

PANEL DATA

GARY CHAMBERLAIN*

University of Wisconsin-Madison and NBER

Contents

1. Introduction and summary	1248
2. Specification and identification: Linear models	1254
2.1. A production function example	1254
2.2. Fixed effects and incidental parameters	1256
2.3. Random effects and specification analysis	1257
2.4. A consumer demand example	1259
2.5. Strict exogeneity conditional on a latent variable	1262
2.6. Lagged dependent variables	1264
2.7. Serial correlation or partial adjustment?	1267
2.8. Residual covariances: Heteroskedasticity and serial correlation	1268
2.9. Measurement error	1269
3. Specification and identification: Nonlinear models	1270
3.1. A random effects probit model	1270
3.2. A fixed effects logit model: Conditional likelihood	1274
3.3. Serial correlation and lagged dependent variables	1278
3.4. Duration models	1282
4. Inference	1285
4.1. The estimation of linear predictors	1286
4.2. Imposing restrictions: The minimum distance estimator	1288
4.3. Simultaneous equations: A generalization of three-stage least squares	1292
4.4. Asymptotic efficiency: A comparison with the quasi-maximum likelihood estimator	1294
4.5. Multivariate probit models	1296
5. Empirical applications	1299
5.1. Linear models: Union wage effects	1299
5.2. Nonlinear models: Labor force participation	1304
6. Conclusion	1311
Appendix	1311
References	1313

*I am grateful to a number of individuals for helpful discussions. Financial support was provided by the National Science Foundation (Grants No. SOC-7925959 and No. SES-8016383), by the University of Wisconsin Graduate School, and by funds granted to the Institute for Research on Poverty at the University of Wisconsin-Madison by the Department of Health, Education and Welfare pursuant to the provisions of the Economic Opportunity Act of 1964.

1. Introduction and summary

The chapter has four parts: the specification of linear models; the specification of nonlinear models; statistical inference; and empirical applications. The choice of topics is highly selective. We shall focus on a few problems and try to develop solutions in some detail.

The discussion of linear models begins with the following specification:

$$y_{it} = \beta x_{it} + c_i + u_{it}, \quad (1.1)$$

$$E(u_{it} | x_{i1}, \dots, x_{iT}, c_i) = 0 \quad (i = 1, \dots, N; t = 1, \dots, T). \quad (1.2)$$

For example, in a panel of farms observed over several years, suppose that y_{it} is a measure of the output of the i th farm in the t th season, x_{it} is a measured input that varies over time, c_i is an unmeasured, fixed input reflecting soil quality and other characteristics of the farm's location, and u_{it} reflects unmeasured inputs that vary over time such as rainfall.

Suppose that data is available on $(x_{i1}, \dots, x_{iT}, y_{i1}, \dots, y_{iT})$ for each of a large number of units, but c_i is not observed. A cross-section regression of y_{i1} on x_{i1} will give a biased estimate of β if c is correlated with x , as we would expect it to be in the production function example. Furthermore, with a single cross section, there may be no internal evidence of this bias. If $T > 1$, we can solve this problem given the assumption in (1.2). The change in y satisfies:

$$E(y_{i2} - y_{i1} | x_{i2} - x_{i1}) = \beta(x_{i2} - x_{i1}),$$

and the least squares regression of $y_{i2} - y_{i1}$ on $x_{i2} - x_{i1}$ provides a consistent estimator of β (as $N \rightarrow \infty$) if the change in x has sufficient variation. A generalization of this estimator when $T > 2$ can be obtained from a least squares regression with individual specific intercepts, as in Mundlak (1961).

The restriction in (1.2) is necessary for this result. For example, consider the following autoregressive specification:

$$y_{it} = \beta y_{i,t-1} + c_i + u_{it},$$

$$E(u_{it} | y_{i,t-1}, c_i) = 0.$$

It is clear that a regression of $y_{it} - y_{i,t-1}$ on $y_{i,t-1} - y_{i,t-2}$ will not provide a consistent estimator of β , since $u_{it} - u_{i,t-1}$ is correlated with $y_{i,t-1} - y_{i,t-2}$. Hence, it is not sufficient to assume that:

$$E(u_{it} | x_{it}, c_i) = 0.$$

Much of our discussion will be directed at testing the stronger restriction in (1.2).

Consider the (minimum mean-square error) linear predictor of c_i conditional on x_{i1}, \dots, x_{iT} :

$$E^*(c_i | x_{i1}, \dots, x_{iT}) = \eta + \lambda_1 x_{i1} + \dots + \lambda_T x_{iT}. \quad (1.3)$$

Given the assumptions that variances are finite and that the distribution of $(x_{i1}, \dots, x_{iT}, c_i)$ does not depend upon i , there are no additional restrictions in (1.3); it is simply notation for the linear predictor. Now consider the linear predictor of y_{it} given x_{i1}, \dots, x_{iT} :

$$E^*(y_{it} | x_{i1}, \dots, x_{iT}) = \zeta_t + \pi_{t1} x_{i1} + \dots + \pi_{tT} x_{iT}.$$

Form the $T \times T$ matrix Π with π_{ts} as the (t, s) element. Then the restriction in (1.2) implies that Π has a distinctive structure:

$$\Pi = \beta I + I\lambda',$$

where I is the $T \times T$ identity matrix, I is a $T \times 1$ vector of ones, and $\lambda' = (\lambda_1, \dots, \lambda_T)$. A test for this structure could usefully accompany estimators of β based on change regressions or on regressions with individual specific intercepts. Moreover, this formulation suggests an alternative estimator for β , which is developed in the inference section.

This test is an exogeneity test and it is useful to relate it to Granger (1969) and Sims (1972) causality. The novel feature is that we are testing for noncausality conditional on a latent variable. Suppose that $t=1$ is the first period of the individual's (economic) life. Within the linear predictor context, a Granger definition of "y does not cause x conditional on a latent variable c" is:

$$E^*(x_{i,t+1} | x_{i1}, \dots, x_{it}, y_{i1}, \dots, y_{it}, c_i) = E^*(x_{i,t+1} | x_{i1}, \dots, x_{it}, c_i) \quad (t=1, 2, \dots).$$

A Sims definition is:

$$E^*(y_{it} | x_{i1}, x_{i2}, \dots, c_i) = E^*(y_{it} | x_{i1}, \dots, x_{it}, c_i) \quad (t=1, 2, \dots).$$

In fact, these two definitions imply identical restrictions on the covariance matrix of $(x_{i1}, \dots, x_{iT}, y_{i1}, \dots, y_{iT})$. The Sims form fits directly into the Π matrix framework and implies the following restrictions:

$$\Pi = B + \gamma\lambda', \quad (1.4)$$

where B is a lower triangular matrix and γ is a $T \times 1$ vector. We show how these nonlinear restrictions can be transformed into linear restrictions on a standard simultaneous equations model.

A Π matrix in the form (1.4) occurs in the autoregressive model of Balestra and Nerlove (1966). The $\gamma\lambda'$ term is generated by the projection of the initial condition onto the x 's. We also consider autoregressive models in which a time-invariant omitted variable is correlated with the x 's.

The methods we shall discuss rely on the measured x_{it} changing over time whereas the unmeasured c_i is time invariant. It seems plausible to me that panel data should be useful in separating the effects of x_{it} and c_i in this case. An important limitation, however, is that measured, time-invariant variables (z_i) can be absorbed into c_i . Their effects are not identified without further restrictions that distinguish them from c_i . Some solutions to this problem are discussed in Chamberlain (1978) and in Hausman and Taylor (1981).

In Section 3 we use a multivariate probit model to illustrate the new issues that arise in models that are nonlinear in the variables. Consider the following specification:

$$\begin{aligned}\tilde{y}_{it} &= \beta x_{it} + c_i + u_{it}, \\ y_{it} &= 1, \quad \text{if } \tilde{y}_{it} \geq 0, \\ &= 0, \quad \text{otherwise} \quad (i = 1, \dots, N; t = 1, \dots, T),\end{aligned}$$

where, conditional on $x_{i1}, \dots, x_{iT}, c_i$, the distribution of (u_{i1}, \dots, u_{iT}) is multivariate normal $(N(\mathbf{0}, \Sigma))$ with mean $\mathbf{0}$ and covariance matrix $\Sigma = (\sigma_{jk})$. We observe $(x_{i1}, \dots, x_{iT}, y_{i1}, \dots, y_{iT})$ for a large number of individuals, but we do not observe c_i . For example, in the reduced form of a labor force participation model, y_{it} can indicate whether or not the i th individual worked during period t , x_{it} can be a measure of the presence of young children, and c_i can capture unmeasured characteristics of the individual that are stable at least over the sample period. In the certainty model of Heckman and MaCurdy (1980), c_i is generated by the single life-time budget constraint.

If we treat the c_i as parameters to be estimated, then there is a severe incidental parameter problem. The consistency of the maximum likelihood estimator requires that $T \rightarrow \infty$, but we want to do asymptotic inference with $N \rightarrow \infty$ for fixed T , which reflects the sample sizes in the panel data sets we are most interested in. So we consider a random effects estimator, which is based on the following specification for the distribution of c conditional on x :

$$c_i = \eta + \lambda_1 x_{i1} + \dots + \lambda_T x_{iT} + v_i, \quad (1.5)$$

where the distribution of v_i conditional on x_{i1}, \dots, x_{iT} is $N(0, \sigma_v^2)$. This is similar to our specification in (1.3) for the linear model, but there is an important difference; (1.3) was just notation for the linear predictor, whereas (1.5) embodies substantive restrictions. We are assuming that the regression function of c on the

x 's is linear and that the residual variation is homoskedastic and normal. Given these assumptions, our analysis runs parallel to the linear case. There is a matrix Π of multivariate probit coefficients which has the following structure:

$$\Pi = \text{diag}\{\alpha_1, \dots, \alpha_T\}[\beta I + \lambda \lambda'],$$

where $\text{diag}\{\alpha_1, \dots, \alpha_T\}$ is a diagonal matrix of normalization factors with $\alpha_t = (\sigma_{it} + \sigma_v^2)^{-1/2}$. We can impose these restrictions to obtain an estimator of $\alpha_t \beta$ which is consistent as $N \rightarrow \infty$ for fixed T . We can also test whether Π in fact has this structure.

A quite different treatment of the incidental parameter problem is possible with a logit functional form for $P(y_{it} = 1 | x_{it}, c_i)$. The sum $\sum_{t=1}^T y_{it}$ provides a sufficient statistic for c_i . Hence we can use the distribution of y_{i1}, \dots, y_{iT} conditional on x_{i1}, \dots, x_{iT} , $\sum_t y_{it}$ to obtain a conditional likelihood function that does not depend upon c_i . Maximizing it with respect to β provides an estimator that is consistent as $N \rightarrow \infty$ for fixed T , and the other standard properties for maximum likelihood hold as well. The power of the procedure is that it places no restrictions on the conditional distribution of c given x . It is perhaps the closest analog to the change regression in the linear model. A shortcoming is that the residual covariance matrix is constrained to be equicorrelated. Just as in the probit model, a key assumption is:

$$P(y_{it} = 1 | x_{i1}, \dots, x_{iT}, c_i) = P(y_{it} = 1 | x_{it}, c_i), \quad (1.6)$$

and we discuss how it can be tested.

It is natural to ask whether (1.6) is testable without imposing the various functional form restrictions that underlie our tests in the probit and logit cases. First, some definitions. Suppose that $t=1$ is the initial period of the individual's (economic) life; an extension of Sims' condition for x to be strictly exogenous is that y_t is independent of x_{t+1}, x_{t+2}, \dots conditional on x_1, \dots, x_t . An extension of Granger's condition for "y does not cause x" is that x_{t+1} is independent of y_1, \dots, y_t conditional on x_1, \dots, x_t . Unlike the linear predictor case, now strict exogeneity is weaker than noncausality. Noncausality requires that y_t be independent of x_{t+1}, x_{t+2}, \dots conditional on x_1, \dots, x_t and on y_1, \dots, y_{t-1} . If x is strictly exogenous and in addition y_t is independent of x_1, \dots, x_{t-1} conditional on x_t , then we shall say that the relationship of x to y is static.

Then our question is whether it is restrictive to assert that there exists a latent variable c such that the relationship of x to y is static conditional on c . We know that this is restrictive in the linear predictor case, since the weaker condition that x be strictly exogenous conditional on c is restrictive. Unfortunately, there are no restrictions when we replace zero partial correlation by conditional independence. It follows that conditional strict exogeneity is restrictive only when combined with specific functional forms—a truly nonparametric test cannot exist.

Section 4 presents our framework for inference. Let $\mathbf{r}_i' = (1, x_{i1}, \dots, x_{iT}, y_{i1}, \dots, y_{iT})$ and assume that \mathbf{r}_i is independent and identically distributed (i.i.d.) for $i = 1, 2, \dots$. Let \mathbf{w}_i be the vector formed from the squares and cross-products of the elements in \mathbf{r}_i . Our framework is based on a simple observation: the matrix Π of linear predictor coefficients is a function of $E(\mathbf{w}_i)$; if \mathbf{r}_i is i.i.d. then so is \mathbf{w}_i ; hence our problem is to make inferences about a function of a population mean under random sampling. This is straightforward and provides an asymptotic distribution theory for least squares that does not require a linear regression function or homoskedasticity.

Stack the columns of Π' into a vector $\boldsymbol{\pi}$ and let $\boldsymbol{\pi} = \mathbf{h}(\boldsymbol{\mu})$, where $\boldsymbol{\mu} = E(\mathbf{w}_i)$. Then the limiting distribution for least squares is normal with covariance matrix:

$$\boldsymbol{\Omega} = \frac{\partial \mathbf{h}}{\partial \boldsymbol{\mu}'} V(\mathbf{w}_i) \frac{\partial \mathbf{h}'}{\partial \boldsymbol{\mu}}.$$

We impose restrictions on Π by using a minimum distance estimator. The restrictions can be expressed as $\boldsymbol{\mu} = \mathbf{g}(\boldsymbol{\theta})$, where $\boldsymbol{\theta}$ is free to vary within some set \mathcal{T} . Given the sample mean $\bar{\mathbf{w}} = \sum_{i=1}^N \mathbf{w}_i / N$, we choose $\hat{\boldsymbol{\theta}}$ to minimize the distance between $\bar{\mathbf{w}}$ and $\mathbf{g}(\boldsymbol{\theta})$, using the following distance function:

$$\min_{\boldsymbol{\theta} \in \mathcal{T}} [\bar{\mathbf{w}} - \mathbf{g}(\boldsymbol{\theta})]' \hat{V}^{-1}(\mathbf{w}_i) [\bar{\mathbf{w}} - \mathbf{g}(\boldsymbol{\theta})],$$

where $\hat{V}(\mathbf{w}_i)$ is a consistent estimator of $V(\mathbf{w}_i)$. This is a generalized least squares estimator for a multivariate regression model with nonlinear restrictions on the parameters; the only explanatory variable is a constant term. The limiting distribution of $\hat{\boldsymbol{\theta}}$ is normal with covariance matrix:

$$\left[\frac{\partial \mathbf{g}'}{\partial \boldsymbol{\theta}} V^{-1}(\mathbf{w}_i) \frac{\partial \mathbf{g}}{\partial \boldsymbol{\theta}'} \right]^{-1}.$$

An asymptotic distribution theory is also available when we use some matrix other than $\hat{V}^{-1}(\mathbf{w}_i)$ in the distance function. This theory shows that $\hat{V}^{-1}(\mathbf{w}_i)$ is the optimal choice. However, by using suboptimal norms, we can place a number of commonly used estimators within this framework.

The results on efficient estimation have some surprising consequences. The simplest example is a univariate linear predictor: $E^*(y_i | x_{i1}, x_{i2}) = \pi_0 + \pi_1 x_{i1} + \pi_2 x_{i2}$. Consider imposing the restriction that $\pi_2 = 0$; we do not want to maintain any other restrictions, such as linear regression, homoskedasticity, or normality. How shall we estimate π_1 ? Let $\hat{\boldsymbol{\pi}}' = (\hat{\pi}_1, \hat{\pi}_2)$ be the estimator obtained from the least squares regression of y on x_1, x_2 . We want to find a vector of the form $(\boldsymbol{\theta}, 0)$ as close as possible to $(\hat{\pi}_1, \hat{\pi}_2)$, using $\hat{V}^{-1}(\hat{\boldsymbol{\pi}})$ in the distance function. Since we are not using the conventional estimator of $V(\hat{\boldsymbol{\pi}})$, the answer to this minimization

problem is not, in general, to set $\hat{\theta} = b_{yx_1}$, the estimator obtained from the least squares regression of y on x_1 . We can do better by using $b_{yx_1} + \tau\hat{\pi}_2$; the asymptotic mean of $\hat{\pi}_2$ is zero if $\pi_2 = 0$, and if b_{yx_1} and $\hat{\pi}_2$ are correlated, then we can choose τ to reduce the asymptotic variance below that of b_{yx_1} .

This point has a direct counterpart in the estimation of simultaneous equations. The restrictions on the reduced form can be imposed using a minimum distance estimator. This is more efficient than conventional estimators since it is using the optimal norm. In addition, there are generalizations of two- and three-stage least squares that achieve this efficiency gain at lower computational cost.

A related application is to the estimation of restricted covariance matrices. Here the assumption to be relaxed is multivariate normality. We show that the conventional maximum likelihood estimator, which assumes normality, is asymptotically equivalent to a minimum distance estimator. But that minimum distance estimator is not, in general, using the optimal norm. Hence, there is a feasible minimum distance estimator that is at least as good as the maximum likelihood estimator; it is strictly better in general for non-normal distributions.

The minimum distance approach has an application to the multivariate probit model of Section 3. We begin by estimating T separate probit specifications in which all leads and lags of x are included in the specification for each y_{it} :

$$P(y_{it} = 1 | x_{i1}, \dots, x_{iT}) = F(\pi_{i0} + \pi_{i1}x_{i1} + \dots + \pi_{iT}x_{iT}),$$

where F is the standard normal distribution function. Each of the T probit specifications is estimated using a maximum likelihood program for univariate probit analysis. There is some sacrifice of efficiency here, but it may be outweighed by the advantage of avoiding numerical integration. Given the estimator for Π , we derive its asymptotic covariance matrix and then impose and test restrictions by using the minimum distance estimator.

Section 5 presents two empirical applications, which implement the specifications discussed in Sections 2 and 3 using the inference procedures from Section 4. The linear example is based on the panel of Young Men in the National Longitudinal Survey (Parnes); y_i is the logarithm of the individual's hourly wage and x_i includes variables to indicate whether or not the individual's wage is set by collective bargaining; whether or not he lives in an SMSA; and whether or not he lives in the South. We present unrestricted least squares regressions of y_i on x_1, \dots, x_T , and we examine the form of the Π matrix. There are significant leads and lags, but there is evidence in favor of a static relationship conditional on a latent variable; the leads and lags could be interpreted as just due to c , with $E(y_i | x_1, \dots, x_T, c) = \beta x_i + c$. The estimates of β that control for c are smaller in absolute value than the cross-section estimates. The union coefficient declines by 40%, with somewhat larger declines for the SMSA and region coefficients.

The second application presents estimates of a model of labor force participation. It is based on a sample of married women in the Michigan Panel Study of

Income Dynamics. We focus on the relationship between participation and the presence of young children. The unrestricted Π matrix for the probit specification has significant leads and lags; but, unlike the wage example, there is evidence here that the leads and lags are not generated just by a latent variable. If we do impose this restriction, then the resulting estimator of β indicates that the cross-section estimates overstate the negative effect of young children on the woman's participation probability.

The estimates for the logit functional form present some interesting contrasts to the probit results. The cross-section estimates, as usual, are in close agreement with the probit estimates. But when we use the conditional maximum likelihood estimator to control for c , the effect of an additional young child on participation becomes substantially more negative than in the cross-section estimates; so the estimated sign of the bias is opposite to that of the probit results. Here the estimation method is having a first order effect on the results. There are a variety of possible explanations. It may be that the unrestricted distribution for c in the logit form is the key. Or, since there is evidence against the restriction that:

$$P(y_{it}|x_{i1}, \dots, x_{iT}, c_i) = P(y_{it}|x_{it}, c_i),$$

perhaps we are finding that imposing this restriction simply leads to different biases in the probit and logit estimates.

2. Specification and identification: Linear models

2.1. A production function example

We shall begin with a production function example, due to Mundlak (1961).¹ Suppose that a farmer is producing a product with a Cobb–Douglas technology:

$$y_{it} = \beta x_{it} + c_i + u_{it} \quad (0 < \beta < 1; i = 1, \dots, N; t = 1, \dots, T),$$

where y_{it} is the logarithm of output on the i th farm in season t , x_{it} is the logarithm of a variable input (labor), c_i represents an input that is fixed over time (soil quality), and u_{it} represents a stochastic input (rainfall), which is not under the farmer's control. We shall assume that the farmer knows the product price (P) and the input price (W), which do not depend on his decisions, and that he knows c_i . The factor input decision, however, is made before knowing u_{it} , and we shall assume that x_{it} is chosen to maximize expected profits. Then the factor demand equation is:

$$x_i = \{ \ln \beta + \ln [E(e^{u_i} | \mathcal{J}_i)] + \ln (P_i / W_i) + c \} / (1 - \beta), \quad (2.1)$$

¹ This example is also discussed in Mundlak (1963) and in Zellner, Kmenta, and Dr ze (1966).

where \mathcal{J}_i is the information set available to the farmer when he chooses x_i , and we have suppressed the i subscript.

Assume first that u_i is independent of \mathcal{J}_i , so that the farmer cannot do better than using the unconditional mean. In that case we have:

$$E(y_i|x_1, \dots, x_T, c) = \beta x_i + c.$$

So if c is observed, only one period of data is needed; the least squares regression of y_1 on x_1 , c provides a consistent estimator of β as $N \rightarrow \infty$.

Now suppose that c is not observed by the econometrician, although it is known to the farmer. Consider the least squares regression of y_1 on x_1 , using just a single cross-section of the data. The population counterpart is:

$$E^*(y_1|x_1) = \pi_0 + \pi x_1,$$

where E^* is the minimum mean-square error linear predictor (the wide-sense regression function):

$$\pi = \text{cov}(y_1, x_1)/V(x_1), \quad \pi_0 = E(y_1) - \pi E(x_1).$$

We see from (2.1) that c and x_1 are correlated; hence $\pi \neq \beta$ and the least squares estimator of β does not converge to β as $N \rightarrow \infty$. Furthermore, with a single cross section, there may be no internal evidence of this omitted-variable bias.

Now the panel can help to solve this problem. Mundlak's solution was to include farm specific indicator variables: a least squares regression of y_{it} on x_{it} , d_{it} ($i=1, \dots, N$; $t=1, \dots, T$), where d_{it} is an $N \times 1$ vector of zeros except for a one in the i th position. So this solution treats the c_i as a set of parameters to be estimated. It is a "fixed effects" solution, which we shall contrast with "random effects". The distinction is that under a fixed effects approach, we condition on the c_i , so that their distribution plays no role. A random effects approach invokes a distribution for c . In a Bayesian framework, β and the c_i would be treated symmetrically, with a prior distribution for both. Since I am only going to use asymptotic results on inference, however, a "gentle" prior distribution for β will be dominated. That this need not be true for the c_i is one of the interesting aspects of our problem.

We shall do asymptotic inference as N tends to infinity for fixed T . Since the number of parameters (c_i) is increasing with sample size, there is a potential "incidental parameters" problem in the fixed effects approach. This does not, however, pose a deep problem in our example. The least squares regression with the indicator variables is algebraically equivalent to the least squares regression of $y_{it} - \bar{y}_i$ on $x_{it} - \bar{x}_i$ ($i=1, \dots, N$; $t=1, \dots, T$), where $\bar{y}_i = \sum_{t=1}^T y_{it}/T$, $\bar{x}_i =$

$\sum_{t=1}^T x_{it}/T$. If $T=2$, this reduces to a least squares regression of $y_{i2} - y_{i1}$ on $x_{i2} - x_{i1}$. Since

$$E(y_{i2} - y_{i1} | x_{i2} - x_{i1}) = \beta(x_{i2} - x_{i1}),$$

the least squares regression will provide a consistent estimator of β if there is sufficient variation in $x_{i2} - x_{i1}$.²

2.2. Fixed effects and incidental parameters

The incidental parameters can create real difficulties. Suppose that u_{it} is independently and identically distributed (i.i.d.) across farms and periods with $V(u_{it}) = \sigma^2$. Then under a normality assumption, the maximum likelihood estimator of σ^2 converges (almost surely) to $\sigma^2(T-1)/T$ as $N \rightarrow \infty$ with T fixed.³ The failure to correct for degrees of freedom leads to a serious inconsistency when T is small. For another example, consider the following autoregression:

$$y_{i1} = \beta y_{i0} + c_i + u_{i1},$$

$$y_{i2} = \beta y_{i1} + c_i + u_{i2}.$$

Assume that u_{i1} and u_{i2} are i.i.d. conditional on y_{i0} and c_i , and that they follow a normal distribution ($N(0, \sigma^2)$). Consider the likelihood function corresponding to the distribution of (y_{i1}, y_{i2}) conditional on y_{i0} and c_i . The log-likelihood function is quadratic in β, c_1, \dots, c_N (given σ^2), and the maximum likelihood estimator of β is obtained from the least squares regression of $y_{i2} - y_{i1}$ on $y_{i1} - y_{i0}$ ($i=1, \dots, N$). Since u_{i1} is correlated with y_{i1} , and

$$y_{i2} - y_{i1} = \beta(y_{i1} - y_{i0}) + u_{i2} - u_{i1},$$

it is clear that

$$E(y_{i2} - y_{i1} | y_{i1} - y_{i0}) \neq \beta(y_{i1} - y_{i0}),$$

and the maximum likelihood estimator of β is not consistent. If the distribution of y_{i0} conditional on c_i does not depend on β or c_i , then the likelihood function based on the distribution of (y_{i0}, y_{i1}, y_{i2}) conditional on c_i gives the same inconsistent maximum likelihood estimator of β . If the distribution of (y_{i0}, y_{i1}, y_{i2})

²We shall not discuss methods for eliminating omitted-variable bias when x does not vary over time ($x_{it} = x_i$). See Chamberlain (1978) and Hausman and Taylor (1981).

³This example is discussed in Neyman and Scott (1948).

is stationary, then the estimator obtained from the least squares regression of $y_{i2} - y_{i1}$ on $y_{i1} - y_{i0}$ converges, as $N \rightarrow \infty$, to $(\beta - 1)/2$.⁴

2.3. Random effects and specification analysis

We have seen that the success of the fixed effects estimator in the production function example must be viewed with some caution. The incidental parameter problem will be even more serious when we consider nonlinear models. So we shall consider next a random effects treatment of the production function example; this will also provide a convenient framework for specification analysis.⁵

Assume that there is some joint distribution for $(x_{i1}, \dots, x_{iT}, c_i)$, which does not depend upon i , and consider the regression function that does not condition on c :

$$E(y_{it}|x_{i1}, \dots, x_{iT}) = \beta x_{it} + E(c_i|x_{i1}, \dots, x_{iT}).$$

The regression function for c_i given $\mathbf{x}_i = (x_{i1}, \dots, x_{iT})$ will generally be some nonlinear function. But we can specify a minimum mean-square error linear predictor:⁶

$$E^*(c_i|x_{i1}, \dots, x_{iT}) = \psi + \lambda_1 x_{i1} + \dots + \lambda_T x_{iT} = \psi + \boldsymbol{\lambda}' \mathbf{x}_i, \quad (2.2)$$

where $\boldsymbol{\lambda} = V^{-1}(\mathbf{x}_i) \text{cov}(\mathbf{x}_i, c_i)$. No restrictions are being imposed here—(2.2) is simply giving our notation for the linear predictor.

Now we have:

$$E^*(y_{it}|\mathbf{x}_i) = \psi + \beta x_{it} + \boldsymbol{\lambda}' \mathbf{x}_i.$$

Combining these linear predictors for the T periods gives the following multivariate linear predictor:⁷

$$\begin{aligned} E^*(\mathbf{y}_i|\mathbf{x}_i) &= \boldsymbol{\pi}_0 + \boldsymbol{\Pi} \mathbf{x}_i, \\ \boldsymbol{\Pi} &= \text{cov}(\mathbf{y}_i, \mathbf{x}_i') V^{-1}(\mathbf{x}_i) = \beta \mathbf{I} + \mathbf{I} \boldsymbol{\lambda}', \end{aligned} \quad (2.3)$$

where $\mathbf{y}_i' = (y_{i1}, \dots, y_{iT})$, \mathbf{I} is the $T \times T$ identity matrix, and \mathbf{I} is a $T \times 1$ vector of ones.

⁴See Chamberlain (1980) and Nickell (1981).

⁵In our notation, Kiefer and Wolfowitz (1956) invoke a distribution for c to pass from the distribution of (y, x) conditional on c to the marginal distribution of (y, x) . Note that they did not assume a parametric form for the distribution of c .

⁶Mundlak (1978) uses a similar specification, but with $\lambda_1 = \dots = \lambda_T$. The appropriateness of these equality constraints is discussed in Chamberlain (1980, 1982a).

⁷We shall not discuss the problems caused by attrition. See Griliches, Hall and Hausman (1978) and Hausman and Wise (1979).

The Π matrix is a useful tool for analyzing this model. Consider first the estimation of β ; if $T = 2$ we have:

$$\Pi = (\pi_{jk}) = \begin{pmatrix} \beta + \lambda_1 & \lambda_2 \\ \lambda_1 & \beta + \lambda_2 \end{pmatrix}.$$

Hence,

$$\beta = \pi_{11} - \pi_{21} = \pi_{22} - \pi_{12}.$$

So given a consistent estimator for Π , we can obtain a consistent estimator for β . The estimation of Π is almost a standard problem in multivariate regression; but, due to the nonlinearity in $E(c_i|x_i)$, we are estimating only a wide-sense regression function, and some care is needed. It turns out that there is a way of looking at the problem which allows a straightforward treatment, under very weak assumptions. We shall develop this in the section on inference.

We see in (2.3) that there are restrictions on the Π matrix. The off-diagonal elements within the same column of Π are all equal. The T^2 elements of Π are functions of the $T + 1$ parameters $\beta, \lambda_1, \dots, \lambda_T$. This suggests an obvious specification test. Or, backing up a bit, we could begin with the specification that $\Pi = \beta I$. Then passing to (2.3) would be a test for whether there is a time-invariant omitted variable that is correlated with the x 's. The test of $\Pi = \beta I + I\lambda'$ against an unrestricted Π would be an omnibus test of a variety of misspecifications, some of which will be considered next.⁸

Suppose that there is serial correlation in u , with $u_t = \rho u_{t-1} + w_t$, where w_t is independent of \mathcal{J}_t and we have suppressed the i subscripts. Now we have:

$$E(e^{u_t}|\mathcal{J}_t) = e^{\rho u_{t-1}} E(e^{w_t}).$$

So the factor demand equation becomes:

$$x_t = \{ \ln \beta + \ln[E(e^{w_t})] + \ln(P_t/W_t) + \rho u_{t-1} + c \} / (1 - \beta).$$

Suppose that there is no variation in prices across the farms, so that the P_t/W_t term is captured in period specific intercepts, which we shall suppress. We can solve for u_t in terms of x_{t+1} and c , and substitute this solution into the y_t equation. Then we have:

$$E^*(y_t|x_1, \dots, x_T) = \beta x_t + (1 - \rho^{-1})(\lambda_1 x_1 + \dots + \lambda_T x_T) + \varphi x_{t+1},$$

⁸This specification test was proposed in Chamberlain (1978a, 1979). The restrictions are similar to those in the MIMIC model of Jöreskog and Goldberger (1975); also see Goldberger (1974a), Griliches (1974), Jöreskog and Sörbom (1977), Chamberlain (1977), and Jöreskog (1978). There are also connections with the work on sibling data, which is surveyed in Griliches (1979).

where $\varphi = \rho^{-1}(1 - \beta)$. So the Π matrix would indicate a distributed lead, even after controlling for c . If instead there is a first order moving average, $u_t = w_t + \rho w_{t-1}$, then:

$$E(e^{u_t} | \mathcal{J}_t) = e^{\rho w_{t-1}} E(e^{w_t}),$$

and a bit of algebra gives:

$$E(y_t | x_1, \dots, x_T) = x_t - \rho^{-1}(\lambda_1 x_1 + \dots + \lambda_T x_T) + \varphi x_{t+1}.$$

Once again there is a distributed lead, but now β is not identified from the Π matrix.

2.4. A consumer demand example

2.4.1. Certainty

We shall follow Ghez and Becker (1975), Heckman and MaCurdy (1980), and MaCurdy (1981) in presenting a life-cycle model under certainty. Suppose that the consumer is maximizing

$$V = \sum_{t=1}^{\tau} \rho^{(t-1)} U_t(C_t)$$

subject to

$$\sum_{t=1}^{\tau} \gamma^{-(t-1)} P_t C_t \leq B, \quad C_t \geq 0 \quad (t=1, \dots, \tau),$$

where $\rho^{-1} - 1$ is the rate of time preference, $\gamma - 1$ is the (nominal) interest rate, C_t is consumption in period t , P_t is the price of the consumption good in period t , and B is the present value in the initial period of lifetime income. In this certainty model, the consumer faces a single lifetime budget constraint.

If the optimal consumption is positive in every period, then

$$U'_t(C_t) = (\gamma\rho)^{-(t-1)} (P_t/P_1) U'_1(C_1).$$

A convenient functional form is $U_t(C) = A_t C^\delta / \delta$ ($A_t > 0$, $\delta < 1$); then we have:

$$y_t = \beta x_t + \varphi(t-1) + c + u_t, \quad (2.4)$$

where $y_t = \ln C_t$, $x_t = \ln P_t$, $c = (\delta - 1)^{-1} \ln[U_1'(C_1)/P_1]$, $u_t = (1 - \delta)^{-1} \ln A_t$, $\beta = (\delta - 1)^{-1}$, and $\varphi = (1 - \delta)^{-1} \ln(\gamma\rho)$. Note that c is determined by the marginal utility of initial wealth: $U_1'(C_1)/P_1 = \partial V/\partial B$.

We shall assume that A_t is not observed by the econometrician, and that it is independent of the P 's. Then the model is similar to the production function example if there is price variation across consumers as well as over time. There will generally be correlation between c and (x_1, \dots, x_T) . As before we have the prediction that $\Pi = \beta I + I\lambda'$, which is testable. A consistent estimator of β can be obtained with only two periods of data since

$$y_t - y_{t-1} = \beta(x_t - x_{t-1}) + \varphi + u_t - u_{t-1}.$$

We shall see next how these results are affected when we allow for some uncertainty.

2.4.2. Uncertainty

We shall present a highly simplified model in order to obtain some explicit results in the uncertainty case. The consumer is maximizing

$$E \left[\sum_{t=1}^{\tau} \rho^{t-1} U_t(C_t) \right]$$

subject to

$$\begin{aligned} P_1 C_1 + S_1 &\leq B, \\ P_t C_t + S_t &\leq \gamma S_{t-1}, \quad C_t \geq 0, \quad S_t \geq 0 \quad (t=1, \dots, \tau). \end{aligned}$$

The only source of uncertainty is the future prices. The consumer is allowed to borrow against his future income, which has a present value of B in the initial period. The consumption plan must have C_t a function only of information available at date t .

It is convenient to set $\tau = \infty$ and to assume that P_{t+1}/P_t is i.i.d. ($t=1, 2, \dots$). If $U_t(C) = A_t C^\delta / \delta$, then we have the following optimal plan:⁹

$$\begin{aligned} C_1 &= d_1 B / P_1, \quad S_1 = (1 - d_1) B, \\ C_t &= d_t \gamma S_{t-1} / P_t, \quad S_t = (1 - d_t) \gamma S_{t-1} \quad (t=2, 3, \dots), \end{aligned} \tag{2.5}$$

⁹We require $\rho\kappa g' < 1$, where $A_t \leq M g'$ for some constant M . Phelps (1962) obtained explicit solutions for models of this type. The derivation of (2.5) can be obtained by following Levhari and Srinivasan (1969) or Dynkin and Yushkevich (1979, Ch. 6.9).

where

$$d_t = [1 + f_{t+1} + (f_{t+1}f_{t+2}) + \dots]^{-1},$$

$$f_t = (\rho\kappa A_t/A_{t-1})^{[1/(1-\delta)]}, \quad \kappa = \gamma^\delta E[(P_{t-1}/P_t)^\delta].$$

It follows that:

$$y_t - y_{t-1} = (-1)(x_t - x_{t-1}) + \zeta + u_t - u_{t-1},$$

where y , x , u are defined as in (2.4) and $\zeta = (1 - \delta)^{-1} \ln(\rho\kappa) + \ln \gamma$.

We see that, in this particular example, the appropriate interpretation of the change regression is very sensitive to the amount of information available to the consumer. In the uncertainty case, a regression of $(\ln C_t - \ln C_{t-1})$ on $(\ln P_t - \ln P_{t-1})$ does not provide a consistent estimator of $(\delta - 1)^{-1}$; in fact, the estimator converges to -1 , with the implied estimator of δ converging to 0.

2.4.3. Labor supply

We shall consider a certainty model in which the consumer is maximizing

$$V = \sum_{t=1}^{\tau} \rho^{(t-1)} U_t(C_t, L_t) \quad (2.6)$$

subject to

$$\sum_{t=1}^{\tau} \gamma^{-(t-1)} (P_t C_t + W_t L_t) \leq B + \sum_{t=1}^{\tau} \gamma^{-(t-1)} W_t \bar{L},$$

$$C_t \geq 0, \quad 0 \leq L_t \leq \bar{L} \quad (t = 1, \dots, \tau),$$

where L_t is leisure, W_t is the wage rate, B is the present value in the initial period of nonlabor income, and \bar{L} is the time endowment. We shall assume that the inequality constraints on L are not binding; the participation decision will be discussed in the section on nonlinear models. If U_t is additively separable:

$$U_t(C, L) = U_t^*(C) + \tilde{U}_t(L),$$

and if $\tilde{U}_t(L) = A_t L^\delta / \delta$, then we have:

$$y_t = \beta x_t + \varphi(t-1) + c + u_t, \quad (2.7)$$

where $y_t = \ln L_t$, $x_t = \ln W_t$, $c = (\delta - 1)^{-1} \ln[\tilde{U}_1'(L_1)/W_1]$, $u_t = (1 - \delta)^{-1} \ln A_t$, $\beta =$

$(\delta - 1)^{-1}$, and $\varphi = (1 - \delta)^{-1} \ln(\gamma\rho)$. Once again c is determined by the marginal utility of initial wealth: $\tilde{U}'_1(L_1)/W_1 = \partial V/\partial B$.

We shall assume that A_t is not observed by the econometrician. There will generally be a correlation between c and (x_1, \dots, x_T) , since L_1 depends upon wages in all periods. If A_t is independent of the W 's, then we have the prediction that $\Pi = \beta I + \lambda \lambda'$. If, however, wages are partly determined by the quantity of previous work experience, then there will be lags and leads in addition to those generated by c , and Π will not have this simple structure.¹⁰

It would be useful at this point to extend the uncertainty model to incorporate uncertainty about future wages. Unfortunately, a comparably simple explicit solution is not available. But we may conjecture that the correct interpretation of a regression of $(\ln L_t - \ln L_{t-1})$ on $(\ln W_t - \ln W_{t-1})$ is also sensitive to the amount of information available to the consumer.

2.5. *Strict exogeneity conditional on a latent variable*

We shall relate the specification analysis of Π to the causality definitions of Granger (1969) and Sims (1972). Consider a sample in which $t=1$ is the first period of the individual's (economic) life.¹¹ A Sims definition of " x is strictly exogenous" is:

$$E^*(y_t | x_1, x_2, \dots) = E^*(y_t | x_1, \dots, x_t) \quad (t=1, 2, \dots).$$

In this case Π is lower triangular: the elements above the main diagonal are all zero. This fails to hold in the models we have been considering, due to the omitted variable c . But, in some cases, we do have the following property:

$$E^*(y_t | x_1, x_2, \dots, c) = E^*(y_t | x_1, \dots, x_t, c) \quad (t=1, 2, \dots). \quad (2.8)$$

It was stressed by Granger (1969) that the assessment of noncausality depends crucially on what other variables are being conditioned on. The novel feature of (2.8) is that we are asking whether there exists some latent variable (c) such that x is strictly exogenous conditional on c . The question is not vacuous since c is restricted to be time invariant.

¹⁰See Blinder and Weiss (1976) and Heckman (1976) for life-cycle labor supply models with human capital accumulation.

¹¹We shall not discuss the problems that arise from truncating the lag distribution. See Griliches and Pakes (1980).

Let us examine what restrictions are implied by (2.8). Define the following linear predictors:¹²

$$y_t = \beta_{t1}x_1 + \cdots + \beta_{tT}x_T + \gamma_t c + u_t, \\ E^*(u_t|x_1, \dots, x_T, c) = 0 \quad (t=1, \dots, T).$$

Then (2.8) is equivalent to $\beta_{ts} = 0$ for $s > t$. If $\gamma_1 \neq 0$, we can choose a scale normalization for c such that $\gamma_1 \equiv 1$. Then we can rewrite the system with $\beta_{ts} = 0$ ($s > t$) as follows:

$$y_t = \tilde{\beta}_{t1}x_1 + \beta_{t2}x_2 + \cdots + \beta_{tt}x_t + \gamma_t y_1 + \tilde{u}_t, \\ \tilde{\beta}_{t1} = \beta_{t1} - \gamma_t \beta_{11}, \quad \tilde{u}_t = u_t - \gamma_t u_1, \\ E(x_s \tilde{u}_t) = 0 \quad (s=1, \dots, T; t=2, \dots, T). \quad (2.9)$$

Consider the “instrumental variable” orthogonality conditions implied by $E(x_s \tilde{u}_t) = 0$. In the y_T equation, we have $T+1$ unknown coefficients: $\tilde{\beta}_{T1}, \beta_{T2}, \dots, \beta_{TT}, \gamma_T$, and T orthogonality conditions. So these coefficients are not identified. In the y_{T-1} equation, however, we have just enough orthogonality conditions; and in the y_{T-j} equation ($j \leq T-2$), we have $j-1$ more than we need since there are $T-j+1$ unknown coefficients: $\tilde{\beta}_{T-j,1}, \beta_{T-j,2}, \dots, \beta_{T-j,T-j}, \gamma_{T-j}$, and T orthogonality conditions: $E(x_s \tilde{u}_{T-j}) = 0$ ($s=1, \dots, T$). It follows that, subject to a rank condition, we can identify β_{ts} , γ_t , and $\tilde{\beta}_{t1}$ for $2 \leq s \leq t \leq T-1$. In addition, the hypothesis in (2.8) implies that if $T \geq 4$, there are $(T-3)(T-2)/2$ over identifying restrictions.

Consider next a Granger definition of “ y does not cause x conditional on c ”:

$$E^*(x_{t+1}|x_1, \dots, x_t, y_1, \dots, y_t, c) = E^*(x_{t+1}|x_1, \dots, x_t, c) \quad (t=1, \dots, T-1). \quad (2.10)$$

Define the following linear predictors:

$$x_{t+1} = \psi_{t1}x_1 + \cdots + \psi_{tt}x_t + \varphi_{t1}y_1 + \cdots + \varphi_{tt}y_t + \zeta_{t+1}c + v_{t+1}, \\ E^*(v_{t+1}|x_1, \dots, x_t, y_1, \dots, y_t, c) = 0 \quad (t=1, \dots, T-1).$$

Then (2.10) is equivalent to $\varphi_{ts} = 0$. We can rewrite the system, imposing $\varphi_{ts} = 0$, as follows:

$$x_{t+1} = \tilde{\psi}_{t1}x_1 + \cdots + \tilde{\psi}_{t,t-1}x_{t-1} + \tau_t x_t + \tilde{v}_{t+1}, \\ \tilde{\psi}_{ts} = \psi_{ts} - (\zeta_{t+1}/\zeta_t)\psi_{t-1,s}, \quad \tau_t = \psi_{tt} + (\zeta_{t+1}/\zeta_t), \\ \tilde{v}_{t+1} = v_{t+1} - (\zeta_{t+1}/\zeta_t)v_t, \quad E(x_s \tilde{v}_{t+1}) = E(y_s \tilde{v}_{t+1}) = 0 \\ (s \leq t-1; t=2, \dots, T-1). \quad (2.11)$$

¹²We are suppressing the period specific intercepts.

In the equation for x_{t+1} , there are t unknown parameters, $\tilde{\psi}_{t1}, \dots, \tilde{\psi}_{t,t-1}, \tau_t$, and $2(t-1)$ orthogonality conditions. Hence, there are $t-2$ restrictions ($3 \leq t \leq T-1$).

It follows that the Granger condition for “ y does not cause x conditional on c ” implies $(T-3)(T-2)/2$ restrictions, which is the same number of restrictions implied by the Sims condition. In fact, it is a consequence of Sims’ (1972) theorem, as extended by Hosoya (1977), that the two sets of restrictions are equivalent; this is not immediately obvious from a direct comparison of (2.9) and (2.11).

In terms of the Π matrix, conditional strict exogeneity implies that:

$$\Pi = B + \gamma\lambda',$$

$$B = \begin{bmatrix} \beta_{11} & 0 & \cdots & \cdots & 0 \\ \beta_{21} & \beta_{22} & 0 & \cdots & 0 \\ \vdots & & & & \\ \beta_{T1} & \beta_{T2} & \cdots & & \beta_{TT} \end{bmatrix}, \quad \gamma = \begin{pmatrix} \gamma_1 \\ \vdots \\ \gamma_T \end{pmatrix}, \quad \lambda = \begin{pmatrix} \lambda_1 \\ \vdots \\ \lambda_T \end{pmatrix}.$$

These nonlinear restrictions can be imposed and tested using the minimum distance estimator to be developed in the inference section. Alternatively, we can use the transformations in (2.9) or in (2.11). These transformations give us “simultaneous equations” systems with linear restrictions; (2.9) can be estimated using three-stage least squares. A generalization of three-stage least squares, which does not require homoskedasticity assumptions, is developed in the inference section. It is asymptotically equivalent to imposing the nonlinear restrictions directly on Π , using the minimum distance estimator.

2.6. Lagged dependent variables

For a specific example, write the labor supply model in (2.7) as follows:

$$y_t = \delta_1 x_t + \delta_2 x_{t-1} + \delta_3 y_{t-1} + v_t,$$

$$E^*(v_t | x_1, \dots, x_T) = 0 \quad (t = 1, \dots, T); \quad (2.12)$$

this reduces to (2.7) if $\delta_2 = -\delta_1$ and $\delta_3 = 1$. If we assume that $v_t = w + e_t$, where w is uncorrelated with the x ’s and e_t is i.i.d. and uncorrelated with the x ’s and w , then we have the autoregressive, variance-components model of Balestra and Nerlove (1966).¹³ In keeping with our general approach, we shall avoid placing

¹³Estimation in variance-components models is discussed in Nerlove (1967, 1971, 1971a), Wallace and Hussain (1969), Amemiya (1971), Madalla (1971), Madalla and Mount (1973), Harville (1977), Mundlak (1978), Mazodier and Trognon (1978), Trognon (1978), Lee (1979), and Taylor (1980).

restrictions on the serial correlation structure of v_t ; our inference procedures will be based on the strict exogeneity condition that $E^*(v_t|x_1, \dots, x_T) = 0$.

We can fit this model into the Π matrix framework by using recursive substitution to obtain the reduced form:

$$y_t = \beta_{t1}x_1 + \dots + \beta_{tt}x_t + \gamma_t c + u_t, \\ E^*(u_t|x_1, \dots, x_T) = 0,$$

where

$$\beta_{ts} = (\delta_2 + \delta_3\delta_1)\delta_3^{t-s-1}, \quad \beta_{tt} = \delta_1, \gamma_t = \delta_3^{t-1}, \\ c = \delta_2x_0 + \delta_3y_0, \quad u_t = v_t + \delta_3v_{t-1} + \dots + \delta_3^{t-1}v_1 \\ (1 \leq s \leq t-1, t=1, \dots, T).$$

[We are assuming that (2.12) holds for $t \geq 1$, but data on (x_0, y_0) are not available.] Hence, this model satisfies the conditional strict exogeneity restrictions:

$$\Pi = B + \gamma\lambda',$$

where B is lower triangular. The $\gamma\lambda'$ term is generated by the projection of the initial condition $(\delta_2x_0 + \delta_3y_0)$ on x_1, \dots, x_T .¹⁴

Estimation can proceed by using the minimum distance procedure to impose the nonlinear restrictions on Π . Alternatively, we can complete the system in (2.12) with:

$$y_1 = \varphi_1x_1 + \dots + \varphi_Tx_T + v_1;$$

this is just notation for the identity:

$$y_1 = E^*(y_1|x_1, \dots, x_T) + [y_1 - E^*(y_1|x_1, \dots, x_T)].$$

Then we can apply the generalized three-stage least squares estimator to be developed in the inference section. It achieves the same limiting distribution at lower computational cost, since the restrictions in this form are linear and can be imposed without requiring iterative optimization techniques.

Now consider a second-order autoregression:

$$y_t = \delta_1x_t + \delta_2x_{t-1} + \delta_3y_{t-1} + \delta_4y_{t-2} + v_t, \\ E^*(v_t|x_1, \dots, x_T) = 0 \quad (t=1, \dots, T).$$

¹⁴The treatment of initial conditions in linear models is also discussed in Anderson and Hsiao (1981, 1982) and MaCurdy (1982). The difficulties that arise in nonlinear models are discussed in Heckman (1981a).

Recursive substitution gives:

$$y_t = \beta_{1t}x_1 + \cdots + \beta_{it}x_t + \gamma_{1t}c_1 + \gamma_{2t}c_2 + u_t, \\ E^*(u_t|x_1, \dots, x_T) = 0 \quad (t=1, \dots, T),$$

where

$$c_1 = \delta_2x_0 + \delta_3y_0 + \delta_4y_{-1}, \quad c_2 = y_0,$$

and there are nonlinear restrictions on the parameters. The Π matrix has the following form:

$$\Pi = B + \gamma_1\lambda'_1 + \gamma_2\lambda'_2, \quad (2.13)$$

where B is lower triangular, $\gamma'_j = (\gamma_{1j}, \dots, \gamma_{Tj})$, and $E^*(c_j|x) = \lambda'_jx$ ($j=1, 2$).

This specification suggests a natural extension of the conditional strict exogeneity idea, with the conditioning set indexed by the number of latent variables. We shall say that " x is strictly exogenous conditional on c_1, c_2 " if:

$$E^*(y_t|\dots, x_{t-1}, x_t, x_{t+1}, \dots, c_1, c_2) = E^*(y_t|x_t, x_{t-1}, \dots, c_1, c_2).$$

We can also introduce a Granger version of this condition and generalize the analysis in Section 2.5.

Finally, consider an autoregressive model with a time-invariant omitted variable that is correlated with x :

$$y_t = \delta_1x_t + \delta_2y_{t-1} + c + v_t,$$

where $E^*(v_t|x_1, \dots, x_T) = 0$. Recursive substitution gives:

$$y_t = \beta_{1t}x_1 + \cdots + \beta_{it}x_t + \gamma_{1t}c_1 + \gamma_{2t}c_2 + u_t, \\ E^*(u_t|x_1, \dots, x_T) = 0 \quad (t=1, \dots, T),$$

where $c_1 = y_0$, $c_2 = c$, and there are nonlinear restrictions on the parameters. So y is strictly exogenous conditional on c_1, c_2 , and setting $E^*(c_j|x) = \psi_j + \lambda'_jx$ ($j=1, 2$) gives a Π matrix in the (2.13) form.

We can impose the restrictions on Π directly, using a minimum distance estimator. There is, however, a transformation of the model that allows a

computationally simpler instrumental variable estimator:

$$\begin{aligned} y_t - y_{t-1} &= \delta_1(x_t - x_{t-1}) + \delta_2(y_{t-1} - y_{t-2}) + v_t - v_{t-1}, \\ E^*(v_t - v_{t-1} | \mathbf{x}) &= 0 \quad (t = 3, \dots, T); \\ y_2 &= \varphi_{21}x_1 + \dots + \varphi_{2T}x_T + w_2, \\ y_1 &= \varphi_{11}x_1 + \dots + \varphi_{1T}x_T + w_1, \\ E^*(w_j | \mathbf{x}) &= 0 \quad (j = 1, 2), \end{aligned}$$

where $E^*(y_j | \mathbf{x}) = \boldsymbol{\varphi}'_j \mathbf{x}$ is unrestricted since $E^*(c_j | \mathbf{x})$ is unrestricted ($j = 1, 2$). Now we can apply the generalized three-stage least squares estimator. This is computationally simple since the parameter restrictions are linear. The estimator is asymptotically equivalent to applying the minimum distance procedure directly to Π . Since the linear predictor equations for y_1 and y_2 are unrestricted, the limiting distribution of $\hat{\delta}_1$ and $\hat{\delta}_2$ is not affected if we drop these equations when we form the generalized three-stage least squares estimator. (See the Appendix.)

2.7. Serial correlation or partial adjustment?

Griliches (1967) considered the problem of distinguishing between the following two models: a partial adjustment model,¹⁵

$$y_t = \beta x_t + \gamma y_{t-1} + v_t, \quad (2.14)$$

and a model with no structural lagged dependent variable but with a residual following a first-order Markov process:

$$\begin{aligned} y_t &= \beta x_t + u_t, \\ u_t &= \rho u_{t-1} + e_t, \quad e_t \text{ i.i.d.}; \end{aligned} \quad (2.15)$$

in both cases x is strictly exogenous:

$$E^*(v_t | x_1, \dots, x_T) = E^*(u_t | x_1, \dots, x_T) = 0 \quad (t = 1, \dots, T).$$

In the serial correlation case, we have:

$$y_t = \beta x_t - \rho \beta x_{t-1} + \rho y_{t-1} + e_t;$$

¹⁵See Nerlove (1972) for distributed lag models based on optimizing behavior in the presence of uncertainty and costs of adjustment.

as Griliches observed, the least squares regression will have a distinctive pattern—the coefficient on lagged x equals (as $N \rightarrow \infty$) minus the product of the coefficients on current x and lagged y .

I want to point out that this prediction does not rest on the serial correlation structure of u . It is a direct implication of the assumption that u is uncorrelated with x_1, \dots, x_T :

$$\begin{aligned} E^*(y_t | x_t, x_{t-1}, y_{t-1}) &= \beta x_t + E^*(u_t | x_t, x_{t-1}, y_{t-1}) \\ &= \beta x_t + E^*(u_t | u_{t-1}) \\ &= \beta x_t + \varphi_t u_{t-1} \\ &= \beta x_t - \varphi_t \beta x_{t-1} + \varphi_t y_{t-1}. \end{aligned}$$

Here $\varphi_t u_{t-1}$ is simply notation for the linear predictor. In general u_t is not a first-order process ($E^*(u_t | u_{t-1}, u_{t-2}) \neq E^*(u_t | u_{t-1})$), but this does not affect our argument.

Within the Π matrix framework, the distinction between the two models is that (2.15) implies a diagonal Π matrix, with no distributed lag, whereas the partial adjustment specification in (2.14) implies that $\Pi = B + \gamma\lambda'$, with a distributed lag in the lower triangular B matrix and a rank one set of lags and leads in $\gamma\lambda'$.

We can generalize the serial correlation model to allow for an individual specific effect that may be correlated with x :

$$y_t = \beta x_t + c + u_t, \quad E^*(u_t | x_1, \dots, x_T) = 0.$$

Now both the serial correlation and the partial adjustment models have a rank one set of lags and leads in Π , but we can distinguish between them because only the partial adjustment model has a distributed lag in the B matrix. So the absence of structural lagged dependent variables is signalled by the following special case of conditional strict exogeneity:

$$E^*(y_t | x_1, \dots, x_T, c) = E^*(y_t | x_t, c).$$

In this case the relationship of x to y is “static” conditional on c . We shall pursue this distinction in nonlinear models in Section 3.3.

2.8. Residual covariances: Heteroskedasticity and serial correlation

2.8.1. Heteroskedasticity

If $E(c_i | x_i) \neq E^*(c_i | x_i)$, then there will be heteroskedasticity, since the residual will contain $E(c_i | x_i) - E^*(c_i | x_i)$. Another source of heteroskedasticity is random

coefficients:

$$\begin{aligned} y_{it} &= b_i x_{it} + c_i + u_{it}, \\ b_i &= \beta + w_i, \quad E(w_i) = 0, \\ y_{it} &= \beta x_{it} + c_i + (w_i x_{it} + u_{it}). \end{aligned}$$

If w is independent of x , then $\Pi = \beta I + I\lambda'$, and our previous discussion is relevant for the estimation of β . We shall handle the heteroskedasticity problem in the inference section by allowing $E[(y_i - \Pi x_i)(y_i - \Pi x_i)' | x_i]$ to be an arbitrary function of x_i .¹⁶

2.8.2. Serial correlation

It may be of interest to impose restrictions on the residual covariances, such as a variance-components structure together with an autoregressive-moving average scheme.¹⁷ Consider the homoskedastic case in which

$$\Omega = E[(y_i - \Pi x_i)(y_i - \Pi x_i)' | x_i]$$

does not depend upon x_i . Then the restrictions can be expressed as $\Omega_{jk} = g_{jk}(\theta)$, where the g 's are known functions and θ is an unrestricted parameter vector. We shall discuss a minimum distance procedure for imposing such restrictions in Section 4.

2.9. Measurement error

Suppose that

$$\begin{aligned} y_{it} &= \beta x_{it}^* + u_{it}, \\ x_{it} &= x_{it}^* + v_{it} \quad (i = 1, \dots, N; t = 1, \dots, T), \end{aligned}$$

where x_{it}^* is not observed. We assume that the measurement error v_{it} satisfies $E^*(v_{it} | x_i) = 0$. If $E^*(u_{it} | x_i) = 0$, then $E^*(y_i | x_i) = \Pi x_i$, with

$$\Pi = \beta V(x_i^*) V^{-1}(x_i). \quad (2.16)$$

¹⁶Anderson (1969, 1970), Swamy (1970, 1974), Hsiao (1975), and Mundlak (1978a) discuss estimators that incorporate the particular form of heteroskedasticity that is generated by random coefficients.

¹⁷Such models for the covariance structure of earnings have been considered by Hause (1977, 1980), Lillard and Willis (1978), Lillard and Weiss (1979), MaCurdy (1982), and others.

Since $V(x_i)$ and $V(x_i^*)$ will generally not be diagonal matrices, (2.16) provides an alternative interpretation of lags and leads in the Π matrix. The Π matrix in (2.16) generally does not have the form $\tau_1 I + \tau_2 I\lambda'$; nevertheless, it may be difficult to distinguish between measurement error and a time-invariant omitted variable if T is small. For example, if the covariance matrices of x_i and x_i^* have the form $\varphi_1 I + \varphi_2 W'$ (equicorrelated), then Π has this form also and no distinction is possible. Although $\text{cov}(x_{it}, x_{is})$ generally declines as $|t - s|$ increases, the equicorrelated approximation may be quite good for small T .

It has been noted in other contexts that the bias from measurement error can be aggravated by analysis of covariance techniques.¹⁸ Consider the following example with $T = 2$:

$$y_{i2} - y_{i1} = \beta(x_{i2} - x_{i1}) + u_{i2} - u_{i1} - \beta(v_{i2} - v_{i1}),$$

so that $E^*(y_{i2} - y_{i1} | x_{i2} - x_{i1}) = \tilde{\beta}(x_{i2} - x_{i1})$ with

$$\tilde{\beta} = \beta \left(1 - \frac{V(v_{i2} - v_{i1})}{V(x_{i2} - x_{i1})} \right).$$

If $V(v_{i1}) = V(v_{i2})$ and $V(x_{i1}) = V(x_{i2})$, then we can rewrite this as:

$$\tilde{\beta} = \beta \left(1 - \frac{V(v_{i1})(1 - r_{v_1 v_2}^2)}{V(x_{i1})(1 - r_{x_1 x_2}^2)} \right),$$

where $r_{v_1 v_2}$ denotes the correlation between v_{i1} and v_{i2} . If x_{i1} and x_{i2} are highly correlated but v_{i1} and v_{i2} are not, then a modest bias from measurement error in a cross-section regression can become large when we relate the change in y to the change in x . On the other hand, if $v_{i1} = v_{i2}$, then the change regression eliminates the bias from measurement error. Data from reinterview surveys should be useful in distinguishing between these two cases.

3. Specification and identification: Nonlinear models

3.1. A random effects probit model

Our treatment of individual effects carries over with some important qualifications to nonlinear models. We shall illustrate with a labor force participation example. If the upper bound on leisure is binding in (2.6), then

$$\rho^{(t-1)} \tilde{U}_t'(\bar{L}) > m \gamma^{-(t-1)} W_t,$$

¹⁸See Griliches (1979) for example.

where m is the Lagrange multiplier corresponding to the lifetime budget constraint (the marginal utility of initial wealth) and $\bar{U}_t(L) = A_t L^\delta / \delta$. Let $y_{it} = 1$ if individual i works in period t , $y_{it} = 0$ otherwise. Let:

$$\ln W_{it} = \varphi_1 x_{it} + e_{1it},$$

$$\ln A_{it} = \varphi_2 x_{it} + e_{2it},$$

where x_{it} contains measured variables that predict wages and tastes for leisure. We shall simplify the notation by supposing that x_{it} consists of a single variable. Then, $y_{it} = 1$ if:

$$(\varphi_1 - \varphi_2)x_{it} - (t-1)\ln(\gamma\rho) + \ln m_i + (1-\delta)\ln \bar{L} + e_{1it} - e_{2it} \geq 0,$$

which we shall write as:

$$\beta x_{it} + \varphi(t-1) + c_i + u_{it} \geq 0. \quad (3.1)$$

Now we need a distributional assumption for the u 's. We shall assume that (u_1, \dots, u_T) is independent of c and the x 's, with a multivariate normal distribution $(N(\mathbf{0}, \Sigma))$. So we have a probit model (suppressing the i subscripts and period-specific intercepts):

$$P(y_t = 1 | x_1, \dots, x_T, c) = F[\sigma_t^{-1/2}(\beta x_t + c)],$$

where $F(\cdot)$ is the standard normal distribution function and σ_t is the t th diagonal element of Σ .

Next we shall specify a distribution for c conditional on $\mathbf{x} = (x_1, \dots, x_T)$:

$$c = \psi + \lambda_1 x_1 + \dots + \lambda_T x_T + v,$$

where v is independent of the x 's and has a normal distribution $(N(0, \sigma_v^2))$. There is a very important difference in this step compared with the linear case. In the linear case it was not restrictive to decompose c into its linear projection on \mathbf{x} and an orthogonal residual. Now, however, we are assuming that the regression function $E(c|\mathbf{x})$ is actually linear, that v is independent of \mathbf{x} , and that v has a normal distribution. These are restrictive assumptions and there may be a payoff to relaxing them.

Given these assumptions, the distribution for y_t conditional on x_1, \dots, x_T but marginal on c also has a probit form:

$$P(y_t = 1 | x_1, \dots, x_T) = F[\alpha_t(\beta x_t + \lambda_1 x_1 + \dots + \lambda_T x_T)],$$

$$\alpha_t = (\sigma_{tt} + \sigma_v^2)^{-1/2}.$$

Combining these T specifications gives the following matrix of coefficients:¹⁹

$$\Pi = \text{diag} \{ \alpha_1, \dots, \alpha_T \} [\beta I_T + I \lambda']. \quad (3.2)$$

This differs from the linear case only in the diagonal matrix of normalization factors α_t . There are now nonlinear restrictions on Π , but the identification analysis is still straightforward. We have:

$$\alpha_t \beta = \frac{\alpha_t}{\alpha_1} \pi_{11} - \pi_{t1} = \pi_{tt} - \frac{\alpha_t}{\alpha_1} \pi_{1t},$$

$$\frac{\alpha_t}{\alpha_1} = \frac{(\pi_{tt} + \pi_{t1})}{(\pi_{11} + \pi_{1t})} \quad (t = 2, \dots, T),$$

if $\beta + \lambda_1 + \lambda_t \neq 0$. Then, as in the linear case, we can solve for $\alpha_1 \beta$ and $\alpha_1 \lambda$. Only ratios of coefficients are identified, and so we can use a scale normalization such as $\alpha_1 \equiv 1$.

As for inference, a computationally simple approach is to estimate T cross-sectional probit specifications by maximum likelihood, where x_1, \dots, x_T are included in each of the T specifications. This gives $\hat{\pi}_t$ ($t = 1, \dots, T$) and we can use a Taylor expansion to derive the covariance matrix of the asymptotic normal distribution for $(\hat{\pi}_1, \dots, \hat{\pi}_T)$. Then restrictions can be imposed on Π using a minimum distance estimator, just as in the linear case.

We shall conclude our discussion of this model by considering the interpretation of the coefficients. We began with the probit specification that

$$P(y_i = 1 | x_1, \dots, x_T, c) = F \left[\sigma_{ii}^{-1/2} (\beta x_i + c) \right].$$

So one might argue that the correct measure of the effect of x_t is based on $\sigma_{ii}^{-1/2} \beta$, whereas we have obtained $(\sigma_{ii} + \sigma_v^2)^{-1/2} \beta$, which is then an underestimate. But there is something curious about this argument, since the "omitted variable" v is independent of x_1, \dots, x_T . Suppose that we decompose u_t in (3.1) into $u_{1t} + u_{2t}$ and that measurements on u_{1t} become available. Then this argument implies that the correct measure of the effect of x_t is based on $[V(u_{2t})]^{-1/2} \beta$. As the data collection becomes increasingly successful, there is less and less variance left in the residual u_{2t} , and $[V(u_{2t})]^{-1/2}$ becomes arbitrarily large.

The resolution of this puzzle is that the effect of x_t depends upon the value of c , and the effect evaluated at the average value for c is not equal to the average of the effects, averaging over the distribution for c . Consider the effect on the probability that $y_i = 1$ of increasing x_t from x' to x'' ; using the average value for c

¹⁹This approach to analysis of covariance in probit models was proposed in Chamberlain (1980). For other applications of multivariate probit models to panel data, see Heckman (1978, 1981).

gives:

$$F[\sigma_{it}^{-1/2}(\beta x'' + E(c))] - F[\sigma_{it}^{-1/2}(\beta x' + E(c))].$$

The problem with this measure is that it may be relevant for only a small fraction of the population. I think that a more appropriate measure is the mean effect for a randomly drawn individual:

$$\int [P(y_t = 1|x_t = x'', c) - P(y_t = 1|x_t = x', c)] \mu(dc),$$

where $\mu(dc)$ gives the population probability measure for c .

We shall see how to recover this measure within our framework. Let $z = \lambda_1 x_1 + \dots + \lambda_T x_T$; let $\mu(dz)$ and $\mu(dv)$ give the population probability measures for the independent random variables z and v . Then:

$$\begin{aligned} P(y_t = 1|x_t, c) &= P(y_t = 1|x_1, \dots, x_T, c) \\ &= P(y_t = 1|x_t, z, v); \\ \int P(y_t = 1|x_t, z, v) \mu(dz) \mu(dv) \\ &= \int P(y_t = 1|x_t, z, v) \mu(dv|x_t, z) \mu(dz) \\ &= \int P(y_t = 1|x_t, z) \mu(dz), \end{aligned}$$

where $\mu(dv|x_t, z)$ is the conditional probability measure, which equals the unconditional measure since v is independent of x_t and z . [It is important to note that the last integral does *not*, in general, equal $P(y_t = 1|x_t)$. For if x_t and z are correlated, as they are in our case, then

$$\begin{aligned} P(y_t = 1|x_t) &= \int P(y_t = 1|x_t, z) \mu(dz|x_t) \\ &\neq \int P(y_t = 1|x_t, z) \mu(dz). \end{aligned}$$

We have shown that:

$$\begin{aligned} \int [P(y_t = 1|x_t = x'', c) - P(y_t = 1|x_t = x', c)] \mu(dc) \\ = \int [P(y_t = 1|x_t = x'', z) - P(y_t = 1|x_t = x', z)] \mu(dz). \end{aligned} \quad (3.3)$$

The integration with respect to the marginal distribution for z can be done using the empirical distribution function, which gives the following consistent (as

$N \rightarrow \infty$) estimator of (3.3):

$$\frac{1}{N} \sum_{i=1}^N \left\{ F \left[\alpha_i (\beta x'' + \lambda_1 x_{i1} + \cdots + \lambda_T x_{iT}) \right] - F \left[\alpha_i (\beta x' + \lambda_1 x_{i1} + \cdots + \lambda_T x_{iT}) \right] \right\}. \quad (3.4)$$

3.2. A fixed effects logit model: Conditional likelihood

A weakness in the probit model was the specification of a distribution for c conditional on x . A convenient form was chosen, but it was only an approximation, perhaps a poor one. We shall discuss a technique that does not require us to specify a particular distribution for c conditional on x ; it will, however, have its own weaknesses.

Consider the following specification:

$$P(y_i = 1 | x_1, \dots, x_T, c) = G(\beta x_i + c), \quad G(z) = e^z / (1 + e^z), \quad (3.5)$$

where y_1, \dots, y_T are independent conditional on x_1, \dots, x_T, c . Suppose that $T = 2$ and compute the probability that $y_2 = 1$ conditional on $y_1 + y_2 = 1$:

$$P(y_2 = 1 | x_1, x_2, c, y_1 + y_2 = 1) = G[\beta(x_2 - x_1)], \quad (3.6)$$

which does not depend upon c . Given a random sample of individuals, the conditional log-likelihood function is:

$$L = \sum_{i \in B} \left\{ w_i \ln G[\beta(x_{i2} - x_{i1})] + (1 - w_i) \ln G[-\beta(x_{i2} - x_{i1})] \right\},$$

where

$$w_i = \begin{cases} 1, & \text{if } (y_{i1}, y_{i2}) = (0, 1), \\ 0, & \text{if } (y_{i1}, y_{i2}) = (1, 0), \end{cases}$$

$$B = \{i | y_{i1} + y_{i2} = 1\}.$$

This conditional likelihood function does not depend upon the incidental parameters. It is in the form of a binary logit likelihood function in which the two outcomes are (0,1) and (1,0) with explanatory variables $x_2 - x_1$. This is the analog of differencing in the two period linear model. The conditional maximum likelihood (ML) estimate of β can be obtained simply from a ML binary logit

program. This conditional likelihood approach was used by Rasch (1960, 1961) in his model for intelligence tests.²⁰

The conditional ML estimator of β is consistent provided that the conditional likelihood function satisfies regularity conditions, which impose mild restrictions on the c_i . These restrictions, which are satisfied if the c_i are a random sample from some distribution, are discussed in Andersen (1970). Furthermore, the inverse of the information matrix based on the conditional likelihood function provides a covariance matrix for the asymptotic ($N \rightarrow \infty$) normal distribution of the conditional ML estimator of β .

These results should be contrasted with the inconsistency of the standard fixed effects ML estimator, in which the likelihood function is based on the distribution of y_1, \dots, y_T conditional on x_1, \dots, x_T, c . For example, suppose that $T = 2$, $x_{i1} = 0$, $x_{i2} = 1$ ($i = 1, \dots, N$). The following limits exist with probability one if the c_i are a random sample from some distribution:

$$\lim_{N \rightarrow \infty} \frac{1}{N} \sum_{i=1}^N E[y_{i1}(1 - y_{i2})|c_i] = \varphi_1,$$

$$\lim_{N \rightarrow \infty} \frac{1}{N} \sum_{i=1}^N E[(1 - y_{i1})y_{i2}|c_i] = \varphi_2,$$

where

$$E[y_{i1}(1 - y_{i2})|c_i] = G(c_i)G(-\beta - c_i),$$

$$E[(1 - y_{i1})y_{i2}|c_i] = G(-c_i)G(\beta + c_i).$$

Andersen (1973, p. 66) shows that the ML estimator of β converges with probability one to 2β as $N \rightarrow \infty$. A simple extension of his argument shows that if G is replaced by any distribution function (\tilde{G}) corresponding to a symmetric, continuous, nonzero probability density, then the ML estimator of β converges

²⁰In Rasch's model, the probability that person i gives a correct answer to item number t is $\exp(\beta_t + c_i)/[1 + \exp(\beta_t + c_i)]$; this is a special case in which x_{it} is a set of dummy indicator variables. An algorithm for maximum likelihood estimation in this case is described in Andersen (1972). The use of conditional likelihood in incidental parameter problems is discussed in Andersen (1970, 1973), Kalbfleisch and Sprott (1970), and Barndorff-Nielsen (1978). The conditional likelihood approach in the logit case is closely related to Fisher's (1935) exact test for independence in a 2×2 table. This exact significance test has been extended by Cox (1970) and others to the case of several contingency tables. Additional references are in Cox (1970) and Bishop et al. (1975). Chamberlain (1979) develops a conditional likelihood estimator for a point process model based on duration data, and Griliches, Hall and Hausman (1981) apply conditional likelihood techniques to panel data on counts.

with probability one to:

$$2\tilde{G}^{-1}\left(\frac{\varphi_2}{\varphi_1 + \varphi_2}\right).$$

The logit case is special in that $\varphi_2/\varphi_1 = e^\beta$ for any distribution for c . In general the limit depends on this distribution; but if all of the $c_i = 0$, then once again we obtain convergence to 2β as $N \rightarrow \infty$.

For general T , conditioning on $\sum_t y_{it}$ ($i=1, \dots, N$) gives the following conditional log-likelihood function:

$$L = \sum_{i=1}^N \ln \left[\exp \left(\beta \sum_{t=1}^T x_{it} y_{it} \right) / \sum_{d \in B_i} \exp \left(\beta \sum_{t=1}^T x_{it} d_t \right) \right],$$

$$B_i = \left\{ d = (d_1, \dots, d_T) \mid d_t = 0 \text{ or } 1 \text{ and } \sum_{t=1}^T d_t = \sum_{t=1}^T y_{it} \right\}.$$

L is in the conditional logit form considered by McFadden (1974), with the alternative set (B_i) varying across the observations. Hence, it can be maximized by standard programs. There are $T+1$ distinct alternative sets corresponding to $\sum_t y_{it} = 0, 1, \dots, T$. Groups for which $\sum_t y_{it} = 0$ or T contribute zero to L , however, and so only $T-1$ alternative sets are relevant. The alternative set for the group with $\sum_t y_{it} = s$ has $\binom{T}{s}$ elements, corresponding to the distinct sequences of T trials with s successes. For example, with $T=3$ and $s=1$ there are three alternatives with the following conditional probabilities:

$$P\left(1, 0, 0 \mid x_i, c_i, \sum_t y_{it} = 1\right) = \exp[\beta(x_{i1} - x_{i3})] / D,$$

$$P\left(0, 1, 0 \mid x_i, c_i, \sum_t y_{it} = 1\right) = \exp[\beta(x_{i2} - x_{i3})] / D,$$

$$P\left(0, 0, 1 \mid x_i, c_i, \sum_t y_{it} = 1\right) = 1/D,$$

$$D = \exp[\beta(x_{i1} - x_{i3})] + \exp[\beta(x_{i2} - x_{i3})] + 1.$$

A weakness in this approach is that it relies on the assumption that the y_t are independent conditional on x , c , with an identical form for the conditional probability each period: $P(y_t = 1 \mid x, c) = G(\beta x_t + c)$. In the probit framework, these assumptions translate into $\Sigma = \sigma^2 I$, so that $v + u_t$ generates an equicorrelated matrix: $\sigma_v^2 W' + \sigma^2 I$. We have seen that it is straightforward to allow Σ to be unrestricted in the probit framework; that is not true here.

An additional weakness is that we are limited in the sorts of probability statements that can be made. We obtain a clean estimate of the effect of x_t on the log odds:

$$\ln \left[\frac{P(y_t = 1 | x_t = x'', c)}{P(y_t = 0 | x_t = x'', c)} \bigg/ \frac{P(y_t = 1 | x_t = x', c)}{P(y_t = 0 | x_t = x', c)} \right] = \beta(x'' - x');$$

the special feature of the logistic functional form is that this function of the probabilities does not depend upon c ; so the problem of integrating over the marginal distribution of c (instead of the conditional distribution of c given x) does not arise. But this is not the only function of the probabilities that one might want to know. In the probit section we considered

$$P(y_t = 1 | x_t = x'', c) - P(y_t = 1 | x_t = x', c),$$

which depends upon c for probit or logit, and we averaged over the marginal distribution for c :

$$\int [P(y_t = 1 | x_t = x'', c) - P(y_t = 1 | x_t = x', c)] \mu(dc). \quad (3.7)$$

This requires us to specify a marginal distribution for c , which is what the conditioning argument tries to avoid. We cannot estimate (3.7) if all we have is the conditional ML estimate of β .

Our specification in (3.5) asserts that y_t is independent of $x_1, \dots, x_{t-1}, x_{t+1}, \dots, x_T$ conditional on x_t, c . This can be relaxed somewhat, but the conditional likelihood argument certainly requires more than

$$P(y_t = 1 | x_t, c) = G(\beta x_t + c);$$

to see this, try to derive (3.6) with $x_2 = y_1$. We can, however, implement the following specification (with $x' = (x_1, \dots, x_T)$):

$$P(y_t = 1 | x, c) = G(\beta_{t0} + \beta_{t1}x_1 + \dots + \beta_{tt}x_t + c), \quad (3.8)$$

where y_1, \dots, y_T are independent conditional on x, c . This corresponds to our specification of “ x is strictly exogenous conditional on c ” in Section 2.5, except that $\gamma_t = 1$ in the term $\gamma_t c$ —it is not straightforward to allow a time-varying coefficient on c in the conditional likelihood approach. The extension of (3.6) is:

$$P(y_t = 1 | x, c, y_1 + y_t = 1) = G(\tilde{\beta}_{t0} + \tilde{\beta}_{t1}x_1 + \beta_{t2}x_2 + \dots + \beta_{tt}x_t) \\ (t = 2, \dots, T), \quad (3.9)$$

where $\tilde{\beta}_{1j} = \beta_{1j} - \beta_{10}$ ($j = 0, 1$). So if x has sufficient variation, we can obtain consistent estimates of $\tilde{\beta}_{10}$, $\tilde{\beta}_{11}$, and β_{1s} ($s = 2, \dots, t$). Only these parameters are identified, since we can transform the model replacing c by $\tilde{c} = \beta_{10} + \beta_{11}x_1 + c$ without violating any restrictions.

The restrictions in (3.5) or in (3.8) can be tested against the following alternative:

$$P(y_t = 1 | x, c) = G(\pi_{t0} + \pi_{t1}x_1 + \dots + \pi_{tT}x_T + c). \quad (3.10)$$

We can identify only $\pi_{tj} - \pi_{1j}$ and so we can normalize $\pi_{1j} = 0$ ($j = 0, \dots, T$; $t = 2, \dots, T$). The maximized values of the conditional log-likelihoods can be used to form χ^2 statistics.²¹ There are $(T-2)(T-1)/2$ restrictions in passing from (3.10) to (3.8), and (3.5) imposes an additional $(T-1)(T+4)/2 - 1$ restrictions.

3.3. Serial correlation and lagged dependent variables

Consider the following two models:

$$y_t = \begin{cases} 1, & \text{if } y_t^* = u_t \geq 0, \\ 0, & \text{otherwise; } u_t = \rho u_{t-1} + e_t; \end{cases} \quad (3.11b)$$

$$y_t = \begin{cases} 1, & \text{if } y_t^* = u_t \geq 0, \\ 0, & \text{otherwise; } u_t = \rho u_{t-1} + e_t; \end{cases} \quad (3.11b)$$

in both cases e_t is i.i.d. $N(0, \sigma^2)$. Heckman (1978) observed that we can distinguish between these two models.²² In the first model,

$$P(y_t = 1 | y_{t-1}, y_{t-2}, \dots) = P(y_t = 1 | y_{t-1}) = F(\gamma y_{t-1} / \sigma),$$

where $F(\cdot)$ is the standard normal distribution function. In the second model, however, $P(y_t = 1 | y_{t-1}, y_{t-2}, \dots)$ depends upon the entire history of the process. If we observed u_{t-1} , then previous outcomes would be irrelevant. In fact, we observe only whether $u_{t-1} \geq 0$; hence conditioning in addition on whether $u_{t-2} \geq 0$ affects the distribution of u_{t-1} and y_t . So the lagged y implies a Markov chain whereas the Markov assumption for the probit residual does not imply a Markov chain for the binary sequence that it generates.

²¹ Conditional likelihood ratio tests are discussed in Andersen (1971).

²² Also see Heckman (1981, 1981b).

There is an analogy with the following linear models:

$$y_t = \gamma y_{t-1} + e_t, \quad (3.12a)$$

$$y_t = u_t, u_t = e_t + \rho e_{t-1}, \quad (3.12b)$$

where e_t is i.i.d. $N(0, \sigma^2)$. We know that if $u_t = \rho u_{t-1} + e_t$, then no distinction would be possible, without introducing more structure, since both models imply a linear Markov process. With the moving average residual, however, the serial correlation model implies that the entire past history is relevant for predicting y . So the distinction between the two models rests on the order of the dependence on previous realizations of y_t .

We can still distinguish between the two models in (3.11) even when (u_1, \dots, u_T) has a general multivariate normal distribution $(N(\mu, \Sigma))$. Given normalizations such as $V(u_t) = 1$ ($t = 1, \dots, T$), the serial correlation model has $T(T+1)/2$ free parameters. Hence, if $T \geq 3$, there are restrictions on the $2^T - 1$ parameters of the multinomial distribution for (y_1, \dots, y_T) . In particular, the most general multivariate probit model cannot generate a Markov chain. So we can add a lagged dependent variable and identify γ .

This result relies heavily on the restrictive nature of the multivariate probit functional form. A more robust distinction between the two models is possible when there is variation over time in x_t . We shall pursue this after first presenting a generalization of strict exogeneity and noncausality for nonlinear models.

Let $t = 1$ be the first period of the individual's (economic) life. An extension of Granger's definition of "y does not cause x" is that x_{t+1} is independent of y_1, \dots, y_t conditional on x_1, \dots, x_t . An extension of Sims' strict exogeneity condition is that y_t is independent of x_{t+1}, x_{t+2}, \dots conditional on x_1, \dots, x_t . In contrast to the linear predictor case, these two definitions are no longer equivalent.²³ For consider the following counterexample: let y_1, y_2 be independent Bernoulli random variables with $P(y_t = 1) = P(y_t = -1) = 1/2$ ($t = 1, 2$). Let $x_3 = y_1 y_2$. Then y_1 is independent of x_3 and y_2 is independent of x_3 . Let all of the other random variables be degenerate (equal to zero, say). Then x is strictly exogenous but x_3 is clearly not independent of y_1, y_2 conditional on x_1, x_2 . The counterexample works for the following reason: if a random variable is uncorrelated with each of two other random variables, then it is uncorrelated with every linear combination of them; but if it is independent of each of the other random variables, it need not be independent of every function of them.

Consider the following modification of Sims' condition: y_t is independent of x_{t+1}, x_{t+2}, \dots conditional on $x_1, \dots, x_t, y_1, \dots, y_{t-1}$ ($t = 1, 2, \dots$). Chamberlain (1982) shows that, subject to a regularity condition, this is equivalent to our

²³See Chamberlain (1982) and Florens and Mouchart (1982).

extended definition of Granger noncausality. The regularity condition is trivially satisfied whenever y_t has a degenerate distribution prior to some point. So it is satisfied in our case since y_0, y_{-1}, \dots have degenerate distributions.

It is straightforward to introduce a time-invariant latent variable into these definitions. We shall say that " y does not cause x conditional on a latent variable c " if either:

x_{t+1} is independent of y_1, \dots, y_t conditional on x_1, \dots, x_t, c ($t = 1, 2, \dots$),

or

y_t is independent of x_{t+1}, x_{t+2}, \dots conditional on $x_1, \dots, x_t, y_1, \dots, y_{t-1}, c$ ($t = 1, 2, \dots$);

they are equivalent. We shall say that " x is strictly exogenous conditional on a latent variable c " if:

y_t is independent of x_{t+1}, x_{t+2}, \dots conditional on x_1, \dots, x_t, c ($t = 1, 2, \dots$).

Now let us return to the problem of distinguishing between serial correlation and structural lagged dependent variables. Assume throughout the discussion that x_t and y_t are not independent. We shall say that the relationship of x to y is *static* if:

x is strictly exogenous and y_t is independent of x_1, \dots, x_{t-1} conditional on x_t .

Then I propose the following distinctions:

There is residual serial correlation if y_t is not independent of y_1, \dots, y_{t-1} conditional on x_1, \dots, x_t ;

If the relationship of x to y is static, then there are no structural lagged dependent variables.

Suppose that y_t and x_t are binary and consider the probability that $y_2 = 1$ conditional on $(x_1, x_2) = (0, 0)$ and conditional on $(x_1, x_2) = (1, 0)$. Since y_t and x_t are assumed to be dependent, the distribution of y_1 is generally different in the two cases. If y_1 has a structural effect on y_2 , then the conditional probability of $y_2 = 1$ should differ in the two cases, so that y_2 is not independent of x_1 conditional on x_2 .

Note that this condition is one-sided: I am only offering a condition for there to be no structural effect of y_{t-1} on y_t . There can be distributed lag relationships in which we would not want to say that y_{t-1} has a structural effect on y_t . Consider the production function example with serial correlation in rainfall; assume for the moment that there is no variation in c . If the serial correlation in rainfall is not incorporated in the farmer's information set, then our definitions assert that there is residual serial correlation but no structural lagged dependent variables, since the relationship of x to y is static. Now suppose that the farmer does use previous rainfall to predict future rainfall. Then the relationship of x to y is not static since

x is not strictly exogenous. But we may not want to say that the relationship between y_{t-1} and y_t is structural, since the technology does not depend upon y_{t-1} .

How are these distinctions affected by latent variables? It should be clear that a time-invariant latent variable can produce residual serial correlation. A major theme of the paper has been that such a latent variable can also produce a failure of strict exogeneity. So consider conditional versions of these properties:

There is residual serial correlation conditional on a latent variable c if y_t is not independent of y_1, \dots, y_{t-1} conditional on x_1, \dots, x_t, c ;

The relationship of x to y is static conditional on a latent variable c if x is strictly exogenous conditional on c and if y_t is independent of x_1, \dots, x_{t-1} conditional on x_t, c ;

If the relationship of x to y is static conditional on a latent variable c , then there are no structural lagged dependent variables.

A surprising feature of the linear predictor definition of strict exogeneity is that it is restrictive to assert that there exists some time-invariant latent variable c such that x is strictly exogenous conditional on c . This is no longer true when we use conditional independence to define strict exogeneity. For a counterexample, suppose that x_t is a binary variable and consider the conditional strict exogeneity question: "Does there exist a time-invariant random variable c such that y_t is independent of x_1, \dots, x_T conditional on x_1, \dots, x_t, c ?" The answer is "yes" since we can order the 2^T possible outcomes of the binary sequence (x_1, \dots, x_T) and set $c = j$ if the j th outcome occurs ($j = 1, \dots, 2^T$). Now y_t is independent of x_1, \dots, x_T conditional on c !

For a nondegenerate counterexample, let y and x be binary random variables with:

$$P(y = \alpha_j, x = \alpha_k) = \tau_{jk} > 0, \quad \sum_{j,k=1}^2 \tau_{jk} = 1,$$

where $\alpha_1 = 1, \alpha_2 = 0$. Let $\gamma' = (\tau_{11}, \tau_{12}, \tau_{21}, \tau_{22})$. Then we can set:

$$\gamma = \sum_{m=1}^4 \gamma_m e_m, \quad \gamma_m > 0, \quad \sum_{m=1}^4 \gamma_m = 1,$$

where e_m is a vector of zeros except for a one in the m th component. Hence γ is in the interior of the convex hull of $\{e_m, m = 1, \dots, 4\}$. Now consider the vector:

$$y(\delta, \lambda) = \begin{bmatrix} \delta\lambda \\ \delta(1-\lambda) \\ (1-\delta)\lambda \\ (1-\delta)(1-\lambda) \end{bmatrix}$$

The components of $y(\delta, \lambda)$ give the probabilities $P(y = \alpha_j, x = \alpha_k)$ when y and x are independent with $P(y = 1) = \delta$, $P(x = 1) = \lambda$. Set $e_m^* = y(\delta_m, \lambda_m)$ with $0 < \delta_m < 1, 0 < \lambda_m < 1$. Then γ will be in the interior of the convex hull of $\{e_m^*, m = 1, \dots, 4\}$ if we choose δ_m, λ_m so that e_m^* is sufficiently close to e_m . Hence:

$$\gamma = \sum_{m=1}^4 \gamma_m^* e_m^*, \quad \gamma_m^* > 0, \quad \sum_{m=1}^4 \gamma_m^* = 1.$$

Let the components of e_m^* be $(\tau_{11}^m, \tau_{12}^m, \tau_{21}^m, \tau_{22}^m)$. Let c be a random variable with $P(c = m) = \gamma_m^* (m = 1, \dots, 4)$, and set

$$P(y = \alpha_j, x = \alpha_k | c = m) = \tau_{jk}^m.$$

Now y is independent of x conditional on c , and the conditional distributions are nondegenerate.

If $(x_1, \dots, x_T, y_1, \dots, y_T)$ has a general multinomial distribution, then a straightforward extension of this argument shows that there exists a random variable c such that (y_1, \dots, y_T) is independent of (x_1, \dots, x_T) conditional on c , and the conditional distributions are nondegenerate.

A similar point applies to factor analysis. Consider a linear one-factor model. The specification is that there exists a latent variable c such that the partial correlations between y_1, \dots, y_T are zero given c . This is restrictive if $T > 3$. But we now know that it is not restrictive to assert that there exists a latent variable c such that y_1, \dots, y_T are independent conditional on c .

It follows that we cannot test for conditional strict exogeneity without imposing functional form restrictions; nor can we test for a conditionally static relationship without restricting the functional forms.

This point is intimately related to the fundamental difficulties created by incidental parameters in nonlinear models. The labor force participation example is assumed to be static conditional on c . We shall present some tests of this in Section 5, but we shall be jointly testing that proposition and the functional forms—a truly nonparametric test cannot exist. We stressed in the probit model that the specification for the distribution of c conditional on x is restrictive; we avoided such a restrictive specification in the logit model but only by imposing a restrictive functional form on the distribution of y conditional on x, c .

3.4. Duration models

In many problems the basic data is the amount of time spent in a state. For example, a complete description of an individual's labor force participation history is the duration of the first spell of participation and the date it began, the

duration of the following spell of nonparticipation, and so on. This complete history will generate a binary sequence when it is cut up into fixed length periods, but these periods may have little to do with the underlying process.²⁴

In particular, the measurement of serial correlation depends upon the period of observation. As the period becomes shorter, the probability that a person who worked last period will work this period approaches one. So finding significant serial correlation may say very little about the underlying process. Or consider a spell that begins near the end of a period; then it is likely to overlap into the next period, so that previous employment raises the probability of current employment.

Consider the underlying process of time spent in one state followed by time spent in the other state. If the individual's history does not help to predict his future given his current state, then this is a Markov process. Whereas serial independence in continuous time has the absurd implication that mean duration of a spell is zero, the Markov property does provide a fruitful starting point. It has two implications: the individual's history prior to the current spell should not affect the distribution of the length of the current spell; and the amount of time spent in the current state should not affect the distribution of remaining time in that state.

So the first requirement of the Markov property is that durations of the spells be independent of each other. Assuming stationarity, this implies an alternating renewal process. The second requirement is that the distribution of duration be exponential, so that we have an alternating Poisson process. We shall refer to departures from this model as duration dependence.

A test of this Markov property using binary sequences will depend upon what sampling scheme is being used. The simplest case is point sampling, where each period we determine the individual's state at a particular point in time, such as July 1 of each year. Then if an individual is following an alternating Poisson process, her history prior to that point is irrelevant in predicting her state at the next interview. So the binary sequence generated by point sampling should be a Markov chain.

It is possible to test this in a fixed effects model that allows each individual to have her own two exponential rate parameters (c_{i1}, c_{i2}) in the alternating Poisson process. The idea is related to the conditional likelihood approach in the fixed effects logit model. Let s_{ijk} be the number of times that individual i is observed making a transition from state j to state k ($j, k = 1, 2$). Then the initial state and these four transition counts are sufficient statistics for the Markov chain. Sequences with the same initial state and the same transition counts should be equally likely. This is the Markov form of de Finetti's (1975) partial exchangeabil-

²⁴This point is discussed in Singer and Spilerman (1974, 1976).

ity.²⁵ So we can test whether the Markov property holds conditional on c_{i1}, c_{i2} by testing whether there is significant variation in the sample frequencies of sequences with the same transition counts.

This analysis is relevant if, for example, each year the survey question is: "Did you have a job on July 1?" In the Michigan Panel Study of Income Dynamics, however, the most commonly used question for generating participation sequences is: "Did your wife do any work for money last year?" This interval sampling leads to a more complex analysis, since even if the individual is following an alternating Poisson process, the binary sequence generated by this sampling scheme is not a Markov chain. Suppose that $y_{t-1} = 1$, so that we know that the individual worked at some point during the previous period. What is relevant, however, is the individual's state at the end of the period, and y_{t-2} will affect the probability that the spell of work occurred early in period $t-1$ instead of late in the period.

Nevertheless, it is possible to test whether the underlying process is alternating Poisson. The reason is that if $y_{t-1} = 0$, we know that the individual never worked during period $t-1$, and so we know the state at the end of that period; hence y_{t-2}, y_{t-3}, \dots are irrelevant. So we have:

$$\begin{aligned} P(y_t = 1 | c_1, c_2, y_{t-1}, y_{t-2}, \dots) \\ &= P(y_t = 1 | c_1, c_2, y_{t-1} = \dots = y_{t-d} = 1, y_{t-d-1} = 0) \\ &= P(y_t = 1 | c_1, c_2, d), \end{aligned}$$

where d is the number of consecutive preceding periods that the individual was in state 1.

Let s_{01} be the number of times in the sequence that 1 is preceded by 0; let s_{011} be the number of times that 1 is preceded by 0, 1; etc. Then sufficient statistics are s_{01}, s_{011}, \dots , as well as the number of consecutive ones at the beginning (n_1) and at the end (n_T) of a sequence.²⁶ For an example with $T = 5$, let $n_1 = 0, n_5 = 0, s_{01} = 1, s_{011} = 1, s_{0111} = \dots = 0$; then we have

$$\begin{aligned} P(0, 1, 1, 0, 0 | c) \\ &= P(y_1 = 0 | c) P(1 | 0, c) P(1 | 0, 1, c) P(0 | 0, 1, 1, c) P(0 | 0, c); \\ P(0, 0, 1, 1, 0 | c) \\ &= P(y_1 = 0 | c) P(0 | 0, c) P(1 | 0, c) P(1 | 0, 1, c) P(0 | 0, 1, 1, c), \end{aligned}$$

²⁵We are using the fact that partial exchangeability is a necessary condition for the distribution to be a mixture of Markov chains. Diaconis and Freedman (1980) study the sufficiency of this condition. Heckman (1978) used exchangeability to test for serial independence in a fixed effects model.

²⁶This test was presented in Chamberlain (1978a, 1979). It has been applied to unemployment sequences by Corcoran and Hill (1980). For related tests and extensions, see Lee (1980).

where $\mathbf{c} = (c_1, c_2)$. Thus these two sequences are equally likely conditional on \mathbf{c} , and letting μ be the probability measure for \mathbf{c} gives:

$$\begin{aligned} P(0, 1, 1, 0, 0) &= \int P(0, 1, 1, 0, 0 | \mathbf{c}) \mu(d\mathbf{c}) \\ &= \int P(0, 0, 1, 1, 0 | \mathbf{c}) \mu(d\mathbf{c}) = P(0, 0, 1, 1, 0). \end{aligned}$$

So the alternating Poisson process implies restrictions on the multinomial distribution for the binary sequence.

These tests are indirect. The duration dependence question is clearly easier to answer using surveys that measure durations of spells. Such duration data raises a number of new econometric problems, but we shall not pursue them here.²⁷ I would simply like to make one connection with the methods that we have been discussing.

Let us simplify to a one state process; for example, y_{it} can be the duration of the time interval between the starting date of the i th individual's t th job and his $(t+1)$ th job. Suppose that we observe $T > 1$ jobs for each of the N individuals, a not innocuous assumption. Impose the restriction that $y_{it} > 0$ by using the following specification:

$$\begin{aligned} y_{it} &= \exp(\beta x_{it} + c_i + u_{it}), \\ E^*(u_{it} | \mathbf{x}_i) &= 0 \quad (t = 1, \dots, T), \end{aligned}$$

where $\mathbf{x}'_i = (x_{i1}, \dots, x_{iT})$. Then:

$$E^*(\ln y_{it} | \mathbf{x}_i) = \beta x_{it} + \lambda x_i,$$

and our Section 2 analysis applies. The strict exogeneity assumption has a surprising implication in this context. Suppose that x_{it} is the individual's age at the beginning of the t th job. Then $x_{it} - x_{i,t-1} = y_{i,t-1}$ —age is not strictly exogenous.²⁸

4. Inference

Consider a sample $\mathbf{r}'_i = (\mathbf{x}'_i, \mathbf{y}'_i)$, $i = 1, \dots, N$, where $\mathbf{x}'_i = (x_{i1}, \dots, x_{iK})$, $\mathbf{y}'_i = (y_{i1}, \dots, y_{iM})$. We shall assume that \mathbf{r}_i is independent and identically distributed (i.i.d.) according to some multivariate distribution with finite fourth moments

²⁷See Tuma (1979, 1980), Lancaster (1979), Nickell (1979), Chamberlain (1979), Lancaster and Nickell (1980), Heckman and Borjas (1980), Kiefer and Neumann (1981), and Flinn and Heckman (1982, 1982a).

²⁸This example is based on Chamberlain (1979).

and $E(\mathbf{x}_i \mathbf{x}_i')$ nonsingular. Consider the minimum mean-square error linear predictors,²⁹

$$E^*(y_{im}|\mathbf{x}_i) = \pi'_m \mathbf{x}_i \quad (m=1, \dots, M),$$

which we can write as:

$$E^*(y_i|\mathbf{x}_i) = \Pi \mathbf{x}_i, \quad \Pi = E(y_i \mathbf{x}_i') [E(\mathbf{x}_i \mathbf{x}_i')]^{-1}.$$

We want to estimate Π subject to restrictions and to test those restrictions. For example, we may want to test whether a submatrix of Π has the form $\beta \mathbf{I} + \lambda \mathbf{X}$.

We shall not assume that the regression function $E(y_i|\mathbf{x}_i)$ is linear. For although $E(y_i|\mathbf{x}_i, c_i)$ may be linear (indeed, we hope that it is), there is generally no reason to insist that $E(c_i|\mathbf{x}_i)$ is linear. So we shall present a theory of inference for linear predictors. Furthermore, even if the regression function is linear, there may be heteroskedasticity—due to random coefficients, for example. So we shall allow $E[(y_i - \Pi \mathbf{x}_i)(y_i - \Pi \mathbf{x}_i)'|\mathbf{x}_i]$ to be an arbitrary function of \mathbf{x}_i .

4.1. The estimation of linear predictors

Let \mathbf{w}_i be the vector formed from the distinct elements of $\mathbf{r}_i \mathbf{r}_i'$ that have nonzero variance.³⁰ Since $\mathbf{r}_i' = (\mathbf{x}_i', y_i')$ is i.i.d., it follows that \mathbf{w}_i is i.i.d. This simple observation is the key to our results. Since Π is a function of $E(\mathbf{w}_i)$, our problem is to make inferences about a function of a population mean, under random sampling.

Let $\boldsymbol{\mu} = E(\mathbf{w}_i)$ and let $\boldsymbol{\pi}$ be the vector formed from the columns of Π' ($\boldsymbol{\pi} = \text{vec}(\Pi')$). Then $\boldsymbol{\pi}$ is a function of $\boldsymbol{\mu}$: $\boldsymbol{\pi} = \mathbf{h}(\boldsymbol{\mu})$. Let $\bar{\mathbf{w}} = \sum_{i=1}^N \mathbf{w}_i / N$; then $\hat{\boldsymbol{\pi}} = \mathbf{h}(\bar{\mathbf{w}})$ is the least squares estimator:

$$\hat{\boldsymbol{\pi}} = \text{vec} \left[\left(\sum_{i=1}^N \mathbf{x}_i \mathbf{x}_i' \right)^{-1} \sum_{i=1}^N \mathbf{x}_i y_i' \right].$$

By the strong law of large numbers, $\bar{\mathbf{w}}$ converges almost surely to $\boldsymbol{\mu}^0$ as $N \rightarrow \infty$ ($\bar{\mathbf{w}} \xrightarrow{\text{a.s.}} \boldsymbol{\mu}^0$), where $\boldsymbol{\mu}^0$ is the true value of $\boldsymbol{\mu}$. Let $\boldsymbol{\pi}^0 = \mathbf{h}(\boldsymbol{\mu}^0)$. Since $\mathbf{h}(\boldsymbol{\mu})$ is continuous at $\boldsymbol{\mu} = \boldsymbol{\mu}^0$, we have $\hat{\boldsymbol{\pi}} \xrightarrow{\text{a.s.}} \boldsymbol{\pi}^0$. The central limit theorem implies that:

$$\sqrt{N}(\bar{\mathbf{w}} - \boldsymbol{\mu}^0) \xrightarrow{D} N(\mathbf{0}, V(\mathbf{w}_i)).$$

²⁹This agrees with the definition in Section 2 if \mathbf{x}_i includes a constant.

³⁰Sections 4.1–4.4 are taken from Chamberlain (1982a).

Since $\mathbf{h}(\boldsymbol{\mu})$ is differentiable at $\boldsymbol{\mu} = \boldsymbol{\mu}^0$, the δ -method gives

$$\sqrt{N}(\hat{\boldsymbol{\pi}} - \boldsymbol{\pi}^0) \xrightarrow{D} N(\mathbf{0}, \boldsymbol{\Omega}),$$

where

$$\boldsymbol{\Omega} = \frac{\partial \mathbf{h}(\boldsymbol{\mu}^0)}{\partial \boldsymbol{\mu}'} V(\mathbf{w}_i) \frac{\partial \mathbf{h}'(\boldsymbol{\mu}^0)}{\partial \boldsymbol{\mu}}.^{31}$$

We have derived the limiting distribution of the least squares estimator. This approach was used by Cramér (1946) to obtain limiting normal distributions for sample correlation and regression coefficients (p. 367); he presents an explicit formula for the variance of the limiting distribution of a sample correlation coefficient (p. 359). Kendall and Stuart (1961, p. 293) and Goldberger (1974) present the formula for the variance of the limiting distribution of a simple regression coefficient.

Evaluating the partial derivatives in the formula for $\boldsymbol{\Omega}$ is tedious. That calculation can be simplified since $\hat{\boldsymbol{\pi}}$ has a “ratio” form. In the case of simple regression with a zero intercept, we have $\boldsymbol{\pi} = E(y_i x_i) / E(x_i^2)$ and

$$\sqrt{N}(\hat{\boldsymbol{\pi}} - \boldsymbol{\pi}^0) = \left(\sum_{i=1}^N y_i x_i - \boldsymbol{\pi}^0 \sum_{i=1}^N x_i^2 \right) / \left[\sqrt{N} \left(\sum_{i=1}^N x_i^2 / N \right) \right].$$

Since $\sum_{i=1}^N x_i^2 / N \xrightarrow{\text{a.s.}} E(x_i^2)$, we obtain the same limiting distribution by working with

$$\sum_{i=1}^N [(y_i - \pi^0 x_i) x_i] / [\sqrt{N} E(x_i^2)].$$

The definition of $\boldsymbol{\pi}^0$ gives $E[(y_i - \pi^0 x_i) x_i] = 0$, and so the central limit theorem implies that:

$$\sqrt{N}(\hat{\boldsymbol{\pi}} - \boldsymbol{\pi}^0) \xrightarrow{D} N\left\{0, E\left[(y_i - \pi^0 x_i)^2 x_i^2\right] / [E(x_i^2)]^2\right\}.$$

This approach was used by White (1980) to obtain the limiting distribution for univariate regression coefficients.³² In the Appendix (Proposition 6) we follow

³¹See Billingsley (1979, example 29.1, p. 340) or Rao (1973, p. 388).

³²Also see White (1980a,b).

White's approach to obtain:

$$\Omega = E\left[(y_i - \Pi^0 x_i)(y_i - \Pi^0 x_i)' \otimes \Phi_x^{-1}(x_i x_i') \Phi_x^{-1}\right], \quad (4.1)$$

where $\Phi_x = E(x_i x_i')$. A consistent estimator of Ω is readily available from the corresponding sample moments:

$$\begin{aligned} \hat{\Omega} &= \frac{1}{N} \sum_{i=1}^N \left[(y_i - \hat{\Pi} x_i)(y_i - \hat{\Pi} x_i)' \otimes S_x^{-1}(x_i x_i') S_x^{-1} \right] \\ &\xrightarrow{\text{a.s.}} \Omega, \end{aligned} \quad (4.2)$$

where $S_x = \sum_{i=1}^N x_i x_i' / N$.

If $E(y_i | x_i) = \Pi x_i$, so that the regression function is linear, then:

$$\Omega = E[V(y_i | x_i) \otimes \Phi_x^{-1}(x_i x_i') \Phi_x^{-1}].$$

If $V(y_i | x_i)$ is uncorrelated with $x_i x_i'$, then:

$$\Omega = E[V(y_i | x_i)] \otimes \Phi_x^{-1}.$$

If the conditional variance is homoskedastic, so that $V(y_i | x_i) = \Sigma$ does not depend on x_i , then:

$$\Omega = \Sigma \otimes \Phi_x^{-1}.$$

4.2. Imposing restrictions: The minimum distance estimator

Since Π is a function of $E(w_i)$, restrictions on Π imply restrictions on $E(w_i)$. Let the dimension of $\mu = E(w_i)$ be q .³³ We shall specify the restrictions by the condition that μ depends only on a $p \times 1$ vector θ of unknown parameters: $\mu = g(\theta)$, where g is a known function and $p \leq q$. The domain of θ is T , a subset of p -dimensional Euclidean space (R^p) that contains the true value θ^0 . So the restrictions imply that $\mu^0 = g(\theta^0)$ is confined to a certain subset of R^q .

We can impose the restrictions by using a minimum distance estimator: choose $\hat{\theta}$ to

$$\min_{\theta \in T} \sum_{i=1}^N [w_i - g(\theta)]' A_N [w_i - g(\theta)],$$

³³ If there is one element in $r_i r_i'$ with zero variance, then $q = [(K + M)(K + M + 1)/2] - 1$.

where $A_N \xrightarrow{\text{a.s.}} \Psi$ and Ψ is positive definite.³⁴ This minimization problem is equivalent to the following one: choose $\hat{\theta}$ to

$$\min_{\theta \in T} [\bar{w} - g(\theta)]' A_N [\bar{w} - g(\theta)].$$

The properties of $\hat{\theta}$ are developed, for example, in Malinvaud (1970, ch. 9). Since g does not depend on any exogenous variables, the derivation of these properties can be simplified considerably, as in Chiang (1956) and Ferguson (1958).³⁵

For completeness, we shall state a set of regularity conditions and the properties that they imply:

Assumption 1

$a_N \xrightarrow{\text{a.s.}} g(\theta^0)$; T is a compact subset of R^p that contains θ^0 ; g is continuous on T , and $g(\theta) = g(\theta^0)$ for $\theta \in T$ implies that $\theta = \theta^0$; $A_N \xrightarrow{\text{a.s.}} \Psi$, where Ψ is positive definite.

Assumption 2

$\sqrt{N}[a_N - g(\theta^0)] \xrightarrow{D} N(\theta, \Delta)$; T contains a neighborhood Ξ_0 of θ^0 in which g has continuous second partial derivatives; $\text{rank}(G) = p$, where $G = \partial g(\theta^0)/\partial \theta'$.

Choose $\hat{\theta}$ to

$$\min_{\theta \in T} [a_N - g(\theta)]' A_N [a_N - g(\theta)].$$

Proposition 1

If Assumption 1 is satisfied, then $\hat{\theta} \xrightarrow{\text{a.s.}} \theta^0$.

Proposition 2

If Assumptions 1 and 2 are satisfied, then $\sqrt{N}(\hat{\theta} - \theta^0) \xrightarrow{D} N(\theta, \Lambda)$, where

$$\Lambda = (G' \Psi G)^{-1} G' \Psi \Delta \Psi G (G' \Psi G)^{-1}.$$

If Δ is positive definite, then $\Lambda - (G' \Delta^{-1} G)^{-1}$ is positive semi-definite; hence an optimal choice for Ψ is Δ^{-1} .

³⁴This application of nonlinear generalized least squares was proposed in Chamberlain (1980a).

³⁵Some simple proofs are collected in Chamberlain (1982a).

Proposition 3

If Assumptions 1 and 2 are satisfied, if Δ is a $q \times q$ positive-definite matrix, and if $A_N \xrightarrow{\text{a.s.}} \Delta^{-1}$, then:

$$N[a_N - g(\hat{\theta})]'A_N[a_N - g(\hat{\theta})] \xrightarrow{D} \chi^2(q-p).$$

Now consider imposing additional restrictions, which are expressed by the condition that $\theta = f(\alpha)$, where α is $s \times 1$ ($s \leq p$). The domain of α is T_1 , a subset of R^s that contains the true value α^0 . So $\theta^0 = f(\alpha^0)$ is confined to a certain subset of R^p .

Assumption 2'

T_1 is a compact subset of R^s that contains α^0 ; f is a continuous mapping from T_1 into T ; $f(\alpha) = \theta^0$ for $\alpha \in T_1$ implies $\alpha = \alpha^0$; T_1 contains a neighborhood of α^0 in which f has continuous second partial derivatives; $\text{rank}(F) = s$, where $F = \partial f(\alpha^0) / \partial \alpha'$.

Let $h(\alpha) = g[f(\alpha)]$. Choose $\hat{\alpha}$ to

$$\min_{\alpha \in T_1} [a_N - h(\alpha)]'A_N[a_N - h(\alpha)].$$

Proposition 3'

If Assumptions 1, 2, and 2' are satisfied, if Δ is positive definite, and if $A_N \xrightarrow{\text{a.s.}} \Delta^{-1}$, then $d_1 - d_2 \xrightarrow{D} \chi^2(p-s)$, where

$$d_1 = N[a_N - h(\hat{\alpha})]'A_N[a_N - h(\hat{\alpha})],$$

$$d_2 = N[a_N - g(\hat{\theta})]'A_N[a_N - g(\hat{\theta})].$$

Furthermore, $d_1 - d_2$ is independent of d_2 in their limiting joint distribution.

Suppose that the restrictions involve only II . We specify the restrictions by the condition that $\pi = f(\delta)$, where δ is $s \times 1$ and the domain of δ is T_1 , a subset of R^s that includes the true value δ^0 . Consider the following estimator of δ^0 : choose $\hat{\delta}$ to

$$\min_{\delta \in T_1} [\hat{\pi} - f(\delta)]'\hat{\Omega}^{-1}[\hat{\pi} - f(\delta)],$$

where $\hat{\Omega}$ is given in (4.2), and we assume that Ω in (4.1) is positive definite. If T_1

and f satisfy Assumptions 1 and 2, then $\hat{\delta} \xrightarrow{\text{a.s.}} \delta^0$,

$$\sqrt{N}(\hat{\delta} - \delta^0) \xrightarrow{D} N(\theta, [F' \Omega^{-1} F]^{-1}),$$

and

$$N[\hat{\pi} - f(\hat{\delta})]' \hat{\Omega}^{-1} [\hat{\pi} - f(\hat{\delta})] \xrightarrow{D} \chi^2(MK - s),$$

where $F = \partial f(\delta^0) / \partial \delta'$.

We can also estimate δ^0 by applying the minimum distance procedure to \bar{w} instead of to $\hat{\pi}$. Suppose that the components of w_i are arranged so that $w_i' = (w_{i1}', w_{i2}')$, where w_{i1} contains the components of $x_i y_i'$. Partition $\mu = E(w_i)$ conformably: $\mu' = (\mu_1', \mu_2')$. Set $\theta' = (\theta_1', \theta_2') = (\delta', \mu_2')$. Assume that $V(w_i)$ is positive definite. Now choose $\hat{\theta}$ to

$$\min_{\theta \in T} [\bar{w} - g(\theta)]' A_N [\bar{w} - g(\theta)],$$

where $A_N \xrightarrow{\text{a.s.}} V^{-1}(w_i)$,

$$g(\theta) = \begin{bmatrix} g_1[f(\delta), \mu_2] \\ \mu_2 \end{bmatrix},$$

and $g_1(\pi, \mu_2) = \mu_1$. Then $\hat{\theta}_1$ gives an estimator of δ^0 ; it has the same limiting distribution as the estimator $\hat{\delta}$ that we obtained by applying the minimum distance procedure to $\hat{\pi}$.³⁶

This framework leads to some surprising results on efficient estimation. For a simple example, we shall use a univariate linear predictor model,

$$E^*(y_i | x_{i1}, x_{i2}) = \pi_0 + \pi_1 x_{i1} + \pi_2 x_{i2}.$$

Consider imposing the restriction $\pi_2 = 0$. Then the conventional estimator of π_1 is b_{yx_1} , the slope coefficient in the least squares regression of y on x_1 . We shall show that this estimator is generally less efficient than the minimum distance estimator if the regression function is nonlinear or if there is heteroskedasticity.

Let $\hat{\pi}_1, \hat{\pi}_2$ be the slope coefficients in the least squares multiple regression of y on x_1, x_2 . The minimum distance estimator of π_1 under the restriction $\pi_2 = 0$ can be obtained as $\hat{\delta} = \hat{\pi}_1 + \tau \hat{\pi}_2$, where τ is chosen to minimize the (estimated)

³⁶See Chamberlain (1982a, proposition 9).

variance of the limiting distribution of $\hat{\delta}$; this gives:

$$\hat{\delta} = \hat{\pi}_1 - \frac{\hat{\omega}_{12}}{\hat{\omega}_{22}} \hat{\pi}_2,$$

where $\hat{\omega}_{jk}$ is the estimated covariance between $\hat{\pi}_j$ and $\hat{\pi}_k$ in their limiting distribution. Since $\hat{\pi}_1 = b_{yx_1} - \hat{\pi}_2 b_{x_2x_1}$, we have:

$$\hat{\delta} = b_{yx_1} - \left(b_{x_2x_1} + \frac{\hat{\omega}_{12}}{\hat{\omega}_{22}} \right) \hat{\pi}_2.$$

If $E(y_i|x_{i1}, x_{i2})$ is linear and if $V(y_i|x_{i1}, x_{i2}) = \sigma^2$, then $\omega_{12}/\omega_{22} = -\text{cov}(x_{i1}, x_{i2})/V(x_{i1})$ and $\hat{\delta} = b_{yx_1}$. But in general $\hat{\delta} \neq b_{yx_1}$ and $\hat{\delta}$ is more efficient than b_{yx_1} . The source of the efficiency gain is that the limiting distribution of $\hat{\pi}_2$ has a zero mean (if $\pi_2 = 0$), and so we can reduce variance without introducing any bias if $\hat{\pi}_2$ is correlated with b_{yx_1} . Under the assumptions of linear regression and homoskedasticity, b_{yx_1} and $\hat{\pi}_2$ are uncorrelated; but this need not be true in the more general framework that we are using.

4.3. Simultaneous equations: A generalization of three-stage least squares

Given the discussion on imposing restrictions, it is not surprising that two-stage least squares is not, in general, an efficient procedure for combining instrumental variables. Also, three-stage least squares, viewed as a minimum distance estimator, is using the wrong norm in general.

Consider the standard simultaneous equations model:

$$y_i = \Pi x_i + u_i, \quad E(u_i x_i') = 0,$$

$$\Gamma y_i + B x_i = v_i,$$

where $\Gamma \Pi + B = 0$ and $\Gamma u_i = v_i$. We are continuing to assume that y_i is $M \times 1$, x_i is $K \times 1$, $r_i' = (x_i', y_i')$ is i.i.d. according to a distribution with finite fourth moments ($i=1, \dots, N$), and that $E(x_i x_i')$ is nonsingular. There are restrictions on Γ and B : $m(\Gamma, B) = 0$, where m is a known function. Assume that the implied restrictions on Π can be specified by the condition that $\pi = \text{vec}(\Pi') = f(\delta)$, where the domain of δ is T_1 , a subset of R^s that includes the true value δ^0 ($s \leq MK$). Assume that T_1 and f satisfy assumptions 1 and 2; these properties could be derived from regularity conditions on m , as in Malinvaud (1970, proposition 2, p. 670).

Choose $\hat{\delta}$ to

$$\min_{\delta \in T_1} [\hat{\pi} - f(\delta)]' \hat{\Omega}^{-1} [\hat{\pi} - f(\delta)],$$

where $\hat{\Omega}$ is given in (4.2) and we assume that Ω in (4.1) is positive definite. Let $F = \partial f(\delta^0) / \partial \delta'$. Then we have $\sqrt{N}(\hat{\delta} - \delta^0) \xrightarrow{D} N(0, \Lambda)$, where $\Lambda = (F' \Omega^{-1} F)^{-1}$. This generalizes Malinvaud's minimum distance estimator (p. 676); it reduces to his estimator if $u_i^0 u_i^{0'}$ is uncorrelated with $x_i x_i'$, so that $\Omega = E(u_i^0 u_i^{0'}) \otimes [E(x_i x_i')]^{-1}$ ($u_i^0 = y_i - \Pi^0 x_i$).

Now suppose that the only restrictions on Γ and B are that certain coefficients are zero, together with the normalization restrictions that the coefficient of y_{im} in the m th structural equation is one. Then we can give an explicit formula for Λ . Write the m th structural equation as:

$$y_{im} = \delta_m' z_{im} + v_{im},$$

where the components of z_{im} are the variables in y_i and x_i that appear in the m th equation with unknown coefficients. Let there be M structural equations and assume that the true value Γ^0 is nonsingular. Let S_{zx} be the following block-diagonal matrix:

$$S_{zx} = \text{diag} \left\{ \frac{1}{N} \sum_{i=1}^N z_{i1} x_i', \dots, \frac{1}{N} \sum_{i=1}^N z_{iM} x_i' \right\},$$

and $s_{xy} = N^{-1} \sum_{i=1}^N y_i \otimes x_i$. Let $v_i^{0'} = (v_{i1}^0, \dots, v_{iM}^0)$, where $v_{im}^0 = y_{im} - \delta_m^{0'} z_{im}$ and δ_m^0 is the true value; let $\Phi_{zx} = E(S_{zx})$ and $\Phi_x = E(x_i x_i')$. Let $\delta' = (\delta_1', \dots, \delta_M')$. Then we have:

$$\Lambda = \left\{ \Phi_{zx} \left[E(v_i^0 v_i^{0'} \otimes x_i x_i') \right]^{-1} \Phi_x' \right\}^{-1}.^{37} \quad (4.3)$$

If $u_i^0 u_i^{0'}$ is uncorrelated with $x_i x_i'$, then this reduces to:

$$\Lambda = \left\{ \Phi_{zx} \left[E^{-1}(v_i^0 v_i^{0'}) \otimes \Phi_x^{-1} \right] \Phi_x' \right\}^{-1},$$

which is the conventional asymptotic covariance matrix for three-stage least squares [Zellner and Theil (1962)].

There is a generalization of three-stage least squares that has the same limiting distribution as the generalized minimum distance estimator. Let $\hat{\Psi} = N^{-1} \sum_{i=1}^N (\hat{v}_i \hat{v}_i' \otimes x_i x_i')$, where $\hat{v}_i = \hat{\Gamma} y_i + \hat{B} x_i$ and $\hat{\Gamma} \xrightarrow{\text{a.s.}} \Gamma^0$, $\hat{B} \xrightarrow{\text{a.s.}} B^0$. Define:

$$\hat{\delta}_{G3} = (S_{zx} \hat{\Psi}^{-1} S_{zx}')^{-1} (S_{zx} \hat{\Psi}^{-1} s_{xy}).$$

³⁷See Chamberlain (1982a).

The limiting distribution of this estimator is derived in the Appendix (Proposition 6). We record it as:

Proposition 4

$\sqrt{N}(\hat{\delta}_{G3} - \delta^0) \xrightarrow{D} N(\theta, \Lambda)$, where Λ is given in (4.3). This generalized three-stage least squares estimator is asymptotically efficient within the class of minimum distance estimators.

Our derivation of the limiting distribution of $\hat{\delta}_{G3}$ relies on linearity. For a generalized nonlinear three-stage least squares estimator, see Hansen (1982).³⁸

4.4. Asymptotic efficiency: A comparison with the quasi-maximum likelihood estimator

Assume that r_i is i.i.d. ($i = 1, 2, \dots$) from a distribution with $E(r_i) = \tau$, $V(r_i) = \Sigma$, where Σ is a $J \times J$ positive-definite matrix; the fourth moments are finite. Suppose that we wish to estimate functions of Σ subject to restrictions. Let $\sigma = \text{vec}(\Sigma)$ and express the restrictions by the condition that $\sigma = g(\theta)$, where g is a function from T into R^q with a domain $T \subset R^p$ that contains the true value θ^0 ($q = J^2$; $p \leq J(J+1)/2$). Let

$$\bar{S} = \frac{1}{N} \sum_{i=1}^N (r_i - \bar{r})(r_i - \bar{r})',$$

and let $\bar{s} = \text{vec}(\bar{S})$.

If the distribution of r_i is multivariate normal, then the log-likelihood function is:

$$L = \frac{N}{2} \ln |\Sigma^{-1}| - \frac{N}{2} \text{tr} \left\{ \Sigma^{-1} \left[\bar{S} + (\bar{r} - \tau)(\bar{r} - \tau)' \right] \right\}.$$

If there are no restrictions on τ , then the maximum likelihood estimator of θ^0 is a solution to the following problem: Choose $\hat{\theta}$ to solve:

$$\frac{\partial g'(\theta)}{\partial \theta} [\Sigma^{-1}(\theta) \otimes \Sigma^{-1}(\theta)] (\bar{s} - g(\theta)) = 0.$$

We shall derive the properties of this estimator when the distribution of r_i is not necessarily normal; in that case we shall refer to the estimator as a quasi-maximum likelihood estimator ($\hat{\theta}_{QML}$).³⁹

³⁸There are generalizations of two-stage least squares in Chamberlain (1982a) and White (1982a).

³⁹The quasi-maximum likelihood terminology was used by the Cowles Commission; see Malinvaud (1970, p. 678).

MaCurdy (1979) considered a version of this problem and showed that, under suitable regularity conditions, $\sqrt{N}(\hat{\theta}_{\text{QML}} - \theta^0)$ has a limiting normal distribution; the covariance matrix, however, is not given by the standard information matrix formula. We would like to compare this distribution with the distribution of the minimum distance estimator.

This comparison can be readily made by using theorem 1 in Ferguson (1958). In our notation, Ferguson considers the following problem: Choose $\hat{\theta}$ to solve

$$W(\bar{s}, \theta)[\bar{s} - g(\theta)] = 0.$$

He derives the limiting distribution of $\sqrt{N}(\hat{\theta} - \theta^0)$ under regularity conditions on the functions W and g . These regularity conditions are particularly simple in our problem since W does not depend on \bar{s} . We can state them as follows:

Assumption 3

$\Xi_0 \subset R^p$ is an open set containing θ^0 ; g is a continuous, one-to-one mapping of Ξ_0 into R^q with a continuous inverse; g has continuous second partial derivatives in Ξ_0 ; $\text{rank} [\partial g(\theta)/\partial \theta'] = p$ for $\theta \in \Xi_0$; $\Sigma(\theta)$ is nonsingular for $\theta \in \Xi_0$.

In addition, we shall need $\bar{s} \xrightarrow{\text{a.s.}} g(\theta^0)$ and the central limit theorem result that

$$\sqrt{N}(\bar{s} - g(\theta^0)) \xrightarrow{D} N(\theta, \Delta), \text{ where } \Delta = V[(r_i - \tau^0) \otimes (r_i - \tau^0)].$$

Then Ferguson's theorem implies that the likelihood equations almost surely have a unique solution within Ξ_0 for sufficiently large N , and $\sqrt{N}(\hat{\theta}_{\text{QML}} - \theta^0) \xrightarrow{D} N(\theta, \Delta)$, where

$$\Delta = (G' \Psi G)^{-1} G' \Psi \Delta \Psi G (G' \Psi G)^{-1},$$

and $G = \partial g(\theta^0)/\partial \theta'$, $\Psi = (\Sigma^0 \otimes \Sigma^0)^{-1}$. It will be convenient to rewrite this, imposing the symmetry restrictions on Σ . Let σ^* be the $J(J+1)/2 \times 1$ vector formed by stacking the columns of the lower triangle of Σ . We can define a $J^2 \times [J(J+1)/2]$ matrix T such that $\sigma = T\sigma^*$. The elements in each row of T are all zero except for a single element which is one; T has full column rank. Let $\bar{s} = T\bar{s}^*$, $g(\theta) = Tg^*(\theta)$, $G^* = \partial g^*(\theta^0)/\partial \theta'$, $\Psi^* = T' \Psi T$; then $\sqrt{N}[\bar{s}^* - g^*(\theta^0)] \xrightarrow{D} N(\theta, \Delta^*)$, where Δ^* is the covariance matrix of the vector formed from the columns of the lower triangle of $(r_i - \tau^0)(r_i - \tau^0)'$. Now we can set

$$\Delta = (G'^* \Psi^* G^*)^{-1} (G'^* \Psi^* \Delta^* \Psi^* G^*) (G'^* \Psi^* G^*)^{-1}.$$

Consider the following minimum distance estimator. Choose $\hat{\theta}_{\text{MD}}$ to

$$\min_{\theta \in T} [\bar{s}^* - g^*(\theta)]' A_N [\bar{s}^* - g^*(\theta)],$$

where \mathcal{T} is a compact subset of Ξ_0 that contains a neighborhood of θ^0 and $A_N \xrightarrow{\text{a.s.}} \Psi^*$. Then the following result is implied by Proposition 2.

Proposition 5

If Assumption 3 is satisfied, then $\sqrt{N}(\hat{\theta}_{\text{QML}} - \theta^0)$ has the same limiting distribution as $\sqrt{N}(\hat{\theta}_{\text{MD}} - \theta^0)$.

If Δ^* is nonsingular, an optimal minimum distance estimator has $A_N \xrightarrow{\text{a.s.}} \zeta \Delta^{*-1}$, where ζ is an arbitrary positive real number. If the distribution of r_i is normal, then $\Delta^{*-1} = (1/2)\Psi^*$; but in general Δ^{*-1} is not proportional to Ψ^* , since Δ^* depends on fourth moments and Ψ^* is a function of second moments. So in general $\hat{\theta}_{\text{QML}}$ is less efficient than the optimal minimum distance estimator that uses

$$A_N = \left[\frac{1}{N} \sum_{i=1}^N (s_i^* - \bar{s}^*)(s_i^* - \bar{s}^*)' \right]^{-1}, \quad (4.4)$$

where s_i^* is the vector formed from the lower triangle of $(r_i - \bar{r})(r_i - \bar{r})'$.

More generally, we can consider the class of consistent estimators that are continuously differentiable functions of \bar{s}^* : $\hat{\theta} = \hat{\theta}(\bar{s}^*)$. Chiang (1956) shows that the minimum distance estimator based on Δ^{*-1} has the minimal asymptotic covariance matrix within this class. The minimum distance estimator based on A_N in (4.4) attains this lower bound.

4.5. Multivariate probit models

Suppose that

$$\begin{aligned} y_{im} &= 1, & \text{if } \pi'_m x_i + u_{im} \geq 0, \\ &= 0, & \text{otherwise} \quad (i = 1, \dots, N; m = 1, \dots, M), \end{aligned}$$

where the distribution of $u'_i = (u_{i1}, \dots, u_{iM})$ conditional on x_i is multivariate normal, $N(\theta, \Sigma)$. There may be restrictions on $\pi' = (\pi'_1, \dots, \pi'_M)$, but we want to allow Σ to be unrestricted, except for the scale normalization that the diagonal elements of Σ are equal to one. In that case, the maximum likelihood estimator has the computational disadvantage of requiring numerical integration over $M - 1$ dimensions.

Our strategy is to avoid numerical integration. We estimate π_m by maximizing the marginal likelihood function that is based on the distribution of y_{im} conditional on x_i :

$$P(y_{im} = 1 | x_i) = F(\pi'_m x_i),$$

where F is the standard normal distribution function. Then under standard assumptions we have $\hat{\pi}_m \xrightarrow{\text{a.s.}} \pi_m^0$, the true value. If $\sqrt{N}(\hat{\pi} - \pi^0) \xrightarrow{D} N(0, \Omega)$, then we can impose the restriction that $\pi = f(\delta)$ by choosing $\hat{\delta}$ to minimize

$$[\hat{\pi} - f(\hat{\delta})]' \hat{\Omega}^{-1} [\hat{\pi} - f(\hat{\delta})].$$

We only need to derive a formula for Ω .⁴⁰

Our estimator of π is solving the following equation:

$$s(\hat{\pi}) = \frac{\partial Q(\hat{\pi})}{\partial \pi} = 0,$$

where

$$Q(\pi) = \sum_{i=1}^N \left\{ \sum_{m=1}^M y_{im} \ln F(\pi'_m x_i) + (1 - y_{im}) \ln [1 - F(\pi'_m x_i)] \right\}.$$

Hence, the asymptotic distribution of $\hat{\pi}$ can be obtained from the theory of “ M -estimators”. Huber (1967) provides general results, which do not impose differentiability restrictions on $s(\pi)$. His results cover, for example, regression estimators based on minimizing the residual sum of absolute deviations. We shall not need this generality here and shall sketch the derivation for the simpler, differentiable case. This case has been considered by Hansen (1982), MaCurdy (1981a), and White (1982).⁴¹

Let z_i be i.i.d. according to a distribution with support $Z \subset R^q$. Let Θ be an open, convex subset of R^p and let $\psi(z, \theta)$ be a function from $Z \times \Theta$ into R^p ; its k th component is $\psi_k(z, \theta)$. For each $\theta \in \Theta$, ψ is a measurable function of z , and there is a $\theta^0 \in \Theta$ with:

$$E[\psi(z_1, \theta^0)] = 0, \quad E[\psi(z_1, \theta^0) \psi'(z_1, \theta^0)] = \Delta < \infty.$$

For each $z \in Z$, ψ is a twice continuously differentiable function of θ . In addition:

$$J = E \left[\frac{\partial \psi(z_1, \theta^0)}{\partial \theta'} \right]$$

is nonsingular, and

$$\left| \frac{\partial \psi_k^2(z, \theta)}{\partial \theta_l \partial \theta_m} \right| \leq h(z) \quad (k, l, m = 1, \dots, p)$$

for $\theta \in \Theta$, where $E[h(z_1)] < \infty$.

⁴⁰For an alternative approach to multivariate probit models, see Avery, Hansen and Hotz (1981).

⁴¹Also see Rao (1973, problem 9, p. 378).

Suppose that we have a (measurable) estimator $\hat{\theta}_N \in \Theta$ such that $\hat{\theta}_N \xrightarrow{\text{a.s.}} \theta^0$ and

$$\sum_{i=1}^N \psi(z_i, \hat{\theta}_N) = \mathbf{0}$$

for sufficiently large N a.s. By Taylor's theorem:

$$\frac{1}{\sqrt{N}} \sum_{i=1}^N \psi_k(z_i, \theta^0) + \left[j'_{Nk} + \frac{1}{2}(\hat{\theta}_N - \theta^0)' C_{Nk} \right] [\sqrt{N}(\hat{\theta}_N - \theta^0)] = \mathbf{0},$$

where

$$j_{Nk} = \frac{1}{N} \sum_{i=1}^N \frac{\partial \psi_k(z_i, \theta^0)}{\partial \theta}, \quad C_{Nk} = \frac{1}{N} \sum_{i=1}^N \frac{\partial^2 \psi_k(z_i, \theta_{Nk}^*)}{\partial \theta \partial \theta'},$$

and θ_{Nk}^* is on the line segment joining $\hat{\theta}_N$ and θ^0 ($k=1, \dots, p$). [The measurability of θ_{Nk}^* follows from lemma 3 of Jennrich (1969).] By the strong law of large numbers, j'_{Nk} converges a.s. to the k th row of J , and

$$\left| \frac{1}{N} \sum_{i=1}^N \frac{\partial^2 \psi_k(z_i, \theta_{Nk}^*)}{\partial \theta_l \partial \theta_m} \right| \leq \frac{1}{N} \sum_{i=1}^N h(z_i) \xrightarrow{\text{a.s.}} E[h(z_1)]$$

($k, l, m=1, \dots, p$). Hence $(\hat{\theta}_N - \theta^0)' C_{Nk} \rightarrow \mathbf{0}$ a.s. and

$$\sqrt{N}(\hat{\theta}_N - \theta^0) = -D_N^{-1} \left[\frac{1}{\sqrt{N}} \sum_{i=1}^N \psi(z_i, \theta^0) \right]$$

for N sufficiently large a.s. where $D_N \xrightarrow{\text{a.s.}} J$. By the central limit theorem,

$$\frac{1}{\sqrt{N}} \sum_{i=1}^N \psi(z_i, \theta^0) \xrightarrow{D} N(\mathbf{0}, \Delta).$$

Hence:

$$\sqrt{N}(\hat{\theta}_N - \theta^0) \xrightarrow{D} N(\mathbf{0}, J^{-1} \Delta J'^{-1}).$$

Applying this result to our multivariate probit estimator gives:

$$\sqrt{N}(\hat{\pi} - \pi^0) \xrightarrow{D} N(\mathbf{0}, J^{-1} \Delta J^{-1}),$$

where $J = \text{diag}\{J_1, \dots, J_M\}$ is a block-diagonal matrix with:

$$J_m = E\left[\left\{(F')^2/[F(1-F)]\right\}x_1x_1'\right]$$

(F and its derivative F' are evaluated at $\pi_m^0x_1$); and

$$\Delta = E[H \otimes x_1x_1'],$$

where the m, n element of the $M \times M$ matrix H is $h_{mn} = e_me_n$ with

$$e_m = \frac{y_{1m} - F}{F(1-F)}F' \quad (m=1, \dots, M)$$

(F and F' are evaluated at $\pi_m^0x_1$). We obtain a consistent estimator ($\hat{\Omega}$) of $J^{-1}\Delta J^{-1}$ by replacing expectations by sample means and using $\hat{\pi}$ in place of π^0 . Then we can apply the minimum distance theory of Section 4.2 to impose restrictions on π .

5. Empirical applications

5.1. Linear models: Union wage effects

We shall present an empirical example that illustrates some of the preceding results.⁴² The data come from the panel of Young Men in the National Longitudinal Survey (Parnes). The sample consists of 1454 young men who were not enrolled in school in 1969, 1970, or 1971, and who had complete data on the variables listed in Table 5.1. Table 5.2(a) presents an unrestricted least squares regression of the logarithm of wage in 1969 on the union, *SMSA*, and region variables for all three years. The regression also includes a constant, schooling, experience, experience squared, and race. This regression is repeated using the 1970 wage and the 1971 wage.

In Section 2 we discussed the implications of a random intercept (c). If the leads and lags are due just to c , then the submatrices of Π corresponding to the union, *SMSA*, or region coefficients should have the form $\beta I + I\lambda$. Consider, for example, the 3×3 submatrix of union coefficients—the off-diagonal elements in each column should be equal to each other. So we compare 0.048 to 0.046, 0.042 to 0.041, and -0.009 to 0.010; not bad.

⁴² This application is taken from Chamberlain (1982a).

Table 5.1
 Characteristics of National Longitudinal Survey Young Men,
 not enrolled in school in 1969, 1970, 1971: Means and
 standard deviations, $N=1454$.

Variable	Mean	Standard deviation
<i>LW1</i>	5.64	0.423
<i>LW2</i>	5.74	0.426
<i>LW3</i>	5.82	0.437
<i>U1</i>	0.336	—
<i>U2</i>	0.362	—
<i>U3</i>	0.364	—
<i>U1U2</i>	0.270	—
<i>U1U3</i>	0.262	—
<i>U2U3</i>	0.303	—
<i>U1U2U3</i>	0.243	—
<i>SMSA1</i>	0.697	—
<i>SMSA2</i>	0.627	—
<i>SMSA3</i>	0.622	—
<i>RNS1</i>	0.409	—
<i>RNS2</i>	0.404	—
<i>RNS3</i>	0.410	—
<i>S</i>	11.7	2.64
<i>EXP69</i>	5.11	3.71
<i>EXP69</i> ²	39.8	46.6
<i>RACE</i>	0.264	—

Notes:

LW1, *LW2*, *LW3*—logarithm of hourly earnings (in cents) on the current or last job in 1969, 1970, 1971; *U1*, *U2*, *U3*—1 if wages on current or last job set by collective bargaining, 0 if not, in 1969, 1970, 1971; *SMSA1*, *SMSA2*, *SMSA3*—1 if respondent in *SMSA*, 0 if not, in 1969, 1970, 1971; *RNS1*, *RNS2*, *RNS3*—1 if respondent in South, 0 if not, in 1969, 1970, 1971; *S*—years of schooling completed; *EXP69*—(age in 1969–*S*–6); *RACE*—1 if respondent black, 0 if not.

In Table 5.2(b) we add a complete set of union interactions, so that, for the union variables at least, we have a general regression function. Now the submatrix of union coefficients is 3×7 . If it equals $(\beta I_3, 0) + I\lambda'$, then in the first three columns, the off-diagonal elements within a column should be equal; in the last four columns, all elements within a column should be equal.

I first imposed the restrictions on the *SMSA* and region coefficients, using the minimum distance estimator. Ω is estimated using the formula in (4.2), and $A_N = \hat{\Omega}^{-1}$. The minimum distance statistic (Proposition 3) is 6.82, which is not a surprising value from a $\chi^2(10)$ distribution. If we impose the restrictions on the union coefficients as well, then the 21 coefficients in Table 5.2(b) are replaced by 8: one β and seven λ 's. This gives an increase in the minimum distance statistic (Proposition 3') of $19.36 - 6.82 = 12.54$, which is not a surprising value from a $\chi^2(13)$ distribution. So there is no evidence here against the hypothesis that all the

lags and leads are generated by c . In the terminology of Section 3.3, the (linear predictor) relationship of x to y appears to be static conditional on c .

Consider a transformation of the model in which the dependent variables are $LW1$, $LW2-LW1$, and $LW3-LW2$. Start with a multivariate regression on all of the lags and leads (and union interactions); then impose the restriction that U , $SMSA$, and RNS appear in the $LW2-LW1$ and $LW3-LW2$ equations only as contemporaneous changes ($E(y_t - y_{t-1} | x_1, x_2, x_3) = \beta(x_t - x_{t-1})$). This is equivalent to the restriction that c generates all of the lags and leads, and we have seen that it is supported by the data. I also considered imposing all of the restrictions with the single exception of allowing separate coefficients for entering and leaving

Table 5.2
Unrestricted least squares regressions.
(a)

Dependent variable	Coefficients (and standard errors) of:								
	$U1$	$U2$	$U3$	$SMSA1$	$SMSA2$	$SMSA3$	$RNS1$	$RNS2$	$RNS3$
$LW1$	0.171 (0.025)	0.042 (0.026)	-0.009 (0.025)	0.135 (0.028)	-0.001 (0.055)	0.032 (0.054)	-0.016 (0.081)	-0.020 (0.081)	-0.108 (0.070)
$LW2$	0.048 (0.023)	0.150 (0.028)	0.010 (0.026)	0.086 (0.027)	0.053 (0.065)	0.020 (0.061)	0.065 (0.099)	-0.039 (0.109)	-0.155 (0.092)
$LW3$	0.046 (0.023)	0.041 (0.030)	0.132 (0.030)	0.083 (0.031)	0.003 (0.058)	0.088 (0.056)	0.074 (0.079)	0.056 (0.093)	-0.232 (0.078)

Notes:

All regressions include $(1, S, EXP69, EXP69^2, RACE)$. The standard errors are calculated using $\hat{\Omega}$ in (4.2).

(b)

Dependent variable	Coefficients (and standard errors) of:						
	$U1$	$U2$	$U3$	$U1U2$	$U1U3$	$U2U3$	$U1U2U3$
$LW1$	0.127 (0.044)	-0.047 (0.042)	-0.072 (0.041)	0.128 (0.072)	0.092 (0.075)	0.156 (0.070)	-0.182 (0.104)
$LW2$	-0.019 (0.040)	0.014 (0.045)	-0.085 (0.040)	0.181 (0.074)	0.118 (0.092)	0.227 (0.066)	-0.229 (0.116)
$LW3$	-0.050 (0.037)	-0.072 (0.053)	-0.022 (0.052)	0.110 (0.079)	0.264 (0.081)	0.246 (0.079)	-0.256 (0.113)

Notes:

All regressions include $(SMSA1, SMSA2, SMSA3, RNS1, RNS2, RNS3, 1, S, EXP69, EXP69^2, RACE)$. The standard errors are calculated using $\hat{\Omega}$ in (4.2).

union coverage in the wage change equations. The estimates (standard errors) are 0.097 (0.019) and -0.119 (0.022). The standard error on the sum of the coefficients is 0.024, so again there is no evidence against the simple model with $E(y_i|x_1, x_2, x_3, c) = \beta x_i + c$.⁴³

Table 5.3(a) exhibits the estimates that result from imposing the restrictions using the optimal minimum distance estimator.⁴⁴ We also give the conventional generalized least squares estimates. They are minimum distance estimates in which the weighting matrix (A_N) is the inverse of

$$\hat{\Omega}_s = \frac{1}{N} \sum_{i=1}^N (y_i - \hat{\Pi}x_i)(y_i - \hat{\Pi}x_i)' \otimes \left(\frac{1}{N} \sum_{i=1}^N x_i x_i' \right)^{-1}. \quad (5.1)$$

We give the conventional standard errors based on $(F'\hat{\Omega}_s^{-1}F)^{-1}$ and the standard errors calculated according to Proposition 2, which do not require an assumption of homoskedastic linear regression. These standard errors are larger than the conventional ones, by about 30%. The estimated gain in efficiency from using the appropriate metric is not very large; the standard errors calculated according to Proposition 2 are about 10% larger when we use conventional GLS instead of the optimum minimum distance estimator.

Table 5.3(a) also presents the estimated λ 's. Consider, for example, an individual who was covered by collective bargaining in 1969. The linear predictor of c increases by 0.089 if he is also covered in 1970, and it increases by an additional 0.036 if he is covered in all three years. The predicted c for someone who is always covered is higher by 0.102 than for someone who is never covered.

Table 5.3(b) presents estimates under the constraint that $\lambda = 0$. The increment in the distance statistic is $89.08 - 19.36 = 69.72$, which is a surprisingly large value to come from a $\chi^2(13)$ distribution. If we constrain only the union λ 's to be zero, then the increment is $57.06 - 19.36 = 37.7$, which is surprisingly large coming from a $\chi^2(7)$ distribution. So there is strong evidence for heterogeneity bias.

The union coefficient declines from 0.157 to 0.107 when we relax the $\lambda = 0$ restriction. The least squares estimates for the separate cross sections, with no

⁴³Using May–May CPS matches for 1977–1978, Mellow (1981) reports coefficients (standard errors) of 0.087 (0.018) and -0.069 (0.020) for entering and leaving union membership in a wage change regression. The sample consists of 6602 males employed as nonagricultural wage and salary workers in both years. He also reports results for 2177 males and females whose age was ≤ 25 . Here the coefficients on entering and leaving union membership are quite different: 0.198 (0.031) and -0.035 (0.041); it would be useful to reconcile these numbers with our results for young men. Also see Stafford and Duncan (1980).

⁴⁴We did not find much evidence for nonstationarity in the slope coefficients. If we allow the union β to vary over the three years, we get 0.105, 0.103, 0.114. The distance statistic declines to 18.51, giving $19.36 - 18.51 = 0.85$; this is not a surprising value from a $\chi^2(2)$ distribution. If we also free up β for SMSA and RNS, then the decline in the distance statistic is $18.51 - 13.44 = 5.07$, which is not a surprising value from a $\chi^2(4)$ distribution.

Table 5.3
Restricted estimates.
(a)

	Coefficients (and standard errors) of:						
	<i>U</i>	<i>SMSA</i>	<i>RNS</i>				
$\hat{\beta}$:	0.107 (0.016)	0.056 (0.020)	-0.082 (0.045)				
$\hat{\beta}_{GLS}$:	0.121 (0.013) (0.018)	0.050 (0.017) (0.021)	-0.085 (0.040) (0.052)				
	<i>U1</i>	<i>U2</i>	<i>U3</i>	<i>U1U2</i>	<i>U1U3</i>	<i>U2U3</i>	<i>U1U2U3</i>
$\hat{\lambda}$:	-0.023 (0.030)	-0.067 (0.040)	-0.082 (0.037)	0.156 (0.057)	0.152 (0.062)	0.195 (0.059)	-0.229 (0.085)
	<i>SMSA1</i>	<i>SMSA2</i>		<i>SMSA3</i>	<i>RNS1</i>	<i>RNS2</i>	<i>RNS3</i>
	0.086 (0.025)	-0.008 (0.046)		0.032 (0.046)	0.100 (0.072)	-0.021 (0.077)	-0.128 (0.068)

$\chi^2(23) = 19.36$

(b) Restrict $\lambda = 0$

	Coefficients (and standard errors) of:		
	<i>U</i>	<i>SMSA</i>	<i>RNS</i>
$\hat{\beta}$:	0.157 (0.012)	0.120 (0.013)	-0.150 (0.016)

$\chi^2(36) = 89.08$

Notes:

$E^*(y|x) = \Pi x = \Pi_1 x_1 + \Pi_2 x_2$; $x'_1 = (U1, U2, U3, U1U2, U1U3, U2U3, U1U2U3, SMSA1, SMSA2, SMSA3, RNS1, RNS2, RNS3)$; $x'_2 = (1, S, EXP69, EXP69^2, RACE)$. $\Pi_1 = (\beta_u I_3, 0, \beta_{SMSA} I_3, \beta_{RNS} I_3) + t\lambda'$; Π_2 is unrestricted. The restrictions are expressed as $\pi = F\delta$, where δ is unrestricted. $\hat{\beta}$ and $\hat{\lambda}$ are minimum distance estimates with $A_N^{-1} = \hat{\Omega}$ in (4.2); $\hat{\beta}_{GLS}$ and $\hat{\lambda}_{GLS}$ are minimum distance estimates with $A_N^{-1} = \hat{\Omega}_s$ in (5.1) ($\hat{\lambda}_{GLS}$ is not shown in the table). The first standard error for $\hat{\beta}_{GLS}$ is the conventional one based on $(F'\hat{\Omega}_s^{-1}F)^{-1}$; the second standard error for $\hat{\beta}_{GLS}$ is based on $(F'\hat{\Omega}_s^{-1}F)^{-1}F'\hat{\Omega}_s^{-1}\hat{\Omega}\hat{\Omega}_s^{-1}F(F'\hat{\Omega}_s^{-1}F)^{-1}$ (Proposition 2). The χ^2 statistics are computed from $N[\hat{\pi} - F\hat{\delta}]'\hat{\Omega}^{-1}[\hat{\pi} - F\hat{\delta}]$ (Proposition 3).

leads or lags, give union coefficients of 0.195, 0.189, and 0.191 in 1969, 1970, and 1971.⁴⁵ So the decline in the union coefficient, when we allow for heterogeneity bias, is 32% or 44% depending on which biased estimate (0.16 or 0.19) one uses. The *SMSA* and region coefficients also decline in absolute value. The least squares estimates for the separate cross sections give an average *SMSA* coefficient of 0.147 and an average region coefficient of -0.131 . So the decline in the *SMSA* coefficient is either 53% or 62%, and the decline in absolute value of the region coefficient is either 45% or 37%.

5.2. *Nonlinear models: Labor force participation*

We shall illustrate some of the results in Section 3. The sample consists of 924 married women in the Michigan Panel Study of Income Dynamics. The sample selection criteria and the means and standard deviations of the variables are in Table 5.4. Participation status is measured by the question: "Did _____ do any work for money last year?" We shall model participation in 1968, 1970, 1972, and 1974.

In terms of the model described in Section 3.1, the wage predictors are schooling, experience, and experience squared, where experience is measured as age minus schooling minus six; the tastes for nonmarket time are predicted by these variables and by children. The specification for children is a conventional one that uses the number of children of age less than six (*YK*) and the total number of children in the family unit (*K*).⁴⁶ Variables that affect only the lifetime budget constraint in this certainty model are captured by *c*. In particular, nonlabor income and the husband's wage are assumed to affect the wife's participation only through the lifetime budget constraint. The individual effect (*c*) will also capture unobserved permanent components in wages or in tastes for nonmarket time.

Table 5.5 presents maximum likelihood (ML) estimates of cross-section probit specifications for each of the four years. Table 5.6 presents unrestricted ML estimates for all lags and leads in *YK* and *K*. If the residuals (u_{it}) in the latent

⁴⁵ Using the NLS Young Men in 1969 ($N = 1362$), Griliches (1976) reports a union membership coefficient of 0.203. Using the NLS Young Men in a pooled regression for 1966–1971 and 1973 ($N = 470$), Brown (1980) reports a coefficient of 0.130 on a variable measuring the probability of union coverage. (The union coverage question was asked only in 1969, 1970, and 1971; so this variable is imputed for the other four years.) The coefficient declines to .081 when individual intercepts are included in the regression. His regressions also include a large number of occupation and industry specific job characteristics.

⁴⁶ Some of the work on participation and fertility is in Mincer (1963), Willis (1973), Gronau (1973, 1976, 1977), Hall (1973), Ben-Porath (1973), Becker and Lewis (1973), Mincer and Polachek (1974), Heckman (1974, 1980), Heckman and Willis (1977), Cain and Dooley (1976), Schultz (1980), Hanoch (1980), and Rosenzweig and Wolpin (1980).

variable model (3.1) have constant variance, then $\alpha_1 = \dots = \alpha_4$ in (3.2), and the submatrices of Π corresponding to YK and K should have the form $\beta I + I\lambda$. There may be some indication of this pattern in Table 5.6, but it is much weaker than in the wage regressions in Table 5.2.

We allow for unequal variances and provide formal tests by using the minimum distance estimator developed in Section 4.5. In Table 5.7(a) we impose the restrictions that

$$\Pi = \text{diag}\{\alpha_1, \dots, \alpha_4\} [\beta_{YK} I_4 + I\lambda_{YK}, \beta_K I_4 + I\lambda_K].$$

The minimum distance statistic is 53.8, which is a very surprising value coming from a $\chi^2(19)$ distribution. So the latent variable c does not appear to provide an

Table 5.4
Characteristics of Michigan Panel Study of Income
Dynamics married women: Means and standard deviations, $N = 924$.

Variable	Mean	Standard deviation
<i>LFP1</i>	0.499	—
<i>LFP2</i>	0.530	—
<i>LFP3</i>	0.529	—
<i>LFP4</i>	0.566	—
<i>YK1</i>	0.969	1.200
<i>YK2</i>	0.764	1.069
<i>YK3</i>	0.551	0.895
<i>YK4</i>	0.363	0.685
<i>K1</i>	2.38	1.69
<i>K2</i>	2.30	1.64
<i>K3</i>	2.11	1.61
<i>K4</i>	1.84	1.52
<i>S</i>	12.1	2.1
<i>EXP68</i>	17.2	8.5
<i>EXP68</i> ²	368.	301.

Notes:

LFP1, ..., *LFP4*—1 if answered "yes" to "Did _____ work for money last year?", 0 otherwise, referring to 1968, 1970, 1972, 1974; *YK1*, ..., *YK4*—number of children of age less than six in 1968, 1970, 1972, 1974; *K1*, ..., *K4*—number of children of age less than eighteen living in the family unit in 1968, 1970, 1972, 1974; *S*—years of schooling completed; *EXP68*—(age in 1968—*S*—6). The sample selection criteria required that the women be married to the same spouse from 1968 to 1976; not part of the low income subsample; between 20 and 50 years old in 1968; white; out of school from 1968 to 1976; not disabled. We required complete data on the variables in the table, and that there be no inconsistency between reported earnings and the answer to the participation question.

Table 5.5
ML probit cross-section estimates.

Dependent variable	Coefficients (and standard errors) of:							
	<i>YK1</i>	<i>YK2</i>	<i>YK3</i>	<i>YK4</i>	<i>K1</i>	<i>K2</i>	<i>K3</i>	<i>K4</i>
<i>LFP1</i>	-0.246 (0.046)	—	—	—	-0.063 (0.031)	—	—	—
<i>LFP2</i>	—	-0.293 (0.055)	—	—	—	-0.075 (0.031)	—	—
<i>LFP3</i>	—	—	-0.342 (0.067)	—	—	—	-0.077 (0.032)	—
<i>LFP4</i>	—	—	—	-0.366 (0.081)	—	—	—	-0.069 (0.034)

Notes:

Separate ML estimates each year. All specifications include (1, *S*, EXP68, EXP68²).

adequate interpretation of the unrestricted leads and lags. It may be that the distributed lag relationship between current participation and previous births is more general than the one implied by summing over the previous six years (*YK*) and over the previous eighteen years (*K*). It may be fruitful to explore this in more detail in future work. Perhaps strict exogeneity conditional on *c* will hold when we use a more general specification for lagged births. But we must keep in mind that this question is intrinsically tied to the functional form restrictions—we saw in Section 3.3 that there always exist specifications in which y_t is independent of x_1, \dots, x_T conditional on *c*.

Table 5.6
Unrestricted ML probit estimates.

Dependent variable	Coefficients (and standard errors) of:							
	<i>YK1</i>	<i>YK2</i>	<i>YK3</i>	<i>YK4</i>	<i>K1</i>	<i>K2</i>	<i>K3</i>	<i>K4</i>
<i>LFP1</i>	-0.205 (0.081)	-0.017 (0.119)	-0.160 (0.141)	0.420 (0.144)	0.176 (0.076)	-0.142 (0.100)	-0.196 (0.110)	0.063 (0.090)
<i>LFP2</i>	-0.047 (0.079)	-0.238 (0.117)	-0.047 (0.140)	0.093 (0.142)	0.320 (0.077)	-0.278 (0.102)	-0.250 (0.110)	0.177 (0.090)
<i>LFP3</i>	-0.254 (0.080)	0.214 (0.116)	-0.190 (0.139)	-0.209 (0.141)	0.204 (0.077)	-0.210 (0.102)	-0.045 (0.112)	0.030 (0.090)
<i>LFP4</i>	-0.195 (0.079)	0.252 (0.118)	-0.211 (0.139)	-0.282 (0.138)	0.020 (0.075)	0.083 (0.100)	-0.181 (0.110)	0.058 (0.090)

Notes:

Separate ML estimates each year. All specifications include (1, *S*, EXP68, EXP68²).

Table 5.7
Restricted estimates.
(a)

Coefficients (and standard errors) of:								
	<i>YK</i>	<i>K</i>						
$\alpha_4 \hat{\beta}$	-0.121 (0.046)	-0.058 (0.029)						
	<i>YK1</i>	<i>YK2</i>	<i>YK3</i>	<i>YK4</i>	<i>K1</i>	<i>K2</i>	<i>K3</i>	<i>K4</i>
$\alpha_4 \hat{\lambda}$	-0.042 (0.041)	0.038 (0.060)	-0.050 (0.070)	0.087 (0.077)	0.194 (0.056)	-0.118 (0.062)	-0.146 (0.073)	0.090 (0.056)
	$\hat{\alpha}_1$	$\hat{\alpha}_2$	$\hat{\alpha}_3$	$\hat{\alpha}_4$				
$\hat{\alpha}$	1.585 (0.392)	1.758 (0.375)	1.279 (0.231)	1.0 (-)				

$\chi^2(19) = 53.8$

(b) Restrict $\lambda = 0$

	<i>YK</i>	<i>K</i>						
$\alpha_4 \hat{\beta}$	-0.273 (0.065)	-0.073 (0.023)						
	$\hat{\alpha}_1$	$\hat{\alpha}_2$	$\hat{\alpha}_3$	$\hat{\alpha}_4$				
$\hat{\alpha}$	0.821 (0.198)	0.930 (0.205)	0.920 (0.191)	1.0 (-)				

$\chi^2(27) = 78.4$

(c) Restrict $\alpha_t = 1$ ($t = 1, \dots, 4$)

Coefficients (and standard errors) of:								
	<i>YK</i>	<i>K</i>						
$\hat{\beta}$	-0.193 (0.043)	-0.070 (0.031)						
	<i>YK1</i>	<i>YK2</i>	<i>YK3</i>	<i>YK4</i>	<i>K1</i>	<i>K2</i>	<i>K3</i>	<i>K4</i>
$\hat{\lambda}$	-0.077 (0.062)	0.082 (0.082)	-0.098 (0.102)	0.102 (0.110)	0.203 (0.063)	-0.108 (0.083)	-0.157 (0.098)	0.072 (0.081)

$\chi^2(22) = 61.6$

(d) Restrict $\alpha_t = 1$; β_t unrestricted ($t = 1, \dots, 4$)

Coefficients (and standard errors) of:								
	<i>YK1</i>	<i>YK2</i>	<i>YK3</i>	<i>YK4</i>	<i>K1</i>	<i>K2</i>	<i>K3</i>	<i>K4</i>
$\hat{\beta}$	-0.107 (0.054)	-0.216 (0.059)	-0.198 (0.067)	-0.277 (0.086)	-0.107 (0.040)	-0.047 (0.035)	-0.046 (0.039)	-0.017 (0.043)
	<i>YK1</i>	<i>YK2</i>	<i>YK3</i>	<i>YK4</i>	<i>K1</i>	<i>K2</i>	<i>K3</i>	<i>K4</i>
$\hat{\lambda}$	-0.111 (0.063)	0.085 (0.083)	-0.102 (0.102)	0.126 (0.111)	0.213 (0.064)	-0.113 (0.083)	-0.155 (0.099)	0.052 (0.082)

$\chi^2(16) = 52.7$

Notes:

$\Pi_1 = \text{diag}\{\alpha_1, \dots, \alpha_4\}[\beta_{YK}I_4 + I\lambda_{YK}, \beta_K I_4 + I\lambda_K]$; Π_2 is unrestricted. In Table 5.7(d) $\beta_{YK}I_4$ and $\beta_K I_4$ are replaced by diagonal matrices with no restrictions on the diagonal elements. All restrictions are imposed by applying the minimum distance procedure to the unrestricted estimates of Π_1 in Table 5.6. The asymptotic covariance matrix of $\hat{\Pi}_1$ is obtained as in Section 4.5. α_4 is normalized to equal one.

If we do impose the restrictions in Table 5.7(a), then there is strong evidence that $\lambda \neq 0$. Constraining $\lambda = 0$ in Table 5.7(b) gives an increase in the distance statistic of $78.4 - 53.8 = 24.6$, which is surprisingly large to come from a $\chi^2(8)$ distribution.

In Table 5.7(c) we constrain all of the residual variances to be equal ($\alpha_i = 1$). An alternative interpretation of the time varying coefficients is provided in Table 5.7(d), where β_{YK} and β_K vary freely over time and $\alpha_i = 1$. In principle, we could also allow the α_i to vary freely, since they can be identified from changes over time in the coefficients of c . In fact that model gives very imprecise results and it is difficult to ensure numerical accuracy.

We shall interpret the coefficients on YK and K by following the procedure in (3.4). Table 5.8 presents estimates of the expected change in the participation probability when we assign an additional young child to a randomly chosen family, so that YK and K increase by one. We compute this measure for the models in Tables 5.7(a), 5.7(c) and 5.7(d). The average change in the participation probability is -0.096 . We can get an indication of omitted variable bias by comparing these estimates with the ones based on Table 5.7(b), where λ is constrained to be zero. Now the average change in the participation probability is -0.122 , so that the decline in absolute value when we control for c is 21%. An alternative comparison can be based on the cross-section estimates, with no leads or lags, in Table 5.5. Now the average change in the participation probability is -0.144 , giving an omitted variable bias of 33%.

Next we shall consider estimates from the logit framework of Section 3.2. Table 5.9 presents (standard) maximum likelihood estimates of cross-section logit specifications for each of the four years. We can use the cross-section probit results in Table 5.5 to construct estimates of the expected change in the log odds of participation when we add a young child to a randomly chosen family. Doing this in each of the four years gives -0.502 , -0.598 , -0.683 , and -0.703 . With the logit estimates, we simply add together the coefficients on YK and K in Table 5.9; this gives -0.507 , -0.612 , -0.691 , and -0.729 . The average over the four years is -0.621 for probit and -0.635 for logit. So at this point there is little difference between the two functional forms.

Now allow for the latent variable (c). Table 5.10 presents the conditional maximum likelihood estimates for the fixed effects logit model. The striking result here is that, unlike the probit case, allowing for c leads to an *increase* in the absolute value of the children coefficients. If we constrain β_{YK} and β_K to be constant over time (Table 5.10(a)), the estimated change in the log odds of participation when we add an additional young child is -0.909 . If we allow β_{YK} and β_K to vary freely over time (Table 5.10(b)), the average of the estimated changes is -0.879 . So the absolute value of the estimates increases by about 40% when we control for c using the logit framework. The estimation method is having a first order effect on the results.

It is commonly found that probit and logit specifications, when properly interpreted, give very similar results; our cross-section estimates are an example of this. But our attempt to incorporate latent variables has turned up marked differences between the probit and logit specifications. There are a number of possible explanations for this. The probit specification restricts c to have a normal distribution conditional on x with a linear regression function and constant variance. The conditional likelihood approach in the logit model does not impose this possibly false restriction. On the other hand, the probit model has a more general specification for the residual covariance matrix.

Table 5.3
Estimated effects of an additional young child.

8.7(a)	Unrestricted λ and α_t ($t = 1, \dots, 4$)			
	$E(P_{1t} - P_{0t})$			
	1968	1970	1972	1974
	-0.105	-0.116	-0.087	-0.069
8.7(b)	Restrict $\lambda = 0$			
	-0.108	-0.123	-0.122	-0.134
8.7(c)	Restrict $\alpha_t = 1$ ($t = 1, \dots, 4$)			
	-0.098	-0.099	-0.099	-0.101
8.7(d)	Restrict $\alpha_t = 1$; β_t unrestricted ($t = 1, \dots, 4$)			
	-0.081	-0.099	-0.092	-0.112
8.5	Cross-section estimates			
	-0.116	-0.139	-0.157	-0.166

Notes:

$$\hat{P}_{0it} = F[\hat{\alpha}_t(\hat{\beta}'_t x_{1it} + \lambda' x_{1t}) + \hat{\pi}'_{2t} x_{2it}],$$

$$\hat{P}_{1it} = F[\hat{\alpha}_t(\hat{\beta}'_t x_{1it}^* + \lambda' x_{1t}) + \hat{\pi}'_{2t} x_{2it}],$$

where $F(\cdot)$ is the standard normal distribution function; $x'_{1it} = (YKt, Kt)_i$;

$$x_{1it}^* = (YKt + 1, Kt + 1)_i \quad (t = 1, \dots, 4);$$

$$x'_{1t} = (YK1, YK2, YK3, YK4, K1, K2, K3, K4)_i;$$

$x'_{2t} = (1, S, \text{EXP68}, \text{EXP68}^2)_i$. The estimate of $E(P_{1t} - P_{0t})$ is:

$$\frac{1}{N} \sum_{i=1}^N (\hat{P}_{1it} - \hat{P}_{0it}).$$

The estimates of $\alpha_t, \beta_t, \lambda, \pi_{2t}$ used in Tables 8.7(a), ..., 8.7(d) are based on the specifications in Tables 5.7(a)–5.7(d); the estimates in Table 8.5 are based on the specification in Table 5.5.

Table 5.9
ML logit cross-section estimates.

Dependent variable	Coefficients (and standard errors) of:							
	<i>YK1</i>	<i>YK2</i>	<i>YK3</i>	<i>YK4</i>	<i>K1</i>	<i>K2</i>	<i>K3</i>	<i>K4</i>
<i>LFP1</i>	-0.404 (0.077)	—	—	—	-0.103 (0.051)	—	—	—
<i>LFP2</i>	—	-0.494 (0.095)	—	—	—	-0.118 (0.035)	—	—
<i>LFP3</i>	—	—	-0.568 (0.114)	—	—	—	-0.123 (0.051)	—
<i>LFP4</i>	—	—	—	-0.617 (0.138)	—	—	—	-0.112 (0.055)

Notes:

Separate ML estimates each year. All specifications include (1, *S*, EXP68, EXP68²).

We have seen that the restrictions on the probit Π matrix, which underlie our estimate of β , appear to be false. An analogous test in the logit framework is based on (3.10). We use conditional ML to estimate a model that includes $YK_s \cdot D_t$, $K_s \cdot D_t$ ($s = 1, \dots, 4$; $t = 2, 3, 4$), where D_t is a dummy variable that is one in period t and zero otherwise. It is not restrictive to exclude $YK_s \cdot D_1$ and $K_s \cdot D_1$, since they can be absorbed in c . We include also D_t , $S \cdot D_t$, EXP68 $\cdot D_t$, and EXP68² $\cdot D_t$ ($t = 2, 3, 4$). Then comparing the maximized conditional likelihoods

Table 5.10
Conditional ML estimates of the fixed effects logit model.

(a)		
Coefficients (and standard errors) of:		
	<i>YK</i>	<i>K</i>
$\hat{\beta}$	-0.573 (0.115)	-0.336 (0.120)

(b) $\hat{\beta}_t$ unrestricted ($t = 1, \dots, 4$)								
Coefficients (and standard errors) of:								
	<i>YK1</i>	<i>YK2</i>	<i>YK3</i>	<i>YK4</i>	<i>K1</i>	<i>K2</i>	<i>K3</i>	<i>K4</i>
$\hat{\beta}$	-0.336 (0.144)	-0.679 (0.172)	-0.780 (0.205)	-0.967 (0.242)	-0.315 (0.135)	-0.178 (0.145)	-0.141 (0.155)	-0.120 (0.165)

Notes:

A conditional likelihood ratio test of $\beta_1 = \dots = \beta_4$ gives $\chi^2(6) = 8.7$. The specifications in Tables 10(a) and 10(b) include dummy variables for 1970, 1972, 1974 (D_t , $t = 2, 3, 4$) and the interactions $S \cdot D_t$, EXP68 $\cdot D_t$, EXP68² $\cdot D_t$ ($t = 2, 3, 4$). (Due to the presence of the fixed effect c_t , it is not restrictive to exclude D_1 , S , EXP68, EXP68², $S \cdot D_1$, EXP68 $\cdot D_1$, EXP68² $\cdot D_1$.)

for this specification and the specification in Table 5.10(b) gives a conditional likelihood ratio statistic of 53.9, which is a very surprising value to come from a $\chi^2(16)$ distribution. So the restrictions underlying our logit estimates of β also appear to be false. It may be that the false restrictions simply imply different biases in the probit and logit specifications.

6. Conclusion

Our discussion has focused on models that are static conditional on a latent variable. The panel aspect of the data has primarily been used to control for the latent variable. Much work needs to be done on models that incorporate uncertainty and interesting dynamics. Exploiting the martingale implications of time-additive utility seems fruitful here, as in Hall (1978) and Hansen and Singleton (1982). There is, however, a potentially important distinction between time averages and cross-section averages. A time average of forecast errors over T periods should converge to zero as $T \rightarrow \infty$. But an average of forecast errors across N individuals surely need not converge to zero as $N \rightarrow \infty$; there may be common components in those errors, due to economy-wide innovations. The same point applies when we consider covariances of forecast errors with variables that are in the agents' information sets. If those conditioning variables are discrete, we can think of averaging over subsets of the forecast errors; as $T \rightarrow \infty$, these averages should converge to zero, but not necessarily as $N \rightarrow \infty$.

As for controlling for latent variables, I think that future work will have to address the lack of identification that we have uncovered. It is not restrictive to assert that (y_1, \dots, y_T) and (x_1, \dots, x_T) are independent conditional on some latent variable c .

Appendix

Let $\mathbf{r}'_i = (\mathbf{x}'_i, \mathbf{y}'_i)$, $i = 1, \dots, N$, where $\mathbf{x}'_i = (x_{i1}, \dots, x_{iK})$ and $\mathbf{y}'_i = (y_{i1}, \dots, y_{iM})$. Write the m th structural equation as:

$$y_{im} = \delta'_m \mathbf{z}_{im} + v_{im} \quad (m = 1, \dots, M),$$

where the components of \mathbf{z}_{im} are the variables in \mathbf{y}_i and \mathbf{x}_i that appear in the m th equation with unknown coefficients. Let \mathbf{S}_{zx} be the following block-diagonal matrix:

$$\mathbf{S}_{zx} = \text{diag} \left\{ \frac{1}{N} \sum_{i=1}^N \mathbf{z}_{i1} \mathbf{x}'_i, \dots, \frac{1}{N} \sum_{i=1}^N \mathbf{z}_{iM} \mathbf{x}'_i \right\},$$

and

$$s_{xy} = \frac{1}{N} \sum_{i=1}^N y_i \otimes x_i.$$

Let $v_i^0 = (v_{i1}^0, \dots, v_{iM}^0)$, where $v_{im}^0 = y_{im} - \delta_m^0 z_{im}$ and δ_m^0 is the true value of δ_m ; let $\Phi_{zx} = E(S_{zx})$. Let $\delta' = (\delta'_1, \dots, \delta'_M)$ be $s \times 1$ and set

$$\hat{\delta} = (S_{zx} D^{-1} S'_{zx})^{-1} (S_{zx} D^{-1} s_{xy}).$$

Proposition 6

Assume that (1) r_i is i.i.d. according to some distribution with finite fourth moments; (2) $E[x_1(y_{1m} - \delta_m^0 z_{1m})] = 0$ ($m=1, \dots, M$); (3) $\text{rank}(\Phi_{zx}) = s$; (4) $D \xrightarrow{\text{a.s.}} \Psi$ as $N \rightarrow \infty$, where Ψ is a positive-definite matrix. Then $\sqrt{N}(\hat{\delta} - \delta^0) \xrightarrow{D} N(0, \Lambda)$, where

$$\Lambda = (\Phi_{zx} \Psi^{-1} \Phi'_{zx})^{-1} \Phi_{zx} \Psi^{-1} [E(v_1^0 v_1^0 \otimes x_1 x_1')] \Psi^{-1} \Phi'_{zx} (\Phi_{zx} \Psi^{-1} \Phi'_{zx})^{-1}.$$

Proof

$$\sqrt{N}(\hat{\delta} - \delta^0) = (S_{zx} D^{-1} S'_{zx})^{-1} S_{zx} D^{-1} \sum_{i=1}^N (v_i^0 \otimes x_i) / \sqrt{N}.$$

By the strong law of large numbers, $S_{zx} \xrightarrow{\text{a.s.}} \Phi_{zx}$; $\Phi_{zx} \Psi^{-1} \Phi'_{zx}$ is an $s \times s$ positive-definite matrix since $\text{rank}(\Phi_{zx}) = s$. So we obtain the same limiting distribution by considering

$$(\Phi_{zx} \Psi^{-1} \Phi'_{zx})^{-1} \Phi_{zx} \Psi^{-1} \sum_{i=1}^N (v_i^0 \otimes x_i) / \sqrt{N}.$$

Note that $v_i^0 \otimes x_i$ is i.i.d. with $E(v_1^0 \otimes x_1) = 0$, $V(v_1^0 \otimes x_1) = E(v_1^0 v_1^0 \otimes x_1 x_1')$. Then applying the central limit theorem gives $\sqrt{N}(\hat{\delta} - \delta^0) \xrightarrow{D} N(0, \Lambda)$. Q.E.D.

This result includes as special cases a number of the commonly used estimators. If $z_{im} = x_i$ ($m=1, \dots, M$) and $D = I$, then $\hat{\delta}$ is the least squares estimator and Λ reduces to the formula for Ω given in (4.1). If $\Psi = E(v_1^0 v_1^0) \otimes E(x_1 x_1')$, then Λ is the asymptotic covariance matrix for the three-stage least squares estimator. If

$\Psi = E(v_1^0 v_1^0 \otimes x_1 x_1')$, then Λ is the asymptotic covariance matrix for the generalized three-stage least squares estimator (4.3).

Consider applying the generalized three-stage least squares estimator to the first J equations ($J < M$). If $E(z_{1j} x_1')$ is nonsingular for $j = J+1, \dots, M$, then this estimator for $(\delta'_1, \dots, \delta'_J)$ has the same asymptotic covariance matrix as the estimator obtained by applying the generalized three-stage least squares estimator to the full set of M equations. This follows from examining the partitioned inverse of (4.3).

References

- Amemiya, T. (1971) "The Estimation of Variances in a Variance-Components Model", *International Economic Review*, 12, 1–13.
- Andersen, E. B. (1970) "Asymptotic Properties of Conditional Maximum Likelihood Estimators", *Journal of the Royal Statistical Society, Series B*, 32, 283–301.
- Andersen, E. B. (1971) "Asymptotic Properties of Conditional Likelihood Ratio Tests", *Journal of the American Statistical Association*, 66, 630–633.
- Andersen, E. B. (1972) "The Numerical Solution of a Set of Conditional Estimation Equations", *Journal of the Royal Statistical Society, Series B*, 34, 42–54.
- Andersen, E. B. (1973) *Conditional Inference and Models for Measuring*. Copenhagen: Mentalhygiejnisk Forlag.
- Anderson, T. W. (1969) "Statistical Inference for Covariance Matrices with Linear Structure", in P. R. Krishnaiah, ed., *Proceedings of the Second International Symposium on Multivariate Analysis*, Academic Press, New York.
- Anderson, T. W. (1970) "Estimation of Covariance Matrices Which are Linear Combinations or Whose Inverses are Linear Combinations of Given Matrices", *Essays in Probability and Statistics*, University of North Carolina Press, Chapel Hill.
- Anderson, T. W. and C. Hsiao (1981) "Estimation of Dynamic Models with Error Components", *Journal of the American Statistical Association*, 76, 598–606.
- Anderson, T. W. and C. Hsiao (1982) "Formulation and Estimation of Dynamic Models Using Panel Data", *Journal of Econometrics*, forthcoming.
- Avery, R. B., L. P. Hansen, and V. J. Hotz (1981) "Multiperiod Probit Models and Orthogonality Condition Estimation", Carnegie-Mellon University, Graduate School of Industrial Administration Working Paper No. 62-80-81.
- Balestra, P. and M. Nerlove (1966) "Pooling Cross Section and Time Series Data in the Estimation of a Dynamic Model: The Demand for Natural Gas", *Econometrica*, 34, 585–612.
- Barndorff-Nielsen, O. (1978) *Information and Exponential Families in Statistical Theory*. New York: Wiley.
- Basmann, R. L. (1965) "On the Application of the Identifiability Test Statistic and Its Exact Finite Sample Distribution Function in Predictive Testing of Explanatory Economic Models", unpublished manuscript.
- Becker, G. S. and H. G. Lewis (1973) "On the Interaction Between the Quantity and Quality of Children", *Journal of Political Economy*, 81, S279–S288.
- Ben-Porath, Y. (1973) "Economic Analysis of Fertility in Israel: Point and Counter-Point", *Journal of Political Economy*, 81, S202–S233.
- Billingsley, P. (1979) *Probability and Measure*, Wiley, New York.
- Bishop, Y. M. M., S. E. Fienberg and P. W. Holland (1975) *Discrete Multivariate Analysis: Theory and Practice*, Cambridge, MA: MIT Press.
- Blinder, A. S. and Y. Weiss (1976) "Human Capital and Labor Supply: A Synthesis", *Journal of Political Economy*, 84, 449–472.

- Brown, C. (1980) "Equalizing Differences in the Labor Market", *Quarterly Journal of Economics*, 94, 113–134.
- Cain, G. G. and M. D. Dooley (1976) "Estimation of a Model of Labor Supply, Fertility, and Wages of Married Women", *Journal of Political Economy*, 84, S179–S199.
- Chamberlain, G. (1977) "An Instrumental Variable Interpretation of Identification in Variance-Components and MIMIC Models", in P. Taubman (ed.), *Kinometrics: The Determinants of Socioeconomic Success Within and Between Families*, North-Holland Publishing Company, Amsterdam.
- Chamberlain, G. (1978) "Omitted Variable Bias in Panel Data: Estimating the Returns to Schooling", *Annales de l'INSEE*, 30/31, 49–82.
- Chamberlain, G. (1978a) "On the Use of Panel Data", unpublished manuscript.
- Chamberlain, G. (1979) "Heterogeneity, Omitted Variable Bias, and Duration Dependence", Harvard Institute for Economic Research Discussion Paper No. 691.
- Chamberlain, G. (1980) "Analysis of Covariance with Qualitative Data", *Review of Economic Studies*, 47, 225–238.
- Chamberlain, G. (1980a) "Studies of Teaching and Learning in Economics: Discussion", *American Economic Review*, Papers and Proceedings, 69, 47–49.
- Chamberlain, G. (1982) "The General Equivalence of Granger and Sims Causality", *Econometrica*, 50, 569–581.
- Chamberlain, G. (1982a) "Multivariate Regression Models for Panel Data", *Journal of Econometrics*, 18, 5–46.
- Chiang, C. L. (1956) "On Regular Best Asymptotically Normal Estimates", *Annals of Mathematical Statistics*, 27, 336–351.
- Corcoran, M. and M. S. Hill (1980) "Persistence in Unemployment Among Adult Men", in G. J. Duncan and J. N. Morgan (eds.), *Five Thousand American Families—Patterns of Economic Progress*, Institute for Social Research, University of Michigan, Ann Arbor.
- Cox, D. R. (1970) *Analysis of Binary Data*, London, Methuen.
- Cramér, H. (1946) *Mathematical Methods of Statistics*, Princeton University Press, Princeton.
- deFinetti, B. (1975) *Theory of Probability*, Vol. 2, Wiley, New York.
- Diaconis, P. and D. Freedman (1980) "deFinetti's Theorem for Markov Chains", *The Annals of Probability*, 8, 115–130.
- Dynkin, E. B. and A. A. Yushkevich (1979) *Controlled Markov Processes*, Springer-Verlag, New York.
- Ferguson, T. S. (1958) "A Method of Generating Best Asymptotically Normal Estimates with Application to the Estimation of Bacterial Densities", *Annals of Mathematical Statistics*, 29, 1046–1062.
- Fisher, R. A. (1935) "The Logic of Inductive Inference", *Journal of the Royal Statistical Society*, Series B, 98, 39–54.
- Flinn, C. J. and J. J. Heckman (1982) "Models for the Analysis of Labor Force Dynamics", in G. Rhodes and R. L. Basmann (eds.), *Advances in Econometrics*, Vol. 1, JAI Press, Greenwich.
- Flinn, C. J. and J. J. Heckman (1982a) "New Methods for Analyzing Structural Models of Labor Force Dynamics", *Journal of Econometrics*, 18, 115–168.
- Florens, J. P. and M. Mouchart (1982) "A Note on Noncausality", *Econometrica*, 50, 583–591.
- Ghez, G. R. and G. S. Becker (1975) *The Allocation of Time and Goods Over the Life Cycle*, Columbia University Press, New York.
- Goldberger, A. S. (1974) "Asymptotics of the Sample Regression Slope", unpublished lecture notes, No. 12.
- Goldberger, A. S. (1974a) "Unobservable Variables in Econometrics", in P. Zarembka (ed.), *Frontiers in Econometrics*, Academic Press, New York.
- Granger, C. W. J. (1969) "Investigating Causal Relations by Econometric Models and Cross-Spectral Methods", *Econometrica*, 37, 424–438.
- Griliches, Z. (1967) "Distributed Lags: A Survey", *Econometrica*, 35, 16–49.
- Griliches, Z. (1974) "Errors in Variables and Other Unobservables", *Econometrica*, 42, 971–998.
- Griliches, Z. (1976) "Wages of Very Young Men", *Journal of Political Economy*, 84, S69–S85.
- Griliches, Z. (1979) "Sibling Models and Data in Economics: Beginnings of a Survey", *Journal of Political Economy*, 87, S37–S64.
- Griliches, Z., B. H. Hall, and J. A. Hausman (1978) "Missing Data and Self-Selection in Large Panels", *Annales de l'INSEE*, 30/31, 137–176.

- Griliches, Z., B. H. Hall, and J. A. Hausman (1981) "Econometric Models for Count Data with an Application to the Patents—R & D Relationship", National Bureau of Economic Research Technical Paper No. 17.
- Griliches, Z. and A. Pakes (1980) "The Estimation of Distributed Lags in Short Panels", National Bureau of Economic Research Technical Paper No. 4.
- Gronau, R. (1973) "The Effect of Children on the Housewife's Value of Time", *Journal of Political Economy*, 81, S168–S199.
- Gronau, R. (1976) "The Allocation of Time of Israeli Women", *Journal of Political Economy*, 84, S201–S220.
- Gronau, R. (1977) "Leisure, Home Production, and Work—the Theory of the Allocation of Time Revisited", *Journal of Political Economy*, 85, 1099–1123.
- Hall, R. E. (1973) "Comment", *Journal of Political Economy*, 81, S200–S201.
- Hall, R. E. (1978) "Stochastic Implications of the Life Cycle—Permanent Income Hypothesis: Theory and Evidence", *Journal of Political Economy*, 86, 971–988.
- Hanoch, G. (1980) "A Multivariate Model of Labor Supply: Methodology and Estimation", in J. P. Smith (ed.), *Female Labor Supply: Theory and Estimation*, Princeton University Press, Princeton.
- Hansen, L. P. (1982) "Large Sample Properties of Generalized Method of Moments Estimators", *Econometrica*, 50, 1029–1054.
- Hansen, L. P. and K. J. Singleton (1982) "Generalized Instrumental Variable Estimation of Nonlinear Rational Expectations Models", *Econometrica*, 50, 1269–1286.
- Harville, D. A. (1977) "Maximum Likelihood Approaches to Variance Component Estimation and to Related Problems", *Journal of the American Statistical Association*, 72, 320–338.
- Hause, J. (1977) "The Covariance Structure of Earnings and the On-the-Job Training Hypothesis", *Annals of Economic and Social Measurement*, 6, 335–365.
- Hause, J. (1980) "The Fine Structure of Earnings and the On-the-Job Training Hypothesis", *Econometrica*, 48, 1013–1029.
- Hausman, J. A. and D. A. Wise (1979) "Attrition Bias in Experimental and Panel Data: The Gary Income Maintenance Experiment", *Econometrica*, 47, 455–473.
- Hausman, J. A. and W. E. Taylor (1981) "Panel Data and Unobservable Individual Effects", *Econometrica*, 49, 1377–1398.
- Heckman, J. J. (1974) "Shadow Prices, Market Wages, and Labor Supply", *Econometrica*, 42, 679–694.
- Heckman, J. J. (1976) "A Life-Cycle Model of Earnings, Learning, and Consumption", *Journal of Political Economy*, 84, S11–S44.
- Heckman, J. J. (1978) "Simple Statistical Models for Discrete Panel Data Developed and Applied to Test the Hypothesis of True State Dependence Against the Hypothesis of Spurious State Dependence", *Annales de l'INSEE*, 30/31, 227–269.
- Heckman, J. J. (1980) "Sample Selection Bias as a Specification Error", in J. P. Smith (ed.), *Female Labor Supply: Theory and Estimation*, Princeton University Press, Princeton.
- Heckman, J. J. (1981) "Statistical Models for Discrete Panel Data", in C. Manski and D. McFadden (eds.), *Structural Analysis of Discrete Data with Econometric Applications*, MIT Press, Cambridge.
- Heckman, J. J. (1981a) "The Incidental Parameters Problem and the Problem of Initial Conditions in Estimating Discrete Time—Discrete Data Stochastic Processes and Some Monte Carlo Evidence", in C. Manski and D. McFadden (eds.), *Structural Analysis of Discrete Data with Econometric Applications*, MIT Press, Cambridge.
- Heckman, J. J. (1981b) "Heterogeneity and State Dependence", in S. Rosen (ed.), *Conference on Labor Markets*, University of Chicago Press, Chicago.
- Heckman, J. J. and G. Borjas (1980) "Does Unemployment Cause Future Unemployment? Definitions, Questions and Answers from a Continuous Time Model of Heterogeneity and State Dependence", *Econometrica*, 47, 247–283.
- Heckman, J. J. and T. E. MaCurdy (1980) "A Life-Cycle Model of Female Labor Supply", *Review of Economic Studies*, 47, 47–74.
- Heckman, J. J. and R. J. Willis (1977) "A Beta-Logistic Model for the Analysis of Sequential Labor Force Participation by Married Women", *Journal of Political Economy*, 85, 27–58.
- Hosoya, Y. (1977) "On the Granger Condition for Non-Causality", *Econometrica*, 45, 1735–1736.
- Hsiao, C. (1975) "Some Estimation Methods for a Random Coefficient Model", *Econometrica*, 43,

- 305–325.
- Huber, P. J. (1967) “The Behavior of Maximum Likelihood Estimates Under Nonstandard Conditions”, *Proceedings of the Fifth Berkeley Symposium on Mathematical Statistics and Probability*, Vol. 1, University of California Press, Berkeley and Los Angeles.
- Jennrich, R. I. (1969) “Asymptotic Properties of Non-Linear Least Squares Estimators”, *The Annals of Mathematical Statistics*, 40, 633–643.
- Jöreskog, K. G. (1978) “An Econometric Model for Multivariate Panel Data”, *Annales de l'INSEE*, 30/31, 355–366.
- Jöreskog, K. G. and A. S. Goldberger (1975) “Estimation of a Model with Multiple Indicators and Multiple Causes of a Single Latent Variable”, *Journal of the American Statistical Association*, 70, 631–639.
- Jöreskog, K. G. and D. Sörbom (1977) “Statistical Models and Methods for Analysis of Longitudinal Data”, in D. J. Aigner and A. S. Goldberger (eds.), *Latent Variables in Socio-economic Models*, North Holland, Amsterdam.
- Kalbfleisch, J. D. and D. A. Sprott (1970) “Application of Likelihood Methods to Models Involving Large Numbers of Parameters”, *Journal of the Royal Statistical Society, Series B*, 32, 175–208.
- Kendall, M. G. and A. Stuart (1961) *The Advanced Theory of Statistics*, Vol. 2, Griffin, London.
- Kiefer, J. and J. Wolfowitz (1956) “Consistency of the Maximum Likelihood Estimator in the Presence of Infinitely Many Incidental Parameters”, *Annals of Mathematical Statistics*, 27, 887–906.
- Kiefer, N. and G. Neumann (1981) “Individual Effects in a Nonlinear Model”, *Econometrica*, 49, 965–980.
- Lancaster, T. (1979) “Econometric Methods for the Duration of Unemployment”, *Econometrica*, 47, 939–956.
- Lancaster, T. and S. Nickell (1980) “The Analysis of Re-Employment Probabilities for the Unemployed”, *Journal of the Royal Statistical Society, Series A*, 143, 141–165.
- Lee, L. F. (1979) “Estimation of Error Components Model with ARMA (P, Q) Time Component—An Exact GLS Approach”, unpublished manuscript.
- Lee, L. F. (1980) “Analysis of Econometric Models for Discrete Panel Data in the Multivariate Log Linear Probability Models”, unpublished manuscript.
- Levhari, D. and T. N. Srinivasan (1969) “Optimal Savings Under Uncertainty”, *Review of Economic Studies*, 36, 153–163.
- Lillard, L. and Y. Weiss (1979) “Components of Variation in Panel Earnings Data: American Scientists 1960–1970”, *Econometrica*, 47, 437–454.
- Lillard, L. and R. J. Willis (1978) “Dynamic Aspects of Earnings Mobility”, *Econometrica*, 46, 985–1012.
- MaCurdy, T. E. (1979) “Multiple Time Series Models Applied to Panel Data: Specification of a Dynamic Model of Labor Supply”, unpublished manuscript.
- MaCurdy, T. E. (1981) “An Empirical Model of Labor Supply in a Life-Cycle Setting”, *Journal of Political Economy*, 89, 1059–1085.
- MaCurdy, T. E. (1981a) “Asymptotic Properties of Quasi-Maximum Likelihood Estimators and Test Statistics”, National Bureau of Economic Research Technical Paper No. 14.
- MaCurdy, T. E. (1982) “The Use of Time Series Processes to Model the Error Structure of Earnings in a Longitudinal Data Analysis”, *Journal of Econometrics*, 18, 83–114.
- Madalla, G. S. (1971) “The Use of Variance Components Models in Pooling Cross Section and Time Series Data”, *Econometrica*, 39, 341–358.
- Madalla, G. S. and T. D. Mount (1973) “A Comparative Study of Alternative Estimators for Variance-Components Models Used in Econometric Applications”, *Journal of the American Statistical Association*, 68, 324–328.
- Malinvaud, E. (1970) *Statistical Methods of Econometrics*, North-Holland, Amsterdam.
- Mazodier, P. and A. Trognon (1978) “Heteroskedasticity and Stratification in Error Components Models”, *Annales de l'INSEE*, 30/31, 451–482.
- McFadden, D. (1974) “Conditional Logit Analysis of Qualitative Choice Behavior”, in P. Zarembka, (ed.), *Frontiers in Econometrics*, New York: Academic Press.
- Mellow, W. (1981) “Unionism and Wages: A Longitudinal Analysis”, *Review of Economics and Statistics*, 63, 43–52.
- Mincer, J. (1963) “Market Prices, Opportunity Costs and Income Effects”, in C. F. Christ, et. al.,

- Measurement in Economics*, Stanford University Press, Stanford.
- Mincer, J. and S. Polachek (1974) "Family Investments in Human Capital: Earnings of Women", *Journal of Political Economy*, 82, S76–S108.
- Mundlak, Y. (1961) "Empirical Production Function Free of Management Bias", *Journal of Farm Economics*, 43, 44–56.
- Mundlak, Y. (1963) "Estimation of Production and Behavioral Functions From a Combination of Time Series and Cross Section Data", in C. F. Christ, et al., *Measurement in Economics*, Stanford University Press, Stanford.
- Mundlak, Y. (1978) "On the Pooling of Time Series and Cross Section Data", *Econometrica*, 46, 69–85.
- Mundlak, Y. (1978a) "Models with Variable Coefficients: Integration and Extension", *Annales de l'INSEE*, 30/31, 483–509.
- Nerlove, M. (1967) "Experimental Evidence on the Estimation of Dynamic Economic Relations from a Time Series of Cross-Sections", *Economic Studies Quarterly*, 18, 42–74.
- Nerlove, M. (1971) "Further Evidence on the Estimation of Dynamic Relations from a Time Series of Cross-Sections", *Econometrica*, 39, 359–382.
- Nerlove, M. (1971a) "A Note on Error-Components Models", *Econometrica*, 39, 383–396.
- Nerlove, M. (1972) "Lags in Economic Behavior", *Econometrica*, 40, 221–251.
- Neyman, J. and E. L. Scott (1948) "Consistent Estimates Based on Partially Consistent Observations", *Econometrica*, 16, 1–32.
- Nickell, S. (1979) "Estimating the Probability of Leaving Unemployment", *Econometrica*, 47, 1249–1266.
- Nickell, S. (1981) "Biases in Dynamic Models with Fixed Effects", *Econometrica*, 49, 1417–1426.
- Phelps, E. S. (1962) "The Accumulation of Risky Capital: A Sequential Utility Analysis", *Econometrica*, 30, 729–743.
- Rao, C. R. (1973) *Linear Statistical Inference and Its Applications*, Wiley, New York.
- Rasch, G. (1960) *Probabilistic Models for Some Intelligence and Attainment Tests*, Denmark's Paedagogiske Institute, Copenhagen.
- Rasch, G. (1961) "On General Laws and the Meaning of Measurement in Psychology", *Proceedings of the Fourth Berkeley Symposium on Mathematical Statistics and Probability*, Vol. 4, University of California Press, Berkeley and Los Angeles.
- Rosenzweig, M. R. and K. I. Wolpin (1980) "Life-Cycle Labor Supply and Fertility", *Journal of Political Economy*, 88, 328–348.
- Rothenberg, T. J. (1973) *Efficient Estimation with A Priori Information*, Yale University Press, New Haven.
- Schultz, T. P. (1980) "Estimating Labor Supply Functions for Married Women", in J. P. Smith (ed.), *Female Labor Supply: Theory and Estimation*, Princeton University Press, Princeton.
- Sims, C. A. (1972) "Money, Income, and Causality", *The American Economic Review*, 62, 540–552.
- Singer, B. and S. Spilerman (1974) "Social Mobility Models for Heterogeneous Populations", in H. L. Costner (ed.), *Sociological Methodology 1973–1974*, Jossey-Bass, Inc., San Francisco.
- Singer, B. and S. Spilerman (1976) "Some Methodological Issues in the Analysis of Longitudinal Surveys", *Annals of Economic and Social Measurement*, 5, 447–474.
- Stafford, F. P. and G. J. Duncan (1980) "Do Union Members Receive Compensating Wage Differentials?", *American Economic Review*, 70, 355–371.
- Swamy, P. A. V. B. (1970) "Efficient Inference in a Random Coefficient Regression Model", *Econometrica*, 38, 311–323.
- Swamy, P. A. V. B. (1974) "Linear Models with Random Coefficients", in P. Zarembka (ed.), *Frontiers in Econometrics*, Academic Press, New York.
- Taylor, W. E. (1980) "Small Sample Considerations in Estimation from Panel Data", *Journal of Econometrics*, 13, 203–223.
- Trognon, A. (1978) "Miscellaneous Asymptotic Properties of Ordinary Least Squares and Maximum Likelihood Estimators in Dynamic Error Components Models", *Annales de l'INSEE*, 30/31, 631–657.
- Tuma, N. B., M. T. Hannan, and L. P. Groeneveld (1979) "Dynamic Analysis of Event Histories", *American Journal of Sociology*, 84, 820–854.
- Tuma, N. B. and P. K. Robins (1980) "A Dynamic Model of Employment Behavior: An Application

- to the Seattle and Denver Income Maintenance Experiments", *Econometrica*, 48, 1031–1052.
- Wallace, T. D. and A. Hussain (1969) "The Use of Error Components Models in Combining Time Series with Cross Section Data", *Econometrica*, 37, 55–72.
- White, H. (1980) "Using Least Squares to Approximate Unknown Regression Functions", *International Economic Review*, 21, 149–170.
- White, H. (1980a) "Nonlinear Regression on Cross Section Data", *Econometrica*, 48, 721–746.
- White, H. (1980b) "A Heteroskedasticity—Consistent Covariance Matrix Estimator and a Direct Test for Heteroskedasticity", *Econometrica*, 48, 817–838.
- White, H. (1982) "Maximum Likelihood Estimation of Misspecified Models", *Econometrica*, 50, 1–25.
- White, H. (1982a) "Instrumental Variable Regression with Independent Observations", *Econometrica*, 50, 483–499.
- Willis, R. J. (1973) "A New Approach to the Economic Theory of Fertility Behavior", *Journal of Political Economy*, 81, S14–S64.
- Zellner, A. and H. Theil (1962) "Three-Stage Least Squares: Simultaneous Estimation of Simultaneous Equations", *Econometrica*, 30, 54–78.
- Zellner, A., J. Kmenta and J. Drèze (1966) "Specification and Estimation of Cobb–Douglas Production Function Models", *Econometrica*, 34, 784–795.