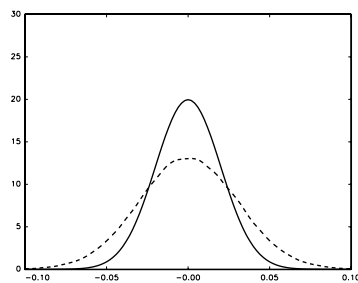
Densities: true  $\beta_i$  (solid) and fixed-effects estimates  $\hat{\beta}_i$  (dashed)



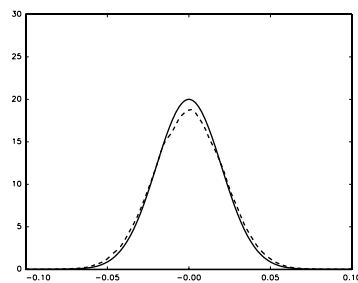
$T = 5$



$T = 10$

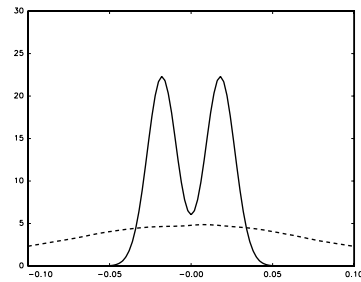


$T = 20$

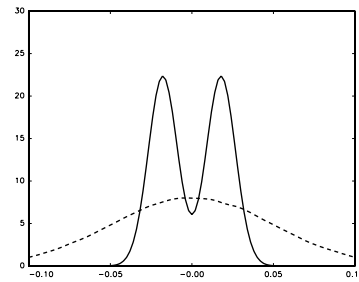


$T = 50$

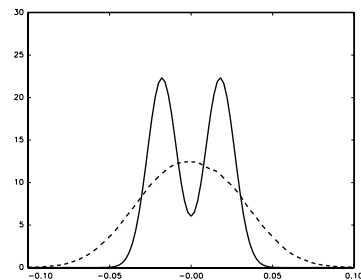
Densities: true  $\beta_i$  (solid) and fixed-effects estimates  $\hat{\beta}_i$  (dashed)



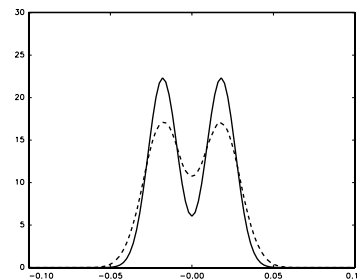
$T = 5$



$T = 10$



$T = 20$



$T = 50$

$\Rightarrow$  Must correct the densities of fixed-effects estimates for the sample noise (for fixed  $T$ ).

### *Estimating variances of effects and distributions*

- Without further restrictions  $\Sigma_0$  is not identified. To see this let  $\Omega_i = E(v_i v_i' | X_i)$  and note that only the variance of  $Q_i v_i$  is identified, which is of reduced rank. In general

$$\Sigma_0 = \text{Var} [H_i (y_i - Z_i \delta_0)] - E (H_i \Omega_i H_i').$$

- If  $\Omega_i = \sigma^2 I_T$  then  $\Sigma_0$  can be estimated as

$$\hat{\Sigma} = \frac{1}{N} \sum_{i=1}^N (\hat{\gamma}_i - \hat{\gamma}) (\hat{\gamma}_i - \hat{\gamma})' - \hat{\sigma}^2 \frac{1}{N} \sum_{i=1}^N (X_i' X_i)^{-1}$$

where

$$\hat{\sigma}^2 = \frac{1}{N(T-q)} \sum_{i=1}^N (y_i - Z_i \hat{\delta})' Q_i (y_i - Z_i \hat{\delta}).$$

- Note that  $E(Q_i v_i v_i' Q_i) = \sigma^2 E(Q_i)$  and  $E(v_i' Q_i v_i) = \sigma^2 (T - q)$ .

### *Estimating variances of effects and distributions (continued)*

- The previous situation can be generalized to less restrictive covariance patterns in  $\Omega_j$ .
- In general

$$E [(y_i - Z_i\delta_0) \otimes (y_i - Z_i\delta_0) \mid Z_i, X_i] = (X_i \otimes X_i) E (\gamma_i \otimes \gamma_i \mid Z_i, X_i) + \text{vec} (\Omega_i).$$

- A WG operator  $M_i = I - G_i (G_i' G_i)^{-1} G_i'$  for the cross-products  $G_i = X_i \otimes X_i$  leads to

$$M_i E [(y_i - Z_i\delta_0) \otimes (y_i - Z_i\delta_0) \mid Z_i, X_i] = M_i \text{vec} (\Omega_i)$$

but since  $M_i$  is singular, (moving-average) restrictions on  $\Omega_j$  are needed:

$$\text{vec} (\Omega_j) = S_2 \omega_j$$

where  $S_2$  is a known selection matrix and  $\omega_j$  is a vector of unrestricted parameters.

- The rank condition for identification of  $\Omega_j$  is

$$\text{rank} (M_i S_2) = \dim (\omega_j).$$

- The variance of  $\gamma_j$  is identified if  $\Omega_j$  is known.
- Moreover, replacing mean independence by full independence assumptions a similar argument can be developed for distributions using second derivatives of log characteristic functions (Arellano and Bonhomme 2012).



## Distributions

- Assume that  $\gamma_i$  and  $v_i$  are independent given  $W_i = (Z_i, X_i)$ .
- Statistical independence leads to functional restrictions on the second derivatives of log characteristic functions, which are formally analogous to the covariance restrictions.
- To derive the identification results, it is convenient to work with characteristic functions.

### *Properties of characteristic functions*

- The conditional characteristic function of  $Y$  (of dimension  $L$ ) given  $X = x$  is defined as:

$$\Psi_{Y|X}(t|x) = E [\exp(jt'Y)|x], \quad t \in R^L$$

where  $j = \sqrt{-1}$ .

- Inverse Fourier transform

$$f_{Y|X}(y|x) = \frac{1}{(2\pi)^L} \int \exp(-jt'y) \Psi_{Y|X}(t|x) dt.$$

- If  $Y_1$  and  $Y_2$  are independent given  $X$  then

$$\Psi_{Y_1+Y_2|X}(t|x) = \Psi_{Y_1|X}(t|x)\Psi_{Y_2|X}(t|x).$$

## Distributions (continued)

- Independence implies that for all  $t$  we have:

$$\Psi_{y_i - Z_i \delta_0 | W_i}(t | W_i) = \Psi_{\gamma_i | W_i}(X_i' t | W_i) \Psi_{v_i | W_i}(t | W_i).$$

- Assuming that the characteristic functions  $\Psi_{\gamma_i | W_i}$  and  $\Psi_{v_i | W_i}$  are nonvanishing we can take logs:

$$\log \Psi_{y_i - Z_i \delta_0 | W_i}(t | W_i) = \log \Psi_{\gamma_i | W_i}(X_i' t | W_i) + \log \Psi_{v_i | W_i}(t | W_i).$$

- If  $\Psi_{v_i | W_i}$  is identified,  $\Psi_{\gamma_i | W_i}$  is also identified.
- Taking second derivatives:

$$\frac{\partial^2 \log \Psi_{y_i - Z_i \delta_0 | W_i}(t | W_i)}{\partial t \partial t'} = X_i \left( \frac{\partial^2 \log \Psi_{\gamma_i | W_i}(X_i' t | W_i)}{\partial t \partial t'} \right) X_i' + \frac{\partial^2 \log \Psi_{v_i | W_i}(t | W_i)}{\partial t \partial t'}.$$

- Evaluating this expression at  $t = 0$  we are back at the variance case.

## Distributions (continued)

- An independent moving-average model implies the following restrictions:

$$\text{vec} \left( \frac{\partial^2 \log \Psi_{v_i|W_i}(t|W_i)}{\partial t \partial t'} \right) = S_2 \omega_i(t), \quad t \in R^T.$$

- So, if  $M_i(X_i \otimes X_i) = 0$  then

$$M_i \text{vec} \left( \frac{\partial^2 \log \Psi_{y_i - Z_i \delta_0 | W_i}(t|W_i)}{\partial t \partial t'} \right) = M_i S_2 \omega_i(t).$$

- The rank and order conditions for identification are the same as for variances.
- $\omega_i(t)$  identified for all  $t$  implies that  $\Psi_{v_i|W_i}$  is identified, because the first derivative of  $\log \Psi_{v_i|W_i}$  at  $t = 0$  vanishes due to mean independence.

## Illustration: the effect of smoking on children outcomes

- Arellano and Bonhomme (2012) apply this methodology to a matched panel dataset of mothers and births constructed in Abrevaya (2006).
- They find that the mean smoking effect on birthweight is significantly negative ( $-160$  grams). Moreover, the effect shows substantial heterogeneity across mothers, the effect being very negative ( $-400$  g) below the 20th percentile.

- The model is

$$y_{ij} = \mathbf{z}_{ij}'\boldsymbol{\delta} + \alpha_i + \beta_i s_{ij} + v_{ij} \quad j = 1, 2, 3$$

$i$ =mother,  $j$ =child.  $y_{ij}$ = weight at birth,  $s_{ij} = 1$  if mother smoked during pregnancy of child  $j$ .

- $v_{ij}$  are assumed i.i.d.
- Production function interpretation. The effect of smoking is mother-specific.
- Abrevaya (2006) estimates a restricted version, where  $\beta_i$  is homogeneous.
- The focus is on mothers with at least 3 children to be able to allow for two heterogeneous quantities.
- Also need  $x_{ij}$  to vary for every mother. So only 1445 mothers who changed smoking status between the three births are considered.
- Under predeterminedness of smoking behavior the moments of  $\beta_i$  are unidentified. However, several interesting average effects can still be identified and estimated when there are no time-varying regressors.

*Estimates of common parameters  $\delta$*

Generalized within-groups		
Variable	Estimate	Standard error
Male	130	22.8
Age	39.0	32.0
Age-sq	-.638	.577
Kessner=2	-82.0	52.7
Kessner=3	-159	81.9
No visit	-18.0	124
Visit=2	83.2	53.9
Visit=3	136	99.2

Regressions of  $\alpha_i$  and  $\beta_i$  on mother-specific characteristics

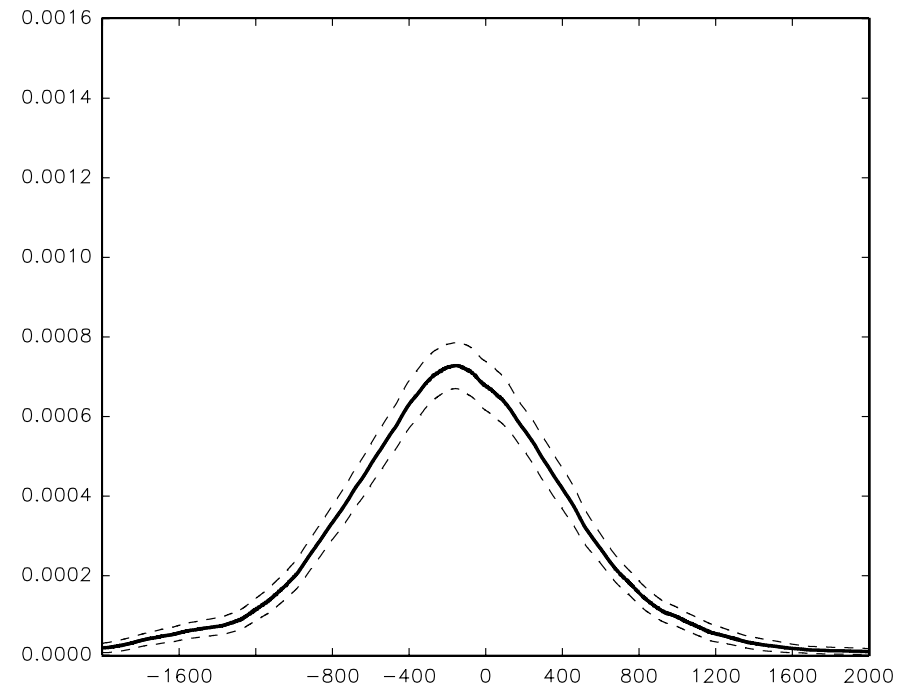
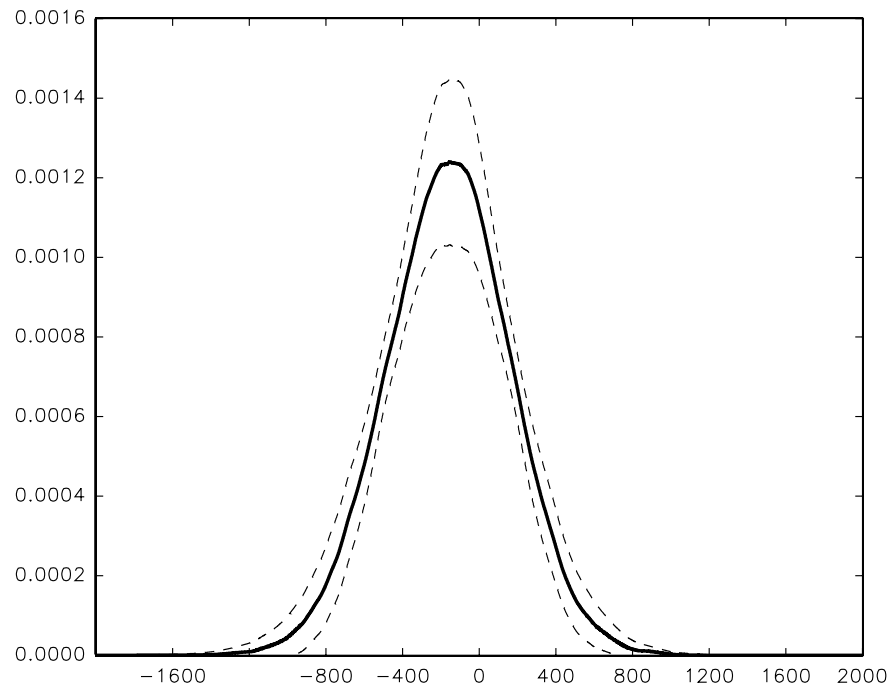
Variable	Estimate	Standard error
$\alpha_i$		
High-school	15.1	42.7
Some college	38.5	55.3
College graduate	58.7	72.1
Married	3.51	34.6
Black	-364	54.0
Mean smoking	-161	83.9
Constant	2879	419
corrected $R^2 = .113$ (instead of .055, uncorrected)		
$\beta_i$		
High-school	-15.9	42.8
Some college	-15.9	42.8
College graduate	64.5	63.8
Married	31.9	41.8
Black	132	60.6
Mean smoking	-49.8	101
Constant	-172	67.1
$R^2 = .021$ (instead of .005)		

Moments of  $\alpha_i$  and  $\beta_i$

Moment	Estimate	Standard error
Mean $\alpha_i$	2782	435
St. Dev. $\alpha_i$	357	21.2
Skewness $\alpha_i$	-1.67	.43
Kurtosis $\alpha_i$	7.12	2.28
Mean $\beta_i$	-161	17.0
St. Dev. $\beta_i$	313	34.6
Skewness $\beta_i$	-1.29	.91
Kurtosis $\beta_i$	-.34	7.84
Correlation ( $\alpha_i, \beta_i$ )	-.47	.07

- Mean effect of smoking is  $-161$  grams, close to Abrevaya's FE estimate of  $-144$  g.
- Density of  $\beta_i$  and  $\hat{\beta}_i$ .
- Quantile function of  $\beta_i$  and  $\hat{\beta}_i$ .

## Density of $\beta_i$ (left) and $\hat{\beta}_i$ (right)





## Quantile function of $\beta_i$ (left) and $\hat{\beta}_i$ (right)

