

NBER TECHNICAL PAPER SERIES

MULTIVARIATE REGRESSION MODELS  
FOR PANEL DATA

Gary Chamberlain

Technical Paper No. 8

NATIONAL BUREAU OF ECONOMIC RESEARCH  
1050 Massachusetts Avenue  
Cambridge MA 02138

December 1980

I am grateful to Arthur Goldberger, Zvi Griliches, and Ariel Pakes for helpful discussions. Financial support was provided by the National Science Foundation (Grant No. SOC 7925959). The research reported here is part of the NBER's research program in Labor Economics and in Productivity. Any opinions expressed are those of the author and not those of the National Bureau of Economic Research.

Multivariate Regression Models For  
Panel Data

ABSTRACT

Under stationarity, the heterogeneous stochastic processes are the non-ergodic ones. We show that if a distributed lag is of finite order, then its coefficients are unconditional means of the underlying random coefficients. This result is applied to linear transformations of the process. The estimation framework is a multivariate wide-sense regression function. The identification analysis requires certain restrictions on the coefficients. The actual regression function is nonlinear, and so we provide a theory of inference for linear approximations. It rests on obtaining the asymptotic distribution of functions of sample moments. Restrictions are imposed by using a minimum distance estimator; it is generally more efficient than the conventional estimators.

Professor Gary Chamberlain  
Department of Economics  
Social Science Building  
1180 Observatory Drive  
University of Wisconsin  
Madison, WI 53706

(608) 262-7789

# MULTIVARIATE REGRESSION MODELS FOR PANEL DATA

by

Gary Chamberlain  
University of Wisconsin-Madison

## 1. INTRODUCTION

We shall consider a stochastic process  $\{(x_t, y_t), t \in T\}$ , where  $T$  is the set of all integers. We are primarily interested in the distribution of  $\{y_t, t \in T\}$  conditional on  $\{x_t, t \in T\}$ . The paper has two main sections: one on identification and one on estimation. The identification problem is that the parameters of interest depend on the entire distribution of  $\{(x_t, y_t), t \in T\}$ , but we are only given the distribution of  $\{(x_t, y_t), t=1, \dots, T\}$ . The estimation problem is to use a sample,  $(x_{i1}, y_{i1}, \dots, x_{iT}, y_{iT}), i=1, \dots, N$ , to make inferences about the distribution of  $\{(x_t, y_t), t=1, \dots, T\}$ .

A typical application will be based on following people or firms through time. We do not assume that observing a long time series on a single individual allows us to infer the distribution of the stochastic process. In the case of a stationary process, the time average  $\sum_{t=1}^T g(x_t, y_t)/T$  converges to a well defined random variable  $(\hat{g})$  as  $T \rightarrow \infty$  if  $E|g(x_t, y_t)| < \infty$ . We do not assume that  $\hat{g} = E[g(x_t, y_t)]$ ; this is true if the variance of  $\hat{g}$  is zero, but in general the variance will not be zero. If there is some  $\hat{g}$  with non-zero variance, then we say that the  $\{x_t, y_t\}$  process is heterogeneous. Then the heterogeneous stationary processes are precisely the non-ergodic ones.

The associated identification problem is that a parameter of interest may be  $E(\hat{b}) = E[\lim_{T \rightarrow \infty} \frac{\sum_{t=1}^T (x_t - \bar{x}_T)(y_t - \bar{y}_T)}{\sum_{t=1}^T (x_t - \bar{x}_T)^2}]$ . If the  $\{x_t, y_t\}$  process is ergodic, then  $E(\hat{b}) = \text{Cov}(x_1, y_1)/V(x_1)$ , and we can identify  $E(\hat{b})$  with  $T = 1$ . With a non-ergodic process, however, we can only say that  $E(\lim_{T \rightarrow \infty} \bar{x}_T) = E(x_1)$ ,  $E(\lim_{T \rightarrow \infty} \bar{y}_T) = E(y_1)$ ,  $E(\lim_{T \rightarrow \infty} \sum_{t=1}^T x_t^2/T) = E(x_1^2)$ , and  $E(\lim_{T \rightarrow \infty} \sum_{t=1}^T x_t y_t/T) = E(x_1 y_1)$ . Since in general  $E(m_1 m_2) \neq E(m_1) E(m_2)$  and  $E(m_1/m_2) \neq E(m_1)/E(m_2)$ , we need more restrictions and may need  $T > 1$ . The population parameter  $E(\hat{b})$  corresponds to the following estimator: accumulate a long time series on each of a large number of randomly sampled individuals; compute a linear regression coefficient for each individual; then average these coefficients across the individuals. This does not, in general, correspond to the regression coefficient obtained from a large cross-section sample at a single point in time.

Our solution to the identification problem takes the following form: we assume that  $E(y_t | \dots, x_{-1}, x_0, x_1, \dots, c, \dots, b_{-1}, b_0, b_1, \dots) = c + \sum_{j=-\infty}^{\infty} b_j x_{t-j}$ , where  $c$  and the  $b_j$  are constant over time but may vary across individuals; then if the  $\{x_t\}$  process satisfies a regularity condition, our theorem asserts that if  $E(y_t | \dots, x_{-1}, x_0, x_1, \dots) = E(y_t | x_t)$ , then  $E(y_t | x_t) = \gamma + \beta x_t$ , where  $\beta = E(b_0) = E(\hat{b})$  and  $E(b_j) = 0$ ,  $j \neq 0$ . There are corresponding results for the case in which  $E(y_t | \dots, x_{-1}, x_0, x_1, \dots)$  is not contemporaneous but displays lags and leads of finite order. We show how the theorem can be applied to linear transformations of the process; for example, if  $\bar{y}_t = y_t - y_{t-1}$ ,  $\bar{x}_t = x_t - x_{t-1}$ , then  $E(\bar{y}_t | \dots, \bar{x}_{-1}, \bar{x}_0, \bar{x}_1, \dots) = E(\bar{y}_t | \bar{x}_t)$  implies that  $E(\bar{y}_t | \bar{x}_t) = \gamma + \beta \bar{x}_t$ , where  $\beta = E(b_0)$  and  $E(b_j) = 0$ ,  $j \neq 0$ .

Alternative linear transformations allow us to treat general rational distributed lag schemes.

Our estimation results focus on linear approximations. We assume that  $(y_{i1}, \dots, y_{iT}, x_{i1}, \dots, x_{iT}) = (\underline{y}_i', \underline{x}_i')$ ,  $i=1, \dots, N$ , are independently and identically distributed (i.i.d). We consider the following minimum mean square error linear predictors:  $E^*(y_{iT} | \underline{x}_i) = \eta_t + \underline{\pi}_t' \underline{x}_i$ , which we write as  $E^*(\underline{y}_i | \underline{x}_i) = \underline{\eta} + \underline{\Pi} \underline{x}_i$ . The identification results suggest examining certain restrictions on  $\underline{\Pi}$ . It is not appropriate, however, to assume that  $E^*(\underline{y}_i | \underline{x}_i) = E(\underline{y}_i | \underline{x}_i)$ ; so we require a theory of inference for linear approximations.

This is provided by noting that  $\underline{\Pi}$  is a function of population moments. The least squares estimator is the corresponding function of the sample moments. Hence the problem is to derive the limiting distribution of a differentiable function of sample moments, under i.i.d. sampling. This is straightforward and gives  $\sqrt{N}(\text{vec } \hat{\underline{\Pi}}' - \text{vec } \underline{\Pi}') \xrightarrow{D} N(0, \underline{\Omega})$  as  $N \rightarrow \infty$ . We simplify the formula for  $\underline{\Omega}$  and provide a consistent estimator  $\hat{\underline{\Omega}}$ .

We can impose restrictions within this framework by using a minimum distance estimator. We find the  $\underline{\Pi}$  matrix satisfying the restrictions that is closest to  $\hat{\underline{\Pi}}$  in the norm provided by  $\hat{\underline{\Omega}}^{-1}$ . This leads to some surprising results. For example, consider a univariate linear predictor ( $T=1$ ):  $E^*(y_i | x_{i1}, x_{i2}) = \eta + \pi_1 x_{i1} + \pi_2 x_{i2}$ . We can impose the restriction that  $\pi_2 = 0$  by using the least squares regression of  $y_i$  on  $x_{i1}$  to provide an estimator of  $\pi_1$ ; however, that estimator is less efficient, in general, than our minimum distance estimator of  $\pi_1$ .

## 2. IDENTIFICATION

Consider a bivariate stochastic process  $\{(x_t, y_t), t \in T\}$ , where  $T$  is the set of all integers. We are interested in the distribution of  $\{y_t, t \in T\}$  conditional on  $\{x_t, t \in T\}$ . A realization of the process is a point  $\omega \in \Omega$  consisting of a doubly infinite sequence of vectors:  $\omega = \{\dots, \omega_{-1}, \omega_0, \omega_1, \dots\} = \{\omega_t, t \in T\}$ , where  $\omega_t = (\omega_{1t}, \omega_{2t})$ ;  $(x_t(\omega), y_t(\omega)) = (\omega_{1t}, \omega_{2t})$  is the  $t^{\text{th}}$  coordinate function.

A distinctive feature is that we shall consider the case in which more than one (partial) realization of the process is observed. Let  $z_i = (x_{i1}, y_{i1}, \dots, x_{iT}, y_{iT})$ . We have a sample  $z_i, i=1, \dots, N$ , where we regard the  $z_i$  as independent and identically distributed (i.i.d) according to the probability law of  $\{(x_t, y_t), t=1, \dots, T\}$ . Typically  $i$  will index people or firms.

We shall consider in some detail the case of a stationary process. This will allow us to state clearly the key definitions and to obtain some sharp results. These results should provide guidance in attempts to deal with the additional identification problems created by non-stationarity. Our estimation results will not rely on any stationarity assumptions.

Consider the sample mean for a single individual (point):

$$\bar{x}_T(\omega) = \frac{1}{T} \sum_{t=1}^T x_t(\omega).$$

We shall not assume that observing a very long time series for a single individual allows us to infer the distribution of the stochastic process. In particular, we shall not assume that  $\lim_{T \rightarrow \infty} \bar{x}_T(\omega) = E(x_t)$  as  $T \rightarrow \infty$ . The

assumption that  $\{(x_t, y_t), t \in T\}$  is a stationary stochastic process implies that  $\lim_{T \rightarrow \infty} \bar{x}_T(\omega)$  exists except for  $\omega$  in a set with probability measure zero, which we shall abbreviate by saying that it exists a.e. We can set  $\hat{x}(\omega) = \lim_{T \rightarrow \infty} \bar{x}_T(\omega)$  and  $\hat{x}$  is a well defined random variable. If it is degenerate, then  $\hat{x} = E(x_t)$  a.e.; but the variance of  $\hat{x}$ , which we shall denote by  $V(\hat{x})$ , need not be zero.

More generally, consider time averages of the following form:

$$\lim_{T \rightarrow \infty} \frac{1}{T} \sum_{t=1}^T g(x_t(\omega), y_t(\omega)) = \lim_{T \rightarrow \infty} \bar{g}_T(\omega).$$

If  $E|g(x_t, y_t)| < \infty$ , then the ergodic theorem asserts that  $\hat{g}(\omega) = \lim_{T \rightarrow \infty} \bar{g}_T(\omega)$  exists a.e. and  $E(\hat{g}) = E[g(x_t, y_t)]$ ; in addition, we can interpret  $\hat{g}$  as a conditional expectation:  $\hat{g} = E[g(x_t, y_t) | J]$ , where  $J$  is the  $\sigma$ -algebra of shift invariant sets. (See Billingsley (1965), theorem 1.3, p. 13 or Rozanov (1967), theorem 5.1, p. 157). If  $V(\hat{g}) = 0$ , then  $\hat{g} = E[g(x_t, y_t)]$  a.e. In general, however,  $V(\hat{g}) \neq 0$ . We shall say that there is heterogeneity if  $V(\hat{g}) \neq 0$  for some  $g$ . Then there is heterogeneity if and only if the process is not ergodic.

What is heterogeneity bias? Consider the following time average:

$$\hat{C}(x_t(\omega), y_t(\omega)) = \lim_{T \rightarrow \infty} \frac{1}{T} \sum_{t=1}^T x_t(\omega) y_t(\omega) - \left( \lim_{T \rightarrow \infty} \frac{1}{T} \sum_{t=1}^T x_t(\omega) \right) \left( \lim_{T \rightarrow \infty} \frac{1}{T} \sum_{t=1}^T y_t(\omega) \right).$$

If  $E(x_t^2) < \infty$  and  $E(y_t^2) < \infty$ , then  $\hat{C}(x_t, y_t)$  is a well defined, integrable random variable. Say that the population parameter of interest is  $E[\hat{C}(x_t, y_t)]$ . Consider measuring  $E[\hat{C}(x_t, y_t)]$  by using  $\text{Cov}(x_t, y_t) = C(x_t, y_t) = E(x_t y_t) - E(x_t) E(y_t)$ . We shall say that there is a heterogeneity bias if  $E[\hat{C}(x_t, y_t)] \neq C(x_t, y_t)$ . Since

$$E[\hat{C}(x_t, y_t)] = E(x_t y_t) - E(\hat{x} \hat{y}), \quad \text{where}$$

$$\hat{x}(\omega) = \lim_{T \rightarrow \infty} \frac{1}{T} \sum_{t=1}^T x_t(\omega), \quad \hat{y}(\omega) = \lim_{T \rightarrow \infty} \frac{1}{T} \sum_{t=1}^T y_t(\omega),$$

and since  $E(\hat{x}) = E(x_t)$  and  $E(\hat{y}) = E(y_t)$ , we see that there is no heterogeneity bias if and only if  $\text{Cov}(\hat{x}, \hat{y}) = 0$ . The absence of heterogeneity in  $x$  or in  $y$  is a sufficient condition for  $\text{Cov}(\hat{x}, \hat{y}) = 0$ , for then  $V(\hat{x}) = 0$  or  $V(\hat{y}) = 0$ .

The associated identification problem is somewhat more general. Clearly  $E[\hat{C}(x_t, y_t)]$  depends on the entire distribution of  $\{(x_t, y_t), t \in T\}$ . We would like to know whether  $E[\hat{C}(x_t, y_t)]$  can be recovered from a knowledge of the marginal distribution of  $\{(x_t, y_t), t = 1, \dots, T\}$ , where  $T$  is small. This problem is relevant if we have an i.i.d. sample of  $(x_{i1}, y_{i1}, \dots, x_{iT}, y_{iT})$ ,  $i = 1, \dots, N$ , where  $N$  is large. The use of  $C(x_t, y_t)$  is an attempt to recover  $E[\hat{C}(x_t, y_t)]$  from the marginal distribution of  $(x_1, y_1)$  when  $T = 1$ .

For a regression example, consider the following:

$$b_0(\omega) = \lim_{T \rightarrow \infty} \frac{\frac{1}{T} \sum_{t=1}^T x_t(\omega) y_t(\omega) - \left( \frac{1}{T} \sum_{t=1}^T x_t(\omega) \right) \left( \frac{1}{T} \sum_{t=1}^T y_t(\omega) \right)}{\frac{1}{T} \sum_{t=1}^T x_t^2(\omega) - \left( \frac{1}{T} \sum_{t=1}^T x_t(\omega) \right)^2} = \frac{\hat{C}(x_t, y_t)}{\hat{V}(x_t)}.$$

This is a well defined random variable if  $\hat{V}(x_t) \neq 0$  a.e. This restriction on  $\hat{V}(x_t)$  is an example of a regularity condition on the  $\{x_t\}$  process; it is violated if there is a set of individuals with non-zero probability for whom  $x_t$  is constant over time. Assume that  $E|b_0| < \infty$  and say that  $E(b_0)$  is the population parameter of interest. Consider measuring  $E(b_0)$  by using



$\beta_0 = C(x_t, y_t)/V(x_t)$ . We shall say that there is a heterogeneity bias if  $E(b_0) \neq \beta_0$ .

In order to analyze the bias, assume that

$$E^*(y_t | x_t, c, b_0) = c + b_0 x_t$$

where

$$b_0 = \frac{\hat{C}(x_t, y_t)}{\hat{V}(x_t)}, \quad c = \hat{y} - b_0 \hat{x},$$

and  $E^*$  is the minimum mean square error linear predictor -- the linear or wide sense regression function. (We would obtain this directly if we used  $E^*(y_t | x_t, J)$  since  $b_0 = \hat{C}(x_t, y_t)/\hat{V}(x_t) = C(x_t, y_t | J)/V(x_t | J)$ .) Since  $\hat{C}(x_t, y_t) = b_0 \hat{V}(x_t)$ , we have

$$\frac{E[\hat{C}(x_t, y_t)]}{E[\hat{V}(x_t)]} = E(b_0)$$

if  $C(b_0, \hat{V}(x_t)) = 0$ , so that there is no correlation between the conditional variance of  $x$  and the slope of the linear predictor. In that case we can use

$$\begin{aligned} C(x_t, y_t) &= E[\hat{C}(x_t, y_t)] + C(\hat{x}, \hat{y}) \\ V(x_t) &= E[\hat{V}(x_t)] + V(\hat{x}) \end{aligned}$$

to obtain

$$\begin{aligned} \beta &= \frac{C(x_t, y_t)}{V(x_t)} = \frac{E[\hat{C}(x_t, y_t)] E[\hat{V}(x_t)]}{E[\hat{V}(x_t)] V(x_t)} + \frac{C(\hat{x}, \hat{y})}{V(\hat{x})} \frac{V(\hat{x})}{V(x_t)} \\ &= \lambda E(b_0) + (1-\lambda) \frac{C(\hat{x}, \hat{y})}{V(\hat{x})}, \end{aligned}$$

where  $\lambda = \frac{E[\hat{V}(x_t)]}{V(x_t)}, \quad 0 < \lambda \leq 1.$

Hence there is no heterogeneity bias if and only if  $V(\hat{x}) = 0$  or

$$\frac{C(x_t, y_t)}{V(x_t)} = \frac{C(\hat{x}, \hat{y})}{V(\hat{x})}$$

-- the cross sectional linear regression must coincide with the linear regression of the long run mean of  $y$  on the long run mean of  $x$ .

This suggests that the linear regression of  $\bar{y}_T = \sum_{t=1}^T y_t / T$  on  $\bar{x}_T = \sum_{t=1}^T x_t / T$

should have the same slope as the linear regression of  $y_t$  on  $x_t$ . If  $T = 2$ , this is equivalent to the linear regression of  $y_2 - y_1$  on  $x_2 - x_1$  having the same slope as the linear regression of  $y_1$  on  $x_1$  or  $y_2$  on  $x_2$ . A sufficient condition is that

$$E^*(y_t | x_1, x_2) = E^*(y_t | x_t) = \gamma + \beta_0 x_t, \quad t=1,2,$$

so that

$$E^*(y_2 - y_1 | x_1, x_2) = \beta_0 (x_2 - x_1).$$

## 2.1. THE THEOREM

A general version of a related argument is presented in appendix A.

We assume that the following regression function is in fact linear:

$$E(y_t | \dots x_{-1}, x_0, x_1, \dots, c, \dots, b_{-1}, b_0, b_1, \dots) \\ = c + \sum_{j=-\infty}^{\infty} b_j x_{t-j}, \quad t \in T,$$

where the series is assumed to be absolutely convergent a.e. and integrable. The coefficients  $c, b_j$  are integrable, invariant random variables. The  $\{x_t\}$  process is assumed to satisfy a regularity condition (C), which will be discussed below. Consider the regression function  $E(y_t | \dots, x_{-1}, x_0, x_1, \dots)$ , which we shall denote by  $E(y_t | x)$ ; it does not condition on the random coefficients. We show that if this regression function displays finite lags and leads, then it must be linear and its coefficients are unconditional means of  $c$  and the  $b_j$ :

$$E(y_t | x) = E(y_t | x_{t-J}, \dots, x_{t+M}) \quad (0 \leq J, M < \infty)$$

implies that

$$E(y_t | x) = \gamma + \sum_{j=-M}^J \beta_j x_{t-j}, \text{ where}$$

$$\gamma = E(c) = E(c | x) \text{ a.e.}$$

$$\beta_j = E(b_j) = E(b_j | x) \text{ a.e. for } j \in [-M, J]$$

$$E(b_j) = E(b_j | x) = 0 \text{ a.e. for } j \notin [-M, J].$$

So in this case there is no heterogeneity bias. If invariant random coefficients are correlated with  $x$ , then they must be correlated with the values of  $x$  in the distant future and distant past.

Next we shall indicate how the theorem can be applied to linear transformations of the model. Then we shall provide some interpretation of the regularity condition on  $\{x_t\}$ .

## 2.2 LINEAR TRANSFORMATIONS

Our assumption that  $E(y_t|x, c, \dots, b_{-1}, b_0, b_1, \dots)$  is a linear regression function implies that

$$\begin{aligned} E(\bar{y}_t|x, c, \dots, b_{-1}, b_0, b_1, \dots) \\ = \bar{c} + \sum_{j=-\infty}^{\infty} b_j \bar{x}_{t-j} = \bar{c} + \sum_{j=-\infty}^{\infty} \bar{b}_j x_{t-j}, \end{aligned}$$

where

$$\bar{y}_t = y_t - \delta_1 y_{t-1} - \dots - \delta_L y_{t-L}$$

$$\bar{x}_t = x_t - \delta_1 x_{t-1} - \dots - \delta_L x_{t-L}$$

$$\bar{b}_j = b_j - \delta_1 b_{j-1} - \dots - \delta_L b_{j-L}, \quad \bar{c} = c(1 - \delta_1 - \dots - \delta_L).$$

Hence there are two ways in which we can apply the theorem, depending on which representation we use for the convolution.

We shall first consider the important case  $\bar{y}_t = y_t - y_{t-1}$ ,  $\bar{x}_t = x_t - x_{t-1}$ ; this is the basis for the analysis of covariance estimator.<sup>1</sup> It will be more fruitful to apply the theorem to the first representation of the convolution. Assume that the  $\{\bar{x}_t\}$  process satisfies the regularity condition (C). Then the theorem asserts that if

$$E(\tilde{y}_t | \tilde{x}) = E(\tilde{y}_t | \tilde{x}_{t-J}, \dots, \tilde{x}_{t+M}) \quad (0 \leq J, M < \infty),$$

then

$$E(\tilde{y}_t | \tilde{x}) = \sum_{j=-M}^J \beta_j \tilde{x}_{t-j} \quad \text{where}$$

$$\beta_j = E(b_j) = E(b_j | x) \text{ a.e. for } j \in [-M, J]$$

$$E(b_j) = E(b_j | x) = 0 \quad \text{a.e. for } j \notin [-M, J].$$

Say that we do not find finite lags and leads with  $E(y_t | x)$  but we do with  $E(\tilde{y}_t | \tilde{x})$ . An interpretation is that  $E(c | x) \neq E(c)$ , which causes a heterogeneity bias in  $E(y_t | x)$  that is eliminated by differencing in  $E(\tilde{y}_t | \tilde{x})$ . Another possibility is that  $E(b_j | x) \neq E(b_j)$  but  $E(b_j | \tilde{x}) = E(b_j)$ . This is quite plausible since it is always true that  $E^*(b_j | \tilde{x}) = E(b_j)$ : by stationarity  $\text{Cov}(b_j, \tilde{x}_t) = E(b_j x_t) - E(b_j x_{t-1}) = 0$  since  $b_j$  is an invariant random variable. So  $E(\tilde{y}_t | x)$  will not necessarily display finite lags and leads, but we do not require that for the identification of  $E(b_j)$ . This is the reason that we chose the first form for the convolution.

Consider next the case  $\tilde{y}_t = y_t - \rho y_{t-1}$ , where  $|\rho| < 1$ . Here it is more fruitful to use the second form for the convolution:

$$E(\tilde{y}_t | x, c, \dots, b_{-1}, b_0, b_1, \dots) = \tilde{c} + \sum_{j=-\infty}^{\infty} \tilde{b}_j x_{t-j}.$$

If  $\{x_t\}$  satisfies the regularity condition (C), then

$$E(\tilde{y}_t | x) = E(\tilde{y}_t | x_{t-J}, \dots, x_{t+M}) \quad (0 \leq J, M < \infty)$$

implies that

$$E(\tilde{y}_t | x) = \tilde{\gamma} + \sum_{j=-M}^J \tilde{\beta}_j x_{t-j}, \quad \text{where}$$

$$\tilde{\gamma} = E(\tilde{c}) = E(\tilde{c} | x) \quad \text{a.e.}$$

$$\tilde{\beta}_j = E(\tilde{b}_j) = E(\tilde{b}_j | x) \quad \text{a.e.} \quad \text{for } j \in [-M, J]$$

$$E(\tilde{b}_j) = E(\tilde{b}_j | x) = 0 \quad \text{a.e.} \quad \text{for } j \notin [-M, J].$$

For example, say that  $E(\tilde{y}_t | x) = E(\tilde{y}_t | x_t)$ . Then  $E(\tilde{y}_t | x) = \tilde{\gamma} + \tilde{\beta}_0 x_t$ ,  $E(\tilde{b}_0) = \tilde{\beta}_0$ ,  $E(\tilde{b}_j) = 0$  for  $j \neq 0$ . The solution is

$$E(b_j) = \rho^j \tilde{\beta}_0, \quad j \geq 0$$

$$E(b_j) = 0, \quad j < 0.$$

So the  $E(b_j)$  follow a geometric distributed lag.

More generally, let  $\tilde{y}_t = y_t - \delta_1 y_{t-1} - \dots - \delta_L y_{t-L}$ , where the roots of the lag polynomial lie outside the unit circle. Then  $E(y_t | x) = E(y_t | x_t, \dots, x_{t-J})$  implies that the  $E(b_j)$  follow a general rational distributed lag scheme. For example, if  $\tilde{y}_t = y_t - (\rho_1 + \rho_2) y_{t-1} + \rho_1 \rho_2 y_{t-2}$ ,  $|\rho_1| < 1$ ,  $|\rho_2| < 1$ , and  $E(\tilde{y}_t | x) = E(\tilde{y}_t | x_t, x_{t-1})$ , then we have

$$E(\tilde{y}_t | x) = \tilde{\gamma} + \tilde{\beta}_0 x_t + \tilde{\beta}_1 x_{t-1}$$

$$E(b_j) = \tilde{\beta}_0 \sum_{m=0}^j \rho_1^m \rho_2^{j-m} + \tilde{\beta}_1 \sum_{m=0}^{j-1} \rho_1^m \rho_2^{j-1-m}, \quad j \geq 1,$$

and  $E(b_0) = \tilde{\beta}_0$ ,  $E(b_j) = 0$  for  $j < 0$ .

### 2.3 REGULARITY CONDITIONS

The regularity condition (C) has two parts. The first part (C.1) rules out mover-stayer type processes for which there is non-zero probability that  $x_t = x_{t-1}$  for all  $t$ . We can see the need for this condition by considering  $\tilde{y}_t = y_t - y_{t-1}$ ,  $\tilde{x}_t = x_t - x_{t-1}$ ,  $E(\tilde{y}_t | \tilde{x}) = \beta_0 \tilde{x}_t$ ; say we also assume that  $b_j = 0$  for  $j \neq 0$  so that  $E(\tilde{y}_t | \tilde{x}, b_0) = b_0 \tilde{x}_t$ . Does this imply that  $\beta_0 = E(b_0)$ ? We can say  $E(\tilde{y}_t | \tilde{x}) = E(b_0 | \tilde{x}) \tilde{x}_t$ , which implies that  $E(b_0 | \tilde{x}) = \beta_0$  if  $x_t \neq x_{t-1}$ ; but  $E(b_0 | \tilde{x})$  can take on any value if  $x_t = x_{t-1}$  for all  $t$ , and so  $E(b_0)$  can be arbitrarily far from  $\beta_0$  if there is non-zero probability that  $x_t = x_{t-1}$  for all  $t$ .

If the probability that  $x_t = x_{t-1}$  for all  $t$  is zero, then for almost all  $\omega$  we can find a  $t$ , which may depend on  $\omega$ , such that  $x_t(\omega) \neq x_{t-1}(\omega)$ ; then for that  $t$ ,  $E(b_0 | x)_\omega \tilde{x}_t(\omega) = \beta_0 \tilde{x}_t(\omega)$  implies that  $E(b_0 | x)_\omega = \beta_0$ . Hence  $E(b_0 | x) = \beta_0$  a.e. and so  $E(b_0) = \beta_0$ .

Condition (C.2) applied to  $\tilde{x}_t = x_t - x_{t-1}$  rules out a process in which  $\tilde{x}_t(\omega) = 0$  or  $u(\omega)$  where  $V(u) \neq 0$ ; in this process,  $x$  changes by the same amount for a given individual whenever it changes, and this amount is a non-degenerate random variable. To see the need for this condition, say that  $E(\tilde{y}_t | \tilde{x}) = E(\tilde{y}_t | \tilde{x}_t) = E(b_0 | \tilde{x}) \tilde{x}_t$ . This does not imply that  $E(b_0 | \tilde{x}) = E(b_0)$  a.e. since we can set  $E(\tilde{y}_t | \tilde{x}_t) = \tilde{x}_t^2$  and  $E(b_0 | \tilde{x}) = u$ . This is a valid solution since  $\tilde{x}_t^2 = u \tilde{x}_t$  for all  $t$ ; but  $E(b_0 | \tilde{x}) = u \neq E(b_0)$  if  $V(u) > 0$ .

### 2.4 IDENTIFICATION UNDER ALTERNATIVE ASSUMPTIONS

An alternative approach is to begin with the assumption that  $E(y_t | x, c, \dots, b_{-1}, b_0, b_1, \dots)$  displays finite lags and leads in addition to

being linear:

$$E(y_t | x, c, \dots, b_{-1}, b_0, b_1, \dots) = c + \sum_{j=-M}^J b_j x_{t-j}.$$

We shall modify this formula by limiting ourselves to probability statements that involve only the finite dimensional distribution of  $\{(x_t, y_t), t = 1, \dots, T\}$ . We shall illustrate by considering the case  $T = 2$  and a purely contemporaneous relationship ( $M = 0, J = 0$ ):

$$E(y_t | x_1, x_2, c, b_0) = c + b_0 x_t, \quad t = 1, 2.$$

Then if  $E(y_2 - y_1 | x_2 - x_1)$  is linear, we have

$$E(y_2 - y_1 | x_2 - x_1) = E(b_0 | x_2 - x_1)(x_2 - x_1) = \beta_0(x_2 - x_1).$$

As noted in our discussion of regularity condition (C.1), this does not by itself imply that  $E(b_0) = \beta_0$ ; but if  $P(x_2 = x_1) = 0$ , then we do obtain  $E(b_0 | x_2 - x_1) = \beta_0$  a.e. and so  $E(b_0) = \beta_0$ . The restriction  $P(x_2 = x_1) = 0$  excludes the case in which  $x$  is discrete; for example, if  $x_t = 0$  or  $1$ , then there will be a non-zero probability that  $x_1 = x_2 = 0$  or  $x_1 = x_2 = 1$ .



If we assume that  $P(x_2 = x_1) = 0$ , then there is a more powerful approach that does not rely on  $E(b_o | x_2 - x_1) = E(b_o)$ . We simply say

$$E(y_2 - y_1 | x_2 - x_1) = E(b_o | x_2 - x_1) (x_2 - x_1) \text{ a.e.}$$

$$E\left(\frac{y_2 - y_1}{x_2 - x_1} | x_2 - x_1\right) = E(b_o | x_2 - x_1) \text{ a.e.}$$

Hence if  $E|(y_2 - y_1)/(x_2 - x_1)| < \infty$ , we have

$$E\left(\frac{y_2 - y_1}{x_2 - x_1}\right) = E(b_o).$$

So if we have an i.i.d. sample of  $z_i = (x_{i1}, y_{i1}, x_{i2}, y_{i2})$ ,  $i=1, \dots, N$ , we can obtain a consistent estimator of  $E(b_o)$  as  $N \rightarrow \infty$  from the sample mean of  $(y_{i2} - y_{i1})/(x_{i2} - x_{i1})$ .

For  $T > 2$  we compute a linear least squares regression of  $y_{it}$  on  $x_{it}$ ,  $t=1, \dots, T$ , for each individual  $i$  to obtain a slope coefficient  $\hat{b}_{oi}$ ; then  $\sum_{i=1}^N \hat{b}_{oi}/N$  converges almost surely to  $E(b_o)$  as  $N \rightarrow \infty$  if  $E|\hat{b}_{oi}| < \infty$ . This works since  $E(\hat{b}_{oi} | x_{i1}, \dots, x_{iT}, b_o) = b_o$ ; hence  $E(\hat{b}_{oi}) = E(b_o)$  if  $E|\hat{b}_{oi}| < \infty$ . So we are averaging unbiased estimators. It would fail miserably if  $x_t = y_{t-1}$  — it is essential that there be strict exogeneity:  $E(y_t | x_1, \dots, x_T, c, b_o) = c + b_o x_t$  and not just  $E(y_t | x_t, c, b_o) = c + b_o x_t$ .

Note that strict exogeneity is usually taken to mean<sup>2</sup>

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

In addition to using regression functions instead of linear predictors, we require finite lags, not just one-sided lags. If we condition on all past

values of  $x$  and then find that the first few future values of  $x$  are irrelevant, this is not evidence against heterogeneity bias. For we show in appendix A that

$$E(c|x) = E(c|x_t, x_{t-1}, \dots) = E(c|x_{t+1}, x_{t+2}, \dots)$$

and the same is true for  $E(b_j|x)$ .

## 3. ESTIMATION

Consider a sample  $\underline{z}_i' = (\underline{x}_i', \underline{y}_i')$ ,  $i=1, \dots, N$ , where  $\underline{x}_i' = (x_{i1}, \dots, x_{iT})$  and  $\underline{y}_i' = (y_{i1}, \dots, y_{iT})$ . We shall assume that the  $\underline{z}_i$ ,  $i=1, \dots, N$ , are independent and identically distributed (i.i.d.) according to some multivariate distribution with  $E(x_{it}^4) < \infty$ ,  $E(y_{it}^4) < \infty$ ,  $t=1, \dots, T$ , and  $V(\underline{x}_i)$  non-singular. The results on estimation will not require stationarity assumptions. We shall focus on the estimation of the linear predictor of  $y_{it}$  conditional on  $x_{i1}, \dots, x_{iT}$ . The key simplification is linearity in the parameters; the variables  $x_{it}$  could be non-linear functions of more basic variables. Also  $x_{it}$  could be a vector of variables, but we simplify notation by discussing the scalar case.

Consider the following minimum mean square error linear predictors:

$$E^*(y_{it} | \underline{x}_i) = \eta_t + \underline{\pi}_t' \underline{x}_i, \quad t=1, \dots, T,$$

which we write as

$$E^*(y_i | \underline{x}_i) = \eta + \underline{\Pi} \underline{x}_i, \quad i=1, \dots, N,$$

where

$$\underline{\Pi}' = V^{-1}(\underline{x}_i) \text{Cov}(\underline{x}_i, \underline{y}_i).$$

Our theorem suggests that we look for finite lags and leads. This may be most promising for  $y_t$  near the end of the sample period, since insufficient lags may result in a full set of leads. (The point is, of course, symmetric;

if  $b_{-j} \neq 0$ , then insufficient leads may result in a full set of lags.)

In general it will not be true that  $E^*(y_i | x_i) = E(y_i | x_i)$ . For small  $t$  there may be omitted lags and  $E(x_{-j} | x_1, \dots, x_T)$  is not linear, in general; throughout the sample period there may be non-linear terms arising from  $E(c | x_1, \dots, x_T) \neq E^*(c | x_1, \dots, x_T)$  and  $E(b_j | x_1, \dots, x_T) \neq E(b_j)$ . So in making inferences about  $\underline{\Pi}$ , we shall be careful to treat  $\underline{\Pi} \underline{x}_i$  as a linear predictor and not assume that it is the multivariate regression function.

We shall also be interested in estimating linear transformations. Given our focus on linear predictors, there are simple relationships that must hold if a linear transformation results in zero coefficients on certain lags or leads. First consider differencing; say that it produces a contemporaneous relationship that is stationary apart from additive period effects:

$$E^*(y_t - y_{t-1} | x_2 - x_1, \dots, x_T - x_{T-1}) = \psi_t + \beta (x_t - x_{t-1}), \quad t=2, \dots, T.$$

Then the original  $\underline{\Pi}$  matrix must have been of the form

$$\underline{\Pi} = \beta \underline{I}_T + \underline{\ell} \underline{\delta}',$$

where  $\underline{\ell}$  is a  $T \times 1$  vector of ones and  $\underline{\delta}$  is an unrestricted  $T \times 1$  vector. So our criterion of looking for transformations that produce finite lags or leads can be stated in terms of restrictions on  $\underline{\Pi}$ .

It is clear that a contemporaneous relationship holds a somewhat special position. It is easier to justify than a claim that two lagged values of  $x$  are necessary but not three lags. For an example where a contemporaneous



$\underline{\delta}$  is an unconstrained  $T \times 1$  vector, and  $\underline{\lambda}$  is a  $T \times 1$  vector with  $\lambda_t = \rho^{t-1}$ ,  $t=1, \dots, T$ . This interpretation of  $\lambda_t$  was pointed out to me by Zvi Griliches and Ariel Pakes. Combining the two cases gives  $\bar{y}_t = (y_t - y_{t-1}) - \rho(y_{t-1} - y_{t-2})$  and

$$\underline{\Pi} = \beta \underline{D} + \underline{\ell} \underline{\delta}'_1 + \underline{\lambda} \underline{\delta}'_2,$$

where  $\underline{\delta}_1$  and  $\underline{\delta}_2$  are unconstrained  $T \times 1$  vectors,  $\underline{\ell}$  is a  $T \times 1$  vector of ones, and  $\underline{\lambda}$  is a  $T \times 1$  vector with  $\lambda_t = \rho^{t-1}$ ,  $t=1, \dots, T$ .

We should stress once again that this is only a wide sense regression function. It may be that  $E^*(\bar{y}_t | x_2 - x_1, \dots, x_T - x_{T-1})$  is in fact a regression function (indeed, we hope that it is); but in integrating back to  $E(y_t | x_1, \dots, x_T)$ , we pick up  $E(c | x_1, \dots, x_T)$  and  $E(\sum_{j=0}^{\infty} \rho^j x_{-j} | x_1, \dots, x_T)$ , which are both non-linear in general. So we cannot rely on an asymptotic theory that begins with the assumption that we are estimating a true regression function.

## 2.1 THE ESTIMATION OF LINEAR APPROXIMATIONS

We are assuming that  $\underline{z}'_i = (x'_i, y'_i)$  is i.i.d.,  $i=1, \dots, N$ . So  $\underline{w}'_i = (\underline{z}'_i, \underline{z}'_i \otimes x'_i)$  is i.i.d. and

$$\underline{\Pi}' = V^{-1}(\underline{x}_i) \text{Cov}(\underline{x}_i, \underline{y}'_i)$$

is a function of  $E(\underline{w}_i)$ . The least squares estimator  $\hat{\underline{\Pi}}'$  is the corresponding function of  $\bar{\underline{w}} = \sum_{i=1}^N \underline{w}_i / N$ . Hence our problem is to derive the limiting distribution of a differentiable function of a sample mean, under i.i.d. sampling.

This is a straightforward application of the  $\delta$ -method and gives

$$\sqrt{N}(\text{vec } \hat{\Pi}' - \text{vec } \Pi') \xrightarrow{D} N(0, \Omega),$$

where

$$\Omega = E[(\underline{y}_i^* - \Pi \underline{x}_i^*)(\underline{y}_i^* - \Pi \underline{x}_i^*)' \otimes V^{-1}(\underline{x}_i)(\underline{x}_i^* \underline{x}_i'^*) V^{-1}(\underline{x}_i)],$$

$\underline{x}_i^* = \underline{x}_i - E(\underline{x}_i)$ ,  $\underline{y}_i^* = \underline{y}_i - E(\underline{y}_i)$ . (See appendix B.) A consistent (as  $N \rightarrow \infty$ ) estimator of  $\Omega$  is readily available from the corresponding sample moments:

$$\hat{\Omega} = \frac{1}{N} \sum_{i=1}^N \{[(\underline{y}_i - \bar{y}) - \hat{\Pi}(\underline{x}_i - \bar{x})][(\underline{y}_i - \bar{y}) - \hat{\Pi}(\underline{x}_i - \bar{x})]'$$

$$\otimes S_x^{-1}(\underline{x}_i - \bar{x})(\underline{x}_i - \bar{x})' S_x^{-1}\},$$

$$\bar{y} = \sum_{i=1}^N y_i / N, \quad \bar{x} = \sum_{i=1}^N \underline{x}_i / N, \quad S_x = \sum_{i=1}^N (\underline{x}_i - \bar{x})(\underline{x}_i - \bar{x})' / N.$$

If  $E(\underline{y}_i | \underline{x}_i) = \eta + \Pi \underline{x}_i$ , then we have

$$\Omega = E[V(\underline{y}_i | \underline{x}_i) \otimes V^{-1}(\underline{x}_i)(\underline{x}_i^* \underline{x}_i'^*) V^{-1}(\underline{x}_i)].$$

If  $V(\underline{y}_i | \underline{x}_i)$  is uncorrelated with  $\underline{x}_i^* \underline{x}_i'^*$ , then

$$\Omega = E[V(\underline{y}_i | \underline{x}_i)] \otimes V^{-1}(\underline{x}_i).$$

If we have homoskedastic variance, so that  $V(\underline{y}_i | \underline{x}_i) = \Sigma$  does not depend on  $\underline{x}_i$ , then  $\Omega = \Sigma \otimes V^{-1}(\underline{x}_i)$ .

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 asymptotic variance of a simple correlation coefficient (p. 359). Kendall and Stuart (1961, p. 293) and

Goldberger (1974) present the formula for the asymptotic sampling variance of a simple regression coefficient. Using a different approach, White (1980) obtains the asymptotic covariance matrix for univariate regression coefficients.

## 2.2 IMPOSING RESTRICTIONS

We can impose restrictions on  $\underline{\Pi}$  by using a minimum distance estimator:

$$\min_{\underline{\theta}} [\text{vec } \hat{\underline{\Pi}}' - \text{vec } \underline{\Pi}'(\underline{\theta})]' \hat{\underline{\Omega}}^{-1} [\text{vec } \hat{\underline{\Pi}}' - \text{vec } \underline{\Pi}'(\underline{\theta})],$$

where the constraints on  $\underline{\Pi}$  are specified by the condition that  $\underline{\Pi}$  depends only on a  $p \times 1$  parameter vector  $\underline{\theta}$ . The domain of  $\underline{\theta}$  is some compact subset of  $R^p$ . The relevant asymptotic theory is developed, for example, in Malinvaud (1970, Chap. 9). Since we already have the asymptotic distribution of  $\hat{\underline{\Pi}}$ , we only need to assume that  $\underline{\Omega}$  is non-singular and that the  $\underline{\Pi}(\underline{\theta})$  functions are continuous and satisfy Malinvaud's assumptions 5 and 6 (p. 349), which impose invertibility and differentiability restrictions in a neighborhood of the true parameter value  $\underline{\theta}^0$ . Then  $\sqrt{N}(\hat{\underline{\theta}} - \underline{\theta}^0) \xrightarrow{D} N(0, \underline{\Lambda})$ ; a consistent (as  $N \rightarrow \infty$ ) estimator of  $\underline{\Lambda}$  is provided by

$$\hat{\underline{\Lambda}} = (\underline{G}' \hat{\underline{\Omega}}^{-1} \underline{G})^{-1},$$

where  $\underline{G} = \partial \text{vec } \underline{\Pi}'(\underline{\theta}) / \partial \underline{\theta}'$  is evaluated at  $\underline{\theta} = \hat{\underline{\theta}}$ . We can test the restrictions by using the limiting distribution of

$$N[\text{vec } \hat{\underline{\Pi}}' - \text{vec } \underline{\Pi}'(\hat{\underline{\theta}})]' \hat{\underline{\Omega}}^{-1} [\text{vec } \hat{\underline{\Pi}}' - \text{vec } \underline{\Pi}'(\hat{\underline{\theta}})],$$

which is  $\chi^2$  with  $MK-p$  degrees of freedom if  $\underline{\Pi}$  is  $M \times K$ .



This framework leads to some surprising results on efficient estimation.

For a simple example, we shall use a univariate linear regression model:

$$E^*(y_i | x_{i1}, x_{i2}) = \eta + \pi_1 x_{i1} + \pi_2 x_{i2} \quad (i=1, \dots, N).$$

Consider imposing the restriction  $\pi_2 = 0$ . If we assume that the distribution of  $y_i$  conditional on  $x_{i1}, x_{i2}$  is  $N(\eta + \pi_1 x_{i1} + \pi_2 x_{i2}, \sigma^2)$ , then under random sampling the efficient, maximum likelihood estimator of  $\pi_1$  is  $b_{yx_1}$ , the slope coefficient in the least squares regression of  $y$  on  $x_1$ . Let  $\hat{\pi}_1, \hat{\pi}_2$  be the coefficients in the multiple regression of  $y$  on  $x_1, x_2$ . If  $\text{Cov}(x_{i1}, x_{i2}) = 0$  and  $\pi_2 = 0$ , then  $\sqrt{N}(b_{yx_1} - \pi_1)$  and  $\sqrt{N}(\hat{\pi}_1 - \pi_1)$  have the same limiting distribution -- but  $\hat{\pi}_1$  is generally less efficient than the minimum distance estimator if the regression function is non-linear or if there is heteroskedasticity.

To see this, consider the asymptotic covariance between  $\hat{\pi}_1$  and  $\hat{\pi}_2$ :

$$\omega_{12} = \frac{E[(y_i^* - \pi_1 x_{i1}^* - \pi_2 x_{i2}^*)^2 x_{i1}^* x_{i2}^*]}{V(x_{i1}) V(x_{i2})}.$$

This is non-zero, in general, even though  $\text{Cov}(x_{i1}, x_{i2}) = 0$ . So if  $\pi_2 = 0$ , we gain efficiency in estimating  $\pi_1$  by considering a linear combination  $\hat{\pi}_1 + \tau \hat{\pi}_2$ , where  $\tau$  is chosen to minimize the variance. (We have consistency for any  $\tau$  since  $\hat{\pi}_2$  almost surely converges to  $\pi_2 = 0$  as  $N \rightarrow \infty$ .) Hence we want  $\tau = -\omega_{12}/\omega_{22}$ . It is a general property of minimum distance estimation that substituting a consistent estimator ( $\hat{\Omega}$ ) for  $\Omega$  does not affect the limiting distribution. So we use

$$\hat{\pi}_1^* = \hat{\pi}_1 - \frac{\hat{\omega}_{12}}{\hat{\omega}_{22}} \hat{\pi}_2.$$

This has a smaller asymptotic variance than  $\hat{\pi}_1$  (or  $b_{yx_1}$ ) unless

$(y_i - \bar{y} - \pi_1 x_{i1} - \pi_2 x_{i2})^2$  is uncorrelated with  $x_{i1}^* x_{i2}^*$ .

If  $\text{Cov}(x_{i1}, x_{i2}) \neq 0$ , then our minimum distance estimator still has the form  $\hat{\pi}_1^* = \hat{\pi}_1 - (\hat{\omega}_{12}/\hat{\omega}_{22}) \hat{\pi}_2$ , where  $\hat{\omega}_{12}$  and  $\hat{\omega}_{22}$  are obtained from our general formula for  $\hat{\Omega}$ . In the special case where  $E(y|x_1, x_2)$  is linear and  $V(y|x_1, x_2) = \sigma^2$ , then  $\omega_{12}/\omega_{22} = -\text{Cov}(x_{i1}, x_{i2})/V(x_{i1})$  and  $\hat{\pi}_1^* = b_{yx_1}$ ; but in general they differ and  $b_{yx_1}$  is not efficient since it is correlated with  $\hat{\pi}_2$ .<sup>3</sup>

### 2.3 AN EXTENSION

Consider a panel of  $N$  farms observed over  $T$  seasons. The technology is Cobb-Douglas:

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

where  $y = \ln$  (output),  $x = \ln$  (input),  $c$  represents an unmeasured input (soil quality) that is constant over the period of the sample, and  $u$  represents a stochastic input such as rainfall. All farms use the same technology so there is no variation in  $\beta$ .

Assume that farmers choose  $x$  to maximize expected profits, conditional on an information set that includes  $c$ . Then  $E(c_i | x_{i1}, \dots, x_{iT}) \neq E(c_i)$ ; under plausible assumptions, farmers on the higher quality soil will use more of the variable factor.<sup>4</sup> The transformation to  $\tilde{y}_{it} = y_{it} - y_{i,t-1}$ ,  $\tilde{x}_{it} = x_{it} - x_{i,t-1}$  is appropriate if  $\text{Cov}(x_{it}, u_{is}) = 0$ ,  $s, t = 1, \dots, T$ . This will not hold, however, if the farmers forecast  $u_t$  on the basis of past values.

Say that  $u_{it} = \rho u_{i,t-1} + v_{it}$ , where  $v_{it}$  is serially independent. Then

$$y_{it} - y_{i,t-1} = \beta (x_{it} - x_{i,t-1}) - \beta \rho (x_{i,t-1} - x_{i,t-2}) \\ + \rho (y_{i,t-1} - y_{i,t-2}) + v_{it} - v_{i,t-1}.$$

Since  $\text{Cov}(x_t, v_s) = 0$  for  $t \leq s$  and  $\text{Cov}(y_t, v_s) = 0$  for  $t < s$ , we can obtain consistent estimates by using instrumental variables based on  $x_{t-1}, x_{t-2}, \dots, x_1$  and  $y_{t-2}, y_{t-3}, \dots, y_1$ .

The asymptotic distribution for the instrumental variable estimator is derived in appendix B. We set

$$\hat{\pi}_m = \left[ \sum_{i=1}^N (q_{im} - \bar{q}_m)(x_{im} - \bar{x}_m)' \right]^{-1} \sum_{i=1}^N (q_{im} - \bar{q}_m)(y_{im} - \bar{y}_m) \\ \pi_m = [\text{Cov}(q_{im}, x_{im}')]^{-1} \text{Cov}(q_{im}, y_{im}) \quad (m=1, \dots, M),$$

where  $x_{im}$  is  $K_m \times 1$ ,  $q_{im}$  is  $K_m \times 1$ , and  $\text{Cov}(q_{im}, x_{im}')$  is assumed to be non-singular. We assume i.i.d. sampling from a distribution with finite fourth moments. Then  $\sqrt{N}(\hat{\pi}_m - \pi_m)$  and  $\sqrt{N}(\hat{\pi}_n - \pi_n)$  have a limiting normal distribution with mean 0 and covariance matrix

$$\Omega_{mn} = E[(y_{im}^* - \pi_m' x_{im}^*)(y_{in}^* - \pi_n' x_{in}^*) \phi_m^{-1} (q_{im}^* \ q_{in}^{*'}) \phi_n^{-1}],$$

where  $\phi_m = \text{Cov}(q_{im}, x_{im}')$ ,  $y_{im}^* = y_{im} - E(y_{im})$ ,  $x_{im}^* = x_{im} - E(x_{im})$ ,  $q_{im}^* = q_{im} - E(q_{im})$ ,  $m, n=1, \dots, M$ .

In our example,  $y_{im}$  is  $y_{it} - y_{i,t-1}$ ;  $x_{im}$  contains  $(x_{it} - x_{i,t-1}), (x_{i,t-1} - x_{i,t-2}), (y_{i,t-1} - y_{i,t-2})$ ; and  $q_{im}$  contains three variables based on  $x_{i,t-1}, x_{i,t-2}, \dots, x_{i1}, y_{i,t-2}, y_{i,t-3}, \dots, y_{i1}$ .

We can impose restrictions on the coefficients, both within and across equations, by using the minimum distance procedure. There is, however, another efficiency issue, since we have more than  $K_m$  instrumental variables for the  $m^{\text{th}}$  equation. Given our previous discussion of imposing restrictions, it is perhaps not surprising that the standard two stage least squares estimator is not, in general, an efficient procedure for combining instrumental variables. We shall demonstrate this by means of a simple example.

Say that  $(y_i, x_i, q_{i1}, q_{i2})$  is i.i.d. according to some distribution with finite fourth moments. Assume that  $E(y_i) = E(x_i) = E(q_{i1}) = E(q_{i2}) = 0$ ,  $E(x_i q_{i1}) \neq 0$ ,  $E(x_i q_{i2}) \neq 0$ , and

$$\frac{E(y_i q_{i1})}{E(x_i q_{i1})} = \frac{E(y_i q_{i2})}{E(x_i q_{i2})} = \pi.$$

Then we have two instrumental variable estimators that both almost surely converge to  $\pi$  as  $N \rightarrow \infty$ :

$$\hat{\pi}_1 = \frac{\sum_{i=1}^N y_i q_{i1}}{\sum_{i=1}^N x_i q_{i1}}, \quad \hat{\pi}_2 = \frac{\sum_{i=1}^N y_i q_{i2}}{\sum_{i=1}^N x_i q_{i2}}.$$

The two stage least squares estimator combines  $\hat{\pi}_1$  and  $\hat{\pi}_2$  by forming  $\hat{q}_3 = \hat{\psi}_1 q_1 + \hat{\psi}_2 q_2$  based on the least squares regression of  $x$  on  $q_1, q_2$  (We shall assume that  $V(q_{i1}, q_{i2})$  is non-singular):

$$\hat{\pi}_{\text{TSL}} = \frac{\sum_{i=1}^N y_i \hat{q}_{i3}}{\sum_{i=1}^N x_i \hat{q}_{i3}} = \hat{\delta} \hat{\pi}_1 + (1-\hat{\delta}) \hat{\pi}_2,$$

where  $\hat{\delta} = \hat{\psi}_1 \frac{\sum_{i=1}^N x_i q_{i1}}{\sum_{i=1}^N x_i q_{i1} + \hat{\psi}_2 \sum_{i=1}^N x_i q_{i2}}$ .

So we have

$$\sqrt{N} (\hat{\pi}_{\text{TSLs}} - \pi) = \hat{\delta} \sqrt{N} (\hat{\pi}_1 - \pi) + (1 - \hat{\delta}) \sqrt{N} (\hat{\pi}_2 - \pi).$$

We can obtain the limiting distribution of  $\sqrt{N}(\hat{\pi}_1 - \pi, \hat{\pi}_2 - \pi)$  from our general result on instrumental variable estimation: set  $M = 2$ ,

$x_{i1} = x_i, x_{i2} = x_i, q_{i1} = q_{i1}, q_{i2} = q_{i2}, y_{i1} = y_i, y_{i2} = y_i$ . (We never assumed that  $V(y_{i1}, y_{i2})$  is non-singular.) So we have

$$\sqrt{N} \left[ \begin{pmatrix} \hat{\pi}_1 \\ \hat{\pi}_2 \end{pmatrix} - \begin{pmatrix} \pi \\ \pi \end{pmatrix} \right] \xrightarrow{D} N(\underline{0}, \underline{\Omega}),$$

where  $\underline{\Omega} = (\omega_{jk})$ :

$$\omega_{jk} = \frac{E[(y_i - \pi x_i)^2 q_{ij} q_{ik}]}{E(q_{ij} x_i) E(q_{ik} x_i)}, \quad j, k = 1, 2.$$

Since

$$\hat{\delta} \xrightarrow{\text{a.s.}} \frac{\psi_1 E(x_i q_{i1})}{\psi_1 E(x_i q_{i1}) + \psi_2 E(x_i q_{i2})} = \delta \text{ as } N \rightarrow \infty,$$

it follows that

$$\sqrt{N} (\hat{\pi}_{\text{TSLs}} - \pi) \xrightarrow{D} N(0, \underline{\delta}' \underline{\Omega} \underline{\delta}) \text{ as } N \rightarrow \infty$$

where  $\underline{\delta}' = (\delta, 1-\delta)$ . (Note that the limiting distribution coincides with what we would obtain using  $q_3 = E^*(x|q_1, q_2) = \psi_1 q_1 + \psi_2 q_2$  instead of  $\hat{q}_3$ .)

We see that the limiting distribution of  $\hat{\pi}_{\text{TSLs}}$  coincides with that of  $\delta \hat{\pi}_1 + (1-\delta) \hat{\pi}_2$ . This suggests finding the  $\tau$  that minimizes the asymptotic variance of  $\tau \hat{\pi}_1 + (1-\tau) \hat{\pi}_2$ ; the answer leads to the minimum distance estimator:

$$\min_{\theta} \left\{ \left[ \begin{pmatrix} \hat{\pi}_1 \\ \hat{\pi}_2 \end{pmatrix} - \begin{pmatrix} \theta \\ \theta \end{pmatrix} \right]' \hat{\Omega}^{-1} \left[ \begin{pmatrix} \hat{\pi}_1 \\ \hat{\pi}_2 \end{pmatrix} - \begin{pmatrix} \theta \\ \theta \end{pmatrix} \right] \right\}$$

gives  $\hat{\pi}^* = \tau \hat{\pi}_1 + (1-\tau) \hat{\pi}_2$ ,

where  $\tau = (\omega^{11} + \omega^{12}) / (\omega^{11} + 2\omega^{12} + \omega^{22})$

and  $\omega^{jk}$  is the  $j, k$  element of  $\hat{\Omega}^{-1}$ . The estimator obtained by using a consistent  $\hat{\Omega}$  instead of  $\hat{\Omega}$  has the same limiting distribution.

In general  $\tau \neq \delta$  since  $\tau$  is a function of fourth moments and  $\delta$  is not.

For example, say that  $x_i = q_{i2}$ . Then  $\delta = 0$  but  $\tau \neq 0$  unless

$$E \left[ (y_i - \pi x_i)^2 \left( \frac{x_i^2}{E(x_i^2)} - \frac{x_i q_{i1}}{E(x_i q_{i1})} \right) \right] = 0.$$

Hence the asymptotic variance of  $\hat{\pi}^*$  is less than that of  $\hat{\pi}_{TSLs}$  unless  $(y_i - \pi x_i)^2$  is uncorrelated with  $[x_i^2/E(x_i^2)] - [(x_i q_{i1})/E(x_i q_{i1})]$ .

If we add another equation then we can consider the conventional three stage least squares estimator. Its limiting distribution can be derived in a straightforward fashion since the estimator is a differentiable function of sample moments; however, viewed as a minimum distance estimator, it is using the wrong norm in general.

Instead we should simply recognize that the restrictions can be expressed as restrictions on population moments. We let  $z_i' = (x_i', y_i')$  and let  $w_i$  be the vector formed from  $z_i$  and the lower triangle of  $z_i z_i'$ . Then the population parameters of interest are functions of  $E(w_i)$ . The restrictions can be expressed by the condition that  $E(w_i)$  depends only on a  $p$ -dimensional parameter vector  $\theta$ :  $E(w_i) = g(\theta)$ , where the domain of  $\theta$  is some compact subset of  $R^p$  that

contains the true parameter vector  $\underline{\theta}^0$ . Under i.i.d. sampling we have  $\sqrt{N} [\bar{w} - E(w_i)] \xrightarrow{D} N(0, V(w_i))$  if  $V(w_i)$  is finite; a consistent estimator of  $V(w_i)$  is provided by

$$\hat{V}(w_i) = \frac{1}{N} \sum_{i=1}^N (w_i - \bar{w})(w_i - \bar{w})'$$

If  $V(w_i)$  is non-singular and if  $g$  is continuous and satisfies Malinvaud's assumptions 5 and 6, then

$$\min_{\underline{\theta}} [\bar{w} - g(\underline{\theta})]' \hat{V}^{-1}(w_i) [\bar{w} - g(\underline{\theta})]$$

gives

$$\sqrt{N}(\hat{\underline{\theta}} - \underline{\theta}^0) \xrightarrow{D} N(0, \underline{\Lambda}) \text{ as } N \rightarrow \infty,$$

where a consistent estimator of  $\underline{\Lambda}$  is provided by

$$\hat{\underline{\Lambda}} = [G' \hat{V}^{-1}(w_i) G]^{-1}$$

with  $G = \partial g / \partial \underline{\theta}'$  evaluated at  $\underline{\theta} = \hat{\underline{\theta}}$ .

#### 4. CONCLUSION

Under stationarity, the heterogeneous stochastic processes are the non-ergodic ones. We showed that if a distributed lag is of finite order, then its coefficients are the unconditional means of the underlying random coefficients. This result can be usefully applied to linear transformations of the process.

The estimation framework is a multivariate wide-sense regression function. The identification analysis suggests that we examine certain restrictions on the coefficients. The actual regression function is non-linear, and so we provide a theory of inference for linear approximations. It rests on obtaining the asymptotic distribution of functions of sample moments. Restrictions are imposed by using a minimum distance estimator; it is generally more efficient than the conventional estimators.



## APPENDIX A

Let  $\Omega$  be a set of points where  $\omega \in \Omega$  is a doubly infinite sequence of vectors of real numbers:  $\omega = \{\dots, \omega_{-1}, \omega_0, \omega_1, \dots\} = \{\omega_t, t \in \mathbb{T}\}$ , where  $\omega_t = (\omega_{1t}, \omega_{2t}) \in \mathbb{R}^2$  and  $\mathbb{T}$  is the set of all integers. Let  $(x_t(\omega), y_t(\omega)) = (\omega_{1t}, \omega_{2t})$  be the  $t^{\text{th}}$  coordinate function. Let  $F$  be the  $\sigma$ -field consisting of sets of the form

$$A = \{\omega: (x_t(\omega), y_t(\omega)) \in B_0, \dots, (x_{t+k}(\omega), y_{t+k}(\omega)) \in B_k\},$$

where  $t, k \in \mathbb{T}$  and the  $B$ 's are two-dimensional Borel sets. Let  $P$  be a probability measure defined on  $F$  such that  $\{(x_t, y_t), t \in \mathbb{T}\}$  is a stationary stochastic process on the probability space  $(\Omega, F, P)$ .

The shift transformation  $S$  is defined by

$$(x_t(S\omega), y_t(S\omega)) = (x_{t+1}(\omega), y_{t+1}(\omega)).$$

It is an invertible, measure preserving transformation. A random variable  $z$  is invariant if  $z(S\omega) = z(\omega)$  a.e. (We shall use "a.e." as an abbreviation for "except on a set of probability measure zero.") A set  $A \in F$  is invariant if its indicator function is an invariant random variable. A  $\sigma$ -field  $G \subset F$  is invariant if  $SG = G$ ; that is, for every  $G \in G$  we have  $SG \in G$  and  $SH = G$  for some  $H \in G$ .

We shall use  $E(z|G)_\omega$  to denote the conditional expectation of the random variable  $z$  with respect to the  $\sigma$ -field  $G$ , evaluated at  $\omega$ . We shall use  $\sigma(x)$  to denote the  $\sigma$ -field generated by  $\dots, x_{-1}, x_0, x_1, \dots$ , and  $E(z|x)$  will denote  $E(z|\sigma(x))$ .

We shall need the following two lemmas.

LEMMA 1. If  $G$  is an invariant  $\sigma$ -field, then

$$E(z|G)_{S\omega} = E(v|G)_{\omega} \quad \text{a.e.}$$

where  $v(\omega) = z(S\omega)$ .

PROOF. A change of variable argument shows that

$$E(z|G)_{S\omega} = E(v|S^{-1}G)_{\omega}.$$

(See Billingsley (1965), example 10.3, p. 109.) The result follows since  $S^{-1}G = G$ .

Q.E.D.

LEMMA 2. If  $G$  is an invariant  $\sigma$ -field and if  $z$  is an invariant random variable, then  $E(z|G)$  is an invariant random variable.

PROOF. By Lemma 1,  $E(z|G)_{S\omega} = E(v|G)_{\omega}$  a.e., where  $v(\omega) = z(S\omega)$ ; since  $z$  is invariant,  $z(S\omega) = z(\omega)$  a.e.; hence  $E(z|G)_{S\omega} = E(z|G)_{\omega}$  a.e.

Q.E.D.

Our theorem will require that the  $\{x_t\}$  process satisfy the following regularity condition:

CONDITION (C). 1. Let  $\{(z_{1t}, z_{2t}), t \in T\}$  be a stationary stochastic process where  $z_{1t}$  is measurable  $\sigma(x_t, \dots, x_{t+k})$  and  $z_{2t}$  is measurable  $\sigma(\dots, x_{t-1}, x_{t+k+1}, \dots)$ ,  $t, k \in T$ . Then

$$P(z_{1t} = z_{2t} \text{ for all } t) = 0$$

unless  $z_{1t}$  is a constant a.e.

2. Let  $A_1 \in \sigma(x_t, \dots, x_{t+k})$ ,  $A_2 \in \sigma(\dots, x_{t-1}, x_{t+k+1}, \dots)$ ,  
 $t, k \in \mathbb{T}$ . Then

$$P(A_1 \cap A_2) = P(A_1)$$

implies  $P(A_1) = 0$  or  $P(A_2) = 1$ .

In addition, we assume that  $\text{Var}(x_t) \neq 0$ .

THEOREM. If  $\{x_t\}$  satisfies the regularity condition (C), if  $E|y_t| < \infty$ ,  
and if

$$E(y_t | \sigma(G \cup \sigma(x))) = c + \sum_{j=-\infty}^{\infty} b_j x_{t-j} \quad \text{a.e., } t \in \mathbb{T},$$

where  $c$  and the  $b_j$  are integrable, invariant random variables and

$$\sum_{j=-\infty}^{\infty} |b_j x_{t-j}| < \infty \quad \text{a.e. and integrable, then}$$

$$E(y_t | x) = E(y_t | x_{t-J}, \dots, x_{t+M}) \quad \text{a.e. } (t \in \mathbb{T}, 0 \leq J, M < \infty)$$

implies that

$$E(y_t | x) = \gamma + \sum_{j=-M}^J \beta_j x_{t-j} \quad \text{a.e., } t \in \mathbb{T},$$

where

$$\gamma = E(c) = E(c | x) \quad \text{a.e.}$$

$$\beta_j = E(b_j) = E(b_j | x) \quad \text{a.e., } j \in [-M, J]$$

$$E(b_j | x) = 0 \quad \text{a.e., } j \notin [-M, J].$$

(We can take  $G$  to be  $J$ , the  $\sigma$ -field of shift invariant sets, but this is not necessary.)

$$\begin{aligned}
\text{PROOF. } E(y_t | x_{t-J}, \dots, x_{t+M}) &= E(y_t | x) = E[E(y_t | \sigma(G \cup \sigma(x))) | x] \\
&= E(c | x) + E\left(\sum_{j=-\infty}^{\infty} b_j x_{t-j} \mid x\right) \\
&= E(c | x) + \sum_{j=-\infty}^{\infty} E(b_j | x) x_{t-j} \quad \text{a.e., } t \in \mathbb{T},
\end{aligned}$$

since

$$\left| \sum_{j=-L}^K b_j x_{t-j} \right| \leq \sum_{j=-L}^K |b_j x_{t-j}| \leq \sum_{j=-\infty}^{\infty} |b_j x_{t-j}|$$

and  $\sum_{j=-\infty}^{\infty} |b_j x_{t-j}|$  is integrable. (See Billingsley (1979), theorem 34.2(v), p. 396.)

Let  $\mathcal{T}_1 = \bigcap_{t=-\infty}^0 \sigma(x_t, x_{t-1}, \dots)$  be the left tail  $\sigma$ -field and  $\mathcal{T}_2 = \bigcap_{t=0}^{\infty} \sigma(x_t, x_{t+1}, \dots)$  be the right tail  $\sigma$ -field. Since  $\sigma(x)$  is an invariant  $\sigma$ -field, lemma 2 implies that  $E(c|x)$  and  $E(b_j|x)$  are invariant random variables. Hence they are measurable  $\mathcal{T}_1$  and measurable  $\mathcal{T}_2$ . (See Rozanov, lemma 6.1, p. 162; more precisely, there exist versions of  $E(c|x)$  and  $E(b_j|x)$  that are measurable  $\mathcal{T}_1$  and there exist versions of  $E(c|x)$  and  $E(b_j|x)$  that are measurable  $\mathcal{T}_2$ .)

Define

$$z_{t-k,k} = [E(y_t | x_{t-J}, \dots, x_{t+M}) - E(c|x) - \sum_{j \neq k} E(b_j | x) x_{t-j}] / E(b_k | x)$$

for  $t \in \mathbb{T}$  and  $\omega \in A_k$ , where

$$A_k = \{\omega : E(b_k | x)_\omega \neq 0\};$$

set  $z_{t-k,k} = 0$  for  $\omega \notin A_k$ . Then  $\{z_{t-k,k}, t \in \mathbb{T}\}$  is a stationary stochastic process since  $E(c|x)$ ,  $E(b_j|x)$ ,  $j \in \mathbb{T}$ , are invariant random variables and  $A_k$  is an invariant set.

If  $k \notin [-M, J]$ , then  $E(y_t | x_{t-J}, \dots, x_{t+M})$  is measurable  $\sigma(\dots, x_{t-k-1}, x_{t-k+1}, \dots)$ . Since  $E(c|x)$  and  $E(b_j|x)$  are measurable  $T_1$ , they are measurable  $(\dots, x_{t-k-2}, x_{t-k-1})$ ; since  $A_k$  is an invariant set, its indicator function is measurable  $T_1$  and so measurable  $(\dots, x_{t-k-2}, x_{t-k-1})$ . Hence  $z_{t-k,k}$  is measurable  $\sigma(\dots, x_{t-k-1}, x_{t-k+1}, \dots)$ ,  $t \in T$ .

Since  $x_{t-k} = z_{t-k,k}$  for all  $t \in T$  for  $\omega \in A_k$ , condition (C.1) implies that  $\text{Var}(x_t) = 0$  if  $P(A_k) > 0$ . We conclude that  $E(b_k|x) = 0$  a.e. for  $k \notin [-M, J]$ .

So we have

$$E(c|x) + \sum_{j=-M}^J E(b_j|x)x_{t-j} = E(y_t | x_{t-J}, \dots, x_{t+M}) \quad \text{a.e.}$$

for  $t \in T$ . If we consider this equation for  $t = J, J+1, \dots, 2J+M+1$ , we can say

$$\underline{Q} \underline{d} = \underline{g} \quad \text{a.e.}$$

where  $\underline{d}' = [E(c|x), E(b_{-M}|x), \dots, E(b_J|x)]$ ,  $\underline{g}' = [g_{2J+M+1}, \dots, g_J]$  where  $g_s = E(y_s | x_{s-J}, \dots, x_{s+M})$ , and  $\underline{Q}$  is a  $(J+M+2) \times (J+M+2)$  matrix:

$$\underline{Q} = \begin{bmatrix} 1 & x_{2J+2M+1} & x_{2J+2M} & \dots & x_{J+M+1} \\ 1 & x_{2J+2M} & x_{2J+2M-1} & \dots & x_{J+M} \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 1 & x_{J+M} & x_{J+M-1} & \dots & x_0 \end{bmatrix}.$$

We show below in lemma 3 that  $P(D) > 0$  where  $D = \{\omega : \underline{Q} \text{ is non-singular}\}$ .

Let  $\underline{h} = \underline{Q}^{-1} \underline{g}$  for  $\omega \in D$ ;  $\underline{h} = \underline{0}$  otherwise. Then

$\underline{h}' = (h_1, \dots, h_{J+M+2}) = \underline{d}'$  for  $\omega \in D$ . We shall show that  $E(c|x) = E(c)$  a.e.

If  $h_1$  is not constant a.e. for  $\omega \in D$ , then we can choose a constant  $r$  such that

$$A_1 = \{\omega: \omega \in D \text{ and } h_1(\omega) \leq r\}, \quad P(A_1) > 0$$

$$A_1' = D - A_1, \quad P(A_1') > 0.$$

Then set

$$A_2 = \{\omega: E(c|x)_\omega \leq r\}, \quad A_2' = \Omega - A_2.$$

Since  $Q$ ,  $g$ , and the indicator function for  $D$  are measurable

$\sigma(x_0, \dots, x_{2J+2M+1})$ , it follows that  $A_1$  and  $A_1'$  are in  $\sigma(x_0, \dots, x_{2J+2M+1})$ .

$E(c|x)$  is measurable  $\mathcal{T}_1$  and so  $A_2$  and  $A_2'$  are in  $\sigma(\dots, x_{-2}, x_{-1})$ .

$E(c|x) = h_1$  for  $\omega \in D$  implies that

$$P(A_1 \cap A_2) = P(A_1), \quad P(A_1' \cap A_2') = P(A_1').$$

Since  $P(A_1) > 0$ ,  $P(A_1') > 0$ , condition (C.2) implies that  $P(A_2) = 1$  and  $P(A_2') = 1$ , a contradiction. Hence  $h_1$  is equal to some constant  $r$  a.e. for  $\omega \in D$ .

Now set  $A_1 = D$ ,  $A_2 = \{\omega: E(c|x)_\omega = r\}$ . Then  $P(A_1 \cap A_2) = P(A_1)$ , and so (C.2) implies  $P(A_2) = 1$ . Hence  $E(c|x) = E(c)$  a.e.

An identical argument establishes that  $E(b_j|x) = E(b_j)$  a.e.,  $j \in [-M, J]$ . So we have shown that

$$E(y_t|x) = E(c) + \sum_{j=-M}^J E(b_j)x_{t-j} \quad \text{a.e.}$$

Q.E.D.

LEMMA 3. Let

$$Q_s = \begin{bmatrix} 1 & x_{2s+1} & x_{2s} & \cdots & x_{s+1} \\ 1 & x_{2s} & x_{2s-1} & \cdots & x_s \\ \vdots & \vdots & \vdots & & \vdots \\ 1 & x_s & x_{s-1} & \cdots & x_0 \end{bmatrix}$$

where  $s$  is a non-negative integer. Then condition (C) implies that

$$P(Q_s \text{ non-singular}) > 0 \quad \text{for all } s.$$

PROOF. We shall use induction on  $s$ . If  $s = 0$ , then  $Q_s$  singular implies  $x_1 = x_0$ . If  $P(x_1 = x_0) = 1$ , then (C.1) or (C.2) implies  $\text{Var}(x_t) = 0$ . So  $P(Q_0 \text{ non-singular}) > 0$ . Now assume that  $P(Q_{s-1} \text{ non-singular}) > 0$  for some  $s \geq 1$  and show that  $P(Q_s \text{ non-singular}) > 0$ .

If  $Q_s$  is singular a.e., then expanding  $\det(Q_s)$  in terms of the cofactors of the first row gives

$$x_{2s+1} \det(Q_{s-1}) = g(x_0, \dots, x_{2s}) \quad \text{a.e.}$$

where  $g$  is a measurable function. By the induction hypothesis,

$P(D_{s-1}) > 0$  where  $D_{s-1} = \{\omega : \det(Q_{s-1}) \neq 0\}$ . Let

$$h = g(x_0, \dots, x_{2s}) / \det(Q_{s-1})$$

for  $\omega \in D_{s-1}$ ;  $h = 0$  for  $\omega \notin D_{s-1}$ . Then  $x_{2s+1} = h$  for  $\omega \in D_{s-1}$ .

If  $h$  is not constant a.e. for  $\omega \in D_{s-1}$ , we can choose a constant  $r$  such that

$$A_1 = \{\omega: \omega \in D_{s-1} \text{ and } h(\omega) \leq r\}, \quad P(A_1) > 0$$

$$A_1' = D_{s-1} - A_1, \quad P(A_1') > 0.$$

Then set

$$A_2 = \{\omega: x_{2s+1}(\omega) \leq r\}, \quad A_2' = \Omega - A_2.$$

Since  $A_1$  and  $A_1'$  are in  $\sigma(x_0, \dots, x_{2s})$  and

$$P(A_1 \cap A_2) = P(A_1), \quad P(A_1' \cap A_2') = P(A_1'),$$

condition (C.2) implies that  $P(A_2) = 1$  and  $P(A_2') = 1$ , a contradiction.

So if  $Q_s$  is singular a.e., there must be a constant  $r$  such that  $h(\omega) = r$  a.e. for  $\omega \in D_{s-1}$ .

Now set  $A_1 = D_{s-1}$ ,  $A_2 = \{\omega: x_{2s+1}(\omega) = r\}$ . Then  $P(A_1 \cap A_2) = P(A_1)$ , and so (C.2) implies that  $P(A_2) = 1$ , which contradicts  $\text{Var}(x_t) \neq 0$ .

We conclude that  $P(Q_s \text{ non-singular}) > 0$ .

Q.E.D.



## APPENDIX B

Let  $\underline{z}'_i = (\underline{x}'_i, \underline{y}'_i)$ ,  $i = 1, \dots, N$ , where  $\underline{x}_i$  is  $K \times 1$ ,  $\underline{y}_i$  is  $M \times 1$ . We shall assume that the  $\underline{z}_i$  are independent and identically distributed (i.i.d.). We shall also assume that  $E(x_{ik}^4) < \infty$  ( $k = 1, \dots, K$ ),  $E(y_{im}^4) < \infty$  ( $m = 1, \dots, M$ ), and  $V(\underline{x}_i)$  non-singular. Let  $\underline{w}'_i = (\underline{w}'_{i1}, \underline{w}'_{i2}) = (\underline{z}'_i, \underline{z}'_i \otimes \underline{x}'_i)$ . Then the  $\underline{w}_i$  are i.i.d. with  $V(\underline{w}_i) < \infty$ . Hence the central limit theorem for i.i.d. random variables implies that

$$\sqrt{N} (\bar{\underline{w}} - E(\underline{w}_i)) \xrightarrow{D} N(0, V(\underline{w}_i)) \text{ as } N \rightarrow \infty,$$

$$\text{where } \bar{\underline{w}} = \sum_{i=1}^N \underline{w}_i / N.$$

Consider the estimator

$$\hat{\underline{\Pi}}' = \left[ \sum_{i=1}^N (\underline{x}_i - \bar{\underline{x}})(\underline{x}_i - \bar{\underline{x}})' \right]^{-1} \sum_{i=1}^N (\underline{x}_i - \bar{\underline{x}})(\underline{y}_i - \bar{\underline{y}})$$

where  $\bar{\underline{x}} = \sum_{i=1}^N \underline{x}_i / N$ ,  $\bar{\underline{y}} = \sum_{i=1}^N \underline{y}_i / N$ .  $\hat{\underline{\Pi}}$  is a function of  $\bar{\underline{w}}$ :

$$\text{vec } \hat{\underline{\Pi}}' = \underline{g}(\bar{\underline{w}}),$$

where  $\text{vec } \hat{\underline{\Pi}}'$  is the  $MK \times 1$  vector formed by stacking the successive columns of  $\hat{\underline{\Pi}}'$ . Evaluating  $\underline{g}$  at  $E(\underline{w}_i)$  gives

$$\text{vec } \underline{\Pi}' = \underline{g}(E(\underline{w}_i)),$$

where

$$\underline{\Pi}' = V^{-1}(\underline{x}_i) \text{Cov}(\underline{x}_i, \underline{y}_i').$$

Since  $g$  is differentiable at  $E(\underline{w}_i)$ , the  $\delta$ -method gives

$$\sqrt{N} (\text{vec } \hat{\underline{\Pi}}' - \text{vec } \underline{\Pi}') \xrightarrow{D} N(0, \underline{\Omega}) \text{ as } N \rightarrow \infty,$$

where

$$\underline{\Omega} = \frac{\partial g}{\partial \underline{r}'} V(\underline{w}_i) \frac{\partial g'}{\partial \underline{r}'}$$

and the partial derivatives are evaluated at  $\underline{r} = E(\underline{w}_i)$ . (See, for example, Billingsley [1979, example 29.1, p. 340] or Rao [1973, p. 388].) We shall evaluate the partial derivatives and simplify to show that

$$\underline{\Omega} = E[(\underline{y}_i^* - \underline{\Pi} \underline{x}_i^*)(\underline{y}_i^* - \underline{\Pi} \underline{x}_i^*)' \otimes V^{-1}(\underline{x}_i)(\underline{x}_i^* \underline{x}_i^{*\prime}) V^{-1}(\underline{x}_i)],$$

where  $\underline{x}_i^* = \underline{x}_i - E(\underline{x}_i)$ ,  $\underline{y}_i^* = \underline{y}_i - E(\underline{y}_i)$ .

Let  $\underline{z}_i^* = (\underline{x}_i^*, \underline{y}_i^*)$ ,  $\underline{w}_i^* = (\underline{w}_{i1}^*, \underline{w}_{i2}^*) = (\underline{z}_i^*, \underline{z}_i^* \otimes \underline{x}_i^*)$ , and

set  $\bar{\underline{w}}^* = \sum_{i=1}^N \underline{w}_i^*/N$ . Now note that  $g(\bar{\underline{w}}^*) = g(\bar{\underline{w}})$  and  $g(E(\underline{w}_i^*)) = g(E(\underline{w}_i))$

since  $\underline{z}_i^* - \bar{\underline{z}}^* = \underline{z}_i - \bar{\underline{z}}$  and  $V(\underline{z}_i^*) = V(\underline{z}_i)$ . Since  $\underline{w}_i^*$  is i.i.d., we have

$$\underline{\Omega} = \frac{\partial g}{\partial \underline{r}'} V(\underline{w}_i^*) \frac{\partial g'}{\partial \underline{r}'},$$

where the partial derivatives are evaluated at  $E(\underline{w}_i^*)$ .

We can write

$$g(\underline{r}) = f[h(\underline{r})], \quad h(\underline{r}) = \underline{r}_2 - \underline{r}_1 \otimes \underline{r}_{11},$$

where the  $\underline{r}' = (\underline{r}_1', \underline{r}_2')$  partition conforms to the  $\underline{w}' = (\underline{w}_1', \underline{w}_2')$  partition

and  $\underline{r}_{11}$  is the sub-vector of  $\underline{r}_1$  containing the first  $K$  elements. When

$\underline{r} = E(\underline{w}_1^*)$ , we have

$$\frac{\partial \underline{h}}{\partial \underline{r}_1'} = \underline{0}, \quad \frac{\partial \underline{h}}{\partial \underline{r}_2'} = \underline{I}_{K(K+M)} ;$$

hence we can simplify to

$$\underline{\Omega} = \frac{\partial \underline{f}}{\partial \underline{h}'} V(\underline{w}_{i2}^*) \frac{\partial \underline{f}'}{\partial \underline{h}},$$

where the partial derivatives are evaluated at  $\underline{h} = E(\underline{w}_{i2}^*)$ .

The function  $\underline{f}$  can be implicitly defined by

$$\underline{S} \underline{f}_m = \underline{d}_m \quad (m = 1, \dots, M),$$

where  $\text{vec}(\underline{s}_1, \dots, \underline{s}_K, \underline{d}_1, \dots, \underline{d}_M) = \underline{h}$ ,  $\underline{S} = (\underline{s}_1, \dots, \underline{s}_K)$ , and

$\text{vec}(\underline{f}_1, \dots, \underline{f}_M) = \underline{f}$ . So we have

$$\underline{S} \frac{\partial \underline{f}_m}{\partial \underline{s}_k'} + f_{mk} \underline{I}_K = \underline{0} \quad (k = 1, \dots, K; m = 1, \dots, M),$$

$$\underline{S} \frac{\partial \underline{f}_m}{\partial \underline{d}_n'} = \begin{cases} \underline{0} & \text{if } n \neq m \\ \underline{I}_K & \text{if } n = m \end{cases} \quad (n, m = 1, \dots, M).$$

If  $\underline{h} = E(\underline{w}_{i2}^*)$ , then  $\underline{S} = V(\underline{x}_i)$ ,  $\underline{d}_m = \text{Cov}(\underline{x}_i, y_{im})$ , and  $\underline{f} = \text{vec } \underline{\Pi}'$ . So the partial derivatives evaluated at  $\underline{h} = E(\underline{w}_{i2}^*)$  are

$$\frac{\partial \underline{f}}{\partial \underline{h}'} = (-\underline{\Pi}, \underline{I}_M) \otimes V^{-1}(\underline{x}_i).$$

Since

$$\begin{aligned} V(\underline{w}_{i2}^*) &= V(\underline{z}_i^* \otimes \underline{x}_i^*) \\ &= E(\underline{z}_i^* \underline{z}_i^{*'} \otimes \underline{x}_i^* \underline{x}_i^{*'}) - E(\underline{z}_i^* \otimes \underline{x}_i^*) E(\underline{z}_i^{*'} \otimes \underline{x}_i^{*'}), \end{aligned}$$

there are two terms to evaluate:

$$\begin{aligned} 1. \quad & \frac{\partial f}{\partial \underline{h}'} [E(\underline{z}_i^* \underline{z}_i^{*'} \otimes \underline{x}_i^* \underline{x}_i^{*'})] \frac{\partial f'}{\partial \underline{h}} \\ &= E\left\{ [(-\underline{\Pi}, \underline{I}_M) \otimes \underline{V}^{-1}(\underline{x}_i)] (\underline{z}_i^* \underline{z}_i^{*'} \otimes \underline{x}_i^* \underline{x}_i^{*'}) [(-\underline{\Pi}, \underline{I}_M)' \otimes \underline{V}^{-1}(\underline{x}_i)] \right\} \\ &= E[(\underline{y}_i^* - \underline{\Pi} \underline{x}_i^*) (\underline{y}_i^* - \underline{\Pi} \underline{x}_i^*)' \otimes \underline{V}^{-1}(\underline{x}_i) (\underline{x}_i^* \underline{x}_i^{*'}) \underline{V}^{-1}(\underline{x}_i)]. \\ 2. \quad & \frac{\partial f}{\partial \underline{h}'} [E(\underline{z}_i^* \otimes \underline{x}_i^*) E(\underline{z}_i^{*'} \otimes \underline{x}_i^{*'})] \frac{\partial f'}{\partial \underline{h}} = 0 \end{aligned}$$

since

$$\begin{aligned} & E\left\{ [(-\underline{\Pi}, \underline{I}_M) \otimes \underline{V}^{-1}(\underline{x}_i)] (\underline{z}_i^* \otimes \underline{x}_i^*) \right\} \\ &= E[(\underline{y}_i^* - \underline{\Pi} \underline{x}_i^*) \otimes \underline{V}^{-1}(\underline{x}_i) \underline{x}_i^*] \end{aligned}$$

and

$$E(\underline{y}_{im}^* - \underline{\pi}_m' \underline{x}_i^*) \underline{x}_i^* = \text{Cov}(\underline{y}_{im}, \underline{x}_i) - \underline{V}(\underline{x}_i) \underline{\pi}_m = 0 \quad (m = 1, \dots, M).$$

We conclude that

$$\underline{\Omega} = E[(\underline{y}_i^* - \underline{\Pi} \underline{x}_i^*) (\underline{y}_i^* - \underline{\Pi} \underline{x}_i^*)' \otimes \underline{V}^{-1}(\underline{x}_i) (\underline{x}_i^* \underline{x}_i^{*'}) \underline{V}^{-1}(\underline{x}_i)].$$

Next we shall apply this procedure to evaluate the limiting distribution of the instrumental variable estimator. Let  $\underline{q}_i$  be  $J \times 1$  and assume that  $(\underline{x}_i, \underline{y}_i, \underline{q}_i)$ ,  $i=1, \dots, N$ , is i.i.d. according to some multivariate distribution with finite fourth moments. Let  $\underline{x}_{im}$  be a  $K_m \times 1$  subvector of  $\underline{x}_i$ ; let  $\underline{q}_{im}$  be a  $K_m \times 1$  subvector of  $\underline{q}_i$ , and consider the instrumental variable estimators:

$$\hat{\underline{\pi}}_m = \left[ \sum_{i=1}^N (\underline{q}_{im} - \bar{\underline{q}}_m)(\underline{x}_{im} - \bar{\underline{x}}_m)' \right]^{-1} \sum_{i=1}^N (\underline{q}_{im} - \bar{\underline{q}}_m)(\underline{y}_{im} - \bar{\underline{y}}_m) \quad (m=1, \dots, M)$$

We shall assume that  $\underline{\phi}_m = \text{Cov}(\underline{q}_{im}, \underline{x}_{im}')$  is non-singular,  $m=1, \dots, M$ . Then  $\hat{\underline{\pi}}_m$  almost surely converges to  $\underline{\pi}_m = [\text{Cov}(\underline{q}_{im}, \underline{x}_{im}')]^{-1} \text{Cov}(\underline{q}_{im}, \underline{y}_{im})$  as  $N \rightarrow \infty$ . Let  $\hat{\underline{\pi}}' = (\hat{\underline{\pi}}_1', \dots, \hat{\underline{\pi}}_M')$  and  $\underline{\pi}' = (\underline{\pi}_1', \dots, \underline{\pi}_M')$ . Since  $\hat{\underline{\pi}}$  is a differentiable function of sample moments, we have

$$\sqrt{N} (\hat{\underline{\pi}} - \underline{\pi}) \xrightarrow{D} N(0, \underline{\Omega}) \quad \text{as } N \rightarrow \infty.$$

In order to simplify  $\underline{\Omega}$ , we shall repeat the procedure used for the multivariate regression case. Now  $\underline{f}' = (f_1', \dots, f_M')$  is implicitly defined by

$$\underline{S}_{m-m} \underline{f} = \underline{d}_m \quad (m=1, \dots, M)$$

where  $\underline{S}_m = (s_{m1}, \dots, s_{mK_m})$  is  $K_m \times K_m$ .

$$\underline{S}_m \frac{\partial \underline{f}_m}{\partial \underline{s}'_{nk}} = \begin{cases} 0 & \text{if } n \neq m \\ -f_{mk} \underline{I}_{K_m} & \text{if } n = m \end{cases}$$

( $k = 1, \dots, K_n$ ;  $n, m = 1, \dots, M$ ).

$$S_{-m} \frac{\partial f_{-m}}{\partial d'_{-n}} = \begin{cases} 0 & \text{if } n \neq m \\ I_{-K_m} & \text{if } n=m \end{cases} \quad (n, m=1, \dots, M).$$

We have  $S_{-m} = \text{Cov}(q_{-im}, x'_{-im})$ ,  $d_{-m} = \text{Cov}(q_{-im}, y_m)$ , and  $f_{-m} = \pi_{-m}$  when these functions are evaluated at the population means. Let  $h_{-m} = \text{vec}(s_{-m1}, \dots, s_{-mK_m}, d_{-m})$ . Then evaluating the partial derivatives at the population means gives

$$\frac{\partial f_{-m}}{\partial h'_{-n}} = \begin{cases} 0 & \text{if } n \neq m \\ (-\pi'_{-m}, 1) \otimes \phi_{-m}^{-1} & \text{if } n=m \end{cases} \quad (n, m=1, \dots, M).$$

The  $m, n$  submatrix of  $\Omega$  gives the asymptotic covariance between  $\hat{\pi}_{-m}$  and  $\hat{\pi}_{-n}$ :

$$\begin{aligned} \Omega_{-mn} &= \frac{\partial f_{-m}}{\partial h'_{-m}} [E(z_{-im}^* z_{-in}^{*'}) \otimes q_{-im}^* q_{-in}^{*'}] \\ &\quad - E(z_{-im}^* \otimes q_{-im}^*) E(z_{-in}^{*'} \otimes q_{-in}^{*'}) \frac{\partial f'_{-n}}{\partial h_{-n}}, \end{aligned}$$

where  $z_{-im}^{*'} = (x_{-im}^{*'}, y_{im}^*)$ ,  $x_{-im}^* = x_{-im} - E(x_{-im})$ ,  $y_{im}^* = y_{im} - E(y_{im})$ , and  $q_{-im}^* = q_{-im} - E(q_{-im})$ . As before there are two terms to evaluate and the second one is zero since  $E(y_{im}^* - \pi'_{-m} x_{-im}^*) q_{-im}^* = 0$ . The first term simplifies to

$$\Omega_{-mn} = E[(y_{im}^* - \pi'_{-m} x_{-im}^*)(y_{in}^* - \pi'_{-n} x_{-in}^*) \phi_{-m}^{-1} q_{-im}^* q_{-in}^{*'} \phi_{-n}^{-1}],$$

where  $\phi_{-m} = \text{Cov}(q_{-im}, x'_{-im})$ ,  $m, n=1, \dots, M$ .

## FOOTNOTES

<sup>1</sup>See Mundlak (1961, 1978).

<sup>2</sup>See Sims (1972).

<sup>3</sup>MaCurdy (1979) considers the conventional maximum likelihood estimator in a homoskedastic multivariate regression model. He shows that its asymptotic distribution can be obtained without normality assumptions by treating it as a minimum distance estimator. In our model, however, that maximum likelihood estimator is generally inefficient since it is using the wrong norm.

<sup>4</sup>See Mundlak (1961), Chamberlain (1980).

## REFERENCES

- Billingsley, P., 1965, Ergodic theory and information (Wiley, New York).
- \_\_\_\_\_, 1979, Probability and measure (Wiley, New York).
- Chamberlain, G., 1980, Analysis of covariance with qualitative data, Review of Economic Studies, 47, 225-238.
- Cramèr, H., 1946, Mathematical methods of statistics (Princeton University Press, Princeton).
- 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.
- Heckman, J. J. and T. E. MaCurdy, 1980, A life cycle model of female labor supply, Review of Economic Studies, 47, 47-74.
- Kendall, M. G. and A. Stuart, 1961, The advanced theory of statistics, Vol. 2 (Griffin, London).
- MaCurdy, T. E., 1979, Multiple time series models applied to panel data: specification of a dynamic model of labor supply, unpublished manuscript.
- \_\_\_\_\_, 1980, An empirical model of labor supply in a life cycle setting, National Bureau of Economic Research Working Paper No. 421.
- Malinvaud, E., 1970, Statistical methods of econometrics (North-Holland, Amsterdam).
- Mundlak, Y., 1961, Empirical production function free of management bias, Journal of Farm Economics, 43, 44-56.
- \_\_\_\_\_, 1978, "On the pooling of time series and cross section data, Econometrica, 46, 69-85.
- Rao, C. R., 1973, Linear statistical inference and its applications (Wiley, New York).
- Rozanov, Y. A., 1967, Stationary random processes (Holden-Day, San Francisco).
- Sims, C. A., 1972, Money, income, and causality, American Economic Review, 62, 540-552.



White, H., 1980, Using least squares to approximate unknown regression functions, International Economic Review 21, 149-170.