

Simulation Based Inference in Econometrics: Motivation and Methods.

Steven Stern

January 19, 1999

1. Introduction

Over the last few years, major advances have occurred in the field of simulation. In particular, McFadden(1989) and Pakes and Pollard(1989) have developed simulation methods to simulate expected values of random functions and have shown how to use those simulators in econometric estimation routines. Important applications where these techniques have been used include patent renewal in Pakes(1986), retirement in Berkovec and Stern(1991), market entry in Berry(1992), dynamic programming problems in Hotz, et al.(1994), exchange rates in Bansal, et al.(1995), and automobile pricing in Berry, Levinsohn, and Pakes(1995). Also, for example, Geweke(1989), Chib(1993), and McCullough and Rossi(1994, 1996) have shown how to use simulation methods to solve previously unsolvable Bayesian econometrics problems. An and Liu(1996) and Diebold and Schuermann(1996) use simulation to solve initial conditions problems in survival models and ARCH models respectively that otherwise seem to have no tractable solution.

Simulation provides an attractive solution for dealing with problems of the following type: Let U be a random variable with density $f(\cdot)$, and let $h(U)$ be some function of U . Then

$$Eh(U) = \int h(u) f(u) du. \quad (1.1)$$

Most econometrics problems including all method of moments problems and many maximum likelihood problems require one to evaluate equation (1.1) as part of an estimation strategy for estimating a set of parameters θ . There are many cases

where $Eh(U)$ can not be evaluated analytically or even numerically with precision. But we usually can simulate $Eh(U)$ on a computer by drawing R “pseudorandom” variables from $f(\cdot)$, u^1, u^2, \dots, u^R , and then constructing

$$\hat{E}h(U) = \frac{1}{R} \sum_{r=1}^R h(u^r). \quad (1.2)$$

Equation (1.2) provides an unbiased simulator of $Eh(U)$ which, for most of the methods discussed later, is enough to provide consistent estimates (or estimates with small bias) of θ .

This chapter provides some examples to motivate the problem. The first example is the multinomial probit problem, the second is a problem with unobserved heterogeneity, and the third is a Monte Carlo experiment. Next, the chapter describes a set of simulators that improve upon the most naive simulator in equation (1.2). Improvement is in terms of variance reduction, increased smoothness, and reduced computation cost. Then the most common simulation estimators are described. Finally, it evaluates the performance of the various simulators and estimation methods.

1.1. Multinomial Probit

The first example is the multinomial probit problem. Consider a model where y_j^* is the value to a person of choosing choice j for $j = 1, 2, \dots, J$ (a person index is suppressed). For example, j might index whether to drive a car, ride in someone else’s car, take a bus, or take a train to get to work ($J = 4$); it might index whether to work full-time, part-time, or retire ($J = 3$); or it might index whether an elderly person lives independently, in a nursing home, with a family member, or with paid help ($J = 4$). It is assumed that the person chooses the choice j with the greatest value; j is chosen iff $y_j^* > y_k^*$ for all $k \neq j$. Furthermore, it is assumed that y_j^* is a linear function of a set of observed variables and an error:

$$y_j^* = X_j\beta + u_j, \quad j = 1, \dots, J. \quad (1.3)$$

Let $u = (u_1, u_2, \dots, u_J)'$ be the vector of errors, and assume that the covariance matrix of u is Ω . The errors sometimes represent variation in values due to unobserved variables, and sometimes they represent variation in β ’s across people. Let $y_j = 1$ if choice j is chosen; $y_j = 1$ iff $y_j^* > y_k^*$ for all $k \neq j$.

Usually in data, we observe the covariates X and $y = (y_1, y_2, \dots, y_J)'$ but not $y^* = (y_1^*, y_2^*, \dots, y_J^*)'$. In order to estimate β and Ω , we need to evaluate the

probability of observing y conditional on X or the moments of y conditional on X . First, note that, since y_j is binary,

$$\begin{aligned} E(y_j | X) &= \Pr[y_j = 1 | X] \\ &= \Pr[y_j^* > y_k^* \forall k \neq j | X]. \end{aligned} \quad (1.4)$$

If we assume that $u_j \sim iid$ Extreme Value, then the probability in equation (1.4) has the analytical form

$$\Pr[y_j = 1 | X] = \exp\{X_j\beta\} / \sum_k \exp\{X_k\beta\}. \quad (1.5)$$

Such a model is called multinomial logit. The problem with multinomial logit is that the independence assumption for the errors is very restrictive. One can read a large literature on the independence of irrelevant alternatives problem caused by the independence of errors assumption. See, for example, Anderson, De Palma, and Thisse(1992).

Alternatively, we could assume that $u \sim N[0, \Omega]$ where Ω can be written in terms of a small number of parameters. When we assume the error distribution is multivariate normal, the resulting choice probabilities are called multinomial probit. For this case, the parameters to estimate are $\theta = (\beta, \Omega)$.¹ The choice probabilities are

$$\Pr[y_j = 1 | X] = \int_{u_1} \cdots \int_{u_J} 1[X_j\beta + u_j > X_k\beta + u_k \forall k \neq j] dF(u | \Omega) \quad (1.6)$$

where $1[\bullet]$ is an indicator function equal to one if the condition inside is true and equal to zero otherwise and $F(u | \Omega)$ is the joint normal distribution of u with covariance matrix Ω (with individual elements ω_{jk}). Let $u_{jk}^* = u_k - u_j$ for all $k \neq j$, and let $u_j^* = (u_{j1}^*, u_{j2}^*, \dots, u_{jj-1}^*, u_{jj+1}^*, \dots, u_{jJ}^*)'$. Then the J -dimensional integral in equation (1.6) can be written as a $J-1$ -dimensional integral:

$$\Pr[y_j = 1 | X] = \int_{u_{1j}^*} \cdots \int_{u_{Jj}^*} 1[X_j\beta - X_k\beta > u_{jk}^* \forall k \neq j] dF^*(u_j^* | \Omega_j^*) \quad (1.7)$$

where $F^*(u_j^* | \Omega_j^*)$ is the joint normal distribution of $u_j^* : u_j^* \sim N[0, \Omega_j^*]$ where $\omega_{jkl}^* = E(u_k - u_j)(u_l - u_j) = \omega_{kl} - \omega_{kj} - \omega_{jl} + \omega_{jj}$ for each element ω_{jkl}^* of Ω_j^* . Equation (1.7) can be written as

$$\Pr[y_j = 1 | X] = \Pr[u_j^* < V_j] \quad (1.8)$$

¹Some restrictions are required for Ω for identification. See, for example, Bunch(1991).

where V_j is a vector with k 'th element equal to $V_{jk} = X_j\beta - X_k\beta$. Note that equation (1.8) can be written as $Eh(U)$ in equation (1.1) with $h(U) = 1 [X_j\beta - X_k\beta > u_{jk}^* \forall k \neq j]$, the integrand in equation (1.7).

In order to make progress in estimating θ , we need to be able to evaluate equation (1.8) for any Ω_j^* and any V_j . For example, the MLE of θ maximizes

$$\frac{1}{N} \sum_i y_{ij} \log \Pr [u_{ij}^* < V_{ij}] \quad (1.9)$$

where i indexes observations, $i = 1, 2, \dots, N$. If $J = 3$, then equation (1.8) involves evaluating a bivariate normal probability; most computers have library routines to perform such a calculation. If $J = 4$, then equation (1.8) involves a 3-dimension integral. One can evaluate such an integral using Gaussian quadrature (see Butler and Moffitt, 1982) or the numerical algorithm in Hausman and Wise(1978). But, if $J > 4$, numerical routines will be cumbersome and frequently imprecise.

Simulation provides an alternative method for evaluating equation (1.8). The simplest simulator of equation (1.8) is

$$\frac{1}{R} \sum_{r=1}^R 1 (u_j^{*r} < V_j) \quad (1.10)$$

where u_j^{*r} is an *iid* draw from $N[0, \Omega_j^*]$. Essentially, the simulator in equation (1.10) draws a random vector from the correct distribution and then checks whether that random vector satisfies the condition, $u_j^* < V_j$. The simulator in equation (1.10) is called a frequency simulator. It is unbiased and bounded between zero and one. But its derivative with respect to θ is either undefined or zero because the simulator is a step function; this characteristic makes it difficult to estimate θ and to compute the covariance matrix of $\hat{\theta}$. Also, especially when $\Pr [y_j = 1 | X]$ is small, the frequency simulator has a significant probability of equaling zero; since MLE requires evaluating $\log \Pr [y_j = 1 | X]$, this is a significant problem. The simulators discussed in Section 2 suggest ways to simulate $\Pr [y_j = 1 | X]$ with small variance, with derivatives, and in computationally efficient ways.

1.2. Unobserved Heterogeneity

The second example involves unobserved heterogeneity in a nonlinear model. Let y_{it} be a random count variable; i.e., $y_{it} = 0, 1, 2, \dots$, with $i = 1, 2, \dots, N$ and $t = 1, 2, \dots, T$. Assume that $y_{it} \sim \text{Poisson}(\lambda_{it})$:

$$f(y_{it} | \lambda_{it}) = \exp\{-\lambda_{it}\} \lambda_{it}^{y_{it}} / y_{it}! \quad (1.11)$$

and that

$$\log \lambda_{it} = X_{it}\beta + u_i + e_{it} \quad (1.12)$$

where $u_i \sim iidG(\cdot | \alpha_G)$, $G(\cdot | \alpha_G)$ is a specified distribution up to a set of parameters α_G ,

$$e_{it} = \rho e_{it-1} + \varepsilon_{it}, \quad (1.13)$$

$\varepsilon_{it} \sim iidH(\cdot | \alpha_H)$, and $H(\cdot | \alpha_H)$ is a specified distribution up to a set of parameters α_H .²For example, y_{it} might be the number of trips person i takes in period t , the number of patents firm i produces in year t , or the number of industrial accidents firm i has in year t . Adding the unobserved heterogeneity u_i and serially correlated error e_{it} allows for richness frequently necessary to explain the data. The goal is to estimate $\theta = (\beta, \rho, \alpha_G, \alpha_H)$. The log likelihood contribution of observation i is

$$L_i = \log \int_{u_i} \int_{\varepsilon_{i1}} \cdots \int_{\varepsilon_{iT}} \prod_{t=1}^T [\exp\{-\lambda_{it}\} \lambda_{it}^{y_{it}} / y_{it}! dH(\varepsilon_{it} | \alpha_H)] dG(u_i | \alpha_G) \quad (1.14)$$

where $\lambda_i = (\lambda_{i1}, \lambda_{i2}, \dots, \lambda_{iT})'$ depends upon $X_{it}\beta$, u_i , and $\varepsilon_i = (\varepsilon_{i1}, \varepsilon_{i2}, \dots, \varepsilon_{iT})'$ through equations (1.12) and (1.13). When there is no serial correlation term e_{it} , the integral in equation (1.14) can be solved analytically for well chosen $G(u_i | \alpha_G)$.³ But for general $G(\cdot | \alpha_G)$ and $H(\cdot | \alpha_H)$, the integral can be evaluated neither analytically nor numerically.

Simulating the integral is quite straightforward. Let ε_i^r be an *iid* pseudorandom draw of ε_i , $r = 1, 2, \dots, R$. Similarly, let u_i^r be an *iid* random draw of u_i , $r = 1, 2, \dots, R$. Then L_i can be simulated by evaluating the integrand for each draw r and taking an average:

$$\hat{L}_i = \log \left\{ \frac{1}{R} \sum_{r=1}^R \left[\prod_{t=1}^T \exp\{-\lambda_{it}^r\} (\lambda_{it}^r)^{y_{it}} / y_{it}! \right] \right\} \quad (1.15)$$

where λ_{it}^r is evaluated using the pseudorandom draws of ε_i and u_i in equation (1.12). The maximum simulated likelihood estimator of θ maximizes $\sum_i \hat{L}_i$. Note

²One might want to specify a different distribution for e_{i0} because of an initial conditions problem.

³See Hausman, Hall, and Griliches(1984).

that even though $\exp\{\hat{L}_i\}$ is unbiased, \hat{L}_i is biased for finite R (because \hat{L}_i is a nonlinear function of $\exp\{\hat{L}_i\}$). This will cause $\hat{\theta}$ to be inconsistent unless $R \rightarrow \infty$ as $NT \rightarrow \infty$. However, Monte Carlo results discussed later show that the asymptotic bias is small as long as “good” simulators are used.

1.3. Monte Carlo Experiments

The last example is a Monte Carlo experiment. Let U be a vector of data and $s(U)$ be a proposed statistic that depends upon U . The statistic $s(U)$ may be an estimator or a test statistic. In general, the user will want to know the distribution of $s(U)$. But, for many statistics $s(\cdot)$, deriving the small sample properties of $s(U)$ is not possible analytically. Simulation can be used to learn about the small sample properties of $s(U)$. All moments of $s(U)$ can be written in the form $Eh(U)$.⁴ Medians and, in fact, the whole distribution of $s(U)$ can be written in the form $Eh(U)$. Monte Carlo experiments are powerful tools to use in evaluating statistical properties of $s(U)$. However care must be taken in conducting such experiments. In particular, one must be careful in generalizing Monte Carlo results to cases not actually simulated; a Monte Carlo experiment really only provides information about the specific case simulated. Also, one must be careful not to attempt simulating objects that do not exist. For example, simulating the expected value of a two stage least squares (2SLS) estimator of a just identified equation would provide an answer (because any particular draw of $s(U)$ is finite) but it would be meaningless because 2SLS estimators of just identified equations have no finite moments. See Hendry(1984) for more on Monte Carlo experiments.

2. Simulators

This section discusses various simulation methods. Throughout, the goal will be to simulate $Eh(U)$ or, in some special cases, $\Pr[y_j = 1 | X]$. The first requirement of a simulation method is to simulate U from its distribution F . In general, if $Z \sim \text{Uniform}(0, 1)$, then $F^{-1}(Z) \sim F$.⁵ For example, the exponential distribution is $F(x) = 1 - \exp\{-\lambda x\}$. Thus, $-\log(1 - Z)/\lambda \sim F$. If F is standard normal, then F^{-1} has no closed form, but most computers have a library routine to approximate

⁴For $Es(U)$, $h(U) = s(U)$, and for $\text{Var}[s(U)]$, $h(U) = [s(U) - Es(U)]^2$.

⁵Most computers have a library routine to generate standard uniform random variables. See, for example, Ripley (1987) for a discussion of standard uniform random number generators.

F^{-1} for the standard normal distribution. Truncated random variables can be simulated in the same way. For example, assume $U \sim N[\mu, \sigma^2]$ but let it be truncated between a and b . Then, since

$$F(u) = \left[\Phi\left(\frac{u-\mu}{\sigma}\right) - \Phi\left(\frac{a-\mu}{\sigma}\right) \right] / \left[\Phi\left(\frac{b-\mu}{\sigma}\right) - \Phi\left(\frac{a-\mu}{\sigma}\right) \right] \quad (2.1)$$

where Φ is the standard normal distribution function, U can be simulated by letting $F(u) = Z$ in equation (2.1) and solving equation (2.1) for u as

$$\sigma\Phi^{-1}\left\{Z\left[\Phi\left(\frac{b-\mu}{\sigma}\right) - \Phi\left(\frac{a-\mu}{\sigma}\right)\right] + \Phi\left(\frac{a-\mu}{\sigma}\right)\right\} + \mu. \quad (2.2)$$

This idea can be applied with a small twist to discrete random variables. Assume $U = i$ with probability p_i for $i = 1, 2, \dots, n$. Let $P_i = \Pr[U \leq i] = \sum_{j=1}^i p_j$. Let $Z \sim \text{Uniform}(0, 1)$, and let $U = i$ iff $P_{i-1} < Z \leq P_i$ (where $P_0 = 0$). Then U is distributed as desired.

Random variables frequently can be simulated by using a composition formula. For example, since a binomial random variable is the sum of independent Bernoulli random variables, we can simulate a binomial random variable by simulating independent Bernoulli's and then adding them up. A more useful example is simulating multivariate $U \sim N[\mu, \Omega]$. Let $Z \sim N[0, I]$, and let C be any matrix such that $CC' = \Omega$ (e.g., the Cholesky decomposition of Ω). Then it is easy to verify that $CZ + \mu \sim N[\mu, \Omega]$. So we can simulate U by simulating Z and then transforming it.

In some cases, it will be necessary to simulate a random variable conditional on some event where the inverse conditional distribution has no analytical form (or good approximation). There are a number of acceptance-rejection methods available for many such cases. Assume (U, Z) have joint distribution $F(u, z)$ and that it is straightforward to draw (U, Z) from its joint distribution. Further, assume we want to draw U conditional on $Z \in S$ where S is a subset of the support of Z . The simplest acceptance-rejection simulation method is:

- (a) Draw (U, Z) from F .
- (b) If $Z \notin S$, go to (a).
- (c) If $Z \in S$, keep.

There are more sophisticated methods that reduce the expected number of draws of (U, Z) needed (see, for example, Devroye(1986), Ripley(1987), or Tierney(1994)), but all acceptance-rejection simulation methods suffer from a) the potentially large

number of draws needed and b) the lack of differentiability of $Eh(U)$ with respect to parameter vector θ .⁶ Thus, for the most part, they should be avoided. For the remainder of the chapter, it will be assumed one can simulate U .

The most straightforward simulator for $Eh(U)$ is

$$\hat{E}h(U) = \frac{1}{R} \sum_{r=1}^R h(u^r) \quad (2.3)$$

where u^r , $r = 1, 2, \dots, R$, are R *iid* pseudorandom draws of U . When simulating $\Pr[y_j = 1 \mid X]$, equation (2.3) becomes equation (1.10). If h is continuous and differentiable with respect to θ , then $\hat{E}h(U)$ will be continuous and differentiable. Equation (2.3) is unbiased, and its variance is $\text{Var}[h(U)]/R$. Note that as $R \rightarrow \infty$, the variance of the simulator \rightarrow zero.

2.1. Importance Sampling

Several methods allow us to improve the performance of a simulator significantly either in terms of reduced variance, better smoothness properties, and/or better computation time properties. For example, the multinomial probit problem in Borsch-Supan, et al.(1992) and the production function estimation problem in Ohanian, et al.(1996) work with simulation only with the use of good importance sampling simulators. The rest of this section describes the most popular simulation methods. The first method is importance sampling. Consider $Eh(U)$ in equation (1.1) where it is either difficult to draw U from F or where h is not smooth. In some cases, one can rewrite equation (1.1) as

$$Eh(U) = \int \frac{h(u) f(u)}{g(u)} g(u) du \quad (2.4)$$

where $g(u)$ is a density with the following properties:

- a) it is easy to draw U from g ,
- b) f and g have the same support,
- c) it is easy to evaluate $h(u) f(u) / g(u)$ given u , and
- d) $h(u) f(u) / g(u)$ is bounded and smooth over the support of U .

⁶Differentiability is important for most estimation procedures. An exception is Gibbs sampling or, more generally, Monte Carlo Markov Chain estimation methods.

Note that equation (2.4) is $E[h(U) f(U) / g(U)]$ where $U \sim g$. Then the importance sampling simulator for $Eh(U)$ is

$$\hat{E}h(U) = \frac{1}{R} \sum_{r=1}^R \frac{h(u^r) f(u^r)}{g(u^r)} \quad (2.5)$$

where u^r , $r = 1, 2, \dots, R$ are R iid draws from g . The purpose of conditions (a) and (c) are to increase computational speed. The purpose of condition (d) is variance bounding and smoothness.

Consider simulating $\Pr[y_j = 1 \mid X]$ for the multinomial probit problem. Equation (1.8) can be written as

$$\int_{u_j^* < V_j} f(u_j^*) du_j^* = \int_{u_j^* < V_j} [f(u_j^*) / g(u_j^*)] g(u_j^*) du_j^* \quad (2.6)$$

for some multivariate density g satisfying Conditions (a) through (d). Consider g where the k th element of u_j^* is distributed independently truncated normal with upper truncation point V_{jk} and variance Ω_{jkk}^* for each k . The candidate g satisfies Conditions (a), (b), and (c), and $h(u) f(u) / g(u)$ is smooth over the support $u_j^* < V_j$. But $h(u) f(u) / g(u)$ is not bounded especially when Ω_j^* has large off-diagonal terms. Thus, this choice of g may be problematic. In fact, in general it is the boundedness condition that is difficult to satisfy. For the multinomial probit problem, the Geweke-Keane-Hajivassiliou (GHK) and decomposition simulators discussed below both can be thought of as importance sampling simulators that satisfy Conditions (a) through (d). The simulators described in Danielsson and Richard(1993) and Richard and Zhang(1996) are more sophisticated importance sampling simulators.

2.2. GHK Simulator

The GHK simulator, developed by Geweke(1991), Hajivassiliou(1990), and Keane(1994), has been found to perform very well in Monte Carlo studies (discussed later) for simulating $\Pr[u_j^* < V_j]$. The GHK algorithm switches back and forth between computing univariate, truncated normal probabilities, simulating draws from univariate normal distributions, and computing normal distributions conditional on previously drawn truncated normal random variables. Since each step is straightforward and fast, the algorithm can decompose the more difficult problem into a series of feasible steps. The algorithm is as follows:

- (a) Set $t = 1$, $\mu = 0$, $\sigma^2 = \Omega_{jtt}^*$, and $\hat{P} = 1$.
- (b) Compute $p = \Pr(u_{jt}^* < V_{jt})$ analytically, and increment $\hat{P} = \hat{P} * p$.
- (c) Draw u_{jt}^* from a truncated normal distribution with mean μ , variance σ^2 , and upper truncation point V_{jt} .
- (d) If $t < J - 1$, increment t by 1; otherwise goto (g).
- (e) Compute (analytically) the distribution of u_{jt}^* conditional on $u_{j1}^*, u_{j2}^*, \dots, u_{jt-1}^*$. Note that this is normal with an analytically computable mean vector μ and variance σ^2 .
- (f) Goto (b).
- (g) \hat{P} is the simulator.

The algorithm relies upon the fact that normal random variables conditional on other normal random variables are still normal. The GHK simulator is strictly bounded between zero and one because each increment to \hat{P} is strictly bounded between zero and one. It is continuous and differentiable in θ because each increment to \hat{P} is continuous and differentiable. Its variance is smaller than the frequency simulator in equation (1.10) because each draw of \hat{P} is strictly bounded between zero and one while each draw of the frequency simulator is either zero or one.

The GHK simulator is an importance sampling simulator. Consider the case where $J = 3$. Then the probability to simulate can be written as

$$\Pr[u < V] = \Phi(V_1, V_2) \int_{-\infty}^{V_1} \int_{-\infty}^{V_2} \Phi(V_3 | u_1, u_2) \frac{\phi(u_2 | u_1) \phi(u_1)}{\Phi(V_1, V_2)} du_2 du_1 \quad (2.7)$$

where $\Phi(V_1, V_2) = \Pr[u_1 < V_1, u_2 < V_2]$, $\Phi(V_3 | u_1, u_2) = \Pr[u_3 < V_3 | u_1, u_2]$, $\phi(u_2 | u_1)$ is the conditional density of u_2 given u_1 , $\phi(u_1)$ is the marginal density of u_1 . Equation (2.7) can be written in the form of equation (2.4) by letting

$$\begin{aligned} h(u) &= \Phi(V_1, V_2) \Phi(V_3 | u_1, u_2), \\ f(u) &= \frac{\phi(u_2 | u_1) \phi(u_1)}{\Phi(V_1, V_2)} 1(u_1 < V_1, u_2 < V_2), \\ g(u) &= \frac{\phi(u_1) \phi(u_2 | u_1)}{\Phi(V_1) \Phi(V_2 | u_1)} \end{aligned} \quad (2.8)$$

where $g(u)$ reflects the GHK algorithm's method of simulation. Because it is an importance sampling simulator, GHK is unbiased.

A minor modification of the algorithm provides draws of normal random variables u_j^* conditional on $u_j^* \leq V_j$. Other minor modifications are useful for related problems.

2.3. Decomposition Simulators

Next, two decomposition simulators are described. The Stern(1992) simulator uses the property that the sum of two normal random vectors is also normal. The goal is to simulate $\Pr [u_j^* < V_j]$. Decompose $u_j^* = Z_1 + Z_2$ where $Z_1 \sim N [0, \lambda]$, $Z_2 \sim N [0, \Omega_j^* - \lambda]$, Z_1 and Z_2 are independent, and λ is chosen to be a diagonal matrix as large as possible such that $\Omega_j^* - \lambda$ is positive definite.⁷ Then equation (1.8) can be written as

$$\begin{aligned} & \int \Pr [Z_1 < V_j - z_2] g(z_2) dz_2 \\ &= \int \prod_k \Phi \left(\frac{V_{jk} - z_{2k}}{\lambda_k} \right) g(z_2) dz_2 \end{aligned} \quad (2.9)$$

where $g(\cdot)$ is the joint normal density of Z_2 . Equation (2.9) can be simulated as

$$\frac{1}{R} \sum_{r=1}^R \prod_k \Phi \left(\frac{V_{jk} - z_{2k}^r}{\lambda_k} \right) \quad (2.10)$$

where z_{2k}^r , $k = 1, 2, \dots, J-1$, are pseudorandom draws of Z_2 . The Stern simulator has all of the properties of the GHK simulator. So which one performs better is an empirical matter left to later discussion.

Another decomposition simulator, suggested by McFadden(1989), changes the specification of equation (1.3) to

$$y_j^* = X_j \beta + u_j + \tau e_j, \quad j = 1, \dots, J \quad (2.11)$$

where τ is a small number and $e_j \sim iid$ Extreme Value. In the limit, as $\tau \rightarrow 0$, $\Pr [y_j = 1 | X]$ converges to a multinomial probit probability. But for any $\tau > 0$,

$$\Pr [y_j = 1 | X] = \int \exp \left\{ \frac{X_j \beta + u_j}{\tau} \right\} / \sum_k \exp \left\{ \frac{X_k \beta + u_k}{\tau} \right\} f(u) du \quad (2.12)$$

which is the multinomial logit probability conditional on $u = (u_1, u_2, \dots, u_J)$ integrated over f . Equation (2.12) can be simulated as

$$\frac{1}{R} \sum_{r=1}^R \left[\exp \left\{ \frac{X_j \beta + u_j^r}{\tau} \right\} / \sum_k \exp \left\{ \frac{X_k \beta + u_k^r}{\tau} \right\} \right] \quad (2.13)$$

⁷An easy way to pick λ is to set each diagonal element of λ equal to the smallest eigenvalue of Ω_j^* minus a small amount.

where u^r are pseudorandom draws of u . The idea in McFadden(1989) is to think of equation (2.11) as a kernel-type approximation of equation (1.3) for small τ . However, assuming equation (2.11) is the true structure (where τ is a parameter that can sometimes be estimated) takes away no flexibility and frequently eases simulation. Multivariate normality is a desirable assumption because of its flexible covariance matrix. But there are very few applications where theory dictates that the error in equation (1.3) should be multivariate normal. Berkovec and Stern(1991) and Berry, Levinsohn, and Pakes(1995) use the McFadden specification as the “true” specification in a structural model of retirement behavior.

2.4. Antithetic Acceleration

Antithetic acceleration is a powerful variance reduction method (see Geweke, 1988). In any simulation method, there is some probability that the pseudorandom draws will be unusually large (or small). Antithetic acceleration prevents such events from occurring and thus reduces the variance of the simulator. Consider the general problem of simulating $Eh(U)$ where $U \sim F$. Let $Z \sim \text{Uniform}(0, 1)$. Then $h(F^{-1}(Z))$ is a simulator of $Eh(U)$. But $h(F^{-1}(1 - Z))$ is also a simulator of $Eh(U)$ (because $1 - Z \sim \text{Uniform}(0, 1)$ also). The antithetic acceleration simulator of $Eh(U)$ is

$$\hat{E}h(u) = \frac{1}{2R} \sum_{r=1}^R [h(F^{-1}(z^r)) + h(F^{-1}(1 - z^r))] \quad (2.14)$$

where z^r is a pseudorandom draw of Z . When F is $N[0, \sigma^2]$, equation (2.14) becomes

$$\hat{E}h(u) = \frac{1}{2R} \sum_{r=1}^R [h(u^r) + h(-u^r)] \quad (2.15)$$

where u^r is a pseudorandom draw of U . For any symmetric F , if h is linear, the variance of $\hat{E}h(U)$ is zero. For monotone h , the variance of $\hat{E}h(U)$ with R draws and antithetic acceleration is smaller than the variance of $\hat{E}h(U)$ with $2R$ draws and no antithetic acceleration. If $Eh(U)$ is being simulated to estimate a parameter θ with N observations and h is monotone, then the increase in $\text{Var}(\hat{\theta})$ due to simulation when antithetic acceleration is used is of order $(1/N)$ times the increase in $\text{Var}(\hat{\theta})$ due to simulation when antithetic acceleration is not used. The value of this is discussed more in the next section.

There are simulation problems where antithetic acceleration does not help. For example, let $U \sim N[0, \sigma^2]$, and let $h(U) = U^2$. Then $\text{Var}[\hat{E}h(U)]$ with antithetic

acceleration and R draws is greater than that without antithetic acceleration and $2R$ draws. This is because $h(-U) = h(U)$ which means that equation (2.15) becomes equation (2.3); the variance is twice as great as with no antithetic acceleration and $2R$ draws. In general, deviations from monotone h will diminish the performance of antithetic acceleration. But Hammersly and Handscomb(1964) suggests generalizations of antithetic acceleration that will reduce variance for more general h .

A related method is the use of control variates. Let $\hat{E}h(U)$ be a simulator of $Eh(U)$, and let $\hat{k}(U)$ be some other simulator with known expected value $Ek(U)$. Then

$$\tilde{E}h(U) = \hat{E}h(U) - \hat{k}(U) + Ek(U) \quad (2.16)$$

has expected value $Eh(U)$ and variance

$$Var[\tilde{E}h(U)] = Var[\hat{E}h(U)] + Var[\hat{k}(U)] - 2Cov[\hat{E}h(U), \hat{k}(U)]. \quad (2.17)$$

If $Cov[\hat{E}h(U), \hat{k}(U)] > Var[\hat{k}(U)]/2$, then $Var[\tilde{E}h(U)] < Var[\hat{E}h(U)]$. This idea can be used effectively in Monte Carlo testing and covariance matrix estimation (where it is easy to find a simulator $\hat{k}(U)$ with known expected value). In fact, it can be used to increase the rate of convergence of such estimators (see, for example, Hendry 1984 or Brown and Newey 1996).

3. Estimation Methods

The goal of this section is to use the simulators developed in the last section in some estimation problems. Four different estimation methods are discussed: method of simulated moments (MSM), maximum simulated likelihood estimation (MSL), method of simulated scores (MSS), and Monte Carlo Markov Chain methods (with emphasis on Gibbs sampling). Each method is described, and its theoretical properties are discussed.

3.1. Method of Simulated Moments

Many estimation problems involve finding a parameter vector θ that solves a set of orthogonality conditions

$$Q'h(y, X | \theta) = 0 \quad (3.1)$$

where Q is a set of instruments with dimension equal to the dimension of θ .⁸ Such estimators are called method of moments (MOM) estimators. All least squares methods are special cases of equation (3.1), and many problems usually estimated as MLE can be recast as MOM estimators. For example, Avery, Hansen, and Hotz(1983) suggest how to recast the the multinomial probit problem as a MOM problem where $h(y, X | \theta)$ is the vector $y - E(y | X)$ in the multinomial probit problem of Section 1 with j th element given by equation (1.4).

In many MOM problems, the orthogonality condition can not be evaluated analytically. For example, in the multinomial probit problem, evaluating $E[y | X]$ involves evaluating equation (1.4). MSM replaces $h(y, X | \theta)$ with an unbiased simulator $\hat{h}(y, X | \theta)$ and then finds the θ that solves

$$Q' \hat{h}(y, X | \theta) = 0. \quad (3.2)$$

The θ that solves equation (3.2) is the MSM estimator of θ , $\hat{\theta}$. McFadden(1989) and Pakes and Pollard(1989) show that, as long as $\hat{h}(y, X | \theta)$ is an unbiased simulator of $h(y, X | \theta)$, deviations between \hat{h} and h will wash out by the Law of Large Numbers because equation (3.2) is linear in \hat{h} and $\text{plim}(\hat{\theta}) = \theta$ as the sample size $N \rightarrow \infty$ even for small R .⁹

Consider the multinomial probit problem in more detail. As in Section 1, let y_i be the vector of dependent variables for observation i , $i = 1, 2, \dots, N$, where $y_{ij} = 1$ iff choice j is chosen by i . The probability of i choosing j conditional on X_i is given in equation (1.8), and its frequency simulator is given in equation (1.10). The frequency simulator should be replaced by one of the simulators discussed in Section 2, but for now we will use the frequency simulator for ease of presentation. As was discussed earlier, $E[y_{ij} | X_i] = \Pr[y_{ij} = 1 | X_i]$. Let P_i be a J -element vector with $\Pr[y_{ij} = 1 | X_i]$ in the j th element of P_i , and let $\varepsilon_i = y_i - P_i$. Then $E[\varepsilon_i | X_i] = 0$, and

$$E \sum_i Q_i' \varepsilon_i = 0 \quad (3.3)$$

for any set of exogenous instruments Q_i . Thus, conditional on a chosen $Q = (Q_1, Q_2, \dots, Q_N)$, the $\theta = (\beta, \Omega)$ that satisfies $\sum_i Q_i' \varepsilon_i = 0$ is the MOM estimator of θ . Let \hat{P}_i be an unbiased simulator of P_i , and let $\hat{\varepsilon}_i = y_i - \hat{P}_i$. Then the θ that solves

$$\sum_i Q_i' \hat{\varepsilon}_i = 0 \quad (3.4)$$

⁸When the dimension of Q is greater than the dimension of θ , the problem can be generalized to a GMM problem.

⁹Extra conditions are found in McFadden(1989) and Pakes and Pollard(1989).

is the MSM estimator of θ .

To find a reasonable Q , consider the log likelihood contribution for the multinomial probit model:

$$L_i = \sum_j y_{ij} \log P_{ij}. \quad (3.5)$$

The score statistics for θ can be written as

$$\begin{aligned} \partial L_i / \partial \theta &= \sum_j y_{ij} \frac{\partial P_{ij} / \partial \theta}{P_{ij}} \\ &= \sum_j \frac{\partial P_{ij} / \partial \theta}{P_{ij}} (y_{ij} - P_{ij} + P_{ij}) \\ &= \sum_j \frac{\partial P_{ij} / \partial \theta}{P_{ij}} (y_{ij} - P_{ij}) + \sum_j \frac{\partial P_{ij}}{\partial \theta} \end{aligned} \quad (3.6)$$

where the last term equals zero because the $\sum_j P_{ij} = 1$. Thus, one can write the score statistics in the form of equation (3.4). With an initial estimate of θ , one can construct $(1/P_{ij}) (\partial P_{ij} / \partial \theta)$ for θ and all j and use it as an instrument matrix Q_i for each i . It is likely that the instruments Q will need to be simulated (e.g., if the elements of Q_i are $(1/P_{ij}) (\partial P_{ij} / \partial \theta)$). This presents no significant problems as long as the pseudorandom variables used to simulate Q_i are independent of those used in the estimation process (to ensure exogeneity). For any exogenous Q , the $\hat{\theta}$ that solves equation (3.4) is a consistent estimate of θ . Thus, once θ is estimated, Q can be updated using $\hat{\theta}$ and then used to find a new $\hat{\theta}$ that solves equation (3.4).

For any exogenous Q , the covariance matrix of $\hat{\theta}$ has two terms: a term due to random variation in the data and a term due to simulation. As long as \hat{P}_i is an exogenous, unbiased simulator of P_i , one can write

$$\hat{P}_i = P_i + \xi_i \quad (3.7)$$

where ξ_i is a random variable caused by simulation with zero mean independent of ϵ_i , the deviation between y_i and P_i . Thus, the covariance matrix of $\hat{\epsilon}_i$ can be written as $E\epsilon\epsilon' + E\xi\xi'$. If \hat{P}_i is the frequency simulator of P_i , then ξ is just an average of R independent pseudorandom variables each with the same covariance matrix as ϵ . Thus, the covariance matrix of $\hat{\epsilon}$ is the covariance matrix of ϵ times $[1 + R^{-1}]$. The asymptotic covariance matrix of $\hat{\theta}$ is a linear function of the covariance matrices for $\hat{\epsilon}_i, i = 1, 2, \dots, N$ (McFadden, 1989, p. 1006). Note that for any $R \geq 1$, $\hat{\theta}$ is consistent; that as $R \rightarrow \infty$, the MSM covariance matrix approaches the MOM covariance matrix (which is efficient when the two-step procedure described above is used); and that the marginal improvement in precision declines

rapidly in R . If an alternative simulator with smaller variance is used, then the loss of precision due to simulation declines. For example, if antithetic acceleration is used, then the loss in precision becomes of order $(1/N)$ (see Geweke 1988) which requires no adjustment to the asymptotic covariance matrix.

Below is a roadmap for using MSM to estimate multinomial probit parameters:

a) Choose an identifiable parameterization for Ω and initial values for $\theta = (\beta, \Omega)$. Make sure that the initial guess results in probabilities reasonably far from zero or one.

b) Choose a simulator.

c) Simulate $2NJR$ ¹⁰ standard normal random variables. Store NJR of them in an instruments random number file and NJR in an estimation random number file. These random numbers will be used throughout the estimation process and never changed.

d) Given the initial guess of θ and the instruments random number file, simulate Q . Store the simulated instruments.

e) Given the initial guess of θ , the simulated Q , and the estimation random number file, solve equation (3.4) for θ . This is an MSM estimator of θ .

f) Given the initial MSM estimator, reperform steps (d) and (e) once.

Solving equation (3.4) requires using an optimization algorithm to find the θ that minimizes

$$\sum_i \hat{\varepsilon}_i' Q_i Q_i' \hat{\varepsilon}_i. \quad (3.8)$$

The derivatives of \hat{P}_i are well behaved, so derivative based optimization routines should be used. At each guess of θ , the standard normal pseudorandom numbers in the estimation random number file are used to create a new set of $N [0, \Omega]$ random numbers using the method described in Section 2. Thus, even though the standard normal random numbers never change, one is always using random numbers from the correct normal distribution.

Consider the unobserved heterogeneity count problem described in equations (1.11) through (1.13). Let y_{it} be the number of events for i at time t . $E[y_{it} | \lambda_{it}]$ is λ_{it} , but the covariance matrix of y_i has no closed form. Let v_i be a vector of residuals with $[T + T(T + 1)/2]$ elements. The first T elements of v_i are $y_{it} - E\lambda_{it}$ for $t = 1, 2, \dots, T$ where the expectation is over e_{it} and u_i in equation (1.12). The last $T(T + 1)/2$ elements correspond to ‘‘covariance residuals.’’ A representative element would be

$$(y_{it} - E\lambda_{it})(y_{is} - E\lambda_{is}) - C_{its} \quad (3.9)$$

¹⁰Remember that N = sample size, J = number of choices, and R = number of draws.

for two periods, t and s , where C_{its} is the $\text{Cov}(y_{it}, y_{is})$. The MOM estimator of $\theta = (\beta, \rho, \sigma_G, \sigma_H)$ solves

$$\sum_i Q_i' v_i = 0 \tag{3.10}$$

given a set of instruments Q . Since both $E\lambda_{it}$ and C_{its} can not be evaluated analytically,¹¹ the MOM estimator is not feasible. But $E\lambda_{it}$ and C_{its} can be simulated. Let \hat{y}_{it}^r be a simulated count variable. We can simulate e_{it} and u_i and therefore λ_{it} . Conditional on the simulated λ_{it} , we can simulate y_{it} either directly or by using the relationship between Poisson random variables and exponential random variables.

Applications using MSM include a retirement problem in Berkovec and Stern(1992), a market entry problem in Berry(1992), a dynamic programming problem in Hotz, et al.(1994), and an automobile pricing model in Berry, Levinsohn, and Pakes(1995).

3.2. Maximum Simulated Likelihood

A common estimation method with good optimality properties is maximum likelihood (ML) estimation. The basic idea is to maximize the log likelihood of the observed data over the vector of estimated parameters. ML estimators are consistent and efficient for a very large class of problems. Their asymptotic distribution is normal for a slightly smaller class of problems. However there are many likelihood functions that can not be evaluated analytically. In many cases, they can be thought of as expected values of some random function that can be simulated.

Consider again the multinomial probit problem. The log likelihood contribution for observation i is defined in equation (3.5). Note that only one element of y_i is not zero, so only one probability needs to be computed. This is a significant advantage of maximum simulated likelihood (MSL) over MSM. Still, to evaluate the log likelihood function, one must be able to evaluate or simulate P_{ij} for the choice chosen. The MSL estimator of θ is the value of θ that maximizes

$$L = \sum_{i=1}^N \sum_j y_{ij} \ln \hat{P}_{ij} \tag{3.11}$$

where \hat{P}_{ij} is the simulated value of P_{ij} .

¹¹Under special assumptions about the distribution of u_i and e_{it} described in Hausman, Hall, and Griliches (1984), the moments have analytical forms.

A significant problem with MSL is that the log likelihood function is not linear in \hat{P} . Thus, unlike MSM, the simulation errors, $\hat{P} - P$, will not wash out asymptotically as $N \rightarrow \infty$ unless $R \rightarrow \infty$ also. Lerman and Manski(1981) suggested using MSL with a frequency simulator. They found that R needed to be quite large to deal with this problem. However, Borsch-Supan and Hajivassiliou(1993) show in Monte Carlo studies that if better simulators are used, in particular smooth, smaller variance simulators bounded away from zero and one, then the bias caused by finite R is small for moderate sized R . In fact, in their study, MSL performs better than MSM.

Consider the unobserved heterogeneity model described in equations (1.11) through (1.13). The log likelihood contribution for observation i is given in equation (1.14). The argument of the log is the expected value of

$$\prod_{t=1}^T [\exp \{-\lambda_{it}\} \lambda_{it}^{y_{it}} / y_{it}!] \quad (3.12)$$

over the distribution of the errors determining λ_{it} . One can simulate λ_{it} for each i and t and therefore the expected value of the term in equation (3.12). Since the simulator of L_i is the log of this term, it is biased, and the bias disappears only as $R \rightarrow \infty$. But the simulator of equation (3.12) is smooth, and antithetic acceleration can be used to significantly reduce the variance. Thus the asymptotic bias associated with simulating the log likelihood function should be small.

Applications of MSL include a patent renewal model in Pakes(1986), a long-term care model in Borsch-Supan, et al.(1992), and the production function model in Ohanian, et al.(1996).

3.3. Method of Simulated Scores

A property of maximum likelihood is that the score statistic, the derivative of the log likelihood function, should have an expected value of zero at the true value of θ . This idea is the motivation behind the method of simulated scores (MSS). Hajivassiliou and McFadden(1990) use MSS in a model of external debt crises. The potential advantage of MSS is to use an estimator with the efficiency properties of ML and the consistency properties of MSM. MSM is asymptotically efficient if the proper weights are used (those that turn the moment condition into a score statistic). MSS ensures that the proper weights are used. The difficulty in this method is to construct an unbiased simulator of the score statistic. The problems this causes will become clear in the multinomial probit example. The

log likelihood contribution of observation i is given in equation (3.5), and its derivative is

$$\begin{aligned}\partial L_i / \partial \theta &= \sum_j y_{ij} \frac{\partial P_{ij} / \partial \theta}{P_{ij}} \\ &= \frac{\partial P_{ij} / \partial \theta}{P_{ij}}\end{aligned}\tag{3.13}$$

for the j corresponding to the chosen alternative. The goal is to construct an unbiased simulator for equation (3.13) so that the problem can be turned into a MSM problem. While it is straightforward to construct an unbiased simulator for both the numerator and denominator in equation (3.13), the ratio will not be unbiased as long as the denominator is random.

Consider constructing an unbiased simulator of the ratio. Suppressing the i subscript, equation (3.13) can be written as

$$\frac{\partial P_j / \partial \theta}{P_j} = \frac{\partial}{\partial \theta} \int_{A_j} f(y^*) dy^* / P_j\tag{3.14}$$

where $y^* = (y_1^*, y_2^*, \dots, y_J^*)$, f is the joint density of y^* , and A_j is the subset of the support of y^* where $y_j^* > y_k^*$ for all $k \neq j$. This equals

$$\begin{aligned}\frac{\partial P_j / \partial \theta}{P_j} &= \int_{A_j} \frac{\partial f(y^*) / \partial \theta}{f(y^*)} f(y^*) dy^* / P_j \\ &= E \left[\frac{\partial}{\partial \theta} \ln f(y^*) \mid y_j = 1 \right]\end{aligned}\tag{3.15}$$

where the expectation is with respect to the joint density of y^* . One usually can simulate the expectation in equation (3.15) (e.g., using the GHK simulator) and thus get an unbiased estimator of the ratio. Hajivassiliou and Ruud (1994) show that this method of simulating the score generalizes for all limited dependent variable problems.

3.4. Monte Carlo Markov Chain Methods

The last estimation procedure discussed is quite different than the others in that it is a Bayesian estimator. In general, we have a model specified up to a set of parameters θ , some data $\{(y_i, X_i)\}_{i=1}^N$, and a prior distribution for θ . The goal is to use the data to update the prior distribution to get a posterior distribution for θ . Computing the posterior involves using Bayes rule which usually involves solving a difficult integral, thus making it an intractable problem. Consider a

general problem where $\pi(z)$ is a known density function and $p(z^{n+1} | z^n)$ is a transition density function such that

$$\pi(z) = \int p(z | z') \pi(z') dz'. \quad (3.16)$$

Then Markov Chain theory tells us that repeated application of the transition density to an arbitrary density $\varphi(z)$ will asymptote to $\pi(z)$:

$$\pi(z) = \int p^n(z | z') \varphi(z') dz' \quad (3.17)$$

where

$$\begin{aligned} p^n(z | z') &= \int p^{n-1}(z | z'') p(z'' | z') \varphi(z'') dz'' \\ p^1(z | z') &= p(z | z') \end{aligned} \quad (3.18)$$

The idea in Monte Carlo Markov Chain (MCMC) methods is to simulate from $p(z^{n+1} | z^n)$ repeatedly and to thus generate a sample from $\pi(z)$. A popular MCMC method is Gibbs sampling (possibly with data augmentation). Continuing with our general notation, assume that there is natural way to partition z into (z_1, z_2, \dots, z_k) such that the conditional densities, $\pi(z_i | z_1, z_2, \dots, z_{i-1}, z_{i+1}, \dots, z_k)$ are easy to simulate from for all i . Then the Gibbs sampling algorithm is:

(a) Initialize $z^0 = (z_1^0, z_2^0, \dots, z_k^0)$ and set $n = 0$.

(b) Simulate $z_1^{n+1} \sim \pi(z_1^{n+1} | z_2^n, z_3^n \dots z_k^n)$;

$z_2^{n+1} \sim \pi(z_2^{n+1} | z_1^{n+1}, z_3^n \dots z_k^n)$;

$z_3^{n+1} \sim \pi(z_3^{n+1} | z_1^{n+1}, z_2^{n+1}, z_4^n, \dots, z_k^n)$;...

$z_k^{n+1} \sim \pi(z_k^{n+1} | z_1^{n+1}, z_2^{n+1}, \dots, z_{k-1}^{n+1})$.

(c) Set $n = n + 1$, and return to (b).

MCMC theory shows that this procedure will generate a sample $\{z^n\}_{n=N_0}^{N_1}$ from $\pi(z)$ where N_0 is chosen to give any effects due to initialization of z^0 an opportunity to die out. Sometimes, while it is difficult to simulate from all of the conditional densities $\pi(z_i | z_1, z_2, \dots, z_{i-1}, z_{i+1}, \dots, z_k)$, there is a way to augment the data with a latent variable so that all of the new conditional densities can be simulated from; such an approach is called Gibbs sampling with data augmentation.

Returning to our Bayesian estimation problem, let's say $\{y_i^*\}_{i=1}^N$ that has the following properties:

- a) the posterior distribution of y_i^* given (y_i, θ) is easy to simulate from, and
- b) the posterior distribution of θ given (y_i^*, y_i) and the prior distribution of θ is easy to compute and simulate from.

Assume there is a $\{y_i^*\}_{i=1}^N$ that satisfies these two conditions. Then the Gibbs sampling with data augmentation algorithm draws $\{y_i^*\}_{i=1}^N$ given $\{y_i\}_{i=1}^N$ and θ , then draws θ given new $\{y_i^*, y_i\}_{i=1}^N$, and repeats this process over and over again. The draws of θ provide information about the posterior distribution of θ . The algorithm is:

- (a) Assume a prior distribution for θ . Choose R_0 such that the first R_0 draws will not count and R_1 such that the process will stop after R_1 draws. Set $r = 0$.
- (b) Simulate one draw of θ from its posterior distribution.
- (c) If $r > R_0$, store the draw of θ as draw $r - R_0$.
- (d) If $r > R_1$, then $r = r + 1$ and goto (g).
- (e) Simulate one draw of $\{y_i^*\}_{i=1}^N$ conditional on $(\{y_i\}_{i=1}^N, \theta)$.
- (f) Evaluate analytically the posterior distribution for θ given $\{(y_i, y_i^*)\}_{i=1}^N$. Increment $r = r + 1$. Goto (b).
- (g) Use the $R_1 - R_0$ draws of θ as a random sample of draws of θ and compute any sample characteristics desired.

Markov chain theory implies that the Gibbs sampling algorithm described above will produce a distribution of draws of θ corresponding to the posterior distribution of θ conditional on $\{(y_i, X_i)\}_{i=1}^N$. See, for example, Casella and George (1992), Gelfand and Smith(1990), Geman and Geman(1984), and Tanner and Wong(1987) for more about Markov chains.

Consider how Gibbs sampling with data augmentation can be applied to the multinomial probit problem.¹² To simplify exposition, assume we know Ω and only need to estimate β . Assume $\beta_1 = 0$ as a normalizing factor. For step (a), we need a prior distribution for $\{\beta_j\}_{j=2}^J$. If we pick R_0 big enough and the prior with a large enough variance, then the choice of prior will become irrelevant. Thus, pick the prior to be diffuse. The diffuse prior makes it easy to compute posterior distributions for $\{\beta_j\}_{j=2}^J$. Next, let y_i^* be the latent variable associated with y_i :

$$y_{ik}^* = X_i \beta_k + u_{ik}, \quad k = 1, 2, \dots, J, \quad i = 1, 2, \dots, N \quad (3.19)$$

where $u_i \sim N [0, \Omega]$.

¹²This is described in much more detail in McCulloch and Rossi (1996).

For step (b), we need to simulate β from its posterior distribution. Since, at any iteration of the algorithm, β is normal, we can simulate β using the method described in Section 2.

For step (e), we need to simulate $\{y_i^*\}_{i=1}^N$ conditional on $(\{y_i\}_{i=1}^N, \beta)$. Since the observations are independent, we need only simulate y_i^* conditional on (y_i, β) for each $i = 1, 2, \dots, N$ separately. Let j be the chosen choice. Then

$$X_i\beta_j + u_{ij} > X_i\beta_k + u_{ik} \quad \forall k \neq j \quad (3.20)$$

or

$$u_{ijk}^* < X_i(\beta_j - \beta_k) \quad \forall k \neq j \quad (3.21)$$

where $u_{ijk}^* = u_{ik} - u_{ij}$. The errors $u_{ijk}^* \quad \forall k \neq j$ can be simulated using the GHK algorithm, and the y_{ik}^* can be constructed $\forall k \neq 1$ ¹³ as

$$y_{ik}^* = X_i\beta_k + u_{ijk}^* - u_{ij1}^*. \quad (3.22)$$

Alternatively, we can use an acceptance-rejection simulator.

For step (f), we need to evaluate the posterior distribution of β given $\{y_i^*\}_{i=1}^N$. Since $u_{ij} \sim N(0, \omega_{jj})$ for each i and j , $y_{ij}^* \sim N[X_i\beta_j, \omega_{jj}]$ which means that computing a posterior distribution for β involves running an OLS regression of y^* on X .

For step (g), the sample of $R_1 - R_0$ draws of β are distributed from the distribution of β conditional on the data (including the dependent variables $\{y_i\}_{i=1}^N$). A few notes of caution are in order here. First, the draws of β are not independent even though any dependence dies out as the number of draws between two draws becomes large. Thus, we must not compute any statistics that depend upon the ordering of the draws. Second, the draws are conditional on $\{y_i\}_{i=1}^N$. This is quite different than what we would expect in classical statistical analysis (where we would condition on only the exogenous variables). The effect of this is that the researcher does not know how the estimator would have behaved had a different realization of the data been observed. This is a fundamental difference between classical estimators and Bayesian estimators. There are other reasonable (and perhaps better) choices for implementing the Gibbs sampler to the multinomial probit problem. The real issues involve also estimating Ω . See McCullough and Rossi(1994, 1996) or Albert and Chib(1993) for a much more extensive discussion.

The unobserved heterogeneity count problem is also easily adaptable to Gibbs sampling. The data should be augmented with $\{\lambda_{it}\}_{t=1}^T, i=1, \dots, N$ and its prior should be

¹³Recall that choice 1 is the base choice.

normal. Steps (b) and (f) are the same as in the multinomial probit problem. Step (e) involves simulating λ_{it} conditional on (y_{it}, β) which is not as straightforward. The density of λ_{it} conditional on (y_{it}, β) is

$$f(\lambda_{it} | y_{it}, \beta) = C(y) e^{-\lambda_{it}} \lambda_{it}^{y_{it}-1} \phi\left(\frac{\log \lambda_{it} - X_{it}\beta}{\sigma_\lambda}\right) \quad (3.23)$$

where σ_λ is the standard deviation of the composite error in equation (1.12), ϕ is the standard normal density function, and $C(y)$ is a proportionality constant chosen so that equation (3.23) integrates to one. One can evaluate the integral of equation (3.23) numerically for each value of $y = 0, 1, \dots$ for a finite number of points: $\delta, 2\delta, \dots, K\delta$ for some small δ . Figure 1 draws the approximate distribution curves for $y = 0, 1, \dots, 5$, $\delta = .01$, and $K = 1000$. Then one can use the discretized distribution as an approximation to draw λ from. This is equivalent to drawing a random point on the vertical axis of Figure 1 (e.g., point A), drawing a horizontal line to the curve corresponding to y (e.g., B when $y = 4$) and choosing λ to be the horizontal component of the curve at that vertical point (e.g., point C).

A generalization of Gibbs sampling is the Metropolis-Hastings (MH) algorithm described in, for example, Chib and Greenberg(1994). Returning to our general notation, let $p(z^{n+1} | z^n)$ be the density we want to simulate from so that we can generate a sample with density $\pi(z)$, but assume it is difficult to simulate from $p(z^{n+1} | z^n)$ directly. Let $q(z^{n+1} | z^n)$ be a ‘‘candidate density’’ chosen according to criteria described in Chib and Greenberg(1994). Then the MH algorithm is:

- (a) Initialize z^0 and set $n = 0$.
- (b) Simulate z^{n+1} from $q(z^{n+1} | z^n)$ and keep it with probability $\alpha(z^{n+1}, z^n)$ where

$$\alpha(z^{n+1}, z^n) = \begin{cases} \min\left[1, \frac{\pi(z^{n+1})q(z^n | z^{n+1})}{\pi(z^n)q(z^{n+1} | z^n)}\right] & \text{if } \pi(z^n) q(z^{n+1} | z^n) > 0 \\ 1 & \text{otherwise.} \end{cases} \quad (3.24)$$

- (c) Set $n = n + 1$ and return to (b).

The MH algorithm is essentially a sophisticated acceptance-rejection method. The acceptance probability $\alpha(z^{n+1}, z^n)$ oversamples transitions where $\pi(z^{n+1})/\pi(z^n)$ is high relative to $q(z^{n+1} | z^n)/q(z^n | z^{n+1})$. Gibbs sampling is a special case of the MH algorithm for $q(z^{n+1} | z^n)$ being the conditional density and $\alpha(z^{n+1}, z^n) = 1$. The difficult part of implementing the MH algorithm is choosing the candidate density $q(z^{n+1} | z^n)$. One wants a $q(z^{n+1} | z^n)$ that moves around fast enough so that the whole support of $\pi(z)$ is sampled but slow enough so that $\alpha(z^{n+1}, z^n)$ is not too small. Also, it is worthwhile to have a candidate density such that

$q(z^{n+1} | z^n) = q(z^n | z^{n+1})$.¹⁴ Chib and Greenberg suggest using candidate densities of the form $q(z^{n+1} | z^n) = q_1(z^{n+1} - z^n)$ or $q(z^{n+1} | z^n) = q_2(z^{n+1})$. They also provide an example estimating an ARMA model with regressors using the MH algorithm.

3.5. Empirical Comparison of Methods

A number of studies have compared the performance of various simulators and estimation methods especially for the multinomial probit problem. This section summarizes the results of four of those studies and presents some new results focusing on questions that are neglected in the other studies.

Borsch-Supan and Hajivassiliou (1993) compare the GHK simulator to the Stern simulator and a frequency simulator. They present convincing evidence that the GHK simulator has a significantly smaller standard deviation than the other two simulators. They further show that the standard deviation of the GHK simulator is small enough so that it can be used in an MSL estimation routine providing parameter estimates with small root mean squared errors (RMSE's). Having a good simulator with a small standard deviation for MSL is important because, unlike MSM, MSL does not provide consistent estimates for fixed R .

Hajivassiliou, McFadden, and Ruud (1994) compare ten different simulators (including the Stern simulator, a Gibbs sampler, and a kernel smoothed simulator) in terms of the RMSE of the multinomial probit probability and its derivatives. They consider a large class of V_j 's and Ω_j^* 's. They find that the GHK simulator performs the best overall. In particular, it performs well relative to the alternatives when Ω_j^* displays high correlation terms. They provide no results concerning parameter estimates.

Geweke, Keane, and Runkle (1994a) compare MSM using GHK, MSL using GHK, Gibbs sampling, and kernel smoothing. In an unrestricted estimation procedure (including covariance parameters), MSM-GHK and Gibbs sampling dominated MSL-GHK. Kernel smoothing was dominated by all methods. In various restricted models, the performance of MSL-GHK improved. In general, as more restrictions were placed on the model, the performance of MSM-GHK, MSL-GHK, and Gibbs sampling converged. But Gibbs sampling seemed to dominate other methods overall.

Geweke, Keane, and Runkle (1994b) compare MSM-GHK, MSL-GHK, and Gibbs sampling in the related multinomial multiple period probit model. They

¹⁴This simplifies evaluation of $\alpha(z^{n+1}, z^n)$ among other things.

find that Gibbs sampling dominates and MSM-GHK is second. Estimated standard errors are good for Gibbs sampling and MSM-GHK but are downward biased for MSL-GHK.

None of these methods compare the computational cost of the alternatives. Computational cost is important because the simulators are essentially a method to reduce computation time; if time was not an issue, we could compute the relevant integrals numerically using arbitrarily precise approximation methods or we could simulate them letting R be an arbitrarily large number. If one method takes twice as much time as another for a given R , then a fair comparison requires using different R for each method to produce comparable times. Also none of the methods considers the effect of using antithetic acceleration (AA) despite Geweke's strong theoretical results.

Table 1 presents the results of a small Monte Carlo study. Its results should be interpreted as suggestive of where more work needs to be done. The methods that are compared are MSM-GHK, MSM-Stern, MSL-GHK, MSL-Stern, Gibbs sampling (with acceptance-rejection), and MSM-KS (kernel smoothing). Three different models are used: a) Ω is diagonal and N (sample size) = 500, b) Ω is diagonal and $N = 1000$, and c) Ω corresponds to an $AR(1)$ process with $\rho = .9$ and $N = 1000$. Except for Gibbs sampling, results are reported with and without AA. RMSE results and average times per estimation procedure are reported.

Kernel smoothing methods performed poorly in terms of RMSE of the simulated multinomial probit probabilities. Also, more importantly, its derivatives with respect to parameters were poorly behaved in that if the bandwidth parameter was small, the derivatives were very volatile (and therefore derivative based optimization algorithms for estimation behaved poorly), and if it was large, parameter bias was very large. Thus kernel smoothing method results are not reported. In terms of RMSE, Gibbs sampling estimators behave reasonably well. But the amount of time involved is an order of magnitude greater than for the MSM and MSL estimates.¹⁵ Thus, there are only limited results reported for the Gibbs samplers.

The remainder of the discussion focuses on MSM, MSL, GHK, Stern, and AA. First, it is clear that MSL dominates MSM in these examples. It provides smaller RMSE's and it requires less computation time. GHK dominates Stern in terms of RMSE, but Stern is significantly faster. One might consider using Stern

¹⁵It should be noted that in these Monte Carlo experiments, I am conditioning on the true value of Ω . It might be the case that the Gibbs sampler performs better relative to the other methods when Ω also is estimated.

with twice as large R . Unreported Monte Carlo experiments suggest that for the examples used here the standard deviation of the multinomial probit probabilities is about twice as large for the Stern simulator as for the GHK simulator when $R = 10$. This would suggest that doubling R for the Stern simulator (relative to the GHK simulator) would make the GHK simulator more efficient by a factor of $\sqrt{2}$. Thus, these results are consistent with Borsch-Supan and Hajivassiliou, suggesting that MSL-GHK provides estimates with the smallest RMSE's even after controlling for variation in computation time. Based on results in Borsch-Supan and Hajivassiliou and Hajivassiliou, McFadden, and Ruud, it probably performs even better for pathological cases with highly correlated errors or small multinomial probit probabilities.

The poor performance of AA is striking. AA almost uniformly improves the performance of the Stern simulator. But it behaves poorly for the GHK simulator. However, Table 2 shows that AA significantly reduces the standard deviation of the simulated multinomial probit probabilities for GHK, Stern, and kernel smoothing. This apparent paradox occurs because of the small sample properties of method of moments (MOM) and maximum likelihood (MLE). In other words, the RMSE of MOM and MLE dominate any extra randomness caused by simulation. This is verified by unreported results showing that when R is increased to 50, MSL-GHK and MSL-Stern RMSE's converge to each other with or without AA and they are similar to the RMSE's for the case when $R = 5$ with AA or $R = 10$ without AA. The bottom line is that for MSM and MSL, the choice of simulation method has a second order effect on RMSE relative to RMSE caused by the underlying estimation method. This further suggests that computation time issues should be given high priority.

Table 1
Monte Carlo Estimation Results

Results for Diagonal Covariance Matrix				
N = 500				
Method	w/ Antithetic Acceleration		wo/ Antithetic Acceleration	
	Avg RMSE	Avg Time	Avg RMSE	Avg Time
MSM-GHK	0.299	3559.0	0.257	3373.8
MSM-Stern	0.270	1047.0	0.288	1097.0
MSL-GHK	0.247	1571.0	0.246	1598.8
MSL-Stern	0.254	654.9	0.252	674.7
Gibbs			0.263	16119.9

Results for Diagonal Covariance Matrix				
N = 1000				
Method	w/ Antithetic Acceleration		wo/ Antithetic Acceleration	
	Avg RMSE	Avg Time	Avg RMSE	Avg Time
MSM-GHK	0.181	6470.9	0.167	6283.1
MSM-Stern	0.173	1911.4	0.186	2006.0
MSL-GHK	0.158	1951.5	0.161	1889.9
MSL-Stern	0.161	802.0	0.163	853.0
Gibbs			0.170	29746.9

Table 1 cont'd

Results for Non-Diagonal Covariance Matrix				
N = 1000				
Method	w/ Antithetic Acceleration		wo/ Antithetic Acceleration	
	Avg RMSE	Avg Time	Avg RMSE	Avg Time
MSM-GHK	0.267	7192.1	0.201	6782.8
MSM-Stern	0.358	2422.5	0.420	2565.2
MSL-GHK	0.175	2194.8	0.192	2010.0
MSL-Stern	0.180	1114.3	0.195	1174.0

Notes:
 There are 200 Monte Carlo draws per experiment.
 There are 6 choices and 5 explanatory variables per choice.
 For experiments with AA, $R = 5$, and for experiments without AA, $R = 10$.
 All experiments are performed on an IBM RS6000 Model 390.
 Gibbs sampling results are based on 10000 draws after skipping 2000 draws;
 i.e., $R_0 = 2000$ and $R_1 = 12000$.

Table 2
Probability Simulations

	wo/ AA	w/ AA
GHK	-0.00050 (0.033)	-0.00269 (0.021)
Stern	-0.00266 (0.059)	-0.00086 (0.023)
Kernel Smoothing	0.00000 (0.077)	0.00000 (0.063)

Notes:
AA is antithetic acceleration.
First row for each simulation method is a sample mean, and second row (in parentheses) is a sample standard deviation.
There are 3000 Monte Carlo draws per experiment.
There are 6 choices and 5 explanatory variables per choice.
For experiments with AA, $R = 5$, and for experiments without AA, $R = 10$.
All experiments are performed on an IBM RS6000 Model 390.

References

- [1] Albert, J. and S. Chib (1993). "Bayesian Analysis of Binary and Polychotomous Data." JASA. 88: 669-679.
- [2] An, Mark Y. and Ming Liu (1996). "Structural Analysis of Labor Market Transitions Using Indirect Inference." Unpublished manuscript, Duke University.
- [3] Anderson, S.P., A. de Palma, and J.F. Thisse (1992). Discrete Choice Theory of Product Differentiation. Cambridge: MIT
- [4] Avery, Robert, Lars Hansen, and V. Joseph Hotz (1983). "Multiperiod Probit Models and Orthogonality Condition Estimation." International Economic Review. 24(1): 21-35.
- [5] Bansal, Ravi, A. Ronald Gallant, Robert Hussey, and George Tauchen (1995). "Nonparametric Estimation of Structural Models for High Frequency Currency Market Data." Journal of Econometrics. 66: 251-287.
- [6] Berkovec, J. and S. Stern (1991). "Job Exit Behavior of Older Men." Econometrica. 59(1): 189-210.
- [7] Berry, Steven (1992). "Estimation of a Model of Entry in the Airline Industry." Econometrica. 60(4): 889-917.
- [8] Berry, Steven, James Levinsohn, and Ariel Pakes (1995). "Automobile Prices in Market Equilibrium." Econometrica. 63(4): 841-890.
- [9] Borsch-Supan, Axel and Vassilis A. Hajivassiliou (1993). "Smooth Unbiased Multivariate Probability Simulators for Maximum Likelihood Estimation of Limited Dependent Variable Models." Journal of Econometrics. 58: 347-368.
- [10] Borsch-Supan, Axel, Vassilis A. Hajivassiliou, Lawrence Kotlikoff, and John Morris (1992). "Health, Children, and Elderly Living Arrangements: A Multiperiod-Multinomial Probit Model with Unobserved Heterogeneity and Autocorrelated Errors." in Topics in the Economics of Aging. (ed.) David Wise. Chicago: University of Chicago Press.

- [11] Brown, Brian and Whitney Newey (1996). “Simulation-Based Inference in Semiparametric Procedures.” in Simulation-Based Inference in Econometrics: Methods and Applications. (ed.) Roberto S. Mariano, Melvyn Weeks, and Til Schuermann. Cambridge: Cambridge University Press.
- [12] Bunch, David (1991). “Estimability in the Multinomial Probit Model.” Transportation Research, Part B, Methodological. 25B: 1-12.
- [13] Butler, J.S. and Robert Moffitt (1982). “A Computationally Efficient Quadrature Procedure for the One-Factor Multinomial Probit Model.” Econometrica. 50: 761-764.
- [14] Casella, G. and E. George (1992). “Explaining the Gibbs Sampler.” American Statistician. 46: 167-174.
- [15] Chib, Siddhartha (1993). “Bayes Regression with Autoregressive Errors: A Gibbs Sampling Approach.” Journal of Econometrics. 58(3): 347-368.
- [16] Chib, Siddhartha and Edward Greenberg (1994). “Understanding the Metropolis-Hastings Algorithm.” Washington University, St. Louis, manuscript.
- [17] Danielsson, J. and J. F. Richard (1993). “Accelerated Gaussian Importance Sampler with Application to Dynamic Latent Variable Models.” Journal of Applied Econometrics. 8: 153-173.
- [18] Devroye, L. (1986). Non-Uniform Random Variate Generation. New York: Springer.
- [19] Diebold, Francis X. and Til Schuermann (1996). “Exact Maximum Likelihood Estimation of Observation-Driven Econometric Models.” in Simulation-Based Inference in Econometrics: Methods and Applications. (ed.) Roberto S. Mariano, Melvyn Weeks, and Til Schuermann. Cambridge: Cambridge University Press.
- [20] Gelfand, A. and A. Smith (1990). “Sampling-Based Approaches to Calculating Marginal Densities.” JASA. 85: 398-409.

- [21] Geman, S. and D. Geman (1984). “Stochastic Relaxation, Gibbs Distributions and the Bayesian Restoration of Images.” IEEE Transactions on Pattern Analysis and Machine Intelligence 6. 721-741.
- [22] Geweke, John F. (1988). “Antithetic Acceleration of Monte Carlo Integration in Bayesian Inference.” Journal of Econometrics. 38: 73-89.
- [23] Geweke, John (1989). “Bayesian Inference in Econometric Models Using Monte Carlo Integration.” Econometrica. 57(6): 1317-1339.
- [24] Geweke, John (1991). “Efficient Simulation from the Multivariate Normal and Student-t Distributions Subject to Linear Constraints.” Computer Science and Statistics: Proceedings of the Twenty-Third Symposium on the Interface. 571-578.
- [25] Geweke, J., M. Keane, and D. Runkle (1994a). “Alternative Computational Approaches to Inference in the Multinomial Probit Model.” Federal Reserve Bank of Minneapolis, Staff Report 170.
- [26] Geweke, J., M. Keane, and D. Runkle (1994b). “Statistical Inference in the Multinomial Multiperiod Probit Model.” Federal Reserve Bank of Minneapolis, Staff Report 177.
- [27] Hajivassiliou, Vassilis (1990). “Smooth Simulation Estimation of Panel Data LDV Models.” Unpublished paper.
- [28] Hajivassiliou, Vassilis and Daniel McFadden (1990). “The Method of Simulated Scores witha Application to Models of External Debt.” Cowles Foundation Discussion Paper No. 967.
- [29] Hajivassiliou, V., D. McFadden, and P. Ruud (1994). “Simulation of Multivariate Normal Rectangle Probabilities and their Derivatives: Theoretical and Computational Results.” Cowles Foundation Discussion Paper No. 1021R.
- [30] Hammersly, J.M. and D.C. Handscomb (1964). Monte Carlo Methods. London: Methuen.

- [31] Hausman, Jerry, Bronwyn Hall, and Zvi Griliches (1984). “Econometric Models for Count Data with an Application to the Patents R & D Relationship.” Econometrica. 52: 903-938.
- [32] Hausman, Jerry A. and David A. Wise (1978). “A Conditional Probit Model for Qualitative Choice: Discrete Decisions Recognizing Interdependence and Heterogenous Preferences.” Econometrica. 46(2): 403-426.
- [33] Hendry, David F. (1984). “Monte Carlo Experimentation in Econometrics.” Handbook of Econometrics, Volume II.
- [34] Hotz, V. Joseph, Robert Miller, Seth Sanders, and Jeffrey Smith (1994). “A Simulation Estimator for Dynamic Models of Discrete Choice.” Review of Economic Studies. 61: 265-289.
- [35] Keane, Michael P. (1994). “A Computationally Practical Simulation Estimator for Panel Data.” Econometrica. 62(1): 95-116.
- [36] Lerman, Steven and Charles Manski (1981). “On the Use of Simulated Frequencies to Approximate Choice Probabilities.” in Structural Analysis of Discrete Data with Econometric Applications. ed. by Charles Manski and Daniel McFadden. Cambridge: MIT Press.
- [37] McCulloch, Robert and Peter Rossi (1994). “An Exact Likelihood Analysis of the Multinomial Probit Model.” Journal of Econometrics. 64: 207-240.
- [38] McFadden, Daniel (1989). “A Method of Simulated Moments for Estimation of Discrete Response Models without Numerical Integration.” Econometrica. 57(5): 995-1026.
- [39] Ohanian, Lee, Giovanni L. Violante, Per Krusell, and José-Victor Ríos-Rull. (1996). “Simulation-Based Estimation of a Nonlinear, Latent Factor Aggregate Production Function.” in Simulation-Based Inference in Econometrics: Methods and Applications. (ed.) Roberto S. Mariano, Melvyn Weeks, and Til Schuermann. Cambridge: Cambridge University Press.
- [40] Pakes, Ariel (1986). “Patents as Options: Some Estimates of the Value of Holding European Patent Stocks.” Econometrica. 54(4): 755-784.

- [41] Pakes, Ariel and David Pollard (1989). “Simulation and the Asymptotics of Optimization Estimators.” Econometrica. 57(5): 1027-1057.
- [42] Richard, J. F. and Wei Zhang (1996). “Accelerated Monte Carlo Integration: An Application to Dynamic Latent Variable Models.” in Simulation-Based Inference in Econometrics: Methods and Applications. (ed.) Roberto S. Mariano, Melvyn Weeks, and Til Schuermann. Cambridge: Cambridge University Press.
- [43] Ripley, Brian (1987). Stochastic Simulation. New York: John Wiley and Sons.
- [44] Stern, Steven (1992). “A Method for Smoothing Simulated Moments of Discrete Probabilities in Multinomial Probit Models.” Econometrica. 60(4): 943-952.
- [45] Tanner, T. and W. Wong (1987). “The Calculation of Posterior Distributions by Data Augmentation.” JASA. 82: 528-549.
- [46] Tierney, L. (1994). “Markov Chains for Exploring Posterior Distributions.” Annals of Statistics. ?

