# Efficiency of repeated-cross-section estimators in fixed-effects models

Montezuma Dumangane[a], Nicoletta Rosati[*,a]

[a]*CEMAPRE and ISEG, Technical University of Lisbon, Portugal*

## Abstract

Efficiency of estimators in additive fixed-effects models is investigated. Asymptotically, panel data are more efficient in case of strong residual autocorrelation; in small samples, variances are comparable, but repeated cross-sections show larger bias for some parameter values.

*Key words:* asymptotic lower bound, conditional moment restrictions, fixed-effects model, repeated cross-sections

## 1. Introduction

The value of repeated cross-sections (RCS) to identifiability of structural parameters in fixed-effects models is known since Heckman and Robb (1985). At the same time, Deaton (1985) stressed that a measurement error problem, peculiar to RCS inference, arises in finite samples due to the need to estimate some nuisance parameters. He also showed that the resulting bias can be dealt with exploiting sample information. A line of research originated out of these seminal papers; among several contributions, see, for instance, Verbeek (1992), Moffit (1993), Verbeek and Nijman (1993), Collado (1997), Girma (2000), McKenzie (2004) and Verbeek and Vella (2005).

---

[*]Corresponding author. Department of Mathematics, ISEG, Technical University of Lisbon. Rua do Quelhas 6, 1200 - 781 Lisbon, Portugal

*Email addresses:* `mdumangane@iseg.utl.pt` (Montezuma Dumangane), `nicoletta@iseg.utl.pt` (Nicoletta Rosati)

A still unsettled issue to our knowledge is the derivation of a benchmark for the amount of precision we might expect in an RCS based inference, which would also be useful to compare the potential of RCS to that of genuine panel information. Heckman and Robb (1985) made an informal contention that even if panel estimators are often claimed to be asymptotically more efficient than the RCS ones, in finite samples this might not happen since sample size is by far larger in cross-sectional surveys.

In this paper we derive a lower bound on the asymptotic variance of a RCS estimator by exploiting a set of moment restrictions implied by the fixed-effects assumption. Under suitable regularity conditions, we allow for a general nonlinear model, provided that unobservables enter the model additively. The lower bound is obtained by applying standard results of inference under conditional moment restrictions (see, for instance, Chamberlain, 1987). Small sample behaviour is studied through simulation.

## 2. The model

Applied sciences often deal with panel data, namely a vector of variables $z_{it}$ observed at times $t = 1, \ldots, T$ on the units $i = 1, \ldots, N$. One of the advantages of panel data is the possibility to control for unobserved heterogeneity among individuals due to time-invariant effects, since repeated observations on each sampling unit are available. When panel data are not available, identification and estimation of the parameters is still possible - under suitable conditions - based on repeated cross-sections, where sample averages are computed from individual observations within pre-defined time-invariant classes (cohorts), playing the role of "macro-individuals" in a pseudo-panel dataset. Estimating techniques are basically the same in both cases. In the case of panel data they are applied directly to individual observations, while in the case of repeated cross-sections they are applied to the class sample averages. The sample averages are considered as error-ridden measurements of the true class means, and therefore the estimators are corrected for the presence of response

2

and covariate measurement error. When the sample size for each class is sufficiently large, the measurement error can be neglected, i.e. no correction is necessary.

To simplify notation, from now on let $i$ indicate the sampling unit or the synthetic macro-individual, according to the sampling scheme (panel or pseudo-panel). For instance, $y_i$ would indicate the response for individual $i$ in a panel, but represents the sample average of the response for class $i$ in a pseudo-panel ($i$-th macro-individual), the only difference being the presence of (possibly negligible) measurement error in the second case. Consider models where the unobservable component has the form:

$$u_{it} = g(z_{it}; \beta) \tag{1}$$

where $g(\ .\ )$ is known, $z_{it} = (y_{it}, x_{it})$ where $y_{it}$ is the response variable and $x_{it}$ is a vector of explanatory variables, and $E(u_{it}|x_{it}) = 0$. Our interest is in estimation of $\beta$, which is both time- and individual-invariant. The error $u_{it}$ is assumed to be such that we can represent it as:

$$u_{it} = \eta_i + v_{it}$$
$$E\{u_{it}|w_i, x_{it}\} = 0 \qquad \forall w_i \in \mathcal{W} \tag{2}$$

where $\eta_i$ is an unobservable individual effect, $w_i$ is a vector of observable time-invariant variables and $\mathcal{W}$ is its support. In the case of multiplicative effects other transformations can be used in order to eliminate the individual effect.

Notice that in the case of RCS the individual effects depend on time too, i.e. $\eta_i = \eta_{i(t)}$, being the sample average of the effects for the specific cohort sample at time $t$. However, its conditional expectation remains time-invariant within cohort $i$, i.e. $E\{\eta_{i(r)}|w_i\} = E\{\eta_{i(s)}|w_i\}$, $\forall r \neq s$. This, together with assumption (2), implies among other things that for two time periods $r \neq s$ the following condition holds

$$E\{g(z_{ir}; \beta) - g(z_{is}; \beta)|w_i, x_{ir}, x_{is}\} = 0 \tag{3}$$

This equation specifies a conditional moment restriction upon which inference on $\beta$ can rest. The lower bound we shall derive exploits such restriction.

3

## 3. Asymptotic efficiency

Suppose one is choosing whether to collect a panel data set or a sequence of cross-sections. Since sampling costs are rather different, it is often the case that the panel sample is of much smaller size than the cross-sectional samples (which we assume - for the sake of simplicity - to be all of the same size). Therefore, we aim at answering the following question: to estimate the parameter of interest, is it more convenient to use panel or RCS information?

An answer to this question can be obtained comparing the optimal asymptotic variance attainable in the two cases; such a comparison can be derived directly from the results in Chamberlain (1987), which we briefly recall here.

Consider a random sample of size $N$ from variables $z$ and $w$, a parameter $\beta$ and a function $m(z, \beta)$ such that the following moment condition holds: $E\{m(z, \beta_0)|\, w\} = 0$; suppose also that function $m$ satisfies the following regularity conditions:

(i) $\beta$ is in an open set $\mathcal{B} \subset \mathbb{R}^p$ such that $m(z, \beta)$ and $\partial m(z, \beta)/\partial \beta^\top$ are continuous for $(z, \beta) \in \mathcal{Z} \times \mathcal{B}$.

(ii) $E\{m(z, \beta_0)|w\} = 0 \ \forall \ w \in \mathcal{W}$.

(iii) $\Sigma(w) = E\{\, m(z, \beta)\, m(z, \beta)^\top |w\}$ is positive-definite for all $w \in \mathcal{W}$.

(iv) Let $D(w) = E\{\, \partial m(z, \beta)/\partial \beta^\top |w\}$ for $w \in \mathcal{W}$. Matrix $E\{D^\top(w)\Sigma^{-1}(w)D(w)\}$ is positive-definite.

Then a lower bound on asymptotic variance for any regular consistent asymptotically normal estimator is given by

$$\Lambda = \left\{\ E\ \left[D^\top(w)\Sigma^{-1}(w)D(w)\right]\ \right\}^{-1} \tag{4}$$

where matrices $D(w)$ and $\Sigma(w)$ are evaluated at $\beta = \beta_0$.

To apply Chamberlain's results to our case, consider the following function

$$m(z_i, \beta) = g(z_{i1}; \beta) - g(z_{i2}; \beta) \tag{5}$$

Equation (3) guarantees that this $m(z, \beta_0)$ satisfies the moment condition above: in the case of panel data this is true at individual level, while observing repeated

4

cross-sections this holds at cohort level (macro-individuals). It follows that a lower bound on asymptotic variance of GMM estimators based on this moment condition can be obtained from (4), if the appopriate regularity conditions hold.

It is easy to verify that if the function $m$ defined in (5) satisfies conditions (i)-(iv) above for panel data, then this happens for RCS too, so that the asymptotic lower bound can be applied to our problem. Asumptotics here refers to large $N$, where $N$ represents the number of individuals in the panel data and the number of cohorts in the repeated cross-sectional data, respectively.

Substituting function $m$ into the definition of $D(w)$ and $\Sigma(w)$, we obtain:

$$D(w_i) = E[\partial g(z_{i1}, \beta)/\partial \beta' | w_i] - E[\partial g(z_{i2}, \beta)/\partial \beta' | w_i]$$

and

$$\Sigma(w_i) = \text{Var}\{g(z_{i1}, \beta) - g(z_{i2}, \beta)| \ w_i \ \} = \text{Var}\{u_{i1} - u_{i2}| \ w_i \ \}$$

To specify the expression of the lower bound in the two sampling schemes, recall that in the case of RCS the variables represented the cohort sample average. Making this explicit, this yields:

$$
\begin{aligned}
D(w_i) &= E\left[\frac{\partial \bar{g}(z_{i1}, \beta)}{\partial \beta'}\bigg| w_i\right] - E\left[\frac{\partial \bar{g}(z_{i2}, \beta)}{\partial \beta'}\bigg| w_i\right] \\
&= E\left[\frac{1}{N}\sum_{i=1}^{N}\frac{\partial g(z_{i1}, \beta)}{\partial \beta'}\bigg| w_i\right] - E\left[\frac{1}{N}\sum_{i=1}^{N}\frac{\partial g(z_{i2}, \beta)}{\partial \beta'}\bigg| w_i\right] \\
&= E\left[\frac{\partial g(z_{i1}, \beta)}{\partial \beta'}\bigg| w_i\right] - E\left[\frac{\partial g(z_{i2}, \beta)}{\partial \beta'}\bigg| w_i\right]
\end{aligned}
$$

This expression coincides with the case of panel data, so that it is not affected by the sampling scheme. Consider now the difference $u_{i1} - u_{i2}$ appearing in the definition of $\Sigma(w_i)$. From assumption (2), in the case of panel data $u_{i1} - u_{i2} = v_{i1} - v_{i2}$ so that in case of homoskedasticity of the residuals

$$
\begin{aligned}
\Sigma_P(w_i) &= \text{Var}(v_{i1}|w_i) + \text{Var}(v_{i2}|w_i) - 2\,\text{Cov}(v_{i1}, v_{i2}|w_i) \\
&= 2[\text{Var}(v_i|w_i) - \text{Cov}(v_{i1}, v_{i2}|w_i)]
\end{aligned}
$$

5

In the case of RCS, instead, $\bar{u}_{i1} - \bar{u}_{i2} = \bar{v}_{i1} - \bar{v}_{i2} + \bar{\eta}_{i1} - \bar{\eta}_{i2}$. This implies that, if $v$ is orthogonal to the other variables in the model

$$
\begin{aligned}
\Sigma_{RCS}(w_i) &= \mathrm{Var}(\bar{v}_{i1}|w_i) + \mathrm{Var}(\bar{v}_{i2}|w_i) - 2\,\mathrm{Cov}(\bar{v}_{i1}, \bar{v}_{i2}|w_i) + \\
&\quad + \mathrm{Var}(\bar{\eta}_{i1}|w_i) + \mathrm{Var}(\bar{\eta}_{i2}|w_i) - 2\,\mathrm{Cov}(\bar{\eta}_{i1}, \bar{\eta}_{i2}|w_i) \\
&= \frac{2}{n^*}[\mathrm{Var}(v_i|w_i) + \mathrm{Var}(\eta_i|w_i)]
\end{aligned}
$$

where $n^* = 2n_{i1}n_{i2}/(n_{i1} + n_{i2})$, being $n_{it}$ the sample size at time $t$ for cohort $i$. If a cohort has always the same number of sampled individuals at all times, then this coincides with $n^*$. Notice that all covariances are zero, as the samples in the two time periods are independent.

In the linear case $g(z_{it}, \beta) = y_i - x_{it}'\beta$, so that $m(z_i, \beta) = (y_{i1} - y_{i2}) - (x_{i1} - x_{i2})'\beta$ and $\partial g(z_{it}, \beta)/\partial\beta' = -x_{it}'$. Quantity $\Sigma(w_i)$ is the same as in the general case, while $D(w_i) = E\{x_{i2}'|\ w_i\ \} - E\{x_{i1}'|\ w_i\ \} = E[\Delta x_i|w_i]$, where $\Delta x_i = x_{i2} - x_{i1}$. Note that, as $E[\Delta x_i|w_i]$ gets closer to zero (non-identifiable model), the variance grows larger and larger.

Analysing the expression for $\Sigma(w_i)$ is enough to compare the asymptotic lower bound under the two sampling schemes, since the term $D(w_i)$ is not affected. The lower bound for panel data is larger than that for RCS iff:

$$
\mathrm{Var}(v_i|w_i) - \mathrm{Cov}(v_{i1}, v_{i2}|w_i) \geq \frac{1}{n^*}[\mathrm{Var}(v_i|w_i) + \mathrm{Var}(\eta_i|w_i)]
$$

that is iff

$$
\mathrm{Var}(v_i|w_i)\left(1 - \frac{1}{n^*}\right) \geq \mathrm{Cov}(v_{i1}, v_{i2}|w_i) + \frac{1}{n^*}\mathrm{Var}(\eta_i|w_i)
$$

If the sample size for each cohort is sufficiently large, then the terms multiplied by $1/n^*$ become negligible, showing that panel is superior to RCS only if the serial correlation of the residuals $v_{it}$ is positive and sufficiently large ($\mathrm{Var}(v_i|w_i)/\mathrm{Cov}(v_{i1}, v_{i2}|w_i) < 1$). Notice that this is never the case for both the $AR(1)$ and the $MA(1)$ specifications for $v_{it}$. In case of uncorrelated residuals, RCS is clearly better, showing that there is no potential efficiency gain using panel data.

The cases of small $n^*$ with large $N$, and small $n^*$ with small $N$ are analysed through simulation, as shown in the next Section.

6

## 4. Simulation study for small sample efficiency

A simulation study has been conducted with the following data generating process:

$$y_{it} = \alpha y_{it-1} + \beta x_{it} + \eta_i + v_{it} \quad i = 1, ..., N \ \ t = 1, 2$$

Let $y_{i0} = 0$ and $\eta_i \sim N(0, \sigma_\eta^2)$. The error follows an AR(1) model specified as $v_{it} = \theta_0 + \theta_1 v_{it-1} + \xi_{it}$, where $\xi_{it} \sim N(0, 1)$. The regressors are also AR(1) and are correlated with the individual effects, the generating process being specified as $x_{it} = \rho x_{it-1} + \gamma \eta_i + \varepsilon_{it}$, with $\varepsilon_{it} \sim N(0, \sigma_\varepsilon^2)$ independent of $\eta_i$.

The $j$-th instrument $z_{it}^{(j)}$ is generated as $z_{it}^{(j)} = \rho z_{it-1}^{(j)} + \varepsilon_{it} \omega_{it}^{(j)}$, for $j = 1, \ldots, J$, were $\rho$ and $\varepsilon_{it}$ are the same as in the $x_{it}$ equation above, $z_{i0}^{(j)} = x_{i0}$ and $\omega_{it}^{(j)} \sim N(1, \sigma_\omega^2)$ independent of $\varepsilon_{it}$. The expression of $\rho_{xz}^{(j)} = \text{Corr}(x_{it}, z_{it}^{(j)})$ is the same for all $j$, is approximately constant over time for $t = 2, 3$, and depends on parameters $\rho, \sigma_\varepsilon^2, \sigma_\eta^2, \gamma$ and $\sigma_\omega^2$.

Estimation of $\beta$ is performed using the difference-based moment condition (3), where the $g$ function is now given by $g(z_{it}; \beta) = y_{it} - \alpha y_{it-1} - \beta x_{it}$. The corresponding orthogonality condition is given by

$$E[(\Delta \eta_i + \Delta v_{it}) \Delta x_{it}] = 0$$

In the case of panel data, $i$ represents the individual, and therefore $\Delta \eta_i = 0$, which gives the condition $E[\Delta v_{it} \Delta x_{it}] = 0$. This is always true for any set of parameter values, therefore the first-differences ordinary-least-squares (FD-OLS) estimator is consistent.

For RCS $i$ represents the cohort, and therefore as noted at the beginning $\eta_i = \eta_{i(t)}$, so the expected value above becomes $E[(\Delta \eta_{i(t)} + \Delta v_{it}) \Delta x_{it}] = \sigma_\eta^2 \gamma (2 - \rho)$. It follows that FD-OLS is inconsistent if both $\gamma$ and $\sigma_\eta^2$ are non-zero, although the estimator has a good behaviour for any combination of the parameters for which the product $\sigma_\eta^2 \gamma$ is approximately zero. If $\sigma_\eta^2 \gamma$ is large, consistent estimates can be obtained through first-differences instrumental-variables (FD-IV) estimation.

### 4.1. Experimental design

Denote by $\bar{n} = N^{-1} \sum_{i=1}^{N} n_{i1} = N^{-1} \sum_{i=1}^{N} n_{i2}$ the average number of individuals within each cohort, where $N$ is the number of cohorts (equal in period one and two) and $n_{it}$ is the number of individuals sampled from cohort $i$ in period $t$. Several cases are considered, for $N \in \{50, 100, 150, 200, 300\}$ and $\bar{n} \in \{8, 16, 24\}$.

The results are compared to panel data where $\bar{n} = 1$, that is a panel data-set with a number of individuals equal to the number of cohorts in the RCS data-set. Simulations were not performed for combinations where the total number of RCS observations $\bar{n} * N$ was very large ($> 2400$).

At each simulation, $\bar{n} \cdot N$ observations are generated, split into cohorts according to a time invariant variable, and the cohort means are computed.

The parameter of interest is $\beta$, which is set equal to one. Some of the other parameters are held fixed, with values respectively $\alpha = 0$, $\theta_0 = 1$, $\rho = 0.5$ and $\sigma_\varepsilon^2 = 1$. Various sets of studies are performed, according to different values for $\theta_1$, $\sigma_\eta^2$, and $\gamma$ as reported in Table 1. In Groups 1 and 2 OLS is consistent for both panel and RCS, so the simulations are performed using the FD-OLS estimator illustrated above. Inconsistent FD-OLS estimates for RCS in Group 3 are compared with FD-IV results, obtained with $J$ instruments, $J \in \{1, 5, 10\}$, using a common correlation $\rho_{xz}^{(j)} \in \{0.5, 0.8\}$ for all intruments in a simulation.

### 4.2. Results

All simulation results presented below are based on 200 replications.

Results for FD-OLS estimation are reported in Table 2. For the case of $\gamma = 0$ (Group 1) the panel and RCS estimators show comparable efficiency, with slight superiority of panel data. The top panel in Table 2 shows results for $\theta_1 = 0.1$ and $\sigma_\eta^2 = 0.5$, but the other simulations in this group yielded analogous evidence.

In Group 2, where the product $\sigma_\eta^2 \gamma$ is small, the choice of the value for parameter $\theta_1$ does not affect the results. The middle panel in Table 2 shows the evidence for $\gamma = \sigma_\eta^2 = 0.1$. The RCS estimates show a small bias (around 1-2%), and standard

8

Table 1: Parameter values chosen for simulation and their influence on the variances and covariances of some design variables, for $\rho = 0.5$.

| Group | $\theta_1$ | $\sigma_\eta^2$ | $\gamma$ | Var(x) | Var(y) | Corr(x,$\eta$) | Var(v) | $r_1^*$ | $r_2^*$ | $R^{2*}$ |
|---|---|---|---|---|---|---|---|---|---|---|
| 1 | 0.5 | 0.1 | 0.0 | 1.47 | 2.90 | 0.00 | 1.33 | 0.07 | 0.06 | 0.51 |
|   | 0.1 | 0.1 | 0.0 | 1.47 | 2.58 | 0.00 | 1.01 | 0.09 | 0.06 | 0.57 |
|   | 0.1 | 0.5 | 0.0 | 2.00 | 3.51 | 0.00 | 1.01 | 0.33 | 0.20 | 0.57 |
| 2 | 0.1 | 0.1 | 0.1 | 1.49 | 2.64 | 0.05 | 1.01 | 0.09 | 0.06 | 0.56 |
|   | 0.5 | 0.1 | 0.1 | 1.49 | 2.97 | 0.05 | 1.33 | 0.07 | 0.06 | 0.50 |
| 3 | 0.5 | 0.1 | 0.5 | 4.00 | 8.33 | 0.50 | 1.33 | 0.43 | 0.14 | 0.48 |
|   | 0.1 | 0.1 | 0.5 | 1.60 | 2.91 | 0.25 | 1.01 | 0.09 | 0.05 | 0.55 |
|   | 0.1 | 0.5 | 0.1 | 2.13 | 3.84 | 0.10 | 1.01 | 0.33 | 0.18 | 0.56 |

\* $r_1 = \sigma_\eta^2/(\text{Var}(v) + \sigma_\eta^2)$, $r_2 = \sigma_\eta^2/\text{Var}(\beta x + \eta)$, $R^2 = \text{Var}(\beta x)/\text{Var}(y)$.

errors slightly lager than the panel estimates. The panel standard errors are smaller in Group 2 than in Group 1.

Group 3 considers the case where $\sigma_\eta^2 \gamma$ is larger, namely five times the value in Group 2, being either $\gamma = 0.1$ and $\sigma_\eta^2 = 0.5$ or $\gamma = 0.5$ and $\sigma_\eta^2 = 0.1$. This case yields bias for the RCS case reaching values aroud 7-8%, as shown in the bottom panel in Table 2 for $\gamma = 0.5$ and $\sigma_\eta^2 = 0.1$. Again, the value for parameter $\theta_1$ is irrelevant. Comparison based on MSE shows superiority of panel data for this specific case.

In general, increasing the number of observations within each cohort, for a fixed number of cohorts, does not improve efficiency nor reduces the bias, while increasing the number of cohorts helps reducing the variance.

Consistent IV estimation for larger $\sigma_\eta^2 \gamma$ achieves bias correction, but for most sets of instruments is far from OLS efficiency. For the sake of comparison with the bottom panel of Table 2, results for the case $\gamma = 0.5$ and $\sigma_\eta^2 = 0.1$ are shown in Table 3 with $J = 5$ and $\rho_{xz}^{(j)} = 0.5$. Table 4 compares efficiency for the case $N = 50$ and $\bar{n} = 8$ for different sets of instruments, and shows how standard errors closer to the OLS ones can be obtained by changing the number and strength of the instruments. However, such a large number of instruments or such a high correlation with the covariates might be difficult to achieve in real data applications.

Table 2: FD-OLS simulation results for 200 replications. True value $\beta = 1$; standard errors in parentheses.

### Group 1: $\theta_1 = 0.1$, $\sigma_\eta^2 = 0.5$ and $\gamma = 0$

| | | $\bar{\text{n}}$ | | | |
|---|---|---|---|---|---|
| | | **1** | **8** | **16** | **24** |
| | **50** | 0.9911 (0.1296) | 0.9912 (0.1618) | 1.0028 (0.1562) | 0.9986 (0.1538) |
| | **100** | 0.9981 (0.0832) | 0.9978 (0.1140) | 1.0023 (0.0909) | 0.9922 (0.1153) |
| **N** | **150** | 1.0053 (0.0703) | 0.9995 (0.0827) | 0.9907 (0.0877) | – |
| | **200** | 0.9993 (0.0639) | 0.9978 (0.0796) | – | – |
| | **300** | 1.0016 (0.0518) | 0.9987 (0.0606) | – | – |

### Group 2: $\theta_1 = 0.5$, $\sigma_\eta^2 = 0.1$ and $\gamma = 0.1$

| | | $\bar{\text{n}}$ | | | |
|---|---|---|---|---|---|
| | | **1** | **8** | **16** | **24** |
| | **50** | 0.9908 (0.1099) | 1.0117 (0.1502) | 1.0186 (0.1508) | 1.0219 (0.1472) |
| | **100** | 0.9974 (0.0675) | 1.0145 (0.1023) | 1.0129 (0.1024) | 1.0160 (0.1076) |
| **N** | **150** | 1.0025 (0.0610) | 1.0139 (0.0927) | 1.0107 (0.0844) | – |
| | **200** | 0.9995 (0.0543) | 1.0132 (0.0806) | – | – |
| | **300** | 1.0001 (0.0438) | 1.0155 (0.0576) | – | – |

### Group 3: $\theta_1 = 0.1$, $\sigma_\eta^2 = 0.1$ and $\gamma = 0.5$

| | | $\bar{\text{n}}$ | | | |
|---|---|---|---|---|---|
| | | **1** | **8** | **16** | **24** |
| | **50** | 0.9911 (0.1297) | 1.0850 (0.1468) | 1.0606 (0.1357) | 1.0634 (0.1277) |
| | **100** | 0.9981 (0.0833) | 1.0567 (0.0895) | 1.0600 (0.0901) | 1.0654 (0.0861) |
| **N** | **150** | 1.0053 (0.0702) | 1.0564 (0.0746) | 1.0521 (0.0687) | – |
| | **200** | 0.9994 (0.0639) | 1.0740 (0.0652) | – | – |
| | **300** | 1.0016 (0.0517) | 1.0699 (0.0557) | – | – |

Table 3: Simulation results for 200 replications, with $\theta_1 = 0.1$, $\sigma_\eta^2 = 0.1$ and $\gamma = 0.5$. True value $\beta = 1$; standard errors in parentheses. FD-OLS for $\bar{n} = 1$, FD-IV with $J = 5$ and $\rho_{xz}^{(j)} = 0.5$ for other values of $\bar{n}$.

| | | $\bar{n}$ | | | |
|---|---|---|---|---|---|
| | | **1** | **8** | **16** | **24** |
| | **50** | 0.9911 (0.1297) | 0.9998 (0.1872) | 0.9838 (0.1661) | 0.9996 (0.1613) |
| | **100** | 0.9981 (0.0833) | 0.9725 (0.1185) | 0.9933 (0.1127) | 1.0007 (0.1108) |
| **N** | **150** | 1.0053 (0.0702) | 0.9793 (0.0976) | 0.9925 (0.0862) | – |
| | **200** | 0.9994 (0.0639) | 0.9981 (0.0818) | – | – |
| | **300** | 1.0016 (0.0517) | 0.9950 (0.0708) | – | – |

Table 4: FD-IV simulation results for 200 replications, with $\theta_1 = 0.1$, $\sigma_\eta^2 = 0.1$ and $\gamma = 0.5$. True value $\beta = 1$; standard errors in parentheses. Case of $N = 50$ and $\bar{n} = 8$. For comparison to FD-OLS results: $\hat{\beta} = 1.0850$, SE=0.1468.

| | | J | | |
|---|---|---|---|---|
| | | **1** | **5** | **10** |
| $\rho_{xz}^{(j)}$ | **0.5** | 1.0108 (0.3465) | 0.9998 (0.1872) | 1.0180 (0.1662) |
| | **0.8** | 1.0115 (0.1811) | 1.0091 (0.1560) | 1.0227 (0.1525) |

## 5. Concluding remarks

It is well known that estimation of models with unobservable individual-specific effects is possible not only if repeated measurements on sampling units are available, but also with data collected in repeated cross-sections. This latter alternative, already investigated in the case of linear models, has been extended in this paper to more general models imposing only that unobservables enter the model additively.

The main result we obtain is a lower bound on the asymptotic variance of RCS estimators. This is derived from known results on estimation with conditional moment restrictions, which we build on to take into account that available information comes from two (or more) independent samples. Using the bound, asymptotic efficiency of RCS and panel data estimators is compared, showing that panel data are more efficient only in case of strong residual autocorrelation; a simulation study is performed in order to study the finite sample behaviour, finding comparable vari-

ances, but larger bias in repeated cross-sections for some sets of parameter values. IV estimation achieved bias correction, although variances comprable to OLS are obtained only for many and strong instruments.

## 6. Acknowledgements

## 7. References

Chamberlain, G. (1987). Asymptotic efficiency in estimation with conditional moment restrictions. *Journal of Econometrics*, **34**, 305-334.

Collado, M. D. (1997). Estimating dynamic models from time series of independent cross-sections. *Journal of Econometrics*, **82**, 37-62.

Deaton, A. (1985). Panel data from time series of cross-sections. *Journal of Econometrics*, **30**, 109-126.

Girma, S. (2000). A quasi-differencing approach to dynamic modelling from a time series of independent cross-sections. *Journal of Econometrics*, **98**, 365-383.

Heckman, J. J. and Robb, R. (1985). Alternative methods for evaluating the impact of interventions. In *Longitudinal Analysis of Labor Market Data*, (eds. J. J. Heckman and B. Singer), Cambridge University Press, New York, 156-245.

McKenzie, D. J. (2004). Asymptotic theory for heterogeneous dynamic pseudo-panels. *Journal of Econometrics*, **120**, 235-262.

Moffitt, R. (1993). Identification and estimation of dynamic models with a time series of repeated cross-sections. *Journal of Econometrics*, **59**, 99-123.

Verbeek, M. (1992). Pseudo panel data. In *The Econometrics of Panel Data*, (eds. L. Mátyás and P. Sevestre), Kluwer Academic Publishers, Dordrecht, 303-315.

Verbeek, M. and Nijman, T. (1993). Minimum MSE estimation of a regression model with fixed effects from a series of cross-sections. *Journal of Econometrics*, **59**, 125-136.

Verbeek, M. and Vella, F. (2005). Estimating dynamic models from repeated cross-sections. *Journal of Econometrics*, **127**, 83-102.