

Estimation of Multivariate Sample Selection Models via a Parameter-Expanded Monte Carlo EM Algorithm

Phillip Li

Department of Economics, Office of the Comptroller of the Currency, Washington, DC, USA
Email: Phillip.Li@occ.treas.gov

Received 6 September 2014; revised 5 October 2014; accepted 2 November 2014

Copyright © 2014 by Phillip Li.

This work is licensed under the Creative Commons Attribution International License (CC BY).

<http://creativecommons.org/licenses/by/4.0/>



Open Access

Abstract

This paper develops a parameter-expanded Monte Carlo EM (PX-MCEM) algorithm to perform maximum likelihood estimation in a multivariate sample selection model. In contrast to the current methods of estimation, the proposed algorithm does not directly depend on the observed-data likelihood, the evaluation of which requires intractable multivariate integrations over normal densities. Moreover, the algorithm is simple to implement and involves only quantities that are easy to simulate or have closed form expressions.

Keywords

Multivariate Sample Selection, Heckman Correction, Incidental Truncation, Expectation Maximization

1. Introduction

Sample selection models, pioneered in [1]-[3], are indispensable to researchers who use observational data for statistical inference. Among the many variants of these types of models, there is a growing interest in multivariate sample selection models. These are used to model a system of two or more seemingly unrelated equations, where the outcome variable for each equation may be non-randomly missing or censored according to its own stochastic selection variable. Applications range from modeling systems of demand equations [4] [5] to household vehicle usage [6]-[8]. A common specification is to assume a correlated multivariate normal distribution underlying both the outcomes of interest and the latent variables in the system.

There are two dominant approaches in the current literature to estimate these models. One approach is to use

maximum likelihood (ML) estimation. However, as noted in the literature, a major hurdle in evaluating the likelihood is that it requires computations of multivariate integrals over normal densities, which do not generally have closed form solutions. [9] discusses the ML estimation of these models and proposes to use the popular Geweke, Hajivassiliou, and Keane (GHK) algorithm to approximate these integrals in a simulated ML framework. While this strategy works reasonably well, the GHK algorithm can be difficult to implement. Another popular approach is to use two-step estimation (see [10] for a survey). In general, there is a tradeoff in the statistical properties and the computational simplicity for these estimators. If efficiency and consistency are of primary concern, then ML estimation should be preferred over two-step estimation.

The objective of this paper is to develop a simple ML estimation algorithm for a commonly used multivariate sample selection model. In particular, this paper develops a parameter-expanded Monte Carlo expectation maximization (PX-MCEM) algorithm that differs from [9] in a few important ways. First, the PX-MCEM algorithm does not use the observed-data likelihood directly, so it avoids the aforementioned integrations. Second, the proposed iterative algorithm does not require the evaluations of gradients or Hessians, which become increasingly difficult to evaluate with more parameters and equations. Third, the algorithm is straightforward to implement. It only depends on quantities that are either easy to simulate or have closed form expressions. This last point is especially appealing when estimating the covariance matrix parameter since there are non-standard restrictions imposed onto it for identification.

This paper is organized as follows. The multivariate sample selection model (MSSM) is formulated in Section 2. Section 3 begins with a brief overview of the EM algorithm for the MSSM and continues with the development of the PX-MCEM algorithm. Methods to obtain the standard errors are discussed. Section 4 offers some concluding remarks.

2. Multivariate Sample Selection Model

The MSSM is

$$y_{i,j}^* = x'_{i,j} \beta_j + \epsilon_{i,j} \tag{1}$$

$$s_{i,j}^* = w'_{i,j} \gamma_j + v_{i,j} \tag{2}$$

$$s_{i,j} = \mathbb{I}(s_{i,j}^* > 0) \tag{3}$$

$$y_{i,j} = \begin{cases} y_{i,j}^* & \text{if } s_{i,j} = 1 \\ \text{missing} & \text{if } s_{i,j} = 0 \end{cases} \tag{4}$$

for observations $i = 1, \dots, N$, and equations $j = 1, \dots, J$. In the previous expressions, $y_{i,j}^*$ is the continuous outcome of interest for observation i and equation j . Using similar indexing notation, $s_{i,j}^*$ is the latent variable underlying the binary selection variable $s_{i,j} = \mathbb{I}(s_{i,j}^* > 0)$, where $\mathbb{I}(A)$ denotes an indicator function that equals 1 if event A is true and 0 otherwise. Sample selection is incorporated by assuming that $y_{i,j}^*$ is missing when $s_{i,j} = 0$. Otherwise, $y_{i,j}^*$ is observed and equal to $y_{i,j}$. For later use, define $s_i = (s_{i,1}, \dots, s_{i,J})'$ and $s_i^* = (s_{i,1}^*, \dots, s_{i,J}^*)'$, where the prime symbol in s_i , s_i^* , and in the rest of this paper is used to denote matrix transpose.

Furthermore, $x_{i,j}$ and $w_{i,j}$ are column vectors of exogenous covariates, and β_j and γ_j are conforming vectors of parameters. Define $\beta = (\beta_1', \beta_2', \dots, \beta_J')'$ and $\gamma = (\gamma_1', \gamma_2', \dots, \gamma_J')'$. For identification, $w_{i,j}$ must contain at least one exogenous covariate that does not overlap with $x_{i,j}$ (refer to [11] for these exclusion restrictions). The unobserved errors $\epsilon_i = (\epsilon_{i,1}, \epsilon_{i,2}, \dots, \epsilon_{i,J})'$ and $v_i = (v_{i,1}, v_{i,2}, \dots, v_{i,J})'$ are jointly distributed as a $2J$ -dimensional multivariate normal with a mean vector of zeros and an unknown covariance matrix of Ω .

Formally, $(\epsilon_i', v_i')' \stackrel{iid}{\sim} \mathcal{N}_{2J}(0, \Omega)$ with

$$\Omega = \begin{pmatrix} \Omega_{\epsilon\epsilon} & \Omega_{\epsilon\nu} \\ \Omega'_{\epsilon\nu} & \Omega_{\nu\nu} \end{pmatrix}. \quad (5)$$

The submatrix $\Omega_{\nu\nu}$ is restricted to be in correlation form to identify the parameters corresponding to the latent variables [9]. The other elements of Ω are restricted such that the matrix is symmetric and positive definite.

The covariates and binary selection variables are always observed. Without loss of generality, assume that the outcomes for any observation i are only missing for the first m_i equations, where $0 \leq m_i \leq J$. Define $y_{i,\text{obs}} = (y_{i,m_i+1}, \dots, y_{i,J})'$, and let $y_{\text{obs}} = \{y_{i,\text{obs}}, s_i\}_{i=1}^N$ denote the observed data. The observed-data likelihood derived from (1) through (5) is denoted as $f(y_{\text{obs}} | \beta, \gamma, \Omega)$. See [9] for an exact expression of this likelihood.

3. Estimation

3.1. Overview of the EM Algorithm

The PX-MCEM algorithm is based on the EM algorithm of [12]. The basic idea behind the EM algorithm is to first augment y_{obs} with a set of “missing data” y_{mis} such that the observed-data likelihood is preserved when the missing data are integrated out of the complete-data likelihood. Formally, the missing data must satisfy

$$f(y_{\text{obs}} | \beta, \gamma, \Omega) = \mathbb{E} \left[f(y_{\text{mis}}, y_{\text{obs}} | \beta, \gamma, \Omega) \right], \quad (6)$$

where $f(y_{\text{mis}}, y_{\text{obs}} | \beta, \gamma, \Omega)$ is the complete-data likelihood to be defined later.

The EM algorithm then proceeds iteratively between an expectation step (E-step) and a maximization step (M-step) as follows. In iteration $(t+1)$ of the algorithm, compute in the E-step

$$Q(\beta, \gamma, \Omega | \beta^{(t)}, \gamma^{(t)}, \Omega^{(t)}) = \mathbb{E} \left[\log(f(y_{\text{mis}}, y_{\text{obs}} | \beta, \gamma, \Omega)) \right], \quad (7)$$

where the expectation is taken with respect to the conditional predictive distribution for the missing data, $\pi(y_{\text{mis}}, y_{\text{obs}} | \beta^{(t)}, \gamma^{(t)}, \Omega^{(t)})$, and in the M-step, find

$$\arg \max_{\beta, \gamma, \Omega} Q(\beta, \gamma, \Omega | \beta^{(t)}, \gamma^{(t)}, \Omega^{(t)}). \quad (8)$$

Denote the maximal values as $\beta^{(t+1)}$, $\gamma^{(t+1)}$, and $\Omega^{(t+1)}$, and continue on with the algorithm until convergence. The final maximal values are at least local maxima of the observed-data likelihood function.

For the MSSM, y_{mis} consists of all the missing outcomes and latent variables. Specifically,

$y_{\text{mis}} = \{y_{i,\text{mis}}, s_i^*\}_{i=1}^N$, where $y_{i,\text{mis}} = (y_{i,1}^*, \dots, y_{i,m_i}^*)'$. Furthermore, denote $y_{i,\text{com}} = (y'_{i,\text{mis}}, y'_{i,\text{obs}}, s_i^{*'})'$ as the vector of complete data, X_i as a block-diagonal matrix with the rows of covariates corresponding to the elements of $y_{i,\text{com}}$ on its block diagonals, and $\theta = (\beta', \gamma')$. The complete-data likelihood for the MSSM is given by

$$f(y_{\text{mis}}, y_{\text{obs}} | \beta, \gamma, \Omega) = \prod_{i=1}^N f(y_{i,\text{com}} | \beta, \gamma, \Omega) p(s_i | y_{i,\text{com}}) \quad (9)$$

with $f(y_{i,\text{com}} | \beta, \gamma, \Omega) = \phi_{2J}(y_{i,\text{com}} | X_i \theta, \Omega)$ which is a density function for a $2J$ -dimensional multivariate normal with mean $X_i \theta$ and covariance Ω , and

$$p(s_i | y_{i,\text{com}}) = \prod_{j=1}^J \left\{ \mathbb{I}(s_{i,j} = 1) \mathbb{I}(s_{i,j}^* > 0) + \mathbb{I}(s_{i,j} = 0) \mathbb{I}(s_{i,j}^* \leq 0) \right\}. \quad (10)$$

Equation (10) is a degenerate density since conditioning on s_i^* in $y_{i,\text{com}}$ determines s_i from (3). Note that the observed-data likelihood from [9] is obtained when y_{mis} is integrated out of (9), hence the condition in (6) holds.

3.2. PX-MCEM Algorithm

The standard EM algorithm using (7) and (8) is difficult to implement for the MSSM as the E-step and M-step are intractable. The PX-MCEM algorithm addresses this issue by modifying the E-step in two ways and leads to an M-step that can be evaluated with closed form quantities. Stated succinctly, the PX-MCEM algorithm is as follows.

1. Initialize $\beta^{(0)}$, $\gamma^{(0)}$, $\Omega^{(0)}$, and the number of Gibbs sampling draws G .
- At iteration $t+1$:
2. Draw G sets of missing data, denoted by $y_{\text{mis}}^{(1)}, \dots, y_{\text{mis}}^{(G)}$, from $\pi(y_{\text{mis}} | y_{\text{obs}}, \beta^{(t)}, \gamma^{(t)}, \Omega^{(t)})$ using Gibbs sampling.
3. PX-MC E-step: Estimate $Q(\alpha, \delta, \Sigma | \beta^{(t)}, \gamma^{(t)}, \Omega^{(t)})$ as

$$Q_G(\alpha, \delta, \Sigma | \beta^{(t)}, \gamma^{(t)}, \Omega^{(t)}) = \frac{1}{G} \sum_{g=1}^G \log(f(y_{\text{mis}}^{(g)}, y_{\text{obs}} | \alpha, \delta, \Sigma)). \quad (11)$$
4. PX-MC M-step: Maximize $Q_G(\alpha, \delta, \Sigma | \beta^{(t)}, \gamma^{(t)}, \Omega^{(t)})$ with iterative generalized least squares (IGLS) to obtain the maximizing parameters $\alpha^{(t+1)}$, $\delta^{(t+1)}$, and $\Sigma^{(t+1)}$.
5. Reduction step: Apply reduction functions to $\alpha^{(t+1)}$, $\delta^{(t+1)}$, and $\Sigma^{(t+1)}$ to obtain $\beta^{(t+1)}$, $\gamma^{(t+1)}$, and $\Omega^{(t+1)}$.
6. Repeat Steps 2 through 5 until convergence. The converged values are the ML estimates $\hat{\beta}$, $\hat{\gamma}$, and $\hat{\Omega}$. Each step is described in more detail in the subsequent sections.

3.2.1. PX-MC E-Step

Following [13], the first modification is to expand the parameter space of the complete-data likelihood function from (β, γ, Ω) to (α, δ, Σ) . The expanded parameters play similar roles as the original parameters, however Σ is expanded into a standard covariance matrix without the correlation restrictions. The parameter-expanded complete-data likelihood function is

$$f(y_{\text{mis}}, y_{\text{obs}} | \alpha, \delta, \Sigma) = \prod_{i=1}^N f(y_{i,\text{com}} | \alpha, \delta, \Sigma) p(s_i | y_{i,\text{com}}) \quad (12)$$

with $f(y_{i,\text{com}} | \alpha, \delta, \Sigma) = \phi_{2J}(y_{i,\text{com}} | X_i \Theta, \Sigma)$, where $\Theta = (\alpha', \delta)'$, and $\alpha = (\alpha_1', \dots, \alpha_J)'$ and $\delta = (\delta_1', \dots, \delta_J)'$ are defined analogously to β and γ . The advantage of using (12) instead of (9) is that Σ is easier to work with in the PX-MC M-step.

Second, instead of computing $Q(\alpha, \delta, \Sigma | \beta^{(t)}, \gamma^{(t)}, \Omega^{(t)}) = \mathbb{E}[\log(f(y_{\text{mis}}, y_{\text{obs}} | \alpha, \delta, \Sigma))]$ analytically, it is approximated as (11) with Monte Carlo methods and Gibbs sampling. To draw from

$\pi(y_{\text{mis}} | y_{\text{obs}}, \beta^{(t)}, \gamma^{(t)}, \Omega^{(t)})$, simply draw $y_{i,\text{mis}}$ and s_i^* from the conditional distribution

$\pi(y_{i,\text{mis}}, s_i^* | y_{i,\text{obs}}, s_i, \beta^{(t)}, \gamma^{(t)}, \Omega^{(t)})$ for $i = 1, \dots, N$. From (9), we have that

$$\pi(y_{i,\text{mis}}, s_i^* | y_{i,\text{obs}}, s_i, \beta^{(t)}, \gamma^{(t)}, \Omega^{(t)}) \propto \phi_{2J}(y_{i,\text{com}} | X_i \theta^{(t)}, \Omega^{(t)}) p(s_i | y_{i,\text{com}}), \quad (13)$$

where $\theta^{(t)} = (\beta^{(t)'}, \gamma^{(t)'})'$. For the missing outcomes, it is easy to see from (13) that

$$y_{i,j}^* | y_{i,\text{mis}(-j)}, s_i^*, y_{i,\text{obs}}, s_i, \beta^{(t)}, \gamma^{(t)}, \Omega^{(t)} \sim \mathcal{N}_1(\mu_{i,j(-j)}, \sigma_{i,j(-j)}^2) \quad (14)$$

for $j = 1, \dots, m_i$, where $y_{i,\text{mis}(-j)}$ is equivalent to $y_{i,\text{mis}}$ with $y_{i,j}^*$ removed, and $\mu_{i,j(-j)}$ and $\sigma_{i,j(-j)}^2$ are respectively the conditional mean and variance of $y_{i,j}^*$ given all other elements in $y_{i,\text{com}}$ from

$$\phi_{2J}(y_{i,\text{com}} | X_i \theta^{(t)}, \Omega^{(t)}).$$

Similarly, for the latent variables,

$$s_{i,j}^* \mid s_{i,(-j)}^*, y_{i,mis}, y_{i,obs}, s_i, \beta^{(t)}, \gamma^{(t)}, \Omega^{(t)} \sim \mathcal{TN}_{B_{i,j}} \left(\lambda_{i,j|(-j)}, \omega_{i,j|(-j)}^2 \right) \tag{15}$$

for $j = 1, \dots, J$, where $\mathcal{TN}_A(a, b^2)$ denotes a univariate normal distribution with mean a and variance b^2 truncated to the region A . In (15), $s_{i,(-j)}^*$ is s_i^* with $s_{i,j}^*$ removed, $B_{i,j}$ is the interval $(-\infty, 0]$ if $s_{i,j} = 0$ and $(0, +\infty)$ otherwise, and $\lambda_{i,j|(-j)}$ and $\omega_{i,j|(-j)}^2$ are respectively the conditional mean and variance of $s_{i,j}^*$ given all other elements of $y_{i,com}$ from $\phi_{2J}(y_{i,com} \mid X_i \theta^{(t)}, \Omega^{(t)})$.

The Gibbs sampler recursively samples from the full conditional distributions in (14) and (15) in the usual way. After a sufficient burn-in period, the last G draws are used in (11).

3.2.2. PX-MC M-Step and Reduction Step

By recognizing that (11) is proportional to the log-likelihood function of a seemingly unrelated regression model with NG observations and $2J$ equations, the maximization can be performed with IGLS. IGLS utilizes the quantities

$$\tilde{\Theta} = \left(\sum_{g=1}^G \sum_{i=1}^N X_i' \tilde{\Sigma}^{-1} X_i \right)^{-1} \left(\sum_{g=1}^G \sum_{i=1}^N X_i' \tilde{\Sigma}^{-1} y_{i,com}^{(g)} \right) \tag{16}$$

and

$$\tilde{\Sigma} = \frac{1}{NG} \sum_{g=1}^G \sum_{i=1}^N \left(y_{i,com}^{(g)} - X_i \tilde{\Theta} \right) \left(y_{i,com}^{(g)} - X_i \tilde{\Theta} \right)', \tag{17}$$

where $y_{i,com}^{(g)}$ is equivalent to $y_{i,com}$ with $y_{i,mis}^{(g)}$ and $s_i^{*(g)}$. First evaluate (16) with $\tilde{\Sigma}^{-1}$ removed, which amounts to estimating Θ equation by equation, and then evaluate (17) based on $\tilde{\Theta}$. Proceed by iterating (16) and (17) recursively until convergence. Denote the converged values as $\alpha^{(t+1)}$, $\delta^{(t+1)}$, and $\Sigma^{(t+1)}$.

In the reduction step, set $\beta^{(t+1)} = \alpha^{(t+1)}$, $\gamma_j^{(t+1)} = \delta_j^{(t+1)} / d_{J+j}$ ($1 \leq j \leq J$), and $\Omega^{(t+1)} = D^{-1} \Sigma^{(t+1)} D^{-1}$, where $D = \text{diag}(1, \dots, 1, d_{J+1}, \dots, d_{2J})$ is a $2J \times 2J$ diagonal matrix with the first J diagonals equal to 1 and the remaining J diagonals equal to the square root of the last J diagonals of $\Sigma^{(t+1)}$. The previous transformations are referred to as the reduction functions, and they are needed because (12) is used instead of (9) in the algorithm [13].

3.3. Standard Errors

The observed information matrix is

$$-\mathbb{E} \left\{ \frac{\partial^2 \log(f(y_{mis}, y_{obs} \mid \beta, \gamma, \Omega))}{\partial \Psi \partial \Psi'} \right\} - \mathbb{V} \left\{ \frac{\partial \log(f(y_{mis}, y_{obs} \mid \beta, \gamma, \Omega))}{\partial \Psi} \right\}, \tag{18}$$

where $\Psi = (\beta', \gamma', \Xi')'$, and Ξ is a column vector denoting the unique elements in Ω . Evaluate (18) at the ML estimates, and take the expectation and variance with respect to $\pi(y_{mis} \mid y_{obs}, \hat{\beta}, \hat{\gamma}, \hat{\Omega})$. These moments are estimated by taking additional draws from the Gibbs sampler and constructing their Monte Carlo analogs. The standard errors are the square roots of the diagonals of the inverse estimated quantity in (18).

4. Concluding Remarks

A new and simple ML estimation algorithm is developed for multivariate sample selection models. Roughly speaking, the implementation of this algorithm only involves iteratively drawing sets of missing data from well-known distributions and using IGLS on the complete data, both of which are inexpensive to perform. By using parameter expansion and Monte Carlo methods, the algorithm only depends on quantities with closed form

expressions, even when estimating the covariance matrix parameter with correlation restrictions. This algorithm is readily extendable to other types of selection models, including extensions to various types of outcome and selection variables with an underlying normal structure, and modifications to time-series or panel data.

Acknowledgements

I would like to thank the referee, Alicia Lloro, Andrew Chang, Jonathan Cook, and Sibel Sirakaya for their helpful comments.

References

- [1] Heckman, J. (1974) Shadow Prices, Market Wages, and Labor Supply. *Econometrica*, **42**, 679-694. <http://dx.doi.org/10.2307/1913937>
- [2] Heckman, J. (1976) The Common Structure of Statistical Models of Truncation, Sample Selection and Limited Dependent Variables and a Simple Estimator for Such Models. *Annals of Economic and Social Measurement*, **5**, 475-492.
- [3] Heckman, J. (1979) Sample Selection Bias as a Specification Error. *Econometrica*, **47**, 153-161. <http://dx.doi.org/10.2307/1912352>
- [4] Su, S.J. and Yen, S.T. (2000) A Censored System of Cigarette and Alcohol Consumption. *Applied Economics*, **32**, 729-737. <http://dx.doi.org/10.1080/000368400322354>
- [5] Yen, S.T., Kan, K. and Su, S.J. (2002) Household Demand for Fats and Oils: Two-Step Estimation of a Censored Demand System. *Applied Economics*, **34**, 1799-1806. <http://dx.doi.org/10.1080/00036840210125008>
- [6] Hao, A.F. (2008) A Discrete-Continuous Model of Households' Vehicle Choice and Usage, with an Application to the Effects of Residential Density. *Transportation Research Part B: Methodological*, **42**, 736-758. <http://dx.doi.org/10.1016/j.trb.2008.01.004>
- [7] Li, P. (2011) Estimation of Sample Selection Models with Two Selection Mechanisms. *Computational Statistics & Data Analysis*, **55**, 1099-1108. <http://dx.doi.org/10.1016/j.csda.2010.09.006>
- [8] Li, P. and Rahman, M.A. (2011) Bayesian Analysis of Multivariate Sample Selection Models Using Gaussian Copulas. *Advances in Econometrics*, **27**, 269-288. [http://dx.doi.org/10.1108/S0731-9053\(2011\)000027A013](http://dx.doi.org/10.1108/S0731-9053(2011)000027A013)
- [9] Yen, S.T. (2005) A Multivariate Sample-Selection Model: Estimating Cigarette and Alcohol Demands with Zero Observations. *American Journal of Agricultural Economics*, **87**, 453-466. <http://dx.doi.org/10.1111/j.1467-8276.2005.00734.x>
- [10] Tauchmann, H. (2010) Consistency of Heckman-Type Two-Step Estimators for the Multivariate Sample-Selection Model. *Applied Economics*, **42**, 3895-3902. <http://dx.doi.org/10.1080/00036840802360179>
- [11] Puhani, P.A. (2000) The Heckman Correction for Sample Selection and Its Critique. *Journal of Economic Surveys*, **14**, 53-68. <http://dx.doi.org/10.1111/1467-6419.00104>
- [12] Dempster, A.P., Laird, N.M. and Rubin, D.B. (1977) Maximum Likelihood from Incomplete Data via the EM Algorithm. *Journal of the Royal Statistical Society*, **39**, 1-38. <http://dx.doi.org/10.2307/2984875>
- [13] Liu, C., Rubin, D.B. and Wu, Y.N. (1998) Parameter Expansion to Accelerate EM: The PX-EM Algorithm. *Biometrika*, **85**, 755-770. <http://dx.doi.org/10.1093/biomet/85.4.755>