

UC Berkeley

CUDARE Working Papers

Title

On efficient estimation of some limited dependent variables models

Permalink

<https://escholarship.org/uc/item/6sd75371>

Author

Tsur, Yacov

Publication Date

1983-07-01

University of California, Berkeley
Department of Agricultural &
Resource Economics

CUDARE Working Papers

Year 1983

Paper 281

On Efficient Estimation of Some Limited
Dependent Variables Models

Yacov Tsur

ON EFFICIENT ESTIMATION OF SOME LIMITED DEPENDENT VARIABLES MODELS†

Yacov Tsur‡

INTRODUCTION

In empirical research, it is often the case that the data set under investigation is incomplete. Consequently, a plethora of estimation methods have been developed to deal with the several possible pattern of missing data, such as in Afifi and Elashoff [1966], Griliches, Hall and Hausman [1978], and Dempster, Rubin and Laird [1977]. In general, these methodologies use prior knowledge or information on the process generating the data and on the missing patterns to "complete" the data set. Standard estimation procedures are then employed on this new data set.

The econometric literature distinguishes two types of models corresponding to missing data problems: incomplete data on predetermined or right-hand-side (RHS) variables only; and incomplete data on the dependent, and possibly on the RHS, variables. In the latter models, known as Limited Dependent Variable (LDV) models, the most widely used estimation procedures follow Heckman [1976,1978] and Lee [1978]. These procedures do not use "filled-in" methods but rather employ a two stage procedure where the second stage focuses only on the *observed* part of the sample with additional terms that account for the

† This research was partially supported by Water Resources Center Grant W622 and by the Giannini Foundation of Agricultural Economics.

‡ Department of Agricultural and Resource Economics, University of California, Berkeley.

I wish to thank T. Rothenberg, L. Le Cam, M. Hanemann and P. Ruud for helpful discussions and comments. Any remaining errors are, of course, my responsibility.

selection mechanism (or missing patterns), derived in the first stage.

Maximum likelihood (ML) methods are frequently considered (Amemiya [1973], Hausman and Wise [1978], Heckman [1978], Duncan [1980], Cosslet [1981], among others) but rarely used in practical applications because of computational and implementation difficulties. (See Griliches *et al.* [1978] as an exception.) This is primarily because the log likelihood functions of LDV models is generally highly non-linear in the parameters and often contains multiple roots so that iterations with arbitrary initial values may result in a root that does not correspond to the global maximum of the function (Amemiya [1973]). To overcome the possibility of a "wrong root" requires a "good" initial parameter value. This leads to the consideration of a single iteration procedure of the form $\vartheta^P - R(\vartheta^P)^{-1}L_{\vartheta}(\vartheta^P)$, where ϑ^P is the preliminary estimate, $R(\vartheta^P)$ is an estimate of the information matrix and $L_{\vartheta}(\vartheta^P)$ is the gradient of the log likelihood function evaluated at ϑ^P . The outcome of this procedure, denoted as the "Linearized Maximum Likelihood" (LML) procedure, is asymptotically equivalent to the ML estimator provided that ϑ^P approaches the true parameter at an "appropriate" rate as the sample size increases. (See Le Cam [1960, 1969], and Rothenberg and Leenders [1964] for a theoretical discussion, and Berndt, Hall, Hall and Hausman [1974] for implementation). In fact any estimate of the form $\vartheta^P - R(\vartheta^P)^{-1}\omega_{\vartheta}(\vartheta^P)$, such that $\sqrt{T}(\omega_{\vartheta}(\vartheta^P) - L_{\vartheta}(\vartheta^P))$ goes to zero in probability as the sample size $T \rightarrow \infty$, has the same limiting distribution as the LML estimator. We will use this generalization of LML and denote it as the GML estimator.

Recently, an alternative algorithm for the ML estimation of models with incomplete data has been offered by Dempster *et al.* [1977]. This method, called the EM algorithm, is an iterative procedure with each iteration comprised of two steps: an Expectation (E) step which can be interpreted as a guideline for filling in the missing data; and a Maximization (M) step which proceeds with a

maximization task. Since this seminal paper, an extensive literature has emerged dealing with the convergence properties of the algorithm and its possible applications (e.g., Wu [1983] and the reference cited therein).

The purpose of the research reported here is to investigate efficient estimation methods for some LDV models in the context of the EM algorithm.¹ The use of the term EM *approach* rather than algorithm is adopted because the M step can take on differing formats in the algorithm and we will modify this step substantially in several cases. Two procedures will be considered: an iterative procedure resulting from a straightforward application of the EM algorithm; and a single iteration procedure, corresponding to the LML procedure, resulting from the EM approach. Each method is evaluated under two scenarios: single equation models (Section 2) and two-equation models (Section 3). Each scenario involves two models corresponding to different patterns of missing data.

The conclusions emanating from the analysis include:

1. ML estimators can be derived via iterative Ordinary Least Squares (OLS) in single equations models, or Generalized Least Squares (GLS) in multiple equation models if the EM algorithm is employed.
2. The EM approach provides a unique method for deriving and interpreting the LML estimator which basically involves the application of an Instrumental Variable (IV) regression.
3. The approach lends itself to an integrative framework which unifies the analysis of the many possible LDV models and links it to classical linear model theory.

¹ The concept of efficiency used in this work together with the conditions that assure the efficiency of the ML estimators are given in Appendix C. Consequently, the terms "efficient estimator" and "ML estimator" are used interchangeably.

The basic notation is as follows: the parameter set is denoted by Θ with ϑ representing an arbitrary member. A specific member of Θ is denoted by a super-script on ϑ except for the cases of ML, LML, and GML estimators which are denoted by $\hat{\vartheta}$, ϑ_{LML} , and ϑ_{GML} respectively. Other exceptions are indicated in the text. The true parameter value is given by ϑ^0 . Superscripts in general indicate evaluations of the quantity at the specific parameter value. A superscript "*" indexing a variable indicates that the variable is incompletely observed. The density and distribution functions of a standard normal variate are denoted by $\phi(\cdot)$ and $\Phi(\cdot)$ respectively. The line segment joining two distinct points ϑ' and ϑ'' of Θ is denoted by $\vartheta' \vartheta''$. The log likelihood function is denoted by $L(\vartheta)$ with $L_{\vartheta}(\cdot)$ and $L_{\vartheta\vartheta}(\cdot)$ indicating the gradient vector of first derivatives and the Hessian matrix of second derivatives respectively. An estimate ϑ^P is consistent of order $T^{-\frac{1}{2}}$, or \sqrt{T} -consistent, if $\sqrt{T}(\vartheta^P - \vartheta^0)$ is bounded in probability as $T \rightarrow \infty$.

2. SINGLE EQUATION MODELS

Two basic LDV single equation models, the Tobit and probit models, corresponding to two different patterns of missing data, are considered in this section.

The underlying structure is given by

$$(2.1) \quad y_t^* = x_t' \beta^0 + u_t, \quad t = 1, 2, \dots, T$$

$$(2.2) \quad u_t \stackrel{iid}{\sim} N(0, \sigma^2), \quad t = 1, 2, \dots, T$$

The Tobit model results from the following observation rule

$$(2.3) \quad y_t = \begin{cases} y_t^* & \text{if } y_t^* \geq 0 \\ 0 & \text{otherwise} \end{cases} \quad t = 1, 2, \dots, T$$

and for the probit model

$$(2.3) \quad y_t = \begin{cases} 1 & \text{if } y_t^* \geq 0 \\ 0 & \text{otherwise} \end{cases} \quad t=1,2,\dots,T$$

For both models it is assumed that x_t is observed for all t , where $X = (x_1, x_2, \dots, x_T)'$ is a $T \times k$ matrix, $y^* = (y_1^*, y_2^*, \dots, y_T^*)'$ is a $T \times 1$ vector of quantities for which we observe $y = (y_1, y_2, \dots, y_T)'$ according to (2.3) or (2.3'), β^0 and σ^0 are respectively $k \times 1$ and 1×1 vectors of parameters to be estimated, and $u = (u_1, u_2, \dots, u_T)'$ is a $T \times 1$ vector of unobserved errors distributed according to (2.2). It is assumed that X is of full rank, that $\lim_{T \rightarrow \infty} \frac{1}{T} X'X$ exists and is positive definite and that each column of X is statistically independent of the error vector u .

The missing data pattern corresponding to the probit model creates identification problems in that any structure that agrees in sign with (2.1) results in the same observations. Specifically the structure $\frac{y_t^*}{a} = \frac{x_t' \beta^0}{a} + \frac{u_t}{a}$ for an arbitrary $0 < a < \infty$ is observationally equivalent to (2.1) and so a normalization rule is required to identify any of the parameters. The normalization usually employed and used in this study is $\sigma^0 = 1$. No identification problems arise in the Tobit model.

Let q_t be defined as the indicator variable

$$(2.4) \quad q_t = \begin{cases} 1 & \text{if } y_t^* \geq 0 \\ 0 & \text{otherwise} \end{cases} \quad t=1,2,\dots,T$$

The sample q_t ; $t = 1, 2, \dots, T$ is a realization of T independent Bernoulli trials with the probability of success for the t^{th} trial equal to $1 - \Phi(-x_t' \beta^0 / \sigma^0)$. This

captures the randomness inherent in the observation rule.

In the application of the EM algorithm, discussed below, the following definitions and results will be used

$$(2.5a) \quad \lambda_i^- = \frac{-\varphi(-\mu_i/\sigma)}{\Phi(-\mu_i/\sigma)}$$

$$(2.5b) \quad \lambda_i^+ = \frac{\varphi(-\mu_i/\sigma)}{1-\Phi(-\mu_i/\sigma)}$$

where $\mu_i = x_i'\beta$. It follows that

$$(2.6a) \quad \frac{\partial \lambda_i^-}{\partial \mu_i} \stackrel{def}{=} \Lambda_i^- = -\frac{1}{\sigma} \lambda_i^- (\mu_i/\sigma + \lambda_i^-)$$

$$(2.6b) \quad \frac{\partial \lambda_i^+}{\partial \mu_i} \stackrel{def}{=} \Lambda_i^+ = -\frac{1}{\sigma} \lambda_i^+ (\mu_i/\sigma + \lambda_i^+)$$

Defining $\tilde{\tau}_i^- = var(u_i | u_i < -\mu_i)$ and $\tilde{\tau}_i^+ = var(u_i | u_i \geq -\mu_i)$ it can be demonstrated that (See Appendix A)

$$(2.7a) \quad \tilde{\tau}_i^- = \sigma^2(1 + \sigma \Lambda_i^-)$$

$$(2.7b) \quad \tilde{\tau}_i^+ = \sigma^2(1 + \sigma \Lambda_i^+).$$

For any given parameter value, ϑ^κ , we define $y_t^\kappa = E_{\vartheta^\kappa}\{y_t^* | y_t\}$ where $E_{\vartheta^\kappa}\{\cdot\}$ denotes the expectation when ϑ^κ holds. So

$$(2.8) \quad y_t^\kappa = q_t y_t + (1-q_t)(\mu_t^\kappa + \sigma^\kappa \lambda_i^{-\kappa}) \text{ under Tobit}$$

and

$$(2.8') \quad y_t^\kappa = q_t(\mu_t^\kappa + \lambda_i^{+\kappa}) + (1-q_t)(\mu_t^\kappa + \lambda_i^{-\kappa}) \text{ under probit.}$$

It can be verified that

$$(2.9) \quad E_{y^k} \{ (y_t^* - \mu_t)^2 | y_t \} = (y_t^k - \mu_t)^2 + \tilde{\tau}_t^k$$

where

$$(2.10) \quad \tilde{\tau}_t = (1 - q_t) \tilde{\tau}_t^- \quad \text{under Tobit}$$

and

$$(2.10') \quad \tilde{\tau}_t = (1 - q_t) \tilde{\tau}_t^- + q_t \tilde{\tau}_t^+ \quad \text{under probit .}$$

The log of the joint density function of y_t^* ; $t = 1, 2, \dots, T$, is given by (disregarding the constant term)

$$(2.11) \quad f(y^*; \vartheta) = - \frac{1}{2\sigma^2} \sum_{t=1}^T (y_t^* - \mu_t)^2 + \frac{T}{2} \log \frac{1}{\sigma^2}$$

From (2.9) it follows that

$$(2.12) \quad E_{y^k} \{ f(y^*; \vartheta) | y \} = - \frac{1}{2\sigma^2} \left\{ \sum_{t=1}^T [(y_t^k - \mu_t)^2 + \tilde{\tau}_t^k] \right\} - \frac{T}{2} \log \sigma^2 .$$

Maximizing (2.12) over β and σ^2 yields

$$(2.13a) \quad \beta^{k+1} = (X'X)^{-1} X' y^k$$

$$(2.13b) \quad \sigma^{k+1^2} = \frac{1}{T} \sum_{t=1}^T (y_t^k - x_t' \beta^{k+1})^2 + \frac{1}{T} \sum_{t=1}^T \tilde{\tau}_t^k$$

Equations (2.12) and (2.13) form, respectively, the E step and M step of one iteration of the EM algorithm which was shown (Dempster *et al.* [1977]) to yield a root of the likelihood function. Therefore ML estimate can be derived via iterative OLS procedures.²

² Note that σ^{σ^2} is not an unknown parameter under the Probit model so (2.13b) is irrelevant in this case.

The same results were obtained by Hartley [1976] Fair [1977] (for the Tobit model), and Green [1982b]. The first two by considering the first order conditions of the log likelihood function and the later by using the EM method. Hence, in single equation models the standard approach of maximizing the log likelihood function can yield the same procedure as the EM algorithm of (2.11) to (2.13). However, as is demonstrated in Section 3, in multiple equation models the EM algorithm provides a unique procedure with advantages that have not previously been recognized.

The other estimation procedure considered in this study is the single step procedure. In this procedure an initial "good" estimate (ϑ^P) is first derived and used to construct $R(\vartheta^P)$ and $\omega_{\vartheta}(\vartheta^P)$. The final estimate is given by $\vartheta_{GML} = \vartheta^P - R(\vartheta^P)^{-1} \omega_{\vartheta}(\vartheta^P)$. Following Le Cam [1960, 1969], and Rothenberg and Leenders [1964], ϑ_{GML} is asymptotically efficient provided that ϑ^P is consistent of order $T^{-1/2}$, $R(\vartheta^P) \xrightarrow{P} \Psi(\vartheta^0)$ the information matrix, and $\sqrt{T} (\omega_{\vartheta}(\vartheta^P) - L_{\vartheta}(\vartheta^P)) \xrightarrow{P} 0$. ϑ_{LML} is accepted from ϑ_{GML} by putting $\omega_{\vartheta}(\vartheta^P) \equiv L_{\vartheta}(\vartheta^P)$. ϑ_{GML} is not unique since there are different $R(\cdot)$'s and $\omega_{\vartheta}(\cdot)$'s satisfying the conditions above (e.g., $R(\cdot)$ can be the Hessian matrix of the second derivatives of the log likelihood function or the one suggested by Berndt *et al.* [1974]). Our approach provides another method to construct ϑ_{GML} which produces an outcome that is identical to the standard method of putting $R(\cdot) = L_{\vartheta\vartheta}(\vartheta^P)$ and $\omega_{\vartheta}(\cdot) = L_{\vartheta}(\vartheta^P)$, in single equation models but, as will be shown in Section 3, is unique in multiple equation models.³ This is the subject of the following discussion which begins by considering the case of known σ^2 (which includes the Probit model automatically).

³ It is worth noting that, starting from a \sqrt{T} -consistent estimate, one EM iteration (of the form given in (2.12)-(2.13)) will not yield an efficient estimator. This is true also in multiple equations models.

2.A The Case of Known σ^2

The first step is to characterize the asymptotic properties of an efficient estimator. From (2.13a) it follows that $\beta^{k+1} = \beta^k + (X'X)^{-1}X'u^k$ where u^k is a $T \times 1$ vector with the t^{th} component defined as

$$(2.14) \quad u_t^k = y_t^k - x_t' \beta^k$$

Since the ML estimator $\hat{\beta}$ is a fixed point of (2.13a) the following orthogonality relation holds at $\hat{\beta}$

$$(2.15) \quad X u^{\hat{\beta}} = \sum_{t=1}^T x_t u_t^{\hat{\beta}} = 0$$

Expanding $\lambda_t^{-\hat{\beta}}$, $\lambda_t^{+\hat{\beta}}$ around $\lambda_t^{-\beta^0}$, $\lambda_t^{+\beta^0}$ respectively allows one to express $u_t^{\hat{\beta}}$ as

$$u_t^{\hat{\beta}} = [1 - (1 - q_t)(1 + \sigma \Lambda_t^{-\hat{\beta}})] x_t' (\beta^0 - \hat{\beta}) + \xi_t$$

under Tobit, where the super-script "o" is dropped from σ until otherwise indicated, and under the probit model as

$$u_t^{\hat{\beta}} = [1 - (1 - q_t)(1 + \Lambda_t^{-\hat{\beta}}) - q_t(1 + \Lambda_t^{+\hat{\beta}})] x_t' (\beta^0 - \hat{\beta}) + \xi_t$$

where $\tilde{\beta} \in \beta^0 \hat{\beta}$ and ξ_t is an error term defined as

$$(2.16) \quad \xi_t = u_t - (1 - q_t)(u_t - \sigma \lambda_t^{-\beta^0}) \quad \text{under Tobit}$$

and

$$(2.16') \quad \xi_t = u_t - (1 - q_t)(u_t - \lambda_t^{-\beta^0}) - q_t(u_t - \lambda_t^{+\beta^0}) \quad \text{under probit.}$$

Using (2.7) $u_t^{\hat{\beta}}$ can be expressed as

$$(2.17) \quad u_t^{\hat{\beta}} = \frac{1}{\sigma^2} [\sigma^2 - \tilde{\tau}_t^{\hat{\beta}}] x_t' (\beta^0 - \hat{\beta}) + \xi_t$$

where $\tilde{\tau}_t$ is defined in (2.10) and (2.10') and is rewritten below for convenience

$$(2.10) \quad \tilde{\tau}_t^+ = (1 - q_t) \tilde{\tau}_t^- \quad \text{under Tobit}$$

$$(2.10') \quad \tilde{\tau}_t = (1 - q_t) \tilde{\tau}_t^- + q_t \tilde{\tau}_t^+ \quad \text{under probit.}$$

It follows from (2.15) and (2.17) that

$$(2.18) \quad \sqrt{T}(\hat{\beta} - \beta^0) = \left\{ \frac{1}{T} \sum_{t=1}^T x_t \frac{1}{\sigma^2} (\sigma^2 - \tilde{\tau}_t^2) x_t' \right\}^{-1} \frac{1}{\sqrt{T}} \sum_{t=1}^T x_t \xi_t$$

Accounting for the randomness underlying q_t the ξ_t 's are independent random variables with zero mean and variance equal to:

$$(2.19) \quad \text{var} \{ \xi_t \} = \sigma^0 - \tau_t^0; \quad t=1, 2, \dots, T$$

where

$$(2.20) \quad \tau_t = \Phi(-\mu_t / \sigma) \tilde{\tau}_t^- \quad \text{under Tobit}$$

and

$$(2.20') \quad \tau_t = \Phi(-\mu_t) \tilde{\tau}_t^- + [1 - \Phi(-\mu_t)] \tilde{\tau}_t^+ \quad \text{under probit.}$$

Applying Liapounoff's Central Limit Theorem to $v' \frac{1}{\sqrt{T}} \sum_{t=1}^T x_t \xi_t$ for an arbitrary

K-vector v yields

$$(2.21) \quad \frac{1}{\sqrt{T}} \sum_{t=1}^T x_t \xi_t \xrightarrow{d} N(0, A^0 - D^0)$$

where

$$(2.22a) \quad A = \lim_{T \rightarrow \infty} \frac{1}{T} \sigma^2 \sum_{t=1}^T x_t x_t'$$

$$(2.22b) \quad D = \lim_{T \rightarrow \infty} \frac{1}{T} \sum_{t=1}^T x_t \tau_t x_t'$$

Note that $A-D \geq 0$ (i.e., is positive semi-definite) since $\sigma^2 - \tau_t \geq 0$ for all t .

Let v be a member of an arbitrary vector space (finite dimension) such that the $\lim_{T \rightarrow \infty} \frac{1}{T} \sum_{i=1}^T v_i v_i'$ exists and let d_i be a Bernoulli variate with success probability equal to P_i . Then the following result holds. (See Amemiya [1973]).

$$(2.24) \quad \lim_{T \rightarrow \infty} \frac{1}{T} \sum_{i=1}^T v_i d_i v_i' = \lim_{T \rightarrow \infty} \frac{1}{T} \sum_{i=1}^T v_i P_i v_i'$$

This result together with the fact that $\tilde{\beta}$ is consistent implies that

$$(2.25) \quad \frac{1}{T} \sum_{i=1}^T x_i \frac{1}{\sigma^2} [\sigma^2 - \tilde{\tau}_i^{\tilde{\beta}}] x_i' \xrightarrow{P} \frac{1}{\sigma^2} [A^0 - D^0]$$

It follows from (2.18), (2.22a and b), and (2.25) that

THEOREM:

$$(2.26) \quad \sqrt{T} (\hat{\beta} - \beta^0) \xrightarrow{d} N\{0, \sigma^4 (A^0 - D^0)^{-1}\}.$$

The inverse of the covariance matrix in the RHS of (2.26), denoted as the asymptotic precision (also the information matrix in this case) is equal to $\frac{1}{\sigma^4} A^0 - \frac{1}{\sigma^4} D^0$. The first term is the asymptotic precision in the absence of missing observations. The second term is the asymptotical loss in precision (information) due to the missing data pattern.⁴

Having characterized the asymptotic distribution of an efficient estimator, we now assume that a "good" initial estimate, denoted as β^P , is available. Define

$$(2.27) \quad \tilde{y}_i^P = y_i^P - \frac{1}{\sigma^2} \tilde{\tau}_i^P x_i' \beta^P; \quad \tilde{y}^P = (\tilde{y}_1^P, \tilde{y}_2^P, \dots, \tilde{y}_T^P)' \text{ a } T \times 1 \text{ vector}$$

⁴ The theorem above is, of course, not new (e.g., see Amemiya [1973]), but the process of deriving it, motivated by the fact that the ML estimator, $\hat{\beta}$, is a fixed point of (2.13a) thereby leading to (2.15) - (2.26), is a natural outcome of using the EM algorithm. This generalizes easily to more complicated LDV models and provides another way of deriving asymptotic distributions of ML estimators.

$$(2.28) \quad \tilde{x}_t^P = x_t' - \frac{1}{\sigma^2} \tilde{\tau}_t^P x_t'; \quad \tilde{X}^P = (\tilde{x}_1^P, \tilde{x}_2^P, \dots, \tilde{x}_T^P)' \text{ a } T \times k \text{ matrix}$$

where y_t^P defined in (2.8) or (2.8') by substituting ϑ^P for ϑ^* and $\tilde{\tau}_t^P$ is $\tilde{\tau}_t$ of (2.10) - (2.10') evaluated at ϑ^P . With the above definitions, the following relation is identified

$$(2.29) \quad \tilde{y}_t^P = \tilde{x}_t^P \beta^0 + \xi_t; \quad t=1,2,\dots,T$$

where ξ_t is the same error term defined in (2.16) or (2.16'). Our estimator, denoted by $\bar{\beta}$, is the result of applying IV regression to (2.29) using x' as the instrument. Formally

$$(2.30) \quad \bar{\beta} = (X\tilde{X}^P)^{-1}X'\tilde{y}^P$$

It can be verified that $\bar{\beta}$ is identical to $\beta_{LML} = \beta^P - L_{\beta\beta'}^{-1}L_{\beta}^P$ by writing (2.30) as

$$\bar{\beta} = \beta^P + (X\tilde{X}^P)^{-1}X'u^P$$

and noting that

$$\frac{\partial L(\beta^P)}{\partial \beta} = \frac{1}{\sigma^2} X'u^P$$

and

$$\frac{\partial^2 L(\beta^P)}{\partial \beta \partial \beta'} = -\frac{1}{\sigma^2} (X'\tilde{X}^P)$$

(see Appendix B). For the sake of completeness the asymptotic efficiency of $\bar{\beta}$ when β^P is consistent of order $T^{-1/2}$, is demonstrated below:

$$\begin{aligned} \bar{\beta} &= (X\tilde{X}^P)^{-1}X'\tilde{y}^P = \beta^P + (X\tilde{X}^P)^{-1}X'u^P \\ &= \beta^P + (X\tilde{X}^P)^{-1} \sum_{t=1}^T \{x_t \frac{1}{\sigma^2} [\sigma^2 - \tilde{\tau}_t^P] x_t' (\beta^0 - \beta^P) + x_t \xi_t\} \end{aligned}$$

(by replacing β^P for $\hat{\beta}$ in (2.16) where $\tilde{\beta} \in \beta^P \beta^0$)

$$= \beta^P + (X' \tilde{X}^P)^{-1} \sum_{t=1}^T \{x_t \frac{1}{\sigma^2} [\sigma^2 - \tilde{\tau}_t^P - o_P(1)] x_t' (\beta^0 - \beta^P) + x_t \xi_t\}$$

(where $h_T = o_P(1)$ if $h_T \xrightarrow{P} 0$ as $T \rightarrow \infty$)

$$= \beta^P + (X' \tilde{X}^P)^{-1} X' \{ \tilde{X}^P (\beta^0 - \beta^P) + X [o_P(1) (\beta^P - \beta^0)] + \xi \}$$

$$= \beta^0 + (X' \tilde{X}^P)^{-1} X' \xi + (X' \tilde{X}^P)^{-1} (X' X) [o_P(1) (\beta^P - \beta^0)]$$

Hence

$$(2.31) \quad \sqrt{T} (\bar{\beta} - \beta^0) = \left(\frac{1}{T} X' \tilde{X}^P \right)^{-1} \frac{1}{\sqrt{T}} X' \xi +$$

$$\left[\frac{1}{T} X' \tilde{X}^P \right]^{-1} \left(\frac{1}{T} X' X \right) [o_P(1) \sqrt{T} (\beta^P - \beta^0)]$$

The convergence of $\frac{1}{T} X' \tilde{X}^P$ and $\frac{1}{T} X' X$ to a proper limit and the assumption that $\sqrt{T} (\beta^P - \beta^0)$ is bounded in probability assures that the second term in the RHS of (2.31) goes to zero in probability. By comparing (2.18) to (2.32) it follows that

$$\sqrt{T} (\hat{\beta} - \beta^0) - \sqrt{T} (\bar{\beta} - \beta^0) = \sqrt{T} (\hat{\beta} - \bar{\beta}) \xrightarrow{P} 0$$

This completes the demonstration. A straightforward implication of the above result is that $\text{var}(\bar{\beta}) = \sigma^2 (X' \tilde{X}^P)^{-1}$ which is the usual formula for the variance of an IV regression coefficient.

The major result is that an efficient estimator can be achieved via instrumental variable regression on a suitably transformed data set. In general, however, σ^{0^2} is unknown and must be estimated along with β^0 .

2.B The Case of Unknown σ^{0^2}

The normalization $\sigma^0 = 1$ excludes the probit model from this case and so only the Tobit model is considered. It is assumed that a "good" initial estimates β^P and σ^{P2} are available. Let us redefine \tilde{y}_i^P and \tilde{x}_i^P in (2.27) and (2.28) by replacing σ with σ^P so $\bar{\beta}$ in (2.30) is now evaluated at σ^P . Denote Ψ as the information matrix and partition it according to β and σ^2 ; i.e.,

$$(2.32) \quad \Psi = \begin{bmatrix} \Psi_{11} & \Psi_{12} \\ \Psi_{21} & \Psi_{22} \end{bmatrix}$$

where Ψ_{11} , $\Psi_{12} = \Psi_{21}$, and Ψ_{22} are, respectively, $k \times k$, $k \times 1$, and 1×1 matrices. Finally, let

$$(2.33a) \quad R^P = -\frac{1}{T}(X' \tilde{X}^P) \frac{1}{\sigma^{P2}}$$

and

$$(2.33b) \quad \Delta^P = \beta^P - \bar{\beta}$$

It is verified in Appendix B that

$$(2.34a) \quad R^P = \frac{\partial^2 L(\vartheta^P)}{\partial \beta \partial \beta'} \frac{1}{T}$$

$$(2.34b) \quad \Delta^P = \left[\frac{\partial^2 L(\vartheta^P)}{\partial \beta \partial \beta'} \right]^{-1} \frac{\partial L(\vartheta^P)}{\partial \beta}$$

Consider now the LML estimator

$$\vartheta_{LML} = \begin{bmatrix} \beta_{LML} \\ \sigma_{LML}^2 \end{bmatrix} = \begin{bmatrix} \beta^P \\ \sigma^{P2} \end{bmatrix} - \begin{bmatrix} \Psi_{11}^P & \Psi_{12}^P \\ \Psi_{21}^P & \Psi_{22}^P \end{bmatrix}^{-1} \begin{bmatrix} L_{\beta}(\vartheta^P) \\ L_{\sigma^2}(\vartheta^P) \end{bmatrix}$$

where $\Psi_{ij}^P = \Psi_{ij}(\vartheta^P)$ is any consistent estimate of $\Psi_{ij}(\vartheta^0)$ for $i, j=1, 2$. From an expression for a partitioned inverse, ϑ_{LML} can be written as

$$\sigma_{LML}^2 = \sigma^{P^2} - []^{P-1} \Psi_{21}^P \Psi_{11}^{P-1} L_{\beta}^P - []^{P-1} L_{\sigma^2}^P$$

$$\beta_{LML} = \beta^P - \Psi_{11}^{P-1} L_{\beta}^P - \Psi_{11}^{P-1} \Psi_{12}^P []^{P-1} \Psi_{21}^P \Psi_{11}^{P-1} L_{\beta}^P + \Psi_{11}^{P-1} \Psi_{12}^P []^{P-1} L_{\sigma^2}^P$$

where

$$[]^P = \Psi_{22}^P - \Psi_{21}^P \Psi_{11}^{P-1} \Psi_{12}^P$$

Choosing $\Psi_{11}^P = R^P$ and using (2.34b) gives

$$(2.36a) \quad \sigma_{LML}^2 = \sigma^{P^2} + []^{P-1} (\Psi_{21}^P \Delta^P - L_{\sigma^2}^P)$$

$$(2.36b) \quad \beta_{LML} = \bar{\beta} + R^{P-1} \Psi_{12}^P (\sigma^{P^2} - \sigma_{LML}^2)$$

$\bar{\beta}$, Δ^P and R^P are constructed using the method outlined above. The other terms needed in the derivation of ϑ_{LML} are Ψ_{12}^P , Ψ_{22}^P , and $L_{\sigma^2}^P$. They can be recovered from the first derivatives of the log likelihood function (Berndt, *at. al.* [1974]) which is given in Appendix B. This completes the discussion of our approach in analyzing single equation LDV models. The results generated by this approach are similar to those methods developed in the literature.

3. TWO EQUATION MODELS

The number of possible models corresponding to different patterns of missing data increases dramatically with the number of equations. This highlights the importance of a unified framework of analysis. Two models, resulting from two different patterns of missing data, are considered. Extension to other possible models is straightforward.

The underlying structure is given by

$$(3.1) \quad y_{it}^* = x'_{it} \beta^0 + u_{it} ; t = 1, 2, \dots, T$$

$$(3.2) \quad y_{2t}^* = x_{2t}'\beta^0 + u_{2t}; t=1,2,\dots,T$$

$$(3.3) \quad \begin{bmatrix} u_{1t} \\ u_{2t} \end{bmatrix} \underset{\sim}{iid} N(0, \Sigma^0); \Sigma = \begin{bmatrix} \sigma_1^2 & \sigma_{12} \\ \sigma_{12} & \sigma_2^2 \end{bmatrix}; t=1,2,\dots,T.$$

Model (I) results from the following observation rule

$$(3.4) \quad y_{1t} = \begin{cases} y_{1t}^* & \text{if } y_{2t}^* \geq 0 \\ NA & \text{otherwise} \end{cases} \quad t = 1,2,\dots,T$$

where NA denotes unobserved data

$$(3.5) \quad y_{2t} = \begin{cases} 1 & \text{if } y_{2t}^* \geq 0 \\ 0 & \text{otherwise} \end{cases} \quad t = 1,2,\dots,T$$

Model (II) results from

$$(3.4') \quad y_{1t} = \begin{cases} y_{1t}^* & \text{if } y_{2t}^* \geq 0 \\ NA & \text{otherwise} \end{cases} \quad t = 1,2,\dots,T$$

$$(3.5') \quad y_{2t} = \begin{cases} y_{2t}^* & \text{if } y_{2t}^* \geq 0 \\ NA & \text{otherwise} \end{cases} \quad t = 1,2,\dots,T$$

For both models it is assumed that x_{jt} , $j=1,2$ is observed for all t . For examples of situations that give rise to such models see Heckman [1976], Lee and Trost [1978], Lee [1978], and Hanemann and Tsur [1982]. The specifications of the quantities involved are as follows: $X_j = (x_{j1}, x_{j2}, \dots, x_{jt})'$, $j=1,2$ is a $T \times k$ matrix for which we assume:

- i) each column of X_j is statistically independent of $U_j = (u_{j1}, u_{j2}, \dots, u_{jt})'$, $j=1,2$

ii) $\lim_{T \rightarrow \infty} \frac{1}{T} X_j' X_j$ exists and is positive definite, $j = 1, 2$

iii) X_j is of full rank, $j = 1, 2$.

$Y_j^* = (y_{j1}^*, y_{j2}^*, \dots, y_{jT}^*)'$ is a $T \times 1$ vector for which we observe $Y_j = (y_{j1}, y_{j2}, \dots, y_{jT})$ according to the observation rule (3.4), (3.5) or (3.4'), (3.5'), $j=1, 2$. β^0 and Σ^0 are, respectively, $k \times 1$ vector and 2×2 positive definite symmetric matrix of parameters to be estimated and $U_j = (u_{j1}, u_{j2}, \dots, u_{jT})'$; $j=1, 2$ is an unobserved error vector with $(u_{1t}, u_{2t})'$ distributed according to (3.3) independently for $t = 1, 2, \dots, T$. The coefficient vectors in (3.1) and (3.2) are not necessarily identical. They can be distinct or have (some or all) common elements. However, it is always possible to merge them into one coefficient vector by a suitable redefinition of x_{1t} and x_{2t} . Hence there is no loss in generality in the way the model is presented. The missing data pattern can create identification problems. In model (I), without exogenous restrictions on the parameters, σ_2^0 , σ_{12}^0 , and the part of β^0 corresponding only to (3.2) are not identifiable. Therefore, a normalization rule is required and is given by $\sigma_2^0 = 1$. No identification problems arise in model (II).

Let us define

$$(3.6a) \quad y_t^* = \begin{bmatrix} y_{1t}^* \\ y_{2t}^* \end{bmatrix}; \quad Y^* = (y_{11}^*, y_{21}^*, \dots, y_{1T}^*)' \text{ a } 2T \times 1 \text{ vector}$$

$$(3.6b) \quad x_t = \begin{bmatrix} x_{1t} \\ x_{2t} \end{bmatrix}; \quad X = (x_1, x_2, \dots, x_T)' \text{ a } 2T \times k \text{ matrix}$$

and

$$(3.6c) \quad u_t = \begin{bmatrix} u_{1t} \\ u_{2t} \end{bmatrix}; \quad U = (u_1, u_2, \dots, u_T)' \text{ a } 2T \times 1 \text{ vector}$$

The structural model can be written now as

$$y_t^* = x_t' \beta^0 + u_t ; u_t \text{ iid } N(0, \Sigma^0), t=1,2,\dots,T.$$

With the observed dependent variable given by

$$(3.6d) \quad y_t = \begin{bmatrix} y_{1t} \\ y_{2t} \end{bmatrix}; Y = (y_1', y_2', \dots, y_T)' \text{ a } 2T \times 1 \text{ vector}$$

where y_{jt} , $j = 1,2$ follows from (3.4) and (3.5) in model (I) and from (3.4') and (3.5') in model (II).

In the application of the EM approach, the following definitions and results will be used:

$$(3.7) \quad \mu_{jt} = x_{jt}' \beta, j=1,2$$

$$(3.8a) \quad \lambda_t^- = \frac{-\varphi(-\mu_{2t}/\sigma_2)}{\Phi(-\mu_{2t}/\sigma_2)}$$

$$(3.8b) \quad \lambda_t^+ = \frac{\varphi(-\mu_{2t}/\sigma_2)}{1-\Phi(-\mu_{2t}/\sigma_2)}$$

From (3.8a and b)

$$(3.9a) \quad \Lambda_t^- = \frac{\partial \lambda_t^-}{\partial \mu_{2t}} = -\frac{1}{\sigma_2} \lambda_t^- \left(\frac{\mu_{2t}}{\sigma_2} + \lambda_t^- \right)$$

$$(3.9b) \quad \Lambda_t^+ = \frac{\partial \lambda_t^+}{\partial \mu_{2t}} = -\frac{1}{\sigma_2} \lambda_t^+ \left(\frac{\mu_{2t}}{\sigma_2} + \lambda_t^+ \right).$$

Let $\tilde{\tau}_{ijt}^- = \text{cov}(u_{it}, u_{jt} | u_{2t} \leq -\mu_{2t})$ and $\tilde{\tau}_{ijt}^+ = \text{cov}(u_{it}, u_{jt} | u_{2t} \geq -\mu_{2t})$ for $i, j = 1,2$. Then using properties of truncated bivariate normal (Johnson and Kotz [1972]. See also Appendix A.), it also follows that

$$(3.10a) \quad \tilde{\tau}_{11t}^- = \sigma_1^2 (1 + \rho^2 \sigma_2 \Lambda_t^-)$$

$$(3.10b) \quad \tilde{\tau}_{22t}^- = \sigma_2^2 (1 + \sigma_2 \Lambda_t^-)$$

$$(3.10c) \quad \tilde{\tau}_{12t}^- = \sigma_{12}(1 + \sigma_2 \Lambda_t^-)$$

$$(3.11a) \quad \tilde{\tau}_{11t}^+ = \sigma_1^2(1 + \rho^2 \sigma_2 \Lambda_t^+)$$

$$(3.11b) \quad \tilde{\tau}_{22t}^+ = \sigma_2^2(1 + \sigma_2 \Lambda_t^+)$$

$$(3.11c) \quad \tilde{\tau}_{12t}^+ = \sigma_{12}(1 + \sigma_2 \Lambda_t^+)$$

where $\rho = \frac{\sigma_{12}}{\sigma_1 \sigma_2}$. Let $\tilde{\tau}_t^-, \tilde{\tau}_t^+$ be a 2×2 matrix with ij element equal to $\tilde{\tau}_{ijt}^-, \tilde{\tau}_{ijt}^+$ respectively and define y_{jt}^k to be $E_{y^k}(y_{jt}^* | y_{jt})$ for $j=1,2$. That is, under model (I):

$$(3.12a) \quad y_{1t}^k = q_t y_{1t} + (1-q_t)(\mu_{1t}^k + \sigma_{12}^k \lambda_t^{-k})$$

$$(3.12b) \quad y_{2t}^k = q_t(\mu_{2t}^k + \lambda_t^{-k}) + (1-q_t)(\mu_{2t}^k + \lambda_t^{-k})$$

and under model (II):

$$(3.12a') \quad y_{1t}^k = q_t y_{1t} + (1-q_t)(\mu_{1t}^k + \frac{\sigma_{12}^k}{\sigma_2^k} \lambda_t^{-k})$$

$$(3.12b') \quad y_{2t}^k = q_t y_{2t} + (1-q_t)(\mu_{2t}^k + \sigma_2^k \lambda_t^{-k})$$

where q_t is the indicator variable defined as

$$(3.13) \quad q_t = \begin{cases} 1 & \text{if } y_{2t}^* \geq 0 \\ 0 & \text{otherwise} \end{cases} \quad t = 1, 2, \dots, T.$$

With the definitions in (3.12) it is straightforward to verify that

$$(3.14) \quad E_{y^k}\{(y_t^* - \mu_t)(y_t^* - \mu_t)' | y_t\} = (y_t^k - \mu_t)(y_t^k - \mu_t)' + \tilde{\tau}_t^k$$

where

$$(3.15) \quad \tilde{\tau}_t = (1-q_t)\tilde{\tau}_t^- + q_t \begin{bmatrix} 0 & 0 \\ 0 & \tilde{\tau}_{22t}^+ \end{bmatrix} \quad \text{under model (I)}$$

$$(3.15') \quad \tilde{\tau}_t = (1-q_t)\tilde{\tau}_t^- \quad \text{under model (II)}$$

and

$$\mu_t = (\mu_{1t}, \mu_{2t})'$$

The log of the joint density of y_t^* , $t = 1, 2, \dots, T$ (disregarding the constant term) is

$$(3.16) \quad f(y^*; \vartheta) = \frac{T}{2} \log |\Sigma^{-1}| - \frac{1}{2} \sum_{t=1}^T \text{tr} \{ \Sigma^{-1} (y_t^* - \mu_t)(y_t^* - \mu_t)' \}$$

where "tr" indicates trace. Using (3.14) gives

$$(3.17) \quad E_{y^k} \{ f(y^*; \vartheta) | y \} = - \frac{1}{2} \sum_{t=1}^T \text{tr} \{ \Sigma^{-1} [(y_t^k - \mu_t)(y_t^k - \mu_t)' + \tilde{\tau}_t^k] \} + \frac{T}{2} \log |\Sigma^{-1}|$$

Maximizing (3.17) over β and Σ yields

$$(3.18a) \quad \beta^{\kappa+1} = X'(I \otimes \Sigma^{\kappa+1-1})X)^{-1}X'(I \otimes \Sigma^{\kappa+1-1})Y^{\kappa}$$

$$(3.18b) \quad \Sigma^{\kappa+1} = \frac{1}{T} \sum_{t=1}^T (y_t^{\kappa} - \mu_t^{\kappa+1})(y_t^{\kappa} - \mu_t^{\kappa+1})' + \frac{1}{T} \sum_{t=1}^T \tilde{\tau}_t^{\kappa}$$

where "I" is the identity matrix (of order T) and " \otimes " is the Kronecker product.

Equations (3.17) and (3.18) form, respectively, the E step and the M step of an iteration of the EM algorithm which was shown (Dempster *et al.*) to yield an estimate corresponding to a (local) maximum of the likelihood function. Solving for $\beta^{\kappa+1}$ and $\Sigma^{\kappa+1}$ that satisfy (3.18) is not trivial and in fact requires an iterative procedure itself. This can be simplified by first solving (3.18a) using Σ^{κ} and then solving (3.18b) using $\beta^{\kappa+1}$.

⁴ this simplification corresponds to using a GEM procedure (see Dempster *et al.* p. 7).

Equation (3.18a) clearly reveals the similarity of our approach to the seemingly unrelated GLS estimation technique of Zellner [1962]. In fact, it implies that the ML estimator can be achieved via iterative GLS procedure where the missing data are being "filled-in" in each iteration according to (3.12a)-(3.12b) or (3.12a')-(3.12b'), depending on the observation rule (I) or (II) respectively. From there the analysis is carried out as if no data are missing with some modifications that account for the "filled-in" values.⁵

The advantage of this procedure is its simplicity of implementation which involves data transformations (that require the evaluation of the standard normal density and distribution functions) and the readily available GLS option. Its convergence properties and computational efficiency is yet to be compared to other iterative methods such as that of Berndt *et al.*

We turn now to the single iteration procedure. First, consider

3.A The Case of Known Σ^o

The first step is to characterize the limiting distribution of the ML estimator. Let us define

$$(3.19) \quad u_t^k = y_t^k - \mu_t^k = \begin{bmatrix} y_{1t}^k \\ y_{2t}^k \end{bmatrix} - \begin{bmatrix} x_{1t}^k \\ x_{2t}^k \end{bmatrix} \beta^k; \quad U^k = \begin{bmatrix} u_{1t}^k \\ \dots \\ u_{Tt}^k \end{bmatrix}, \quad \text{a } 2T \times 1 \text{ vector.}$$

From (3.18a) $\beta^{k+1} = \beta^k + (X'(I \otimes \Sigma^{-1})U^k)$. Since $\hat{\beta}$ is a fixed point of (3.18a) the following orthogonality relation holds

$$(3.20) \quad X'(I \otimes \Sigma^{-1})U^{\hat{\beta}} = 0$$

where the superscript "o" is dropped from Σ until otherwise indicated. By

⁵ The identification problem of Model (I) requires the normalization $\sigma_2 = 1$. This constraint should be incorporated into the maximization of (3.18) which is done over β and the identifiable elements of Σ . No such problem arises under model (II).

replacing $\lambda_t^{-\beta}$, $\lambda_t^{+\beta}$ with their first order expansions around λ_t^{-o} , λ_t^{+o} , respectively, and rearranging we can express u_t^{β} under model (I) as

$$u_t^{\beta} = \{\Sigma - (1-q_t)\} \begin{bmatrix} 1 & \frac{\sigma_{12}\Lambda_t^{-\beta}}{\sigma_2} \\ 0 & 1 + \sigma_2\Lambda_t^{-\beta} \end{bmatrix} \Sigma - q_t \begin{bmatrix} 0 & 0 \\ 0 & 1 + \sigma_2\Lambda_t^{+\beta} \end{bmatrix} \Sigma \} \Sigma^{-1} x_t'(\beta^o - \hat{\beta}) + \xi_t$$

and under model (II) as

$$u_t^{\beta} = \{\Sigma - (1-q_t)\} \begin{bmatrix} 1 & \frac{\sigma_{12}\Lambda_t^{-\beta}}{\sigma_2} \\ 0 & 1 + \sigma_2\Lambda_t^{-\beta} \end{bmatrix} \Sigma \} \Sigma^{-1} x_t'(\beta^o - \hat{\beta}) + \xi_t .$$

or by using (3.10) and (3.11) u_t^{β} can be expressed as

$$(3.21) \quad u_t^{\beta} = \{\Sigma - \tau_t^{*\beta}\} \Sigma^{-1} x_t'(\beta^o - \hat{\beta}) + \xi_t$$

where

$$(3.22) \quad \tau_t^* = (1-q_t)\tilde{\tau}_t^+ + q_t \begin{bmatrix} 0 & 0 \\ \tilde{\tau}_{21t}^+ & \tilde{\tau}_{22t}^+ \end{bmatrix} \text{ under model (I)}$$

$$(3.22') \quad \tau_t^* = (1-q_t)\tilde{\tau}_t^- \quad \text{under model (II)}$$

$$(3.23) \quad \xi_t = \begin{bmatrix} u_{1t} \\ u_{2t} \end{bmatrix} - (1-q_t) \begin{bmatrix} u_{1t} - \frac{\sigma_{12}\lambda_t^{-o}}{\sigma_2} \\ u_{2t} - \sigma_2\lambda_t^{-o} \end{bmatrix} - q_t \begin{bmatrix} 0 \\ u_{2t} - \sigma_2\lambda_t^{+o} \end{bmatrix} \text{ under model (I)}$$

$$(3.23') \quad \xi_t = \begin{bmatrix} u_{1t} \\ u_{2t} \end{bmatrix} - (1-q_t) \begin{bmatrix} u_{1t} - \frac{\sigma_{12}\lambda_t^{-o}}{\sigma_2} \\ u_{2t} - \sigma_2\lambda_t^{-o} \end{bmatrix} \text{ under model (II)}$$

and $\tilde{\beta} \in \hat{\beta}\beta^o$.

It follows from (3.20) and (3.21) that

$$(3.24) \quad \sqrt{2T} (\hat{\beta} - \beta^o) = \left\{ \frac{1}{2T} X'(I \otimes \Sigma^{-1}) [(I \otimes \Sigma) - \Gamma^{*\beta}] (I \otimes \Sigma^{-1}) X \right\}^{-1}$$

$$\frac{1}{\sqrt{2T}} X'(I \otimes \Sigma^{-1}) \xi$$

where Γ^* is a block diagonal matrix with the t^{th} block equal to τ_t^* and $\xi = (\xi_1', \dots, \xi_T)'$. When taking the q_t in (3.23) to be a Bernoulli variate with success probability of $1 - \Phi(-\mu_{2t}^0 / \sigma_2^0)$ the ξ_t are found to be 2×1 independent random vectors with zero mean and

$$(3.25) \quad \text{var}(\xi_t) = \Sigma^0 - \tau_t^0 ; t = 1, 2, \dots, T$$

where

$$(3.30) \quad \tau_t = (1 - \Phi_t) \begin{bmatrix} 0 & \tilde{\tau}_{12t}^+ \\ \tilde{\tau}_{12t}^+ & \tilde{\tau}_{22t}^+ \end{bmatrix} + \Phi_t \tilde{\tau}_t^- \text{ under model (I)}$$

$$(3.30') \quad \tau_t = \Phi_t \tilde{\tau}_t^- \quad \text{under model (II)}$$

and $\Phi_t = \Phi(-\mu_{2t} / \sigma_2)$. Let Γ be a $2T \times 2T$ block diagonal matrix with the t^{th} block equal to τ_t and define

$$(3.31a) \quad A = \lim_{T \rightarrow \infty} \frac{1}{2T} X'(I \otimes \Sigma^{-1}) X$$

$$(3.31b) \quad D = \lim_{T \rightarrow \infty} \frac{1}{2T} X'(I \otimes \Sigma^{-1}) \Gamma (I \otimes \Sigma^{-1}) X$$

$$(3.31c) \quad D^* = \lim_{T \rightarrow \infty} \frac{1}{2T} X'(I \otimes \Sigma^{-1}) \Gamma^* (I \otimes \Sigma^{-1}) X$$

The assumed convergence of $\frac{1}{2T} X_j' X_j$, $j = 1, 2$ assures the existence of A, D and D^* . Note that under model (II) $D = D^*$. This can be seen from the definitions of τ_t and τ_t^* in (3.30') and (3.22'), respectively, and from result (2.25). Formally,

$$(3.31d) \quad D = D^* \text{ under model (II)}$$

Since $\tilde{\beta}$ is consistent, it follows that

$$(3.32) \quad \frac{1}{2T} X'(I \otimes \Sigma^{-1})[(I \otimes \Sigma) - \Gamma^{*\beta}](I \otimes \Sigma^{-1})X \xrightarrow{P} A^0 - D^0$$

and from (3.25)

$$(3.33) \quad \frac{1}{2T} X'(I \otimes \Sigma^{-1})\xi \xrightarrow{d} N(0, A^0 - D^0)$$

These two results together with (3.24) yield

THEOREM:

$$(3.34) \quad \sqrt{2T}(\hat{\beta} - \beta^0) \xrightarrow{d} N\{0, (A^0 - D^0)^{-1}(A^0 - D^0)(A^0 - D^0)^{-1}\}.$$

Note that in model (II) $D = D^*$. Hence the asymptotic covariance matrix reduces to $(A^0 - D^0)^{-1}$.

The process of obtaining result (3.34) relies heavily on the EM approach (starting from the fact that $\hat{\beta}$ is a fixed point of (3.18a) hence (3.20), etc.) and is a straightforward extension of the single equation case that can be generalized easily to other multiple equation LDV models. The asymptotical precision (or information in these cases) is given by the inverse of the covariance matrix and is equal to $A - [D - (D - D^*)(A - D)^{-1}(D - D^*)]$ for model (I) and $A - D$ for model (II). In both cases A is the asymptotic precision in the absence of unobservables. Hence the second term can be interpreted as the asymptotic loss of information due to the unobservable pattern.

Suppose now that a "good" initial estimate β^P is available. Let us define

$$(3.35) \quad \tilde{y}_t^P = y_t^P - \tau_t^{*P} \Sigma^{-1} x_t^P \beta^P; \quad \tilde{Y}^P = (\tilde{y}_1^P, \tilde{y}_2^P, \dots, \tilde{y}_T^P)' \text{ a } 2T \times 1 \text{ vector}$$

$$(3.36a) \quad \tilde{x}_t^{*P} = x_t^P - \tau_t^{*P} \Sigma^{-1} x_t^P; \quad \tilde{X}^{*P} = (\tilde{x}_1^{*P}, \tilde{x}_2^{*P}, \dots, \tilde{x}_T^{*P})' \text{ a } 2T \times k \text{ matrix}$$

$$(3.36b) \quad \tilde{x}_t^P = x_t^P - \tau_t^P \Sigma^{-1} x_t^P; \quad \tilde{X}^P = (\tilde{x}_1^P, \tilde{x}_2^P, \dots, \tilde{x}_T^P)' \text{ a } 2T \times k \text{ matrix}$$

where τ_t^* and τ_t defined in (3.22)-(3.22') and (3.30)-(3.30') respectively.⁶ With the

⁶ Note that under model (II) $\tau_t^* = \tau_t$ which implies $\tilde{x}_t^{*P} = \tilde{x}_t^P$.

above definitions the following relation can be identified

$$(3.37) \quad \tilde{y}_t^P = \tilde{x}_t^{*P} \beta^o + \xi_t; t=1,2,\dots,T$$

where ξ_t is the same error defined in (3.23) or (3.23'). Our estimator $\bar{\beta}$ is the outcome of applying an IV regression on (3.37) using $X'(I \otimes \Sigma^{-1})$ as the instrument. That is

$$(3.38) \quad \bar{\beta} = (X'(I \otimes \Sigma^{-1})\tilde{X}^{*P})^{-1}X'(I \otimes \Sigma^{-1})\tilde{Y}^P$$

Unlike in one equation models, in this case it is impossible to show the equivalence of $\bar{\beta}$ and $\beta_{LML} = \beta^P - L_{\beta\beta}^{P-1}L_{\beta}^P$ by identifying the components at the RHS of (3.38) as the appropriate derivatives of the log likelihood function. However, it is proven below that

THEOREM:

$\bar{\beta}$ is asymptotically efficient, provided that β^P is consistent of order $T^{-1/2}$.

PROOF:

$$\begin{aligned} \bar{\beta} &= (X'(I \otimes \Sigma^{-1})\tilde{X}^{*P})^{-1}X'(I \otimes \Sigma^{-1})\tilde{Y}^P \\ &= \beta^P + (X'(I \otimes \Sigma^{-1})\tilde{X}^{*P})^{-1}X'(I \otimes \Sigma^{-1})U^P \end{aligned}$$

$$(where \ U^P = \tilde{Y}^P - \tilde{X}^{*P}\beta^P = Y^P - X\beta^P)$$

$$\begin{aligned} &= \beta^P + (X'(I \otimes \Sigma^{-1})\tilde{X}^{*P})^{-1} \sum_{t=1}^T x_t \Sigma^{-1} u_t^P \\ &= \beta^P + (X'(I \otimes \Sigma^{-1})\tilde{X}^{*P})^{-1} \sum_{t=1}^T x_t \Sigma^{-1} \{ (\Sigma - \tau_t^{*P}) \Sigma^{-1} x_t' (\beta^o - \beta^P) + \xi_t \} \end{aligned}$$

(from (3.21) by replacing $\hat{\beta}$ with β^P and evaluating at $\tilde{\beta} \in \beta^P \beta^o$)

$$= \beta^P + (X'(I \otimes \Sigma^{-1})\tilde{X}^{*P})^{-1} \sum_{t=1}^T x_t \Sigma^{-1} \{ (\Sigma - \tau_t^{*P} - o_P(1)) \Sigma^{-1} x_t' (\beta^o - \beta^P) + \xi_t \}$$

(where the identify matrix multiplying $o_P(1)$ is of order 2)

$$\begin{aligned} &= \beta^P + (X'(I \otimes \Sigma^{-1})\tilde{X}^{*P})^{-1}X'(I \otimes \Sigma^{-1})\{\tilde{X}^{*P}(\beta^o - \beta^P) + \xi + X[o_P(1)(\beta^P - \beta^o)]\} \\ &= \beta^o + (X'(I \otimes \Sigma^{-1})\tilde{X}^{*P})^{-1}X'(I \otimes \Sigma^{-1})\xi + \\ &\quad (X'(I \otimes \Sigma^{-1})\tilde{X}^{*P})^{-1}X'(I \otimes \Sigma^{-1})X[o_P(1)(\beta^P - \beta^o)] \end{aligned}$$

Hence

$$\begin{aligned} \sqrt{2T}(\bar{\beta} - \beta^o) &= \left(\frac{1}{2T}X'(I \otimes \Sigma^{-1})\tilde{X}^{*P}\right)^{-1}\frac{1}{\sqrt{2T}}X'(I \otimes \Sigma^{-1})\xi + \\ &\quad \left(\frac{1}{2T}X'(I \otimes \Sigma^{-1})\tilde{X}^{*P}\right)^{-1}\frac{1}{2T}X'(I \otimes \Sigma^{-1})X[o_P(1)\sqrt{2T}(\beta^P - \beta^o)] \end{aligned}$$

The \sqrt{T} -consistency of β^P and the convergence of $\frac{1}{2T}X'(I \otimes \Sigma^{-1})\tilde{X}^{*P}$ and $\frac{1}{2T}X'(I \otimes \Sigma^{-1})X$ to proper limits assures that the second term in the RHS above approaches zero in probability as $T \rightarrow \infty$, that is

$$\sqrt{2T}(\bar{\beta} - \beta^o) = \left(\frac{1}{2T}X'(I \otimes \Sigma^{-1})\tilde{X}^{*P}\right)^{-1}\frac{1}{\sqrt{2T}}X'(I \otimes \Sigma^{-1})\xi + o_P(1)$$

From (3.24) then

$$\sqrt{2T}(\hat{\beta} - \beta^o) - \sqrt{2T}(\bar{\beta} - \beta^o) = \sqrt{2T}(\hat{\beta} - \bar{\beta}) \xrightarrow{P} 0 \quad Q.E.D.$$

A direct implication of the above theorem is that $\bar{\beta}$ defined in (3.38) above is a GML estimator. To see this from another angle define

$$(3.39a) \quad Q^{*P} = X'(I \otimes \Sigma^{-1})\tilde{X}^{*P}$$

$$(3.39b) \quad Q^P = X'(I \otimes \Sigma^{-1})\tilde{X}^P$$

Then from (3.34) and (3.31a) to (3.31c) it follows that

$$(3.40) \quad R^P = \frac{1}{2T} Q^{*P} Q^{P-1} Q^{*P} \xrightarrow{P} \Psi_{11}(\vartheta^0) \text{ the information matrix.}$$

Now $\bar{\beta}$ can be expressed as

$$(3.41) \quad \bar{\beta} = \beta^P - R^{P-1} V^P$$

where

$$(3.42) \quad V^P = \frac{1}{2T} Q^{*P} Q^{P-1} X'(I \otimes \Sigma^{-1}) U^P.$$

The fact that $\bar{\beta}$ and β_{LML} are both efficient implies:

LEMMA:

Let α be the vector comprised of the distinct elements of Σ , $V(\cdot)$ as in (3.42), and β^P a \sqrt{T} -consistent estimate of β^0 . Then

$$(3.43) \quad \sqrt{T} [V(\beta^P, \alpha^0) - L_{\beta}(\beta^P, \alpha^0)] \xrightarrow{P} 0.$$

Proof is given in Appendix C. The above Lemma together with (3.40) imply that $\bar{\beta}$ is a GML estimator.

The critical point here is that, in the presence of a "good" initial estimate, efficiency is attainable via (single) IV regression on a suitably transformed data set, using filled-in values for the "missing" dependent variables.

Achieving "good" initial estimators is quite feasible (Amemiya [1973], Heckman [1976], Quandt and Ramsey [1978], Green [1981, 1983]) however, it is rarely the case that the covariance matrix is known *a priori*.

3.B The Case of Unknown Σ^0

It is assumed that "good" initial estimators β^P , α^P , of β^0 , α^0 are available where α is the vector of the identifiable parameters of Σ . That is

$$(3.44) \quad \alpha = (\sigma_1^2, \sigma_{01}/\sigma_2) \text{ a } 2 \times 1 \text{ vector, under model (I)}$$

$$(3.44') \quad \alpha = (\sigma_1^2, \sigma_{01}, \sigma_2^2)', \text{ a } 3 \times 1 \text{ vector, under model (II)}$$

The information matrix, partitioned according to β and α , is given by

$$(3.45) \quad \Psi = \begin{bmatrix} \Psi_{11} & \Psi_{12} \\ \Psi_{21} & \Psi_{22} \end{bmatrix}$$

where Ψ_{11} , $\Psi_{12} = \Psi_{21}'$, and Ψ_{22} are, respectively, $k \times k$, $k \times 2$ ($k \times 3$), and 2×2 (3×3) matrices in model (I) (model (II)). The GML estimator is defined as

$$(3.46) \quad \hat{\alpha}_{GML} = \begin{bmatrix} \beta_{GML} \\ \alpha_{GML} \end{bmatrix} = \begin{bmatrix} \beta^P \\ \alpha^P \end{bmatrix} - \begin{bmatrix} \Psi_{11}^P & \Psi_{12}^P \\ \Psi_{21}^P & \Psi_{22}^P \end{bmatrix}^{-1} \begin{bmatrix} \omega_\beta^P \\ \omega_\alpha^P \end{bmatrix}$$

where Ψ_{ij}^P is any consistent estimate of Ψ_{ij}^0 , $i, j = 1, 2$; and $(\omega_\beta^P, \omega_\alpha^P)'$ satisfy

$$\sqrt{T} [(\omega_\beta^P, \omega_\alpha^P) - (L_\beta^P, L_\alpha^P)] \xrightarrow{P} 0.$$

When using expressions for the inverse of a partitioned matrix, β_{GML} and α_{GML} can be written as

$$(3.47a) \quad \alpha_{GML} = \alpha^P - []^{P-1} \Psi_{21}^P \Psi_{11}^{P-1} \omega_\beta - []^{P-1} \omega_\alpha^P$$

$$(3.47b) \quad \beta_{GML} = \beta^P - \Psi_{11}^{P-1} \omega_\beta^P - \Psi_{11}^{P-1} \Psi_{12} []^{P-1} \Psi_{21}^P \Psi_{11}^{P-1} \omega_\beta^P + \Psi_{11}^{P-1} \Psi_{12} []^{P-1} \omega_\alpha^P$$

$$\text{where } []^P = \Psi_{22}^P - \Psi_{21}^P \Psi_{11}^{P-1} \Psi_{12}^P$$

Let us redefine \tilde{y}_t^P , \tilde{x}_t^P , and \tilde{z}_t^P in (3.35), (3.36a), and (3.36b) by substituting Σ^P for Σ . Also redefine τ_t of (3.30) or (3.30') by replacing $\Phi(-\mu_{2t}/\sigma_2)$ with $1 - q_t$. Finally define

$$(3.48) \quad \Delta^P = \beta^P - \bar{\beta} = R^{P-1} V^P$$

That R^P , defined in (3.40) with Σ^P replacing Σ , remains a consistent estimator of

Ψ_{11}^0 is not surprising. What is less obvious is that, under some weak conditions involving the existence of $\Psi(\cdot)$ and $\partial V(\cdot)/\partial \alpha$ at an appropriate vicinity of ϑ^0 , (3.43) implies

LEMMA:

$$(3.48) \quad \sqrt{T} (V(\beta^P, \alpha^P) - L_{\beta}(\beta^P, \alpha^P)) \xrightarrow{P} 0$$

Provided that β^P, α^P are \sqrt{T} -consistent estimates of β^0, α^0 . Proof is given in Appendix C.

The fact that $R^P \xrightarrow{P} \Psi_{11}^0$ together with (3.48) allows one to substitute $\Delta^P = R^{P-1} V^P$ for $\Psi_{11}^{P-1} \omega_{\beta}^P$ and R^P for Ψ_{11}^P in (3.47a) and (3.47b) and to arrive with the following expressions for ϑ_{GML} :

$$(3.49a) \quad \alpha_{GML} = \alpha^P + []^{P-1} \{ \Psi_{21}^P \Delta^P - \omega_{\alpha}^P \}$$

$$(3.49b) \quad \beta_{GML} = \bar{\beta} + R^{P-1} \Psi_{12}^P (\alpha^P - \alpha_{GML})$$

$\bar{\beta}, \Delta^P$ and R^P are constructed using the method outlined above. The other terms needed in the derivation of ϑ_{GML} are Ψ_{12}^P, Ψ_{22}^P , and ω_{α}^P . They can be recovered from the first derivatives of the log likelihood function (berndt *et al* [1974] which are given in Appendix B.

4. SUMMARY

Methods of Maximum Likelihood estimation of Limited Dependent Variables models are investigated in the context of the EM algorithm developed by Dempster *et al* [1977]. It is shown that this approach lends itself to an integrative framework of analysis that unifies various LDV models and link them to classical linear model theory. The computational consequences of this approach are: i) ML estimators can be obtained via iterative OLS or GLS procedure when the EM

algorithm is applied directly; ii) A unique method for deriving an asymptotically efficient single iteration estimator is achieved which produces an outcome identical to that of the standard (Hessian) method under single equation models but differs from those of other methods in multiple equation models.

The extension of this technique to models of more than two equations is straightforward but the use of this approach in simultaneous equations LDV models (eg., Amemiya [1974,1978], Heckman [1978], Lee [1981]) requires further research. Another topic for investigation is the practical question of convergence and computational efficiency in empirical examples. In Tsur [1983] I have estimated a 3-equation LDV model using the EM algorithm outlined above as well as two nonlinear ML algorithms. The results showed that EM was the only method which converged over the whole parameter space, including the identifiable covariance terms, whereas the ML algorithms broke down. However, further computational experience is required before a comprehensive evaluation of the role of the EM method in LDV models can be made.

APPENDIX A

Let $(u_1, u_2)'$ be $N\left\{0, \begin{bmatrix} \sigma_1^2 & \sigma_{12} \\ & \sigma_2^2 \end{bmatrix}\right\}$ and define $\tilde{u}_1 = u_1/\sigma_1$, $\tilde{u}_2 = u_2/\sigma_2$, and $\rho = \sigma_{12}/\sigma_1\sigma_2$. The density function of $\tilde{u}_2 | \tilde{u}_2 > a$ and $\tilde{u}_2 | \tilde{u}_2 \leq a$ are respectively $[1-\Phi(a)]^{-1}\varphi(x)$ for $x > a$ and $[\Phi(a)]^{-1}\varphi(y)$ for $y \leq a$.

$$\begin{aligned} (A.1) \quad E(\tilde{u}_2 | \tilde{u}_2 > a) &= [1-\Phi(a)]^{-1} \int_a^{\infty} x\varphi(x)dx \\ &= [1-\Phi(a)]^{-1} \int_a^{\infty} -\frac{d\varphi(x)}{dx} dx \\ &= [1-\Phi(a)]^{-1}\varphi(a) \stackrel{def}{=} \lambda^+(a) \end{aligned}$$

In the same manner

$$\begin{aligned} (A.2) \quad E(\tilde{u}_2 | \tilde{u}_2 \leq a) &= [\Phi(a)]^{-1} \int_{-\infty}^a x\varphi(x)dx \\ &= -[\Phi(a)]^{-1}\varphi(a) \stackrel{def}{=} \lambda^-(a) \end{aligned}$$

$$E(\tilde{u}_2^2 | \tilde{u}_2 > a) = [1-\Phi(a)]^{-1} \int_a^{\infty} x^2\varphi(x)dx$$

(integrating by parts)

$$= 1 + a\lambda^+(a)$$

Hence

$$(A.3) \quad var(\tilde{u}_2 | \tilde{u}_2 > a) = 1 + a\lambda^+(a) - [\lambda^+(a)]^2 = 1 + \Lambda^+(a)$$

where

$$\Lambda^+(a) \stackrel{def}{=} \frac{-\partial\lambda^+(a)}{\partial a} = -\lambda^+(a)[-a + \lambda^+(a)]$$

Note that in the text $\alpha = -\mu_t / \sigma$ and $\Lambda_t^s = \partial \lambda_t^s / \partial \mu = \frac{1}{\sigma} \frac{\partial \lambda_t^s}{\partial (\mu_t / \sigma)}$ for $s = -$ and $+$.

In the same manner

$$(A.4) \quad \begin{aligned} \text{var}(\tilde{u}_2 \mid \tilde{u}_2 \leq \alpha) &= 1 + \alpha \lambda^-(\alpha) - [\lambda^-(\alpha)]^2 \\ &= 1 + \Lambda^-(\alpha) \end{aligned}$$

where

$$\Lambda^-(\alpha) \stackrel{\text{def}}{=} - \frac{\partial \lambda^-(\alpha)}{\partial \alpha} = \alpha \lambda^-(\alpha) - [\lambda^-(\alpha)]^2.$$

The variate \tilde{u}_1 can be expressed as

$$(A.5) \quad \tilde{u}_1 = \rho \tilde{u}_2 + (1 - \rho^2)^{1/2} z$$

where $z \sim N(0, 1)$ independent of \tilde{u}_2 . Therefore

$$(A.6) \quad E(\tilde{u}_1 \mid \tilde{u}_2 > \alpha) = \rho E(\tilde{u}_2 \mid \tilde{u}_2 > \alpha) = \rho \lambda^+(\alpha)$$

$$(A.7) \quad E(\tilde{u}_1 \mid \tilde{u}_2 \leq \alpha) = \rho E(\tilde{u}_2 \mid \tilde{u}_2 \leq \alpha) = \rho \lambda^-(\alpha)$$

Using (A.5)

$$\begin{aligned} E(\tilde{u}_1^2 \mid \tilde{u}_2 > \alpha) &= E\{\rho^2 \tilde{u}_2^2 + (1 - \rho^2) z^2 + \rho(1 - \rho^2)^{1/2} \tilde{u}_2 z \mid \tilde{u}_2 > \alpha\} \\ &= \rho^2 E(\tilde{u}_2^2 \mid \tilde{u}_2 > \alpha) + 1 - \rho^2 \end{aligned}$$

$$(A.8) \quad \begin{aligned} \text{var}(\tilde{u}_1 \mid \tilde{u}_2 > \alpha) &= \rho^2 E(\tilde{u}_2^2 \mid \tilde{u}_2 > \alpha) + 1 - \rho^2 - \rho^2 [\lambda^+(\alpha)]^2 \\ &= \rho^2 \text{var}(\tilde{u}_2 \mid \tilde{u}_2 > \alpha) + 1 - \rho^2 \\ &= \rho^2 (1 + \Lambda^+(\alpha)) + 1 - \rho^2 \\ &= 1 + \rho^2 \Lambda^+(\alpha) \end{aligned}$$

In the same manner

$$(A.9) \quad \text{var}(\tilde{u}_1 | \tilde{u}_2 \leq a) = 1 + \rho^2 \Lambda^-(a)$$

Using (A.5) again

$$\begin{aligned} E(\tilde{u}_1 \tilde{u}_2 | \tilde{u}_2 > a) &= E\{\rho \tilde{u}_2^2 + (1 - \rho^2)^{1/2} \tilde{u}_2 | \tilde{u}_2 > a\} \\ &= \rho[\text{var}(\tilde{u}_2 | \tilde{u}_2 > a) + [\lambda^+(a)]^2] \\ &= \rho(1 + \Lambda^+(a)) + \rho[\lambda^+(a)]^2 \end{aligned}$$

Hence

$$\begin{aligned} (A.10) \quad \text{cov}(\tilde{u}_1, \tilde{u}_2 | \tilde{u}_2 > a) &= \rho(1 + \Lambda^+(a)) + \rho[\lambda^+(a)]^2 - \\ &\quad \rho[\lambda^+(a)]^2 \\ &= \rho(1 + \Lambda^+(a)) \end{aligned}$$

In the same manner

$$(A.11) \quad \text{cov}(\tilde{u}_1, \tilde{u}_2 | \tilde{u}_2 \leq a) = \rho(1 + \Lambda^-(a))$$

The extension to the corresponding moments of $u_1 = \sigma_1 \tilde{u}_1$ and $u_2 = \sigma_2 \tilde{u}_2$ is straightforward.

APPENDIX B

This appendix provides likelihood functions and their first derivatives for the Tobit model and models (I) and (II) of the two-equation case.

The Tobit Model

The log of the likelihood of the t^{th} observation is denoted by $L_t(\vartheta)$ and is given by

$$(B.1) \quad L_t(\vartheta) = q_t \left\{ -\frac{1}{2} \log \sigma^2 - \frac{1}{2\sigma^2} (y_t - x_t' \beta)^2 \right\} + (1 - q_t) \log \Phi(-x_t' \beta / \sigma) + \text{constant}$$

$$(B.2) \quad \frac{\partial L_t(\vartheta)}{\partial \beta} = x_t \frac{1}{\sigma^2} \{ q_t (y_t - x_t' \beta) + (1 - q_t) \sigma \lambda_t^- \}$$

$$(B.3) \quad \frac{\partial L_t(\vartheta)}{\partial \sigma^2} = \frac{1}{2\sigma^4} \{ q_t [(y_t - x_t' \beta)^2 - \sigma^2] + (1 - q_t) \sigma \lambda_t^- x_t' \beta \}$$

Note from (B.2) that

$$\begin{aligned} \frac{\partial L(\vartheta^P)}{\partial \beta} &= \frac{1}{\sigma^{P^2}} \sum_{t=1}^T x_t \{ q_t (y_t - x_t' \beta^P) + (1 - q_t) \sigma^P \lambda_t^{-P} \} \\ &= \frac{1}{\sigma^{P^2}} \sum_{t=1}^T x_t u_t^P = \frac{1}{\sigma^{P^2}} X' u^P \end{aligned}$$

where

$$u_t^P = y_t^P - x_t' \beta^P; \quad u^P = (u_1^P, u_2^P, \dots, u_T^P)'$$

Likewise

$$\begin{aligned} \frac{\partial^2 L(\vartheta^P)}{\partial \beta \partial \beta'} &= -\frac{1}{\sigma^{P^2}} \sum_{t=1}^T x_t \{ q_t - (1 - q_t) \lambda_t^{-P} \} x_t' \\ &= -\frac{1}{\sigma^{P^2}} \sum_{t=1}^T x_t \left\{ 1 - \frac{1}{\sigma^{P^2} \tilde{\tau}_t^P} \right\} x_t' \\ &= -\frac{1}{\sigma^{P^2}} X' \tilde{X}^P \end{aligned}$$

where $\tilde{\tau}_t$ and \tilde{X}^P are defined in (2.18) and (2.29) respectively.

Two-Equation Model (I)

$$(B.4) \quad L_t(\vartheta) = q_t \left\{ -\frac{1}{2} \log \sigma_1^2 - \frac{1}{2\sigma_1^2} (y_{1t} - x_{1t}'\beta)^2 + \right. \\ \left. \log \Phi(g_t) \right\} + (1-q_t) \log \Phi(-x_{2t}'\beta) + \text{constant}$$

where

$$(B.5) \quad g_t = \left\{ x_{2t}'\beta + \frac{\sigma_{12}}{\sigma_1^2} (y_{1t} - x_{1t}'\beta) \right\} / (1-\rho^2)^{\frac{1}{2}}$$

and recall that σ_2 is normalized to unity.

$$(B.6) \quad \frac{\partial L_t(\vartheta)}{\partial \beta} = q_t \left\{ x_{1t}' \frac{1}{\sigma_1^2} (y_{1t} - x_{1t}'\beta) + \frac{h_t}{(1-\rho^2)^{\frac{1}{2}}} \left[x_{2t}' - \frac{\sigma_{12}}{\sigma_1^2} x_{1t}' \right] \right\} \\ + (1-q_t) \lambda_t^- x_{2t}'$$

where

$$(B.7) \quad h_t = \frac{\varphi(g_t)}{\Phi(g_t)}$$

$$(B.8) \quad \frac{\partial L_t(\vartheta)}{\partial \sigma_1^2} = q_t \frac{1}{2\sigma_1^4} \left\{ -\sigma_1^2 + (y_{1t} - x_{1t}'\beta)^2 - h_t \frac{2\sigma_{12}}{1-\rho^2} \right. \\ \left. [(y_{1t} - x_{1t}'\beta) + \frac{g_t}{2} (1-\rho^2)^{-\frac{1}{2}}] \right\}$$

$$(B.9) \quad \frac{\partial L_t(\vartheta)}{\partial \sigma_{12}} = q_t h_t \frac{1}{\sigma_1^2 (1-\rho^2)^{\frac{1}{2}}} [(y_{1t} - x_{1t}'\beta) + g_t \sigma_{12} (1-\rho^2)^{-\frac{1}{2}}]$$

Two-Equation Model (II)

$$(B.10) \quad L_t(\vartheta) = q_t \left\{ \frac{1}{2} \log |\Sigma^{-1}| - \frac{1}{2} (y_t - x_t'\beta)' \Sigma^{-1} (y_t - x_t'\beta) \right\} + \\ (1-q_t) \log \Phi(-x_{2t}'\beta / \sigma_2) + \text{constant}$$

$$(B.11) \quad \frac{\partial L_t(\vartheta)}{\partial \beta} = q_t \left\{ -2x_t' \Sigma^{-1} y_t + 2x_t' \Sigma^{-1} x_t'\beta \right\} - (1-q_t) \lambda_t^- \frac{x_{2t}'}{\sigma_2}$$

$$(B.12) \quad \frac{\partial L_t(\vartheta)}{\partial \sigma_1^2} = tr \left\{ \frac{\partial L_t(\vartheta)}{\partial \Sigma^{-1}} \frac{\partial \Sigma^{-1}}{\partial \sigma_1^2} \right\}$$

$$= tr \left\{ \frac{q_t}{2} [\Sigma - (y_t - x_t' \beta)(y_t - x_t' \eta)'] \frac{-1}{(\sigma_1^2 \sigma_2^2 - \sigma_{12}^2)^2} \begin{bmatrix} \sigma_2^{11} & -\sigma_{12} \sigma_2^2 \\ -\sigma_{12} \sigma_2^2 & -\sigma_{12}^2 \end{bmatrix} \right\}$$

$$(B.13) \quad \frac{\partial L_t(\vartheta)}{\partial \sigma_{12}} = tr \left\{ \frac{\partial L_t(\vartheta)}{\partial \Sigma^{-1}} \frac{\partial \Sigma^{-1}}{\partial \sigma_{12}} \right\}$$

$$= tr \left\{ \frac{q_t}{2} [\Sigma - (y_t - x_t' \beta)(y_t - x_t' \beta)'] \frac{1}{(\sigma_1^2 \sigma_2^2 - \sigma_{12}^2)^2} \begin{bmatrix} -2\sigma_{12} \sigma_2^2 & -\sigma_1^2 \sigma_2^2 - \sigma_{12}^2 \\ -\sigma_{12} \sigma_2^2 - \sigma_{12}^2 & -2\sigma_{12} \sigma_1^2 \end{bmatrix} \right\}$$

$$(B.14) \quad \frac{\partial L_t(\vartheta)}{\partial \sigma_2^2} = tr \left\{ q_t \frac{1}{2} [\Sigma - (y_t - x_t' \beta)(y_t - x_t' \beta)'] \frac{-1}{(\sigma_1^2 \sigma_2^2 - \sigma_{12}^2)^2} \begin{bmatrix} -\sigma_{12}^2 & -\sigma_{12} \sigma_1^2 \\ -\sigma_{12} \sigma_1^2 & \sigma_1^4 \end{bmatrix} \right\} +$$

$$\frac{(1-q_t)}{2} \sigma_2 \lambda_t - \frac{x_{2t}' \beta}{\sigma_2^4}$$

APPENDIX C

This appendix clarifies what is meant by *efficiency* in this paper and provides proofs for (3.43) and (3.48).

The concept of *efficiency of an estimator* takes on several forms in statistical literature. (See Rao [1973, pp. 346-351] and Hájek [1972] for a discussion of the different approaches.) One such form entails consistency and attainment of the Cramer-Rao lower bound. That is, an estimator $\hat{\vartheta}_T$ is (asymptotically) efficient if

$$(C.1) \quad P_{\vartheta^0} \{ \sqrt{T} (\hat{\vartheta}_T - \vartheta^0) \} \rightarrow N(0, \Psi^{-1}(\vartheta^0))$$

where $P_{\vartheta}\{x\}$ indicates the distribution law of x when ϑ holds and $\Psi(\cdot)$ is the

information matrix. Definition (C.1) suffers from some deficiencies as is demonstrated in Rao [1973, p. 347]. Hence, a slightly modified definition is often employed (e.g., Bickel [1982]) and will be used here:

Definition: An estimator $\hat{\vartheta}_T$ is (asymptotically) efficient if

$$(C.2) \quad P_{\vartheta^T} \{ \sqrt{T} (\hat{\vartheta}_T - \vartheta^T) \} \rightarrow N(0, \Psi^{-1}(\vartheta^0)) \text{ for any sequence } \vartheta^T$$

such that $\sqrt{T} (\vartheta^T - \vartheta^0)$ stay bounded.

Since in the text the ML estimators are shown to satisfy (C.1), the question arises as to whether they also satisfy (C.2), and are therefore efficient.

Le Cam [1960] provided a general answer to this question. He proved that, in general, (C.1) implies (C.2) if:

(i) $\hat{\vartheta}_T$ is a "right" root of the likelihood function in the sense that it maximizes the likelihood at a "small" neighborhood of ϑ^0 . That is, let

$$N_{\vartheta^0}(T) = \{ \vartheta^T \mid \vartheta^T = \vartheta^0 + \frac{h_T}{\sqrt{T}}, h_T \text{ is a bounded sequence} \} \quad \text{then}$$

$$\hat{\vartheta}_T = \underset{\vartheta \in N_{\vartheta^0}}{\text{MAX}}^{-1} \{ L(\vartheta) \}.$$

(ii) The Local Asymptotic Normality (LAN) assumption holds (Le Cam *op cit.*).

the LAN assumption implies:

$$\Lambda(\vartheta^0, h_T, T) \stackrel{\text{def}}{=} L\left(\vartheta^0 + \frac{h_T}{\sqrt{T}}\right) - L(\vartheta^0) = h'_T \Delta_T(\vartheta^0) - \frac{1}{2} h'_T \Psi(\vartheta^0) h_T + o_p(1)$$

where

$$P_{\vartheta^0} \{ \Delta_T(\vartheta^0) \} \rightarrow N\{ 0, \Psi(\vartheta^0) \} \quad \text{and } h_T \text{ stay bounded.}$$

A sketch of the proof is as follow:

maximizing $\Lambda(\cdot)$ over h_T (disregarding the residual term) to get

$$\hat{h}_T = \Psi^{-1}(\vartheta^0) \Delta_T(\vartheta^0)$$

from (i)

$$\hat{\vartheta}_T = \vartheta^0 + \frac{\hat{h}_T}{\sqrt{T}}$$

and from (ii)

$$COV_{\vartheta^0} \{ \sqrt{T} (\hat{\vartheta}_T - \vartheta^0), h'_T \Delta_T(\vartheta^0) \} \rightarrow h_T.$$

From (C.1) and LAN

$$P_{\vartheta^0} \{ \sqrt{T} (\hat{\vartheta}_T - \vartheta^0), h'_T \Delta_T(\vartheta^0) \} \rightarrow N \left\{ 0, \begin{bmatrix} \Psi^{-1}(\vartheta^0) & h_T \\ h'_T & \Psi(\vartheta^0) \end{bmatrix} \right\}$$

which implies, using theorem of Le Cam [1960], that

$$P_{\vartheta^T} \{ \sqrt{T} (\hat{\vartheta}_T - \vartheta^0) \} \rightarrow N \{ h_T, \Psi^{-1}(\vartheta^0) \}$$

where $\vartheta^T = \vartheta^0 + \frac{h_T}{\sqrt{T}}$, hence

$$P_{\vartheta^T} \{ \sqrt{T} (\hat{\vartheta}_T - \vartheta^T) \} \rightarrow N \{ 0, \Psi^{-1}(\vartheta^0) \}$$

For the models considered in the text the ML estimators satisfy (i) by construction, and the assumption that $\frac{1}{T} X'X$ approaches a proper limit assures that the LAN condition holds.

A proof of (3.43) and (3.48):

Let B, A be a partition of Θ so $\vartheta = (\beta, \alpha)$; $\beta \in B, \alpha \in A$, and partition $\Psi(\cdot)$

according to β and α ;

$$(C.3) \quad \Psi(\vartheta) = \begin{bmatrix} \Psi_{11}(\vartheta) & \Psi_{12}(\vartheta) \\ \Psi_{21}(\vartheta) & \Psi_{22}(\vartheta) \end{bmatrix}$$

It is assumed that:

(AC) $\Psi(\cdot)$ exists and is positive definite in a small vicinity N_{ϑ^0} of ϑ^0 .

Let us consider first the case of known α , so α is fixed at α^0 , and take two asymptotically efficient estimates β_j^1 and β_j^2 , that is

$$(C.4) \quad P_{(\beta^T, \alpha^0)} \{ \sqrt{T} (\beta_j^i - \beta^T) \} \rightarrow N(0, \Psi_{11}^{-1}(\vartheta^0)); \quad j=1,2$$

for any sequence β^T such that $\sqrt{T}(\beta^T - \beta^0)$ stays bounded. Then, under (C.4), the following holds:

LEMMA:

$$(C.5) \quad \sqrt{T}(\beta_j^1 - \beta_j^2) \xrightarrow{P} 0$$

PROOF:

See Theorem (4.1) of Hajek [1972].

Let $\beta_j^1 = \beta^P - \Psi_{11}^{P-1} L_\beta(\beta^P, \alpha^0)$, where β^P is \sqrt{T} -consistent estimate of β^0 , Ψ_{11}^P is a consistent estimate of $\Psi_{11}(\vartheta^0)$ and $L_\beta(\beta^P, \alpha^0) = \partial L(\beta^P, \alpha^0) / \partial \beta$ and let $\beta_j^2 = \beta^P - \Psi_{11}^{P-1} V(\beta^P, \alpha^0)$ where $V(\cdot)$ is a $k \times 1$ vector. If β_j^2 is asymptotically efficient then (C.5) implies

$$(C.6) \quad \sqrt{T}(L_\beta(\beta^P, \alpha^0) - V(\beta^P, \alpha^0)) \xrightarrow{P} 0.$$

We now extend (C.6) to the case of unknown α and prove:

THEOREM:

Let β^P, α^P be \sqrt{T} -consistent estimates of β^0, α^0 , and let $V(\beta, \alpha)$ be such that (C.6) holds at (β^P, α^0) and $V_\alpha(\cdot) = \frac{\partial V(\cdot)}{\partial \alpha}$ exists at $N_{\alpha^0} = N_{\alpha^0} \cap A$, finally assume that (Ac) holds. Then

$$(C.7) \quad \sqrt{T}[L_\beta(\beta^P, \alpha^P) - V(\beta^P, \alpha^P)] \xrightarrow{P} 0$$

PROOF:

Let us define

$$(C.8) \quad \Delta(\beta, \alpha) = L_\beta(\beta, \alpha) - V(\beta, \alpha)$$

and note that, from the \sqrt{T} -consistency, α^P can be expressed as

$$(C.9) \quad \alpha^P = \alpha^0 + \frac{1}{\sqrt{T}} S_T \quad \text{where } 0 < |S_T| < \infty$$

Hence

$$\begin{aligned} (C.10) \quad \sqrt{T} \Delta(\beta^P, \alpha^P) &= \sqrt{T} \Delta(\beta^P, \alpha^0) + \sqrt{T} \left[\Delta(\beta^P, \alpha^0 + \frac{S_T}{\sqrt{T}}) - \Delta(\beta^P, \alpha^0) \right] \\ &= \sqrt{T} \Delta(\beta^P, \alpha^0) + \\ &\quad \frac{\Delta[\beta^P, \alpha^0 + (1/\sqrt{T}) S_T] - \Delta(\beta^P, \alpha^0)}{S_T/\sqrt{T}} \sqrt{T} (\alpha^P - \alpha^0) \end{aligned}$$

As $T \rightarrow \infty$, $\sqrt{T} \Delta(\beta^P, \alpha^0) \xrightarrow{P} 0$ by (C.6) and the second term on the R.H.S. of (C.10) tends to $\Delta_\alpha(\beta^P, \alpha^0) \sqrt{T} (\alpha^P - \alpha^0)$. Since $\sqrt{T} (\alpha^P - \alpha^0)$ is bounded in probability, to prove (C.7) it is sufficient to show

$$(C.11) \quad \Delta_\alpha(\beta^P, \alpha^0) \xrightarrow{P} 0.$$

Let us return to the case of known α and fix α at $\alpha = \alpha' \in N_{\alpha^0} \cap A$. All the

results attained under $\alpha = \alpha^0$ are now followed for $\alpha = \alpha'$. Therefore the following can be concluded:

$$\left. \begin{array}{l} \sqrt{T} \Delta(\beta^P, \alpha^0) \xrightarrow{P} 0 \\ \sqrt{T} \Delta(\beta^P, \alpha) \xrightarrow{P} 0 \end{array} \right\} \implies \frac{\Delta(\beta^P, \alpha') - \Delta(\beta^P, \alpha^0)}{1/\sqrt{T}} \xrightarrow{P} 0 \text{ as } T \rightarrow \infty.$$

Now let $\alpha' \rightarrow \alpha^0$ to get (C.11). Q.E.D.

References

- Affi, A.A. and R.M.E. Elashoff; "Missing Observations in Multivariate Statistics--I. Review of the Literature," *Journal of the American Statistical Association*, 61(1966), 595-605.
- Amemiya, T.; "Regression Analysis when the Dependent Variable is Truncated Normal," *Econometrica*, 41(1973), 997-1017.
- _____; "Multivariate Regression and Simultaneous Equation Models when the Dependent Variables are Truncated Normal," *Econometrica*, 42(1974), 999-1012.
- _____; "The Estimation of Simultaneous Equation Generalized Probit Model," *Econometrica*, 46(1978), 1193-1206.
- Berndt, E.B., B. Hall, R. Hall, and J.A. Hausman; "Estimation and Inference in Nonlinear Structural Models," *Annals of Economics and Social Measurement*, 3(1974), 653-665.
- Bickel, P.J.; "On Adaptive Estimation," *Ann. Statist.*, 10(1982), 647-671.
- Cosslett, S.; "Maximum Likelihood Estimator for Choice-Based Samples," *Econometrica*, 49(1981), 1289-1316.
- Dempster, A.P., N.M. Laird, and D.B. Rubin; "Maximum Likelihood from Incomplete Data Via the EM Algorithm," *Journal of the Royal Statistical Society, Series B*, 39(1977), 1-22.
- Duncan, Gregory M.; "Formulation and Statistical Analysis of the Mixed, Continuous/Discrete Dependent Variable Model in Classical Production Theory," *Econometrica*, 48(1980), 839-852.
- Fair, R.C.; "A Note on the Computation of the Tobit Estimator," *Econometrica*, 45(1977), 1723-1727.

- Green, W.H.; "On the Asymptotic Bias of the Ordinary Least Squares Estimator of the Tobit Model," *Econometrica*, 49(1981), 505-513.
- _____; "Estimation of Limited Dependent Variable Models By Ordinary Least Squares and the Method of Moments," *Journal of Econometrics*, 21(1983), 195-212.
- Griliches, Zvi, B.H. Hall, J.A. Hausman; "Missing Data and Self-Selection in Large Panels," *Annales de INSEE*, 30-31(1978), 137-176.
- Hájek, J.; "Local Asymptotic Minimax and Admissibility in Estimation," *Sixth Berkeley Symp. Math. Statist. Prob.*, 1(1972), 175-194. Univ. of California Press, Berkeley.
- Hanemann, W.M., and Yacov Tsur; "Econometric Models of Discrete/Continuous Supply Decisions under Uncertainty," *Dept. of Agricultural and Resource Economics, Univ. of California, Berkeley, Giannini Foundation Working Paper, No. 195*, (1982).
- Hartley, M.J.; "The Tobit and Probit Models: Maximum Likelihood by Ordinary Least Squares," *Dept. of Economics, State Univ. of New York at Buffalo, Discussion Paper, Number 374*, (1976).
- Hausman, J.A., and D.A. Wise; "A Conditional Probit Model for Qualitative Choice: Discrete Data Decisions Recognizing Interdependence and Heterogeneous Preferences," *Econometrica*, 46(1978), 403-420.
- Heckman, J.; "The Common Structure of Statistical Models of Truncation, Sample Selection and Limited Dependent Variables and a Simple Estimator for Such Models," *The Annals of Economic and Social Measurement*, 5(1976), 475-492.
- _____; "Dummy Endogenous Variables in a Simultaneous Equation System," *Econometrica*, 46(1978), 931-961.

- Johnson, N., and S. Kotz; *Distribution in Statistics: Continuous Multivariate Distributions*, New York: John Wiley & Sons, 1972.
- Le Cam, L.; "Locally Asymptotically Normal Families of Distributions," *Univ. California Publ. Statist.* 3(1960), 27-98.
- _____; *Theorie Asymptotique de la Decision Statistique*. Les Presses de l'Universite de Montreal, 1969.
- Lee, Lung-Fei; "Unionism and Wage Rates: A Simultaneous Equations Model with Qualitative and Limited Dependent Variables," *International Economic Review*, 19(1978), 415-433.
- _____; "Identification and Estimation in Binary Choice Models with Limited (Censored) Dependent Variables," *Econometrica*, 47(1979), 977-996.
- _____; "Simultaneous Equations Models with Discrete and Censored Variables," in C.F. Manski and D. McFadden (eds.), *Structural Analysis of Discrete Data with Econometric Applications*, Cambridge: MIT Press, 1981.
- _____, and R.P. Trost; "Estimation of Some Limited Dependent Variable Models with Application to Housing Demand," *Journal of Econometrics*, 8(1978), 357-382.
- Quandt, R.E., and J.B. Ramsey; "Estimating Mixtures of Normal Distributions and Switching Regressions," *Journal of the American Statistical Association*, 73(1978), 730-738.
- Rao, C.R.; *Linear Statistical Inference and its Applications*, New York: Wiley, 1973.
- Rothenberg, T., and C. Leenders; "Efficient Estimation of Simultaneous Equation Systems," *Econometrica*, 32(1964), 57-76.
- Tsur, Y.; "The Formulation and Estimation of Discrete/Continuous Supply Models Under Uncertainty," Ph.D. dissertation, Department of Agricultural and

Resource Economics, University of California, Berkeley, (1983).

Wu, C.F.J.; "On the Convergence Properties of the EM Algorithm," *Ann. Statist.*,
11(1983), 95-103.

Zellner, A.; "An Efficient Method of Estimating Seemingly Unrelated Regressions
and Tests for aggregation Bias," *Journal of the American Statistical Associ-
ation*, 57(1962), 348-368.