# Comments on: Panel data analysis—advantages and challenges

**Jacques Mairesse**

*I have always been struck by the old Hindu saying: There is no door but only a little window that opens upon a great world, Paul Douglas, 1948.*

Cheng Hsiao's 1986 first edition of *Analysis of Panel Data* was practically the first textbook in this very rapidly developing field of the econometric literature; it has been necessary reading and a landmark for 15 years. In spite of the publication of many other good textbooks specifically dedicated to the methods of panel data econometrics (and existing specialized chapters in most general econometric textbooks), his 2003 revised and much expanded second edition will no doubt stay an outstanding competitor and a landmark for more years to come. His present article provides a very good introductory bird-eye view to the whole book, remarkable both by its clarity and concision, and by its depth and focus on most important and specific issues of panel data econometrics. My brief comments here are from the standpoint of an "ap-

J. Mairesse (✉)
CREST, ENSAE, Malakoff, France
e-mail: Jacques.Mairesse@ensae.fr

J. Mairesse
NBER, Cambridge, USA

J. Mairesse
UNU-MERIT, Maastricht University, Maastricht, The Netherlands

plied" econometrician.[1] They should be viewed as stressing some aspects of Hsiao's presentation and adding some complementary remarks.

(1) The "random" effects and "fixed" effects distinction has been (and remains to some extent) a source of much confusion among students of panel data econometrics: "random" effects being (improperly) used as a synonym for unobserved random effects which are uncorrelated with the explanatory variables, while "fixed" effects can be correlated with them. Hsiao's presentation of the issue in the paper (and long discussion in the book) is accurate and thoughtful. It comforts the basic view of Mundlak (1978) who considers that all individual unobserved effects are random and that the basic distinction is between uncorrelated and correlated random effects. In fact, the underlying real distinction here is whether we are analyzing micro-panel data (firms, households, etc.) where usually the number of individuals $N$ is large (with a number of periods $T$ small) and the unobserved individual effects can be treated as random, or analyzing macro-panel data (countries, industries, etc.) where the number of units $N$ is small (and given). The characteristics of these two kinds of panel data are extremely different in terms of the magnitude of the dispersion and heterogeneity of the variables, and hence for their implications on model specification and estimation methods (and corresponding relevant asymptotic properties, that is, for $N$ large and $T$ fixed in the case of micro-panels, and the opposite in the case of macro-panels).[2]

My following remarks concern only usual micro-panel data (respectively "short" and "long" in the time and cross section dimensions), which I am tempted to consider as "true" panel data, without evidently denying the interest of analyzing macro-panel data as best as possible.

(2) A very common feature of simple panel data linear regression "fixed effects" estimates (or "within" group estimates) is that in practice they differ very significantly from the corresponding "random effects" estimates (and from the "between" group estimates and the overall pooled ordinary least squares estimates, so-called "total" estimates). This is not a problem, on the contrary, whenever the "within" estimates, which control for unobserved correlated individual effects, are plausible. But, more than often, unfortunately, they are not. Furthermore, they are in many cases much less plausible than the other estimates ("between," "total," and "random effects"), which do not control for the correlations of unobserved individual effects with included explanatory variables, although these correlations seem themselves a priori very plausible. Such inconsistencies between estimates, which are based on the time dimension of the data (the "within" as well as the "first differences" type estimates) and those based also on the cross section dimension (the "total," "between," and "random effects" type estimates) are not confined to linear models, but they are also frequent in estimating nonlinear panel data models.

(3) Besides unobserved individual correlated effects, reasons for the inconsistencies between the different usual types of panel data estimates are the biases arising from other possible specification errors in the models, such as typically measurement errors in variables, omitted variables (time varying and correlated with the included

---

[1]That is from the standpoint of a practicing econometrician or econometrician *stricto sensu*, econometrics being by definition a part of economic analysis and an increasingly important one, and not mainly a specialization within mathematical statistics, even if very specific, active, and successful.

[2]See also Wooldridge (2002) for a clear-cut assessment of these points.

explanatory variables), and endogeneity of variables due to simultaneity, anticipation and "feed-back" effects, or to sample selectivity. The ideal way to deal with these specification errors will be of course to be able to attack them directly, that is, for example, by improving the measurement of variables (or assessing the nature and extent of major measurement errors), obtaining information on individual unobserved effects and other important omitted variables, considering a broader "structural" simultaneous equations model while possibly imposing a priori relevant restrictions on parameters. In practice, such direct approaches appear generally much too costly and difficult, if not impossible, and are not often attempted.[3]

(4) The usual approaches to treat the potential specification errors fall mainly into two kinds. The first kind of approach, which remains quite common in actual practice, is just to ignore them more or less openly. That option may not be in fact (too) bad, since simple estimates mainly based on the cross section dimension can be more reliable and reasonable than more sophisticated and in principle appropriate estimates. The second kind of approach is to circumvent the more likely potential specification errors by resorting to various instrumental variables (IV) methods or generalized methods of moments (GMM), usually with "internal" instruments (such as lagged values of explanatory variables in first differences or in levels being used respectively as instruments on linear regressions in levels or written in first differences) and rarely with also "external" instruments justified on a priori economic grounds.[4] This route, though in theory preferable, has its own serious problems. Basically, the IV and GMM methods tend to "protect" the estimates against a mix of specification errors biases by relying on smaller and smaller fractions of the total information in the unprocessed data. Unfortunately, in doing so and in trying to solve one type of specification error, one often exacerbates strongly the seriousness of another type of specification error. In trying to solve two types of specification errors, one can magnify yet another type, and typically one may end up in getting estimates too imprecise and fragile, if not more or less implausible.

(5) Perhaps the best (and simplest) example of the difficulty to avoid possible specification errors biases is that of trying to deal with unobserved possibly correlated individual effects, in presence of even modest random measurement errors in variables. This example goes a long way in accounting for the discrepancies between the estimates respectively based (wholly or mainly) on the time dimension and the cross section dimension of the data, as well as in explaining the fact that the first type of estimates are less plausible than the second one (see Remark 2 above). In the case of the simple Cobb–Douglas (or log-linear) production function estimates discussed in Mairesse (1990) and Griliches and Mairesse (1998), the "within firm" and "first differences" transformations reduce the net variance of the log capital variable (explaining the log labor productivity variable) to about 8% and 3% respectively of its magnitude in "total levels", and they magnify the downward biases resulting from a random measurement error in the log levels of capital (if there is one, which seems

---

[3]The very influential Olley and Pakes (1996) approach to estimating a production function equation goes in this direction, by attempting to proxy for the unobserved individual effects on the basis of an additional equation (i.e., "inverting" the firm investment demand equation). See, for example, Griliches and Mairesse (1998) for an assessment of this approach.

[4]See, for example, Blundell and Bond (2000).

most likely) by a factor of about 11 and 70 respectively (!). Such downward biases explain a large part of the observed discrepancies between the three corresponding estimates of the elasticity of capital: 0.30 in "total levels," 0.15 in "within firm levels" and 0.06 in "first differences." They do not explain them fully, however, leaving room for other likely specification errors.[5],[6]

(6) Measurement errors in variables are pervasive in (micro-)panel data and they are likely to be particularly critical among various other types of specification errors. Their consequences (as illustrated by the example summarized in Remark 5) tend to be much aggravated on usual estimators, which take advantage of the panel data structure to control for unobserved individual effects. They are especially susceptible to affect panel data studies in which one wants to control for, and possibly disentangle, two types of unobserved individual effects. These studies, which represent a very promising new development in econometrics, are based on samples integrating two types of micro-panel data such as matched employer-employee data.[7] One may expect, however, that for them the difficulties arising from errors in variables in estimating parameters of interest will be as serious, if not worse, than in the typical studies controlling for only one type of individual effects. When the interest is also in estimating the two types of unobserved individual effects, say, the workers' and firms' effects, one may expect further difficulties. The separate identification of these individual effects relies basically on the observed network of workers' mobility between firms and thus will tend to be poor in situations (the usual ones) where the probability that a worker stays in the same firm is much higher than the probability that he or she moves to another firm (possibly after some period of unemployment). In such situations, the estimation of "overall" effects, defined as sums of workers' individual effects and workers' firm effects, which is based on the observations for stayers, will be more accurate than the estimation of the separate workers' and firms' effects, which is based on the much less frequent observations for movers. In particular, it is not unlikely that in these circumstances the correlation between the two types of effects would appear spuriously negative.

(7) One basic advantage of (micro-)panel data (shared with micro-cross section data) is the considerable variability in the (unprocessed) data. Thanks to such variability, the econometrician is largely relieved from worrying mostly about issues of estimation efficiency. He is allowed to be concerned instead with the important and challenging issues of model specifications, which he can also investigate with panel data much more thoroughly than with only cross section data or with only time-series. Somewhat paradoxically (as already stressed in Remarks 4 and 5), many of the simple or more sophisticated methods in panel data analysis tend to discard very much

---

[5]These results are obtained for a balanced panel data sample of some 450 French manufacturing firms for a 13 year period; they are quite similar for two comparable panel data samples of Japanese and US manufacturing firms over the same 13 year period, and with some 450 firms for the US sample and 850 firms for the Japanese sample; see also the simulations in Crépon and Mairesse (1996).

[6]Our other attempts to rely on different variants of GMM to deal with the likely mix of specification errors in estimating production functions have had mixed success, see in particular Griliches and Mairesse (1998). This is perhaps not that surprising since, as indicated in Remark 4, these different variants rely on the part of the variability in the explanatory variables accounted by the instrumental variables, which can be very small when restricted only to the time dimension of the data.

[7]See, for example, the innovative analysis of Abowd et al. (1999).

of such variability and to rely on very little of it. Otherwise, the most basic limitation in the practice of panel data econometrics lies certainly in the data themselves: their unsatisfactory quality (see Remarks 5 and 6 emphasizing errors in variables issues) or often their nonexistence (or nonavailability when existing).[8] To give just an example, econometric research on industrial organization and firm behavior suffers from a general lack of available information on prices at the firm and product market levels; prices playing an essential role in economics, this is indeed a major shortcoming. One can also say, *a contrario*, that the solutions to many of the problems encountered in panel data analysis will come in practice from progress in the data themselves— which is not to undervalue, of course, the importance of progress in our econometric techniques and know-how, as well as in our economic theories and understanding.

To close this set of remarks, I readily agree with Cheng Hsiao's own conclusion in his paper (also the very last phrase of his book): "*Although panel data offer many advantages, they are not panacea*." Yet, I would rather turn it around and proclaim, in paraphrasing the old Hindu saying given in the epigraph to these comments (and quoted from Paul Douglas' 1947 presidential address to the American Economic Association "Are there laws of production?" Douglas 1948), that, although panel data are "*only a little window that opens upon a great world*," they are nevertheless the best window in econometrics.

## References

Abowd JM, Kramarz F, Margolis D (1999) High wage workers and high wage firms. Econometrica 67:53–74

Blundell R, Bond S (2000) GMMM estimation with persistent data: an application to production functions. Econom Rev 19:321–340

Crépon B, Mairesse J (1996) Chamberlain and GMM estimates: an overview and some simulation experiments. In: Mátyás L, Sevestre P (eds) The econometrics of panel data, 2nd edn. Kluwer Academic, Boston, pp 323–391

Douglas P (1948) Are there laws of production? Am Econ Rev 38:1–41

Griliches Z (1986) Economic data issues. In: Griliches Z, Intriligator MD (eds) Handbook of econometrics, vol 3. North-Holland, Amsterdam, pp 1465–1514, Chap 25

Griliches Z, Mairesse J (1998) Production functions: the search for identification. In: Ström S (ed) Econometrics and economic theory in the 20th century: the Ragnar Frish centennial symposium. Cambridge University Press, New York, pp 169–203

Mairesse J (1990) Time-series and cross-sectional estimates on panel data: why are they different and why should they be equal? In: Hartog J, Ridder G, Theeuwes J (eds) Panel data and labor market studies. North-Holland, Amsterdam, pp 81–95

Mundlak Y (1978) On the pooling of time series and cross section data. Econometrica 46:69–85

Olley S, Pakes A (1996) The dynamics of productivity in the telecommunications equipment industry. Econometrica 64:1263–1297

Wooldridge J (2002) Econometric analysis of cross section and panel data. MIT Press, Cambridge

---

[8]This is the message that Zvi Griliches has stressed in many of his works; see, in particular, Griliches (1986).