

PANEL DATA FROM TIME SERIES OF CROSS-SECTIONS

Angus DEATON*

Woodrow Wilson School, Princeton University, Princeton, NJ 08544, USA

In many countries, there are few or no panel data, but there exists a series of independent cross-sections. For example, in the United Kingdom, there are no panel data on consumers' expenditure or on household labor supply, but there are several large household surveys carried out every year. Samples for these surveys are drawn anew each year, so that it is impossible to track individual households over time. This paper considers the possibility of tracking 'cohorts' through such data. A 'cohort' is defined as a group with fixed membership, individuals of which can be identified as they show up in the surveys. The most obvious example is an age cohort, e.g. all males born between 1945 and 1950, but there are other possibilities (Korean war veterans or founding members of the Econometric Society). Consider any economic relationship of interest that is linear in parameters (but not necessarily in variables). Corresponding to individual behavior, there will exist a cohort version of the relationship of the same form, but with cohort means replacing individual observations. If there are additive individual fixed effects, there will be corresponding additive cohort fixed effects. Further, the sample cohort means from the surveys are consistent but error-ridden estimates of the true cohort means. Hence, provided errors-in-variables techniques are used (and error variances and covariances can be estimated from the surveys), the sample cohort means can be used as panel data for estimating the relationship. Such data are immune to attrition bias and can be extended for long time periods. There is also evidence to suggest that the errors in variables problems may be just as severe for genuine panel data; in the created panels considered here, the problem is controllable. The paper discusses appropriate errors in variables estimators, with and without fixed effects.

1. Introduction

In many countries, there are few or no panel data, but there exists a series of independent cross-sections. For example, in the United Kingdom, there are no panel data on consumers' expenditure or on household labor supply, but there are several large household surveys that are carried out every year. Samples for these surveys are drawn anew each year, so that *individual* households cannot be traced over time. This paper is concerned with the possibility of tracking 'cohorts' through such data. A 'cohort' is defined as a group with fixed membership, individuals of which can be identified as they show up in the surveys. The most obvious example is an age cohort, for example, all males born between 1945 and 1950, but there are many other possibilities; consider Korean war veterans, or founder members of the Econometric Society. For

*I am grateful to Badi Baltagi, Bill Barnett, Martin Browning, Margaret Irish, Whitney Newey, and an anonymous referee for help and comments.

large enough cohorts, or large enough samples, successive surveys will generate successive random samples of individuals from each of the cohorts. Summary statistics from these random samples generate a time series that can be used to infer behavioral relationships for the cohort as a whole just as if panel data were available. Procedures for constructing such cohorts and for estimation using the resulting data are discussed in this paper.

I consider economic relationships that are linear in parameters, though not necessarily in data, and that may or may not contain individual fixed effects. Corresponding to these individual relationships, there will exist averaged versions of the same form for the cohort population, but with unobservable data points. If there are additive individual fixed effects, there will be corresponding additive cohort fixed effects for the cohort population. Furthermore, the *sample* cohort means from the surveys are consistent but error-ridden estimates of the unobservable cohort populations means. Since the micro data are used to construct the means, they can also be used to construct estimates of the variances and covariances of the sample means. It is therefore possible to use errors-in-variable estimators to estimate consistently the population relationships. Sections 3 and 4 of this paper derive appropriate estimators in the absence and in the presence of individual fixed effects.

Section 2, below, presents some of the models for which the technique is designed with particular emphasis on models of consumption and labor supply. I suggest that the estimation procedures discussed here may shed light on some long-standing puzzles in empirical demand analysis, and that they are likely to be useful for the estimation of life-cycle models in the absence of panel data.

Although the methods discussed here are primarily formulated as a response to the absence of panel data, it is not necessarily the case that they will give inferior results. The attrition problem that effectively curtails the useful length of much panel data is absent here. Because new samples are drawn each year, representativeness is constantly maintained. Indeed, to the extent that long-running panels replace respondents that drop out by 'look alike', the resulting data will have many of the features discussed here. Of course, the errors-in-variables nature of the current methodology is absent in genuine panel studies, but I suspect that the difference is more apparent than real. Survey statisticians who collect panel data are in little doubt as to the magnitude of response error, particularly in the differenced data for individual respondents. And as Ashenfelter (1983) has shown, it is extremely difficult to interpret the diversity of labor supply elasticities obtainable from the Michigan PSID data (certainly the most heavily used panel data set among American economists), without assigning a central role to large and persistent errors of measurement. The technique discussed here has the advantage of recognizing measurement error from the outset and explicitly controlling for it.

2. Model formation

The models discussed in this section are of substantive interest in their own right and form the basis for current research using the methods discussed below. The main purpose in presenting them here, however, is to motivate the methodological discussion in the following sections. Further, the provision of concrete economic examples will justify some of the specific issues of estimation that are dealt with later.

In all of this work, my ultimate aim has been to bring to bear the methodology of panel data on problems of consumer demand analysis. Even in the United States, there is little panel information on the details of consumer expenditure, and a few isolated examples apart, the same is true of the rest of the world. For cross-sectional consumer expenditure surveys, however, most countries of the world are better supplied than is the United States. For example, the British Family Expenditure Survey is in continuous operation and surveys some 7000 households annually; questions on income, labor supply, and a very detailed disaggregation of consumers' expenditure provide a mass of high-quality data. There are also many excellent series of household surveys from LDC's; India's National Sample Survey Office has run some twenty nationwide household expenditure surveys since independence, and both Indonesia and Sri Lanka have closely comparable household surveys for two or more separate years. There is therefore a large potential for any method that can 'convert' these data sets into panel data.

To illustrate some of the more important issues, consider the Engel curve model,

$$q_{iht} = f_{ih}(x_{ht}, a_{ht}), \quad (1)$$

for quantity purchased of good (or leisure) i by household h in period t , household total outlay x_{ht} , and vector of socio-economic or demographic characteristics a_{ht} . A convenient functional form is provided by taking the budget share $w_{iht} = p_{it}q_{iht}/x_{ht}$ as dependent variable and writing

$$w_{iht} = \alpha_i + \beta_i \log x_{ht} + \sum_{j=1}^J \gamma_{ij} a_{jht} + \varepsilon_{iht}, \quad (2)$$

where there are J socio-economic characteristics, α , β and γ are parameters, and ε_{iht} represents an error term. Eq. (2) is typically estimated in one of two derived forms. In the first, using a single cross-section, the t subscript is dropped, and systems of Engel curves are estimated. In the second, the equation is aggregated over h to give, for example,

$$\tilde{w}_{it} = \alpha_i + \beta_i \log \tilde{x}_t + \sum_{j=1}^J \gamma_{ij} \tilde{a}_{jt} + \tilde{\varepsilon}_{it}, \quad (3)$$

where \tilde{w}_{it} , \tilde{a}_{jt} and $\tilde{\epsilon}_{it}$ are weighted averages using $x_{ht}/\sum x_h$ as weights, and \tilde{x}_t is a representative budget level defined by the aggregation procedure [see Deaton and Muellbauer (1980) for a full discussion of the model]. The weighting procedure guarantees that \tilde{w}_{it} is the share of good i in the total of consumers expenditures, and, provided the distribution over households of x does not change over time (as measured by Theil's entropy measure of inequality), \tilde{x}_t can be replaced by \bar{x}_t . Consequently, apart from the substitution of demographics for prices, neither of which explains very much in aggregate, (3) is a conventional aggregate demand system.

The point I want to emphasize is that the values of β_i estimated from cross-sections tend to differ substantially from those estimated using time series. Such contradictions were first extensively documented by Kuznets (1962), not only for aggregate consumption, but also for the components of consumption. It is widely known that savings ratios rise more in cross-sections than in aggregate time series. It is less well-known that it is generally true that total expenditure elasticities for many commodities and groups of commodities are further dispersed from unity when estimated on cross-sections than when estimated on time series. For example, the food share in England in 1955 was almost identical to its value a century earlier in spite of a manyfold increase in real income and in spite of the repeated confirmation of Engel's law in every household survey during that century; see Deaton (1975), Stone and Rowe (1966) and Clark (1957) for further details. The presence of such phenomena also poses problems for *forecasting* demands in those situations where only cross-sectional data are available.

In terms of the foregoing model, goods are necessities if $\beta_i < 0$, luxuries if $\beta_i > 0$, and neither if $\beta_i = 0$. The Kuznets' finding, so stylized, is that β_i is closer to zero in time series than in cross-sections. Presumably, the problem lies in inadequate statistical control. Expenditure differences between poor and rich consumers are not likely to be replicated by making a poor man rich unless the poor and rich consumers are otherwise identical. Controlling for even a long list of socio-economic characteristics is not satisfactory compared with the opportunity yielded only by panel data to use individuals as their own controls. Recognizing this, write (2) as

$$w_{iht} = \alpha_i + \beta_i \log x_{ht} + \sum \gamma_{ij} a_{jht} + \theta_{ih} + \epsilon_{iht}, \quad (4)$$

for individual fixed effect θ_{ih} . Since, in general, θ_{ih} will be correlated with the other explanatory variables, such an equation can only be consistently estimated from panel data. Consider, however, the case where h is a member of a well-defined cohort group that can be tracked via its (randomly chosen) representatives through successive surveys. Let h belong to a cohort c , and

take simple *population* averages of (4) over all h belonging to c to obtain

$$w_{ict}^* = \alpha_i + \beta_i (\log x_{ct})^* + \sum \gamma_{ij} a_{jct}^* + \theta_{ic}^* + \varepsilon_{ict}^*, \quad (5)$$

where asterisks denote population (i.e., cohort population) means. If it were possible to observe the true cohort means, eq. (5) would hold for each cohort in each time period rather than for each household in each time period, and could be directly estimated using cohort dummy variables for the cohort fixed effects θ_{ic}^* . This would be feasible since each cohort appears in each time period; it is of course *infeasible* on the individual model (4) since each individual household appears only once. In practice, the other starred variables can only be proxied by cohort means from the *sample*; these will contain sampling errors and if used without appropriate correction will generally lead to inconsistent estimates since the model is effectively one of errors in variables with *all* variables (except dummies) subject to error. However, the sample can be used in the standard way to derive estimates of sampling variances and covariances and these estimates can be used to derive consistent estimators using more or less standard errors in variables procedures; see sections 3 and 4 below.

Two other features of eq. (5) should be noted. First, the total outlay variable is the mean of the logarithms, not the logarithm of the means. There is no need in this context to fudge the issue since the sample can be used just as easily to estimate the mean of a non-linear function as to estimate the non-linear function of the mean. Second, it is usually possible to select cohorts that are more or less broadly defined. Ultimately, the cohort that is all inclusive is the total population and (5) becomes a macroeconomic aggregate time-series model. In consequence, selection of cohort size allows us to move by degrees from micro to macro data; this is ideal for detecting the roots of a contradiction between micro and macro results.

In the foregoing example, the formation of cohorts can be thought of as an instrumentation procedure that removes the inconsistencies associated with the fixed effects. In my second example, the cohort structure arises naturally out of the formulation of the problem. Consider an individual household choosing consumption and labor supply in an intertemporal setting to maximize the expectation

$$E_t \left\{ \sum_{\tau=t}^T u_{\tau}(q_{\tau}) \right\}, \quad (6)$$

subject to an evolving and uncertain budget constraint

$$W_{t+1} = (1 + i_{t+1}) \{ W_t + y_t - p_t \cdot q_t \}, \quad (7)$$

where, as before, q_τ includes leisure demands, u_τ is period τ 's utility function, W_t is assets at t , y_t is income, and i_{t+1} is the money interest rate from t to $t+1$.

In Browning, Deaton and Irish (1985), it is shown that the solution to this problem can be straightforwardly characterized in terms of Frisch demand/supply functions

$$q_{it} = f_{it}(r_t, p_t), \quad (8)$$

the vector of which is the gradient with respect to p_{it} of a period t 'profit function' $\pi_t(r_t, p_t)$, a convex linearly homogeneous function. The quantity r_t , the period t price of utility, evolves stochastically according to

$$E_t \{ (1 + i_{t+1}) / r_{t+1} \} = 1 / r_t. \quad (9)$$

Once again, the discussion is more useful given a specific functional form. Browning, Deaton and Irish show that the following is consistent with the theory:

$$q_{it} = \alpha_{it} + \beta_i \log p_{it} + \sum_{j \neq i} \theta_{ij} \{ p_{jt} / p_{it} \}^{1/2} - \beta_i \log r_t. \quad (10)$$

This model is correct both under certainty and under uncertainty. In the former case, (9) holds without the expectation operator so that r_t is simply proportional to a discount factor $\Pi(1 + i_t)$ relative to some arbitrary date. Re-introducing the household subscript h , the certainty version of (10) is therefore

$$q_{iht} = \alpha_{iht} + \beta_i \log \tilde{p}_{it} + \sum_{j \neq i} \theta_{ij} (p_{jt} / p_{it})^{1/2} - \beta_i \log r_{0h}, \quad (11)$$

where r_{0h} is independent both of time and the commodity under consideration and \tilde{p}_{it} is p_{it} discounted back to the arbitrary date 0. It is therefore an individual fixed effect which is essentially a sufficient statistic for the influence of current and future values of assets, prices, interest rates, and wages; see MaCurdy (1981). Since r_{0h} is the price of life-time utility to h , it is an increasing function of life-time real wealth given concavity of (6). Consequently r_{0h} will vary with h and thus with cohorts in the cohort version of (11). Moreover, since younger cohorts are on average wealthier than older ones, we should expect the cohort dummy variables to be monotonically related to cohort age. The cohort structure here not only has the advantage of linking micro with macro, but also explicitly recognizes the life-cycle nature of consumption and labor supply. Indeed cohort methods have been widely used

in work with life-cycle models; see e.g. Ghez and Becker (1975) and Smith (1977), though these authors work with single cross-sections which lack the panel element introduced here.

Under uncertainty, (9) can be written approximately as

$$\Delta \log r_{t+1} = \log(1 + i_{t+1}) + v_{t+1}, \quad (12)$$

where $E_t(v_{t+1}) = 0$. Taking differences of (10) and substituting

$$\Delta q_{iht} = \Delta \alpha_{iht} - \beta_i \log(1 + \rho_{it}) + \sum \theta_{ij} \Delta (p_{jt}/p_{it})^{1/2} + v_{th}, \quad (13)$$

where ρ_{it} is the real commodity i rate of interest. Note that in this case, even if the shocks to the system, v_t , are stationary, $\log r_t$ will be non-stationary, so that differencing is required to obtain consistent estimates. In general, the v_t can only be guaranteed to be mean stationary, and further assumptions will be required for the consistency of techniques applied to (13). Even so, the differenced version is likely to be a better starting point for estimation than the original version in levels once uncertainty is taken into account. Once again, the aggregation to cohorts provides the repeated observations necessary for differencing, while the microdata provide the estimates of cohort means together with their sampling errors. The differenced versions will have a different measurement error structure than the levels models, and this is discussed below in sections 3 and 4. Note also that in (13), to the extent that the current prices and interest rates contain relevant new information, the innovation v_{th} will be correlated with the explanatory variables necessitating an estimator that can deal with both errors of measurement and simultaneity.

3. Estimation of models in levels

Before presenting the estimator to be discussed, consider an alternative, and perhaps more obvious approach to the estimation of an equation like (4) of section 2. To unify notation, rewrite this in standard form as

$$y_{ht} = x_{ht} \cdot \beta + \theta_h + \varepsilon_{ht}, \quad (14)$$

where the i subscript is no longer required, y_{ht} is the dependent variable for individual h at t , x_{ht} is a vector of explanatory variables, and θ_h is the fixed effect. Aggregate first over those h belonging to cohort c that happen to be observed in the survey taken at t . We then get *observed sample* cohort means which satisfy the relationship

$$\bar{y}_{ct} = \bar{x}_{ct} \cdot \bar{\beta} + \bar{\theta}_{ct} + \bar{\varepsilon}_{ct}. \quad (15)$$

Note that $\bar{\theta}_{ct}$ is the average of the fixed effects for those members of c that show up in the survey; unlike the unobserved fixed effect for the cohort population mean, θ_c , say, $\bar{\theta}_{ct}$ is *not* constant over time. Furthermore, $\bar{\theta}_{ct}$ is unobserved and, in general, is correlated with the \bar{x}_{ct} . Hence, although (15) may be useful for ‘signing’ the bias in regressing \bar{y}_{ct} on \bar{x}_{ct} , it is not an appropriate basis for consistent estimation any more than is (14), unless the cohort sample sizes are so large that $\bar{\theta}_{ct}$ is a very good approximation for θ_c . In this case, (15) can be estimated by replacing $\bar{\theta}_{ct}$ by dummy variables, one for each cohort.

Consider, instead of (15), the cohort *population* version of (14). I write this

$$y_{ct}^* = x_{ct}^* \cdot \beta + \theta_c + \varepsilon_{ct}, \quad (16)$$

where y_{ct}^* and x_{ct}^* are the unobservable cohort population means, and θ_c is the cohort fixed effect. Since the population belonging to the cohort is assumed to be fixed through time, θ_c is a constant for each c and can be replaced in (16) by cohort dummies. The y_{ct}^* and x_{ct}^* cannot be observed, but the cohort sample means \bar{y}_{ct} and \bar{x}_{ct} are error-ridden estimators, with variances that can also be estimated from the micro survey data. Eq. (16) can then be estimated by errors in variables techniques where *all* variables, except the dummies, are measured subject to error.

Eq. (16) can now be written in convenient standard form

$$y_t^* = x_t^* \cdot \beta + \varepsilon_t, \quad (17)$$

where the cohorts and surveys have been ‘stacked’ into a single index t , running from 1 to T where T is the product of the number of surveys and the number of cohorts. The cohort dummies θ_c have been absorbed into the x_t^* ’s; there is no loss of generality since the dummies can be thought of as being measured with an error that has zero mean and variance. To fix ideas, take the British Family Expenditure Survey as an illustration. Currently, there are about ten years of data available, with about 7000 observations per year. In Browning, Deaton and Irish’s (1985) work on consumption and labor supply, various selection criteria (which are always likely to be present in one form or another) reduce this to between 2500 and 3000 observations, which were formed into sixteen cohorts of five-year age bands. Hence, $T = 80$, but the cohorts, with a maximum size of 300, are not large enough for us to ignore the sampling errors in estimating y_t^* and x_t^* by \bar{y}_t and \bar{x}_t . Since, in this context (and in many others) there is a new survey every year, it is sensible to construct estimators that are consistent as T tends to infinity; with sixteen cohorts $T \rightarrow \infty$ sixteen times faster than annual, and four times as fast as quarterly time series data. The cohort size, however, is held fixed as T becomes large.

The error ε_t in (17) is assumed to be normal, independent over t , and homoskedastic; if the cohorts are very different in size, this will require that each observation be weighted by the square root of the cohort size. I shall assume this has been done as necessary. The model is completed by adding to (17) the assumed measurement structure. The cohort means, \bar{y}_t and \bar{x}_t are observed; dropping the overbars – from now on these are the basic data – I assume

$$\begin{pmatrix} \bar{y}_t \\ \bar{x}_t \end{pmatrix} \sim N \begin{pmatrix} y_t^* & \sigma_{00} & \sigma' \\ x_t^* & \sigma & \Sigma \end{pmatrix}. \quad (18)$$

Given the sampling structure, the normality does not seem to be an implausible assumption. However, the error variances σ_{ij} will in general have to be estimated by their sample counterparts s_{ij} based on the micro survey data. Note that, in estimating the σ_{ij} 's, all T observations can be pooled, so that, if there are n_c observations in each cohort, the sampling variance of s_{ij} diminishes in proportion to $(Tn_c)^{-1}$ and that of \bar{y}_t and \bar{x}_t as $(n_c)^{-1}$. The former is (a) smaller, and (b) tends to zero as $T \rightarrow \infty$ instead of remaining fixed, so that it may be reasonable to assume that the σ 's are known in carrying out the estimation. Nevertheless, I shall derive formulae for both cases, when σ_{00} , σ and Σ are known, and when they are estimated by s_{00} , s and S .

The model is now in the form in which I can apply the results of Fuller (1975, 1981); indeed, in the rest of this section, I essentially repeat Fuller's (1975) formulae in the current context and notation. The interested reader is referred to that paper for further details.

Assume that means have been removed from all data and let the sample moments and cross-product matrices of X and y be M_{xx} , m_{xy} and m_{yy} in an obvious notation. Write σ_ε^2 for the variance of ε_t and Ω for the moment matrix of the unobservable x_t^* 's. Hence

$$E(M_{xx}) = \Omega + \Sigma = \Sigma_{xx}, \quad \text{say}, \quad (19)$$

$$E(m_{xy}) = \Omega\beta + \sigma = \sigma_{xy}, \quad \text{say}, \quad (20)$$

$$E(m_{yy}) = \beta'\Omega\beta + \sigma_{00} + \sigma_\varepsilon^2 = \sigma_{yy}, \quad \text{say}. \quad (21)$$

The estimator $\tilde{\beta}$ is then clearly consistent as $T \rightarrow \infty$, where

$$\tilde{\beta} = (M_{xx} - \Sigma)^{-1}(m_{xy} - \sigma) \quad (22)$$

$$= (X'X - T\Sigma)^{-1}(X'y - T\sigma), \quad (23)$$

and, provided $(X'X - T\Sigma)$ is positive definite, will be a MLE under the normality assumptions on \mathbf{x}_t^* that are sometimes made in errors in variables models. I shall not assume normality of \mathbf{x}_t^* here. Note that if Σ and σ are replaced by estimators S and s , $\tilde{\beta}$ is replaced by

$$\beta^* = (M_{xx} - S)^{-1}(\mathbf{m}_{xy} - s) = (X'X - TS)^{-1}(X'y - Ts). \quad (24)$$

The formula for the variances are derived first for the case where Σ , σ and σ_{00} are known, i.e., for $\tilde{\beta}$. Expanding (22) yields

$$\tilde{\beta} - \beta = \Omega^{-1}(\mathbf{m}_{xy} - \sigma_{xy}) - \Omega^{-1}(M_{xx} - \Sigma_{xx})\Omega^{-1}(\sigma_{xy} - \sigma) + O_p(T^{-1}), \quad (25)$$

But, from (20), $\Omega^{-1}(\sigma_{xy} - \sigma) = \beta$, so that

$$\tilde{\beta} - \beta = \Omega^{-1}[(\mathbf{m}_{xy} - M_{xx}\beta) - (\sigma_{xy} - \Sigma_{xx}\beta)] + O_p(T^{-1}). \quad (26)$$

The variance of $\tilde{\beta}$ thus depends asymptotically on the variance-covariance matrix of $\mathbf{m}_{xy} - M_{xx}\beta$. But

$$\{\mathbf{m}_{xy} - M_{xx}\beta\}_i = \frac{1}{T} \sum_t x_{ti} \left(y_t - \sum_k x_{tk} \beta_k \right) \quad (27)$$

$$= \frac{1}{T} \sum_t (x_{ti}^* + u_{ti}) \left(\varepsilon_t + u_{t0} - \sum_k u_{tk} \beta_k \right), \quad (28)$$

where

$$u_{t0} = y_{t0} - y_{t0}^* \quad \text{and} \quad u_{ti} = x_{ti} - x_{ti}^*.$$

Treating the x_{ti}^* as fixed but unknown constants, and using the standard properties of the normal distribution yields

$$\begin{aligned} TV\{\mathbf{m}_{xy} - M_{xx}\beta\} &= \Sigma_{xx}(\sigma_\varepsilon^2 + \sigma_{00} + \beta'\Sigma\beta - 2\sigma'\beta) \\ &\quad + (\sigma - \Sigma\beta)(\sigma - \Sigma\beta)'. \end{aligned} \quad (29)$$

Hence the asymptotic variance-covariance matrix of $\tilde{\beta}$ is given by

$$TV(\tilde{\beta}) = \Omega^{-1}[\Sigma_{xx}\omega^2 + (\sigma - \Sigma\beta)(\sigma - \Sigma\beta)']\Omega^{-1}, \quad (30)$$

where

$$\omega^2 = \sigma_\varepsilon^2 + \sigma_{00} + \beta'\Sigma\beta - 2\sigma'\beta. \quad (31)$$

An estimate of (30) is straightforwardly derived from the observable moment matrices. From (19) to (21)

$$\tilde{\Omega} = M_{xx} - \Sigma, \quad (32)$$

$$\tilde{\sigma}_e^2 = m_{yy} - \sigma_{00} - \tilde{\beta}' \tilde{\Omega} \tilde{\beta}. \quad (33)$$

Substitution in (31) yields an estimate of ω^2 , i.e.,

$$\tilde{\omega}^2 = \frac{1}{T} (y - X\tilde{\beta})'(y - X\tilde{\beta}) - \frac{1}{T} e'e, \quad (34)$$

where

$$e = y - X\tilde{\beta}. \quad (35)$$

But

$$\Sigma\tilde{\beta} - \sigma = (\Sigma\tilde{\beta} + \Omega\tilde{\beta}) - (\sigma + \Omega\tilde{\beta}) \quad (36)$$

$$= M_{xx}\tilde{\beta} - \sigma_{xy} = \frac{1}{T} X'e. \quad (37)$$

Hence, the estimated variance-covariance matrix is given by

$$T\tilde{V}(\tilde{\beta}) = \tilde{\Omega}^{-1} [T^{-1}M_{xx}e'e + T^{-2}X'ee'X] \tilde{\Omega}^{-1}, \quad (38)$$

which is straightforwardly evaluated in practice.

The derivation of the variance-covariance matrix of β^* , the estimator using the *estimated* error variances, requires only minor modification of the above. I assume that the estimates s_{ij} of σ_{ij} , $i, j = 0, 1, \dots, K$, the dimension of x , are based on νT degrees of freedom. If all cohorts are pooled in estimating the s_{ij} 's, $\nu = n_c$, the number of observations per cohort, but clearly other schemes are possible. Some estimate of the variances and covariances of the s_{ij} is also required; to focus attention I shall use that derived from sampling under normality. Hence, I assume that s_{ij} is consistent for σ_{ij} , and that asymptotically,

$$(\nu T)E\{(s_{ij} - \sigma_{ij})(s_{kl} - \sigma_{kl})\} = \sigma_{ik}\sigma_{jl} + \sigma_{il}\sigma_{jk}. \quad (39)$$

The derivation proceeds as before except that the expansion (25) has an additional term corresponding to the stochastic variation in S and s . Hence,

(26) becomes

$$\begin{aligned}\beta^* - \beta &= \Omega^{-1}[(m_{xy} - M_{xx}\beta) - (\sigma_{xy} - \Sigma_{xx}\beta)] \\ &\quad - \Omega^{-1}[(s - S\beta) - (\sigma - \Sigma\beta)] + O_p(T^{-1}).\end{aligned}\quad (40)$$

By the properties of sampling under normality, the first and second terms are independent, so that, asymptotically

$$\begin{aligned}TV(\beta^*) &= \Omega^{-1}[\Sigma_{xx}\omega^2 + (\sigma - \Sigma\beta)(\sigma - \Sigma\beta)']\Omega^{-1} \\ &\quad + \nu^{-1}\Omega^{-1}V(s - S\beta)\Omega^{-1}.\end{aligned}\quad (41)$$

Elementary manipulation yields

$$V(s - S\beta) = \Sigma(\sigma_{00} - 2\sigma'\beta + \beta'\Sigma\beta) + (\sigma - \Sigma\beta)(\sigma - \Sigma\beta)'. \quad (42)$$

Note that if ν is large, (41) reduces to (30), the case of known variances; this latter is likely to be a formula that would normally be adequate in practice.

Eqs. (23), (24), (38) and (41) are the basic results of this section. I conclude with four issues of practical importance. First, the error variances σ_{00} , σ and Σ will generally vary from survey to survey and cohort to cohort. Write σ'_{00} , σ' and Σ' for the values at observation t , so that (19) and (20) become

$$E(M_{xx}) = \Omega + \bar{\Sigma}, \quad (19')$$

$$E(m_{xy}) = \Omega\beta + \bar{\sigma}, \quad (20')$$

where $\bar{\Sigma}$, $\bar{\sigma}$ are the mean values over the t observations. The analysis then goes through with σ_{00} , σ and Σ (or s_{00} , s and S in the case of estimated variances) replaced by their means. Given the nature of the variation with t , an appropriate variance-covariance matrix for the \bar{S}_{ij} 's can be derived and substituted for $V(s - S\beta)$ in (41). Second, it is necessary to allow for the presence of some x variables that are measured without errors. For example, relationships like (5) of section 2 contain cohort dummies that are clearly error-free. Other variables may not be drawn from the surveys but from other data sources; macroeconomic variables that are the same for all cohorts but vary with time (prices) are the obvious examples, and, exceptionally, there may be other relevant data on the cohorts themselves. One way to proceed is to introduce additional error-free variables to the right-hand side of (14) and to track them through the analysis. This turns out to be equivalent to the simpler (and intuitive) procedure of setting the appropriate elements of σ and rows and columns of Σ to zero; the formulae for $\tilde{\beta}$, β^* and their asymptotic

variances then remain unchanged. Third, for the reasons discussed in section 2, it may be necessary to recognize contemporaneous correlations between some of the x_t^* 's and ε_t . Instrumental variables will typically be available; in the example of section 2, in the form of lagged cohort wages or prices. On the assumption that the instrument vector w_t is constructed from a survey prior to that used for x_t so that their errors of measurement are uncorrelated, the appropriate instrumental variable estimator is

$$\tilde{\beta}_{IV} = [W'X(W'W - T\Sigma_w)^{-1}X'W]^{-1}[W'X(W'W - T\Sigma_w)^{-1}W'y], \quad (43)$$

with a variance matrix calculable by the methods given above. If measurement errors in W and X are correlated, the obvious additional corrections can be made. Fourth, and finally, note that there will typically be some flexibility in constructing cohorts. If cohorts are constructed by age, the age bands can be taken broad or narrow (e.g., a five-year window versus a one-year window), and other qualifying characteristics can be left unspecified or tightly defined. Clearly, the construction of cohorts with members that are distinct from one another and internally homogeneous will minimize the errors in variable problem and enhance precision. Beyond that, it is possible to use trial cohorts to estimate $\tilde{\beta}$ and its variance, and to use these consistent estimates to gauge the consequences of combining or separating cohorts.

4. Estimation of models in differences

In this section, I develop the estimators appropriate for the case where the model, like the second model of section 2, requires differencing prior to estimation. The previous results do not go through directly because the measurement errors induced by the sampling scheme now have a MA(1) structure relative to the unobservable true first differences.

I now write the model in the form

$$\Delta y_t^* = \Delta x_t^* \cdot \beta + \varepsilon_t, \quad (44)$$

for the true unobservable first differences. Corresponding to (18), the measurement structure is

$$\Delta y_t = \Delta y_t^* + v_{t0}, \quad (45)$$

$$\Delta x_t = \Delta x_t^* + v_t, \quad (46)$$

and

$$y_{t0} = u_{t0} - u_{t-10}, \quad (47)$$

$$v_{ti} = u_{ti} - u_{t-1i}, \quad (48)$$

with u_{t0} and u_{ti} the original measurement errors on the y_t and x_{ti} variables respectively. The relationship between Δy_t and Δx_t is therefore given by

$$\Delta y_t = \beta \cdot \Delta x_t + (\varepsilon_t + v_{t0} - \beta \cdot v_t). \quad (49)$$

In passing note a tempting but ineffective possible route to estimation. Consider $Nx_t = x_t + x_{t-1}$, the moving average, as a possible instrument for Δx_t , the first-difference, and for simplicity, assume x_t is scalar. Since $Nx_t = Nx_t^* + u_t + u_{t-1}$, and since $u_t + u_{t-1}$ is independent of $v_t = u_t - u_{t-1}$, Nx_t is orthogonal to the compound error in (49). However, in large samples Nx_t is also orthogonal to the regressor Δx_t ; $E(Nx_t \cdot \Delta x_t) = x_t^{*2} - x_{t-1}^{*2}$, so that $E(T^{-1} \Sigma Nx_t \cdot \Delta x_t) = T^{-1}(x_t^{*2} - x_0^{*2})$ with a limit as $T \rightarrow \infty$ of zero. In consequence, instrumental variables estimation of this type will not work. It is therefore necessary to follow through a scheme similar to that of section 3.

To ease notation, write $n_t \equiv \Delta y_t$ and $z_t \equiv \Delta x_t$ so that, corresponding to (19) to (21) the moment matrices are now

$$E(M_{zz}) = W + 2\Sigma = \Sigma_{zz}, \quad \text{say}, \quad (50)$$

$$E(m_{zn}) = W\beta + 2\sigma = \sigma_{zn}, \quad \text{say}, \quad (51)$$

$$E(m_{nn}) = \beta'W\beta + \sigma_\varepsilon^2 + 2\sigma_{00} = \sigma_{nn}, \quad \text{say}, \quad (52)$$

where W is the sample moment matrix of the $\Delta x_t^* \equiv z_t^*$ variables. The doubled role of measurement error comes from the moving average errors in (47) and (48). The first-difference estimator, $\tilde{\beta}_\Delta$, is immediately given as

$$\tilde{\beta}_\Delta = (M_{zz} - 2\Sigma)^{-1}(m_{zn} - 2\sigma), \quad (53)$$

or equivalently

$$\tilde{\beta}_\Delta = (Z'Z - 2T\Sigma)^{-1}(Z'n - 2T\sigma). \quad (54)$$

Expansion, as in section 3, yields

$$\tilde{\beta}_\Delta - \beta = (\Sigma_{zz} - 2\Sigma)^{-1}\{(m_{zn} - M_{zz}\beta) - (\sigma_{zn} - \Sigma_{zz}\beta)\} + O_p(T^{-1}), \quad (55)$$

so that if C is the variance–covariance matrix of $\mathbf{m}_{zn} - M_{zz}\boldsymbol{\beta}$, the asymptotic variance of $\tilde{\boldsymbol{\beta}}_\Delta$ is given by

$$TV(\tilde{\boldsymbol{\beta}}_\Delta) = (\Sigma_{zz} - 2\Sigma)^{-1}C(\Sigma_{zz} - 2\Sigma)^{-1}. \quad (56)$$

Now

$$(\mathbf{m}_{zn} - M_{zz}\boldsymbol{\beta})_i = \frac{1}{T} \sum z_{ti}(n_i - \boldsymbol{\beta} \cdot \mathbf{z}_i),$$

so that

$$(\mathbf{m}_{zn} - M_{zz}\boldsymbol{\beta})_i = \frac{1}{T} \sum (z_{ti}^* + v_{ti}) \left(\varepsilon_i + v_{i0} - \sum_k \beta_k v_{ik} \right). \quad (57)$$

The variance of this is tedious to calculate, particularly given the MA structure (47) and (48). In the appendix it is shown that

$$\begin{aligned} C = & (W + 2\Sigma)(\sigma_\varepsilon^2 + 2\sigma_A^2) - \sigma_A^2(W^+ + W^- - \Sigma) \\ & + 14(\boldsymbol{\sigma} - \Sigma\boldsymbol{\beta})(\boldsymbol{\sigma} - \Sigma\boldsymbol{\beta})', \end{aligned} \quad (58)$$

where

$$\sigma_A^2 = \sigma_{00}^2 - 2\boldsymbol{\sigma} \cdot \boldsymbol{\beta} + \boldsymbol{\beta}'\Sigma\boldsymbol{\beta}, \quad (59)$$

and

$$W^+ = \frac{1}{T} \sum \mathbf{z}_i^* \mathbf{z}_{i-1}^{*'}, \quad (60)$$

$$W^- = \frac{1}{T} \sum \mathbf{z}_{i-1}^* \mathbf{z}_i^{*'}. \quad (61)$$

The presence of W^+ and W^- reflects the autocorrelation in the measurement error. Comparing with (29), and taking the case where $W^+ = W^- = 0$, the measurement errors now play a much larger role in determining the variance. Put differently, given the same amount of variance in the true unobservables under levels and differences, estimation precision will be lower in the latter case. This result, which is not surprising, can be enhanced or modified by positive or negative autocorrelation in the true \mathbf{z}_i^* series.

The asymptotic variance covariance matrix of $\tilde{\boldsymbol{\beta}}_\Delta$ is obtained by substituting (58) into (56). In practice, an estimation formula can be obtained by noting

that, from (52),

$$\tilde{\sigma}_\varepsilon^2 = m_{nn} - 2\sigma_{00} - \tilde{\beta}'_\Delta \tilde{W} \tilde{\beta}_\Delta. \quad (62)$$

Hence, from (59),

$$\tilde{\sigma}_\varepsilon^2 + 2\sigma_A^2 = m_{nn} - 2\tilde{\beta}'_\Delta(2\sigma + \tilde{W}\tilde{\beta}_\Delta) + \tilde{\beta}'_\Delta(\tilde{W} + 2\Sigma)\tilde{\beta}_\Delta \quad (63)$$

$$= m_{nn} - 2\tilde{\beta}'_\Delta m_{zn} + \tilde{\beta}'_\Delta M_{zz}\tilde{\beta}_\Delta \quad (64)$$

$$= T^{-1}e'e, \quad (65)$$

where

$$e = n - Z \cdot \tilde{\beta}_\Delta = \Delta y - \Delta X \cdot \tilde{\beta}_\Delta. \quad (66)$$

Similarly,

$$\begin{aligned} 2\Sigma\tilde{\beta}_\Delta - 2\sigma &= 2\Sigma\tilde{\beta}_\Delta + \tilde{W}\tilde{\beta}_\Delta - \tilde{W}\tilde{\beta}_\Delta - 2\sigma \\ &= M_{zz}\tilde{\beta}_\Delta - m_{zn} = T^{-1}Z'e. \end{aligned} \quad (67)$$

Finally, therefore, writing M_{zz}^+ for $T^{-1}\Sigma z_t z_{t-1}$ and M_{zz}^- similarly, the estimated variance matrix of $\tilde{\beta}_\Delta$ is given by

$$T\hat{V}(\tilde{\beta}_\Delta) = \tilde{W}^{-1} \left[M_{zz} \frac{1}{T} e'e - \tilde{\sigma}_A^2 (M_{zz}^+ + M_{zz}^-) + \frac{7}{2T^2} Z'ee'Z \right] \tilde{W}^{-1}, \quad (68)$$

$$\tilde{\sigma}_A^2 = \sigma_{00}^2 - 2\sigma \cdot \tilde{\beta}_\Delta + \tilde{\beta}'_\Delta \Sigma \tilde{\beta}_\Delta. \quad (69)$$

The modification of these formulae for the case where σ_{00} , σ and Σ are replaced by estimates is straightforward and is left to the reader.

Appendix: Derivation of the variance for the differenced case

Starting from eq. (57), define

$$\theta_i = (m_{zn} - M_{zz}\beta)_i, \quad (A.1)$$

$$\xi_i = \varepsilon_i + v_{i0} - \sum_k \beta_k v_{ik}, \quad (A.2)$$

so that

$$\theta_i = \frac{1}{T} \sum (z_{ti}^* + v_{ti}) \xi_t. \quad (\text{A.3})$$

To obtain the variance-covariance matrix of θ_i , start from

$$E(\theta_i \theta_j) = E \left\{ T^{-2} \sum_t \sum_s (z_{ti}^* + v_{ti}) \xi_t (z_{sj}^* + v_{sj}) \xi_s \right\} \quad (\text{A.4})$$

$$= E \left\{ T^{-2} \sum_t \sum_s (z_{ti}^* z_{sj}^* \xi_t \xi_s + v_{ti} \xi_t v_{sj} \xi_s) \right\}, \quad (\text{A.5})$$

since ξ_t and v_t are jointly normal with zero third moments. By the MA(1) structure of the v 's,

$$\begin{aligned} E(\theta_i \theta_j) &= T^{-1} \left\{ E(z_{ti}^* z_{tj}^* \xi_t^2) + E(z_{ti}^* z_{t+1j}^* \xi_t \xi_{t+1}) \right. \\ &\quad + E(z_{t+1i}^* z_{tj}^* \xi_{t+1} \xi_t) + E(v_{ti} \xi_t v_{tj} \xi_t) \\ &\quad \left. + 2E(v_{ti} \xi_t v_{t+1j} \xi_{t+1}) \right\}, \end{aligned} \quad (\text{A.6})$$

Now

$$E(\xi_t^2) = \sigma_\epsilon^2 + 2\sigma_{00}^2 + 2\beta' \Sigma \beta - 4\beta' \sigma = \sigma_\epsilon^2 + 2\sigma_A^2, \quad \text{say}, \quad (\text{A.7})$$

$$E(\xi_t \xi_{t+1}) = -\sigma_{00} + 2\beta' \sigma - \beta' \Sigma \beta = -\sigma_A^2. \quad (\text{A.8})$$

Hence, evaluating (A.6) term by term using, where necessary, the standard formulae for fourth moments of normals, gives

$$E(z_{ti}^* z_{tj}^* \xi_t^2) = w_{ij} (\sigma_\epsilon^2 + 2\sigma_A^2), \quad (\text{A.9})$$

$$E(z_{ti}^* z_{t+1j}^* \xi_t \xi_{t+1}) = -w_{ij}^+ \sigma_A^2, \quad (\text{A.10})$$

with a similar expression for its transpose,

$$\begin{aligned} E(v_{ti} \xi_t v_{tj} \xi_t^2) &= 2\sigma_{ij} (\sigma_\epsilon^2 + 2\sigma_A^2) + 2 \left(2\sigma_i - 2 \sum_k \beta_k \sigma_{ik} \right) \left(2\sigma_j - 2 \sum_k \beta_k \sigma_{jk} \right) \\ &= 2\sigma_{ij} (\sigma_\epsilon^2 + 2\sigma_A^2) + 8(\sigma - \Sigma \beta)_i (\sigma - \Sigma \beta)_j. \end{aligned} \quad (\text{A.11})$$

$$E(v_{ti} \xi_t v_{t+1j} \xi_{t+1}) = \sigma_{ij} \sigma_A^2 + 5(\sigma - \Sigma \beta)_i (\sigma - \Sigma \beta)_j. \quad (\text{A.12})$$

Hence, collecting terms and subtracting $E(\theta_i)E(\theta_j) = 4(\sigma - \Sigma\beta)_i(\sigma - \Sigma\beta)_j$ yields

$$\begin{aligned} V(\theta) = & (W + 2\Sigma)(\sigma_\epsilon^2 + 2\sigma_A^2) - \sigma_A^2(W^+ + W^- - \Sigma) \\ & + 14(\sigma - \Sigma\beta)(\sigma - \Sigma\beta), \end{aligned} \quad (\text{A.13})$$

which is eq. (58) of the main text.

References

- Ashenfelter, O., 1983, Macroeconomic analyses of labor supply and microeconomic analyses of labor supply, Presented to Carnegie-Rochester Conference, Bal Harbor, FL, Nov. 1983.
- Browning, M., A. Deaton and M. Irish, 1985, A profitable approach to labor supply and commodity demands over the life-cycle, *Econometrica* 53, 503–543.
- Clark, C., 1957, The conditions of economic progress (Macmillan, London).
- Deaton, A., 1975, The structure of demand 1920–1970, in: C.M. Cipolla, ed., *The Fontana economic history of Europe* (Collins-Fontana, London).
- Deaton, A. and J. Muellbauer, 1980, An almost ideal demand system, *American Economic Review* 70, 312–326.
- Fuller, W.A., 1975, Regression analysis for sample survey, *Sankhya: The Indian Journal of Statistics* C37, 117–132.
- Fuller, W.A., 1981, Measurement error models (Department of Statistics, Iowa State University, Ames, IA).
- Ghez, G.R. and G.S. Becker, 1975, The allocation of time and goods over the life-cycle (NBER, New York).
- Kuznets, S., 1962, Quantitative aspects of the economic growth of nations: VII, The share and structure of consumption, *Economic Development and Cultural Change* 10, 1–92.
- MaCurdy, T.E., 1981, An empirical model of labor supply in a life-cycle setting, *Journal of Political Economy* 89, 1059–1085.
- Smith, J., 1977, Family labor supply over the life-cycle, *Explorations in Research* 4, 205–276.
- Stone, R. and D.A. Rowe, 1966, The measurement of consumers' expenditures and behavior in the United Kingdom 1920–1938, Vol. II (Cambridge University Press, Cambridge).