

Federal Reserve Bank of Minneapolis
Research Department Staff Report 192

May 1995

Monte Carlo Simulation and Numerical Integration*

John Geweke

Federal Reserve Bank of Minneapolis
and University of Minnesota

ABSTRACT

This is a survey of simulation methods in economics, with a specific focus on integration problems. It describes acceptance methods, importance sampling procedures, and Markov chain Monte Carlo methods for simulation from univariate and multivariate distributions and their application to the approximation of integrals. The exposition gives emphasis to combinations of different approaches and assessment of the accuracy of numerical approximations to integrals and expectations. The survey illustrates these procedures with applications to simulation and integration problems in economics.

*This is a chapter prepared for the *Handbook of Computational Economics*, edited by Hans Amman, David Kendrick, and John Rust, to be published by North-Holland. Comments from John Rust and two anonymous referees, who bear no responsibility for errors or omissions, are gratefully acknowledged. Kathleen Rolfe provided editorial assistance, and Maureen O'Connor created the charts. The work was supported in part by National Science Foundation Grant SES-9210070. The views expressed herein are those of the author and not necessarily those of the Federal Reserve Bank of Minneapolis or the Federal Reserve System.

1. Introduction

Optimization problems in dynamic, stochastic environments are an increasingly important part of economic theory and applied economics. Inspired by the potential returns to richer and more realistic models of a variety of policy problems and the promise of ever-growing computational power, economists have turned more and more to models that can be simulated but not solved in closed form. Simulation methods can provide solutions for two related integration problems. One integration problem arises in model solution, for agents whose expected utilities cannot be expressed as a closed function of state and decision variables. The other occurs when the investigator combines sources of uncertainty about models to draw conclusions about policy.

This chapter concentrates on simulation methods that are both important and useful in the solution of these integration problems. In mathematics there is a long-standing use of simulation in the solution of integration problems, notably partial differential equations, where the form of the simulation is often suggested by the problem itself. The history of simulation methods to solve integration problems in economics is shorter, but these methods are appealing there for the same reason: integration generally involves probability distributions in the integrand, which thereby suggests the simulation methods to be employed.

This pervasive use of simulation methods in science persists despite the well-known asymptotic advantages of deterministic approaches to integration. This continued use of simulation methods occurs in part because astronomical computing time is often required to realize the promise of deterministic methods. A more important fact is that simulation methods are generally straightforward for the investigator to implement, relying on an understanding of a few principles of simulation and the structure of the problem at hand. By contrast, deterministic methods typically require much larger problem-specific investments in numerical methods. Simulation methods economize the use of that most valuable resource, the investigator's time.

The objective of this chapter is to convey an understanding of principles for the practical application of simulation in economics, with a specific focus on integration problems. It begins with a discussion of circumstances in which deterministic methods are preferred to simulation, in Section 2. The next section takes up general procedures for simulation from univariate and multivariate distributions, including acceptance and adaptive methods. The construction and use of independent, identically distributed random vectors to solve the multidimensional integration problems that typically arise in economic models is taken up in Section 4, with special attention to combination of different approaches and

assessment of the accuracy of numerical approximations to the integral. Section 5 discusses some modifications of these methods to produce identically but not independently distributed random vectors, that often greatly reduce approximation error in applications in economics. Recently developed Markov chain Monte Carlo methods, which make use of samples that are neither independently nor identically distributed, have greatly expanded the scope of integration problems with convenient practical solutions. These procedures are taken up in Section 6. The chapter concludes with some examples of recent applications of simulation to integration problems in economics.

2. Deterministic methods of integration

The evaluation of the integral $I = \int_a^b f(x)dx$ is a problem as old as the calculus itself and is equivalent to solution of the differential equation $dy/dx = f(x)$ subject to the boundary condition $y(a) = 0$. In well-catalogued instances, analytical solutions are available. (Gradshteyn and Ryzhik, 1965, is a useful standard reference.) The literature on numerical approaches to each problem is huge, a review of any small part of which could occupy this entire volume. This section focuses on those procedures that provide the most useful tools in economics and are readily available in commercial software. This means neglecting the classical but dated approaches using equally spaced abscissas, like Newton-Cotes; a useful overview of these methods is provided by Press *et al.* (1986, Chapter 4), and a more extended discussion may be found in Davis and Rabinowitz (1984, Chapter 2).

2.1 Unidimensional quadrature

The principle underlying most state-of-the-art deterministic evaluations of $I = \int_a^b f(x)dx$ is Gaussian quadrature. If $f(x) = p(x)w(x)$, where $p(x)$ is any polynomial of degree $2n - 1$ or lower and $w(x)$ is a chosen basis function, then there exist points $x_i \in [a, b]$ and a weight ω_i associated with each point such that

$$\int_a^b f(x)dx = \int_a^b p(x)w(x)dx = \sum_{i=1}^n \omega_i p(x_i).$$

The points and weights depend only on a, b , and the function $w(x)$, and if they are known for $a=0$ and $b=1$, then it is straightforward to determine their values for any other choices of a and b . If $r(x) = f(x)/w(x)$ is not a polynomial of degree $2n - 1$ or lower, then

$$\sum_{i=1}^n \omega_i r(x_i)$$

may be taken as an approximation to $I = \int_a^b f(x)dx$. If $r(x)$ is smooth relative to a polynomial of degree $2n - 1$, then the approximation should be good. More precisely, one may show that if $r(x)$ is $2n$ -times differentiable, then

$$\int_a^b f(x)dx - \sum_{i=1}^n \omega_i r(x_i) = c_n r^{(2n)}(\xi)$$

for some $\xi \in [a, b]$, where $\{c_n\}$ is a sequence of constants with $\lim_{n \rightarrow \infty} c_n = 0$. For example, if $w(x) = 1$, $a = -1$, $b = +1$, then $c_n = 2^{2n+1}(n!)^4 / \{(2n+1)![2n!]\}^3$ (Judd, 1991, pp. 6-7, 6-8).

This approach can be applied to any subinterval of $[a, b]$ as well. As long as $r(x)$ is $2n$ -times differentiable, one may satisfy prespecified convergence or error criteria through successive bisection. Error criteria are usually specified as the absolute or relative

difference in the computed approximation to $I = \int_a^b f(x)dx$ using n -point and m -point quadrature (Golub and Welsch, 1969).

Infinite and semi-infinite intervals can be treated through appropriate transformation of variable to a finite interval (Piessens *et al.*, 1983). Existence and boundedness of $r^{(2n)}$ depend in part on the choice of basis function $w(x)$. Some of the most useful are indicated in the following table.

$w(x)$	Interval	Name
1	(-1,1)	Legendre
$1/\sqrt{1-x^2}$	(-1,1)	Chebyshev first kind
$\sqrt{1-x^2}$	(-1,1)	Chebyshev second kind
$\exp(-x^2)$	$(-\infty, +\infty)$	Hermite
$(1+x)^\alpha(1-x)^\beta$	(-1,1)	Jacobi
$\exp(-x)x^\alpha$	$(0, \infty)$	Generalized Laguerre
$1/\cosh(x)$	$(-\infty, +\infty)$	Hyperbolic cosine

For many purposes Gauss-Legendre rules are adequate, and there is a substantial stock of commercially supplied software to evaluate one-dimensional integrals up to specified tolerances. These methods have been adapted to include functions having singularities at identified points in the interval of integration (Piessens, *et al.*, 1983).

2.2 Multidimensional quadrature

Some multidimensional integration problems in fact reduce to an integration in a single variable that must be carried out numerically. For example, all but one dimension may be integrable analytically, or the multidimensional integral may in fact be a product of integrals each in a single variable, perhaps after a suitable change of variable. In such cases quadrature for one-dimensional integrals usually provides a neat solution. Such cases are rare in economics and econometrics. If the dimension of the domain of integration is not too high and the integrand is sufficiently smooth, then one-dimensional methods may be extended with practical results. These cases cover a small subset of integration problems in economics, but when they arise they deserve attention because quadrature-based methods are then often efficient and easy to use.

The straightforward extension of quadrature methods to higher dimensions shows both its strengths and weaknesses. Following Davis and Rabinowitz (1984, pp. 354-359), suppose that R is an m -point rule of integration over $B \subseteq \mathfrak{X}'$, leading to the approximation

$$R(f) = \sum_{j=1}^m \omega_j f(\mathbf{x}_j) \approx \int_B f(\mathbf{x})d\mathbf{x}, \mathbf{x}_j \in B,$$

and that S is an n -point rule over $G \subseteq \mathfrak{R}^s$, leading to the approximation

$$S(f) = \sum_{k=1}^n v_k f(\mathbf{y}_k) \approx \int_G f(\mathbf{y}) d\mathbf{y}, \mathbf{y}_k \in G.$$

The product rule of R and S is the mn -point rule applicable to $B \times G$,

$$R \times S(f) = \sum_{j=1}^m \sum_{k=1}^n \omega_j v_k f(\mathbf{x}_j, \mathbf{y}_k) \approx \int_{B \times G} f(\mathbf{x}, \mathbf{y}) d\mathbf{x} d\mathbf{y}, \mathbf{x}_j \in B, \mathbf{y}_k \in G.$$

If $h(\mathbf{x}, \mathbf{y}) = \sum_{i=1}^k f_i(\mathbf{x}) g_i(\mathbf{y})$, and if R integrates $f_i(\mathbf{x})$ exactly over B and S integrates $g_i(\mathbf{y})$ exactly over G ($i = 1, \dots, k$), then $R \times S$ will integrate $h(\mathbf{x}, \mathbf{y})$ exactly over $B \times G$. The obvious extensions to the product of three or more rules can be made. These extensions can be expected to work well when (a) quadrature is adequate in the lower dimensional marginals of the function at hand, (b) $h(\mathbf{x}, \mathbf{y}) \approx f(\mathbf{x})g(\mathbf{y})$, and (c) the product mn is small enough that computation time is reasonable. Condition (c) and perhaps (a) are violated when the support of h is concentrated on a set small relative to the Cartesian boundaries for that support, as illustrated in Figure 1(a). A more common occurrence in economics involves violations of (b) and (c): $B \times G = \mathfrak{R}^r \times \mathfrak{R}^s$, but the function is concentrated on a small subset of its support that cannot be expressed as a Cartesian product, as illustrated in Figure 1(b). Whether these difficulties are present or not, the number of function evaluations and products required in any product rule increases geometrically with the number of arguments of the function, a phenomenon sometimes dubbed “the curse of dimensionality.”

These constitute the dominant problems for quadrature methods in economics. To a point, one may extend quadrature to higher dimensions using extensions more sophisticated than product rules. These extensions are usually specific to functions of a certain type, and for this reason the literature is large, but reliable software for a problem at hand may be hard to come by. For example, there has been considerable attention to monomials (polynomials for which the highest degree in any one product is bounded), e.g., McNamee and Stenger (1967), Genz and Malik (1983), Davis and Rabinowitz (1984, Section 5.7). Compound, or subregion, methods provide the most widely applied extensions of quadrature to higher dimensions. In these procedures, a finer and finer subdivision of the original integration region is dynamically constructed, with smaller subregions concentrated where the integrand is most irregular. Within each subregion, a local rule with a moderate number of points is used to approximate the integral. If, at a given step, a prespecified global convergence criterion is not satisfied, those regions for which the convergence criterion is farthest from being satisfied are subdivided, and the local rule is applied to the new subdivisions (van Dooren and de Ridder, 1976; Genz and Malik, 1980; Genz, 1991). For these procedures to work successfully, it is important to have a scheme for construction of subregions well suited to the problem at hand, as reconsideration of Figure 1(b) will make clear. For

example, Genz (1993) provides an algorithm that copes well with the isolated peaks in high-dimensional spaces often found in Bayesian multiparameter problems.

These extensions of quadrature are routinely successful for integrals through dimension four or five. Beyond four or five, success depends on whether the problem at hand is of a type for which existing subregion methods are well suited. Whereas the application of quadrature to a function of a single variable can be successful as a “black box” procedure, problems of dimensions three and four are more likely to require transformations or other analytical work before quadrature can be applied. There are very few applications of quadrature-based methods to integrals of more than five dimensions in the literature.

2.3 Low discrepancy methods

A low discrepancy method defines a deterministic sequence of points $\{\mathbf{x}_j\}_{j=1}^{\infty}$ and a corresponding m -point integration rule $m^{-1} \sum_{j=1}^m f(\mathbf{x}_j) \approx \int_B f(\mathbf{x}) d\mathbf{x}$. Gaussian quadrature organizes the choice of points to evaluate interactions of polynomials with basis functions exactly. Low discrepancy methods choose the sequence to minimize the difference between the number of points in a set and its measure. (The discussion here closely follows parts of Niederreiter, 1992, Chapters 2 and 3.)

The canonical problem sets $B = \bar{I}^d$, the d -dimensional hypercube. (This stipulation is less restrictive than it might seem, and we shall return to this point in an example in Section 4.4.) For arbitrary $S \subseteq B$ define

$$A(S; \{\mathbf{x}_j\}_{j=1}^m) = \sum_{j=1}^m \chi_S(\mathbf{x}_j),$$

where $\chi_S(\mathbf{x})$ is the characteristic function of S , $\chi_S(\mathbf{x}) = 1$ if $\mathbf{x} \in S$ and $\chi_S(\mathbf{x}) = 0$ if $\mathbf{x} \notin S$.

Thus $A(S; \{\mathbf{x}_j\}_{j=1}^m)$ is the counting function that indicates the number of j with $1 \leq j \leq m$ for which $\mathbf{x}_j \in S$. If \mathcal{S} is a nonempty family of Lebesgue measurable subsets of \bar{I}^d , then the *discrepancy* of the point set $\{\mathbf{x}_j\}_{j=1}^m$ is

$$D_m(\mathcal{S}; \{\mathbf{x}_j\}_{j=1}^m) = \sup_{S \in \mathcal{S}} \left| A(S; \{\mathbf{x}_j\}_{j=1}^m) / m - \lambda_d(S) \right|,$$

where $\lambda_d(\cdot)$ denotes d -dimensional Lebesgue measure. Let \mathcal{S}^* be the family of all subintervals of \bar{I}^d of the form $\prod_{i=1}^d [0, u_i]$. Then the *star discrepancy* of $\{\mathbf{x}_j\}_{j=1}^m$ is

$$D_m^*(\{\mathbf{x}_j\}_{j=1}^m) = D_m(\mathcal{S}^*; \{\mathbf{x}_j\}_{j=1}^m).$$

The star discrepancy of $\{\mathbf{x}_j\}_{j=1}^m$ may be used to bound the error of approximation of $\int_{\bar{I}^d} f(\mathbf{x}) d\mathbf{x}$ by $m^{-1} \sum_{j=1}^m f(\mathbf{x}_j)$. To do so, first define the *variation of f on \bar{I}^d in the sense of Vitali*,

$$V^{(d)}(f) = \int_0^1 \int_0^1 \left| \frac{\partial^d f}{\partial x_1 \partial \dots \partial x_d} \right| dx_1 \dots dx_d$$

for functions f for which the individual partial derivatives are continuous on \bar{I}^d . Next, let $V^{(k)}(f; i_1, \dots, i_k)$ be the variation in the sense of Vitali of the restriction of f to the k -dimensional face $\{(x_1, \dots, x_d) \in \bar{I}^d : x_j = 1 \text{ for } j \neq i_1, \dots, i_k\}$. The *variation of f on \bar{I}^d in the sense of Hardy and Krause* is

$$V(f) = \sum_{k=1}^d \sum_{1 \leq i_1 \leq \dots \leq i_k \leq d} V^{(k)}(f; i_1, \dots, i_k).$$

(See Niederreiter, 1992, Section 2.2, for an extension of this definition to functions f that are not d times continuously differentiable.) For any sequence $\{\mathbf{x}_j\}, \mathbf{x}_j \in \bar{I}^d$,

$$\left| m^{-1} \sum_{j=1}^m f(\mathbf{x}_j) - \int_{\bar{I}^d} f(\mathbf{x}) d\mathbf{x} \right| \leq V(f) D_m^*(\mathbf{x}_1, \dots, \mathbf{x}_m),$$

the Koksma-Hlawka inequality (Hlawka, 1961; Niederreiter, 1992, Theorem 2.11). The bound is strict (Niederreiter, 1992, Theorem 2.12).

Low discrepancy methods choose sequences $\{\mathbf{x}_j\}$ so as to minimize $D_m^*(\{\mathbf{x}_j\}_{j=1}^m)$. Intuitively, the star discrepancy can be kept small by spacing the points \mathbf{x}_j evenly. A naive grid on \bar{I}^d will achieve this, but requires an impractically large number of points for $d \geq 5$ in the same way as quadrature does. Low discrepancy methods substantially extend the range of practical d before succumbing to the curse of dimensionality. To describe two such sequences, begin with the unique base- b expansion of any integer n ,

$$n = \sum_{j=0}^{\infty} a_j(n) b^j,$$

where b is an integer exceeding 1 and $0 \leq a_j(n) < b$. The *radical-inverse function ϕ_b in base b* is defined by

$$\phi_b(n) = \sum_{j=0}^{\infty} a_j(n) b^{-(j+1)}.$$

This function maps the integers $1, \dots, m$ into m distinct points in the unit interval, maintaining a regular spacing between the points: if $m = b^k - 1$, k integer, then there are m evenly spaced points beginning with b^{-k} and ending with $1 - b^{-k}$. Let $\{b_j\}$ be a sequence of relatively prime integers all exceeding 1. (For example, $b_1 = 2, b_2 = 3, b_3 = 5, \dots$.) The *Halton sequence in bases b_1, \dots, b_d* is

$$\{\mathbf{x}_j\}_{j=1}^{\infty}, \quad \mathbf{x}_j = [\phi_{b_1}(j), \dots, \phi_{b_d}(j)]'$$

(Halton, 1960). The m -element *Hammersley sequence in bases* b_1, \dots, b_d is

$$\{\mathbf{x}_j\}_{j=1}^m, \quad \mathbf{x}_j = \left[j/m, \phi_{b_1}(j), \dots, \phi_{b_d}(j) \right]$$

(Hammersley, 1960). (An even earlier, closely related sequence is that of Richtmeyer, 1952, 1958, described in Hammersley and Handscomb, 1964.)

It may be shown (Niederreiter, 1992, Theorem 3.6) that for a Halton sequence in the pairwise relatively prime bases b_1, \dots, b_d ,

$$\begin{aligned} D_m^* \left(\{\mathbf{x}_j\}_{j=1}^m \right) &\leq \frac{d}{m} + \frac{1}{m} \prod_{j=1}^d \left(\frac{b_{j-1}}{2 \log b_j} \log m + \frac{b_j + 1}{2} \right) \\ &\leq \left(\prod_{j=1}^d \frac{b_{j-1}}{2 \log b_j} \right) m^{-1} (\log m)^d + 0 \left[m^{-1} (\log m)^{d-1} \right]. \end{aligned} \quad (2.3.1)$$

For the corresponding Hammersley sequence, there is the somewhat better bound

$$\begin{aligned} D_m^* \left(\{\mathbf{x}_j\}_{j=1}^m \right) &\leq \frac{d}{m} + \frac{1}{m} \prod_{j=1}^{d-1} \left(\frac{b_{j-1}}{2 \log b_j} \log m + \frac{b_j + 1}{2} \right) \\ &\leq \left(\prod_{j=1}^{d-1} \frac{b_{j-1}}{2 \log b_j} \right) m^{-1} (\log m)^{d-1} + 0 \left[m^{-1} (\log m)^{d-2} \right]. \end{aligned} \quad (2.3.2)$$

The second inequalities in (2.3.1) and (2.3.2) imply that the optimal bases are the primes themselves, $b_1 = 2, b_2 = 3, b_3 = 5, \dots$.

If the upper bounds in (2.3.1)-(2.3.2) are used to govern accuracy, then the number of function evaluations increases faster than geometrically with dimension, d , because of the presence of the term $\prod_{i=1}^d (b_i - 1)/2 \log b_i$ or $\prod_{i=1}^{d-1} (b_i - 1)/2 \log b_i$. Table 1 provides the number of evaluations required to assure that $\left| \sum_{j=1}^m f(\mathbf{x}_j) - \int_{\mathbb{I}^d} f(\mathbf{x}) d\mathbf{x} \right| \leq c$ ($c = 10^{-2}$ or 10^{-5}) for a function for which the Hardy-Krause total variation is d . It also provides the actual number of evaluations required to guarantee an approximation error of c or less for the function $f(\mathbf{x}) = \sum_{j=1}^d x_j$. While the upper bound on the number of evaluations required increases faster than exponentially in the dimension d , the actual number required increases not much faster than linearly and is *much* smaller. In general, however, one will not know the value of the actual error of approximation. The difficulty of assessing this error is a major disadvantage of low discrepancy and other deterministic algorithms for integration.

2.4 Other deterministic methods

In specialized settings integration in high dimensions can be made more tractable. The obvious limiting case is the one in which the entire problem may be solved analytically. But there are also classes of problems that cannot be solved analytically, with common features

that suggest specific approximations. An example is provided by Tierney and Kadane (1986) for a class of problems arising in Bayesian statistics and econometrics:

$$E_n(g) = \frac{\int_{\Theta} g(\theta) \exp[\mathfrak{l}(\theta)] \pi(\theta) d\theta}{\int_{\Theta} \exp[\mathfrak{l}(\theta)] \pi(\theta) d\theta} = \frac{\int_{\Theta} \exp[nL^*(\theta)] d\theta}{\int_{\Theta} \exp[nL(\theta)] d\theta},$$

where $\mathfrak{l}(\theta)$ is a log-likelihood function; $\pi(\theta)$ is a prior density kernel; $g(\theta)$ is a strictly positive function of interest; n is the number of observations entering the log-likelihood function; $L(\theta) = [\log \pi(\theta) + \mathfrak{l}(\theta)]/n$; and $L^*(\theta) = [\log g(\theta) + \log \pi(\theta) + \mathfrak{l}(\theta)]/n$.

Let $\hat{\theta}$ denote the mode of L , and let $\Sigma = \partial^2 L(\hat{\theta})/\partial\theta\partial\theta'$. Laplace's approximation is

$$\int_{\Theta} \exp[nL(\theta)] d\theta \approx \int_{\Theta} \exp\left[nL(\hat{\theta}) - \frac{1}{2}n(\theta - \hat{\theta})' \Sigma(\theta - \hat{\theta})\right] d\theta = (2\pi)^{k/2} |\Sigma|^{1/2} \exp[nL(\hat{\theta})].$$

Similarly, if $\hat{\theta}^*$ is the mode of L^* and $\Sigma^* = \partial^2 L^*(\hat{\theta}^*)/\partial\theta\partial\theta'$, then

$$\int_{\Theta} \exp[nL^*(\theta)] d\theta \approx (2\pi)^{k/2} |\Sigma^*|^{1/2} \exp[nL^*(\hat{\theta}^*)].$$

The error of approximation in each case is $O(n^{-1/2})$, but in the corresponding approximation

$$\hat{E}_n(g) = (|\Sigma^*|/|\Sigma|)^{1/2} \exp\left\{n[L^*(\hat{\theta}^*) - L(\hat{\theta})]\right\},$$

the leading terms in the numerator and denominator cancel, and the resulting error of approximation for $\hat{E}_n(g)$ is $O(n^{-1})$ (Tierney and Kadane, 1986).

The approximate solution provided by this method is a substantial improvement on previous approximations of this kind, which worked with a single expansion about $\hat{\theta}$. The method exhibits two attractions shared by most specialized approximations to integration in higher dimensions. First, it avoids the need for specific adaptive subregion analysis required for quadrature, if indeed quadrature can be made to work at all. Second, once function-specific code has been written, the computations involve standard ascent algorithms to find $\hat{\theta}$ and $\hat{\theta}^*$ and are usually extremely fast. This example also shares some limitations of this approach. First, reduction of approximation error through higher order approximation is tedious at best, whereas in quadrature one can increase the number of points or subregions used and in Monte Carlo one can increase the number of iterations. Second, there is no way to evaluate the error of approximation; again, quadrature and Monte Carlo will provide error estimates. Third, there is possibly time intensive analytical work required for each problem in forming derivatives for different g as well as different \mathfrak{l} . And finally, the requirement that g be strictly positive is restrictive. The method may be extended to more general functions at the cost of some increase in complexity (Tierney, Kass, and Kadane, 1989).

3. Pseudorandom number generation

The analytical properties of virtually all Monte Carlo methods for numerical integration, and more generally for simulation, are rooted in the assumption that it is possible to observe sequences of independent random variables, each distributed uniformly on the unit interval. Given this assumption, various methods, described in Section 3.2, may be used to construct random variables and vectors with more complex distributions. Specific transformations from the uniform distribution on the unit interval to virtually all of the classical distributions of mathematical statistics have been constructed using these methods. Some examples are reviewed in Sections 3.3 and 3.4. These distributions, in turn, constitute building blocks for the solutions of integration and simulation problems described subsequently in this chapter.

The assumption that it is possible to observe sequences of independent random variables, distributed uniformly or otherwise, constitutes a model or idealization of what actually occurs. In this regard it plays the same role here with respect to what follows as does the assumption of randomness in much of economic theory with respect to the derived implications for optimizing behavior or does the assumption of randomness with respect to the development of methods of statistical inference in econometrics. In current methods for pseudorandom number generation, the observed sequences of numbers for which the assumption of an i.i.d. uniform distribution on the unit interval is the model, are in fact deterministic. Since the algorithms that produce these observed sequences are known, the properties of the sequences may be studied analytically in a way that events in the real world corresponding to assumptions of randomness in economic models may not. Thus, the adequacy or inadequacy of stochastic independence as a model for these sequences is on a surer footing than is this assumption as a model in economic or econometric theory. We begin this section with an overview of current methods of generating sequences for which the independent uniform assumption should be an adequate model.

3.1 Uniform pseudorandom number generation

Virtually all pseudorandom number generators employed in practice are linear congruential generators and their elaborations. In the linear congruential generator a sequence of integers $\{J_i\}$ is determined by the recursion

$$J_i = (aJ_{i-1} + c) \bmod m. \quad (3.1.1)$$

The parameters a , c , and m determine the qualities of the generator. If $c = 0$, the resulting generator is a pure multiplicative congruential generator. For example, the multiplicative generator with $m = 2^{31} - 1 = 2147483647$ (a prime) and $a = 16807$, $a = 397204094$, or $a = 950706376$ is used in the IMSL scientific library (IMSL, 1994), and the user may

choose between different values of c as well as set the seed J_0 . The sequence $\{J_i\}$ is mapped into the pseudorandom uniform sequence $\{U_i\}$ by the transformation

$$U_i = J_i/m. \quad (3.1.2)$$

If m is prime, the sequence will cycle after producing exactly m distinct values; clearly one can do no better than $m = 2^{31} - 1$ for a sequence of positive integers with 32-bit arithmetic. There are many criteria for evaluating the i.i.d. uniform distribution on the unit interval as a model for the resulting sequences $\{U_i\}$. Informal but useful discussions are provided by Press *et al.* (1986, pp. 192-194) and Bratley, Fox and Schrage (1987, pp. 216-220). More technical and detailed evaluations, including discussion of the choice of c , may be found in Coveyou and McPherson (1967), Marsaglia (1972), Knuth (1981), and Fishman and Moore (1982, 1986).

There are many elaborations on pseudorandom number generation that build on the primitive of the linear or multiplicative congruential generator. In the shuffled generator, a table is initialized with q seeds. The generator is then used in the obvious way to select a table entry pseudorandomly, and J_1 and U_1 are generated as described in the preceding paragraph. Then a new entry is selected pseudorandomly, U_2 is generated from that entry, and so on. If the congruential generator produced i.i.d. uniform random variables, so would the shuffled generator, and shuffled generators extend the upper bound on cycle length to mq ; this option is provided conveniently in IMSL. A shuffled generator described by L'Ecuyer (1986) has cycle length over 10^{19} . However, the analytical properties of the shuffled generator are harder to evaluate. In another elaboration on the basic approach, one may combine two pseudorandom sequences $\{J_i\}$ and $\{K_i\}$ from the congruential generator to produce a third sequence $\{L_i\}$ that is then mapped into U_i , $U_i = L_i/m$, in one of two ways: (a) Let $L_i = (J_i + K_i) \bmod m$, or (b) use $\{K_i\}$ to randomly shuffle $\{J_i\}$ and then set $\{L_i\}$ to the shuffled sequence. Both of these generators extend cycle length, but subtle issues arise in the combination of sequences. For a discussion of these issues and comparison of properties, consult Wichmann and Hill (1982) or L'Ecuyer (1986) for (a), Marsaglia and Bray (1968) or Knuth (1981, p. 32) for (b).

The add with carry generator (Marsaglia and Zaman, 1991) has a base b , lags r and s ($r > s$), and a seed vector $\mathbf{j}' = (j_1, \dots, j_r, c)$ with integer elements $j_i: 0 \leq j_i < b$ ($i = 1, \dots, r$) and carry bit $c = 0$ or 1 . The generated sequence is $\mathbf{j}, f(\mathbf{j}), f[f(\mathbf{j})], \dots$ with

$$f(j_1, \dots, j_r, c) = \begin{cases} (j_2, \dots, j_r, j_{r+1-s} + j_1 + c, 0) & \text{if } j_{r+1-s} + j_1 + c < b \\ (j_2, \dots, j_r, j_{r+1-s} + j_1 + c - b, 1) & \text{if } j_{r+1-s} + j_1 + c \geq b \end{cases}$$

With appropriately chosen base b , lags r and s , and seed vector \mathbf{j} , the generated sequence has period $b^r + b^s - 2$. Marsaglia and Zaman (1991) discuss appropriate choices of these

values. One example is $b = 2^{32} - 5$, $r = 43$, $s = 22$, and seed vector consisting of any 43 integers in $[0, 2^{32} - 6]$. The sequence of vectors has a cycle exceeding 10^{414} , and all possible sequences of 43 integers appear within a cycle. (The add with carry generator is one of a family of closely related generators. Marsaglia and Zaman, 1991, discuss the family.)

Since pseudorandom numbers are in fact deterministic, some consideration must be given to systematic differences between the two. One important quality is the cycle length. Most simulations on personal computers or workstations are unlikely to exceed the cycle length of 2^{31} of typical good linear congruential generators. But a study carried out with vector or parallel processors could well exceed this length, and in such cases the shuffled or add with carry generator should be considered. Another quality is absence of serial correlation. This is easily tested but generally is not a problem. Gre enberger (1961) shows that the first order serial correlation coefficient of any linear congruential generator is bounded above by $a^{-1}[1 - (6c/m) + 6(c/m)^2] + (a + 6)/m$, and Knuth, 1981, p. 84, points out that for nearly all m the serial correlation coefficient is less than $1/\sqrt{m}$.

Evidence of pseudorandomness is usually exhibited in high dimensional spaces. If one plots successive overlapping sequences of n pseudorandom numbers, then the sequences typically lie in a few hyperplanes of dimension $n - 1$ each. For example, in the case of linear congruential generators the number of hyperplanes is no more than $(n!/m)^{1/n}$ (Marsaglia, 1968): e.g., if $m = 2^{31} - 1$, then sequences of length 6 lie on at most 108 distinct hyperplanes. In the add with carry generator, successive overlapping sequences of more than r values lie on hyperplanes with a separating distance is at least $1/\sqrt{3}$ (Tezuka *et al.*, 1993). One can determine the existence of such hyperplanes using the spectral test first proposed in Coveyou and MacPherson (1967). Accessible descriptions of this test are provided in Knuth (1981) and Bratley, Fox, and Schrage (1987). Most simulation methods employ highly nonlinear transformations of $\{U_i\}$, as we shall see subsequently, so the distribution of sequences on hyperplanes does not carry over. (However, new problems can arise: see the discussion below of the Box and Muller transformation to construct normally distributed random variables.)

A few practical steps will avoid most problems. First, use only uniform pseudorandom number generators that are completely documented with references to the academic literature. Second, questions of execution time, often discussed in the academic literature, are irrelevant in computational economics: subsequent computations using pseudorandom uniform random sequences take much longer than the most elaborate variants on linear congruential generators, so that even if execution time for these generators could be driven

to zero, there would be no significant improvement in overall execution time. Third, one should ensure that cycle length is substantially greater than the length of the pseudorandom sequence to be generated. Finally, any publicly reported result based in part on a sequence of pseudorandom numbers should be checked for sensitivity to the choice of generator. This does not imply numerical analysis that takes the investigator far from the problem of interest. A key advantage of Monte Carlo methods, to be discussed in Section 4, is that measures of accuracy are produced as a by-product based on the assumption that successive pseudorandom numbers are independently and identically distributed. Results obtained using variants of methods for producing these sequences should agree within these measures of accuracy. For example, computations can be executed with different seeds, with different values of c in (3.1.1), with or without shuffling, or using an add with carry or related generator. This requires only minor changes in code for most software.

3.2 General methods for nonuniform distributions

Throughout this section, x will denote a random variable with cumulative distribution function (c.d.f.) F and support C , and u will denote a random variable with uniform distribution on the unit interval. If x is continuous, its probability density function (p.d.f.) will be denoted by f . We turn first to several general methods for mapping u into x .

Inverse c.d.f. Suppose x is continuous, and consequently the inverse c.d.f.

$$F^{-1}(p) = \{c: P(x \leq c) = p\}$$

exists. Then x and $F^{-1}(u)$ have the same distribution: $P[F^{-1}(u) \leq d] = P[u \leq F(d)] = F(d)$. Hence pseudorandom drawings $\{x_i\}_{i=1}^N$ of x may be constructed as $F^{-1}(u_i)$, where $\{u_i\}_{i=1}^N$ is a sequence of pseudorandom uniform numbers.

A simple example is provided by the exponential distribution with probability density $f(x) = \lambda \exp(-\lambda x)$, $x \geq 0$. Correspondingly, $F(x) = 1 - \exp(-\lambda x)$, $F^{-1}(p) = -\log(1 - p)/\lambda$, and consequently, $x = -\log(u)/\lambda$.

The inverse c.d.f. method is very easy to apply if an explicit, closed form expression for the inverse c.d.f. is available. Since most inverse c.d.f.'s require the evaluation of transcendental functions, the method may be inefficient relative to others. (That is the case in the foregoing example; see von Neumann, 1951, or Forsythe, 1972, for a more efficient alternative.) In some cases, evaluation of the c.d.f. is superficially closed form to the user of a mathematical software library but in fact involves nontrivial numerical integration of the kind discussed in Section 2. A leading example is provided by the standard normal distribution, for which specialized methods can be applied to the computation of F^{-1} (Hart

et al., 1968; Strecok, 1968), but for which acceptance and composition methods (discussed below) are more efficient.

Discrete distributions. Suppose that the random variable X takes on a finite number of values, without loss of generality the integers $1, \dots, n$ and $P(X = i) = p_i$. The preferred methods will depend (among other things) on the number of draws to be made from the distribution. If only a few draws are to be made (as may be the case with the Markov chain Monte Carlo methods discussed in Section 6), then the obvious inverse mapping from the unit interval to the integers $1, \dots, n$ can be constructed and subsequently used to search for the appropriate integer corresponding to the drawn u . The disadvantage of this method is that the search time can be substantial. If many draws are to be made, then the alias method due to Walker (1974) and refined by Walker (1977) and Kronmal and Peterson (1979) is more efficient. The basic idea is to draw an integer i from an equiprobable distribution on the first n integers, and choose i with probability r_i and its alias a_i with probability $1 - r_i$. If the values of a_i and r_i are chosen correctly, then the resulting choice probabilities are p_i for i ($i = 1, \dots, n$). Setting up the table of r_i and a_i requires $O(n)$ time (see Bratley, Fox, and Schrage, 1987, pp. 158-160, for an accessible discussion); whether this overhead is worthwhile depends on the value of n and the number of draws to be made from the discrete distribution. The aliasing algorithm is implemented in many mathematical software libraries.

Acceptance methods. Suppose that x is continuous with p.d.f. $f(x)$ and support C . Let g be the p.d.f. of a different continuous random variable z with p.d.f. $g(z)$ which has a distribution from which it is possible to draw i.i.d. random variables and for which

$$\sup_{x \in C} [f(x)/g(x)] = a < \infty.$$

The function g is known as an *envelope* or *majorizing density* of f , and the distribution with p.d.f. g is known as the *source distribution*. To generate x_i ,

- (a) Generate u ;
- (b) Generate z ;
- (c) If $u > f(z)/[ag(z)]$, go to (a);
- (d) $x_i = z$.

The unconditional probability of proceeding from step (c) to step (d) in any pass is

$$\int_{-\infty}^{\infty} \{f(z)/[ag(z)]\} g(z) dz = a^{-1},$$

and the unconditional probability of reaching step (d) with value at most c in any pass is

$$\int_{-\infty}^c \{f(z)/[ag(z)]\} g(z) dz = a^{-1} F(c).$$

Hence the probability that x_i is at most c at step (d) is $F(c)$.

The principle of acceptance sampling is illustrated in Figure 2. The two essentials of applying this procedure are the ability to generate z and the finite upper bound on $f(x)/g(x)$. The efficiency of the method depends on the efficiency of generating z and the unconditional probability of acceptance, which is just the inverse of the upper bound on $f(x)/g(x)$. (In this respect, acceptance sampling is closely related to importance sampling discussed in Section 4.3.) The great advantage of acceptance sampling is its ability to cope with arbitrary probability density functions as long as the two essential conditions are met and efficiency is acceptable for the purposes at hand. Notice that the method will work in exactly the same way if $f(x)$ is merely the kernel of the p.d.f. of x (i.e., proportional to the p.d.f.) as long as $a = \sup_{x \in C} [f(x)/g(x)]$ (although in this case a^{-1} no longer provides the unconditional acceptance probability). This property can be exploited to advantage to avoid numerical approximation of unknown constants of integration.

Specific examples providing insight into the method may be found in the family of truncated univariate normal distributions. As a first example, consider the standard normal probability distribution truncated to the interval $(0, .5)$:

$$f(x) = (.19146)^{-1} (2\pi)^{-1/2} \exp(-x^2/2) = 2.0837 \exp(-x^2/2), 0 < x \leq .5.$$

The standard normal distribution itself is a legitimate source distribution, but since $\sup_{0 < x \leq .5} [f(x)/g(x)] = (.19146)^{-1}$, the efficiency of this method is low. However, for a source distribution uniform on $(0, .5]$, $\sup_{0 < x \leq .5} [f(x)/g(x)] = 2.0837/2.0 = 1.0418$: the unconditional probability of acceptance is $(1.0418)^{-1} = .95985$. As a second example, consider the same distribution truncated to the interval $(5, 8]$:

$$f(x) = (2.8665 \times 10^{-7})^{-1} (2\pi)^{-1/2} \exp(-x^2/2) = (1.3917 \times 10^6) \exp(-x^2/2), 5 < x \leq 8.$$

The standard normal fails as a source distribution since the acceptance probability is 2.8665×10^{-7} . A uniform source density yields an acceptance probability of only .064271. An exponential distribution translated to the truncation point is for many purposes an excellent approximation to a severely truncated normal distribution (Marsaglia, 1964; Geweke, 1986), and for the exponential source density, setting the parameter equal to the truncation point is an optimal or near optimal choice (Geweke, 1991). One can readily verify that the acceptance probability for the source density

$$g(x) = 5 \exp[-5(x - 5)], 5 < x \leq 8,$$

is .96406.

Optimizing acceptance sampling. Acceptance methods may readily be extended to multivariate distributions. This topic is taken up in detail in Section 4.2. We turn now to

the question of finding an optimal source distribution for a specified problem and develop results for the general case of univariate or multivariate distributions .

In general, suppose that it is desired to draw i.i.d. variables from a distribution with target density kernel $f(\mathbf{x};\theta), \theta \in \Theta$, having support $C(\theta) \subseteq \mathfrak{R}^m$; the parameter vector θ indexes a family of density kernels $f(\cdot)$. Suppose that a family of source distributions with densities $g(\mathbf{x};\alpha), \alpha \in A \subseteq \mathfrak{R}^p$, having support $D(\alpha)$, has been identified, with the property that for all $\theta \in \Theta$, there exists at least one α for which $\sup_{\mathbf{x} \in C(\theta)} f(\mathbf{x};\theta)/g(\mathbf{x};\alpha) < \infty$. To accomplish the goal of i.i.d. sampling from $f(\mathbf{x};\theta)$, draws from $g(\mathbf{x};\alpha)$ are retained with probability $q(\alpha, \theta) f(\mathbf{x};\theta)/g(\mathbf{x};\alpha)$, where

$$q(\alpha, \theta) \equiv \left[\sup_{\mathbf{x} \in C(\theta)} f(\mathbf{x};\theta)/g(\mathbf{x};\alpha) \right]^{-1}.$$

Suppose the family of source densities $g(\cdot; \cdot)$ has been fixed, but not the value of α , and that the objective is to maximize the unconditional probability of accepting the draw from the source distribution. Just as in the foregoing examples, this unconditional probability is

$$\int_{D(\alpha)} [q(\alpha, \theta) f(\mathbf{x};\theta)/g(\mathbf{x};\alpha)] g(\mathbf{x};\alpha) d\mathbf{x} = q(\alpha, \theta).$$

Hence the problem is to determine the saddle point

$$\min_{\alpha \in A} \left\{ \max_{\theta \in \Theta} [\log f(\mathbf{x};\theta) - \log g(\mathbf{x};\alpha)] \right\}.$$

Given the usual regularity conditions, a necessary condition is that α be part of a solution of the $(m + p)$ -equation system

$$\begin{aligned} \partial [\log f(\mathbf{x};\theta) - \log g(\mathbf{x};\alpha)] / \partial \mathbf{x} &= \mathbf{0} \\ \partial \log g(\mathbf{x};\alpha) / \partial \alpha &= \mathbf{0}. \end{aligned}$$

As an example, consider the target density kernel

$$f(x; T, \eta) = (x/2)^{Tx/2} [\Gamma(x/2)]^{-T} \exp(-\eta x),$$

which arises as a conditional posterior density kernel for the degrees-of-freedom parameter in a Student- t distribution (Geweke, 1992b, Appendix B). For the exponential family of source densities $g(x; \alpha) = \alpha \exp(-\alpha x)$, the regular necessary conditions are that

$$\begin{aligned} (T/2) [\log(x/2) + 1 - \psi(x/2)] + (\alpha - \eta) &= 0, \\ \alpha^{-1} - x &= 0, \end{aligned}$$

where $\psi(\cdot) = \Gamma'(\cdot)/\Gamma(\cdot)$ is the digamma function. The desired value of α is the solution of

$$(T/2) [-\log(2\alpha) + 1 - \psi(1/2\alpha)] + (\alpha - \eta) = 0,$$

which may be found using standard root-finding algorithms. Acceptance rates of about .15 are reported in Geweke (1992b).

Adaptive methods. It may be possible to improve upon a source distribution, using information about the target distribution acquired in the sampling process itself. A very useful application of this idea has been made to the problem of sampling from distributions with log-concave probability density functions. It is especially attractive when it is costly to evaluate the target density kernel at a point or when known source densities are inefficient or nonexistent. The exposition here closely follows Gilks and Wild (1992), who build on some earlier work by Devroye (1986); see Wild and Gilks (1993) for a published algorithm. An application of this algorithm is discussed in Section 7.1.

Let $h(x) = \log f(x)$. The support D of $f(x)$ is connected, and $h(x)$ is differentiable and weakly concave everywhere in D ; i.e., $h'(x)$ is monotonically nonincreasing in x on D .

Suppose that $h(x)$ and $h'(x)$ have been evaluated at k points in D , $x_1 \leq x_2 \leq \dots \leq x_k$, $k \geq 2$. We assume that if D is unbounded below, then $h'(x_1) > 0$ and that if D is unbounded above, then $h'(x_k) < 0$. Let the piecewise linear upper hull $u(x)$ of $h(x)$ be formed from the tangents to $h(\cdot)$ at the x_j , as shown in Figure 3. For $j = 1, \dots, k-1$ the tangents at x_j and x_{j+1} intersect at

$$w_j = \frac{h(x_{j+1}) - h(x_j) - x_{j+1} h'(x_{j+1}) + x_j h'(x_j)}{h'(x_j) - h'(x_{j+1})}.$$

Further let w_0 denote the lower bound of D (possibly $-\infty$) and w_k the upper bound of D (possibly $+\infty$). Then

$$u(x) = h(x_j) + (x - x_j)h'(x_j), \quad x \in (w_{j-1}, w_j].$$

Similarly the piecewise linear lower hull $l(x)$ of $h(x)$ is formed from the chords between the x_j ,

$$l(x) = \frac{(x_{j+1} - x)h(x_j) + (x - x_j)h(x_{j+1})}{x_{j+1} - x_j}, \quad x \in (x_j, x_{j+1}].$$

For subsequent purposes it is useful to extend the definition to include

$$l(x) = -\infty, \quad x < x_1 \text{ or } x > x_k.$$

At the start of an acceptance/rejection iteration, the function $\exp[u(x)]$ forms a source density kernel, and $\exp[l(x)]$ is a squeezing density kernel. The iteration begins by drawing a value z from the distribution with kernel density function $\exp[u(x)]$. This may be done in two steps:

- (a) Compute $p_j = P(w_{j-1} < x \leq w_j) = I_j / I$ ($j = 1, \dots, k$), where

$$I_j = \begin{cases} \exp[h(x_j) - x_j h'(x_j)] \exp[h'(x_j)(w_j - w_{j-1})] / h'(x_j) & \text{if } h'(x_j) \neq 0 \\ h(x_j)(w_j - w_{j-1}) & \text{if } h'(x_j) = 0 \end{cases}$$

and $I = \sum_{j=1}^k I_j$. Choose an interval $(w_{j-1}, w_j]$ from this discrete distribution as described above.

- (b) Conditional on the choice of interval the source distribution is exponential. Draw z from this distribution as previously discussed.

The draw z is accepted or rejected by means of the acceptance sampling algorithm described above, but using the following shortcut. Having drawn u , we know that z will be accepted if $u \leq \exp[\mathbb{1}(z) - u(z)]$, and in this case no further computations are required. If $u > \exp[\mathbb{1}(z) - u(z)]$, then evaluate $h(z)$ and $h'(z)$ and accept z if and only if $u \leq \exp[h(z) - u(z)]$. In the latter case add z to the set of points (x_1, \dots, x_k) , reordering the x_j 's, and update $u(\cdot)$ and $\mathbb{1}(\cdot)$, unless z is accepted and no more draws from the target distribution are needed. This completes the acceptance iteration.

Notice that this algorithm is more likely to update the source and squeezing densities the more discordant are these functions at a point. As the algorithm proceeds, the probability of acceptance of any draw increases toward 1, and the probability that an evaluation of h will be required for any draw falls to 0.

Composition algorithms. Formally, composition arises from a p.d.f. representation

$$f(x) = \int_{-\infty}^{\infty} g_y(x) dH(y).$$

A random variable Y from distribution H is generated, followed by a random variable X with p.d.f. g_y . This method goes back at least to Marsaglia (1961), who used it to generate normal random variables. It is also the natural method to use for mixture distributions. For example, suppose that x is drawn from a $N(0, .1^2)$ distribution with probability .95 and a $N(0, 10^2)$ distribution with probability .05. The probability density

$$.95(2\pi)^{-1/2} (.1)^{-1} \exp(-x^2/.02) + .05(2\pi)^{-1/2} (10)^{-1} \exp(-x^2/200)$$

is strongly leptokurtic and not well suited to acceptance sampling. But the construction of the random variable in fact corresponds to a composition with

$$P(Y = 0) = .95, \quad P(Y = 1) = .05, \\ g_{Y=0}(x) = (2\pi)^{-1/2} (.1)^{-1} \exp(-x^2/.02), \quad g_{Y=1}(x) = .05(2\pi)^{-1/2} (10)^{-1} \exp(-x^2/200).$$

3.3 Selected univariate distributions

In most cases there is associated with each of the classical univariate distributions a substantial literature on the generation of corresponding pseudorandom variables. Good mathematical and statistical software libraries have drawn on this literature and are widely available. In many cases the most efficient and accurate routines are not simply implementations of the constructions that appear in the mathematical statistics literature, and

the user is well-advised to take advantage of the capital embodied in good libraries. The discussion here is limited to illustrating how the techniques discussed in Section 3.2 are used in specific cases. More thorough surveys in the literature are provided by Bratley, Fox, and Schrage (1987, pp.164-189) and Devroye (1986). All of the methods discussed here are implemented in good software libraries, which should always be used. This discussion is not intended to form the basis of reliable code.

Binomial distribution. The binomial distribution indicates the probability of k successes in n independent trials if p is the probability of success in any given trial:

$$p(k) = \binom{n}{k} p^k (1-p)^{(n-k)}.$$

The definition provides a direct method for generating the random variable k , but is acceptably rapid only if n is small. For small values of np , the inverse c.d.f. method is practical since $p(k)$ will typically require evaluation for only a few values of k . In all other cases, however, composition algorithms with acceptance methods are more efficient. Examples are given by Ahrens and Dieter (1980) and Kachitvichyanukul (1982).

Univariate normal distributions. Inverse c.d.f methods for the standard normal have already been mentioned. Acceptance sampling methods are not hard to design, especially if one exploits the exponential source distribution as first noted by Marsaglia (1964). Related and succeeding work by Marsaglia and Bray (1964); Marsaglia, MacLaren, and Bray (1964); and Kinderman and Ramage (1976) combining acceptance sampling and composition form the basis for the generation of standard normal variables in most software libraries.

Box and Muller (1958) showed that if U_1 and U_2 are mutually independent standard uniform random variables, then

$$X = \cos(2\pi U_1) \sqrt{-2 \log U_2}, \quad Y = \sin(2\pi U_1) \sqrt{-2 \log U_2}$$

are independent standard normal random variables. (The key to the demonstration lies in a transformation to polar coordinates.) The combination of this method with the linear congruential random number generator produces a pathology, however. If U_i and U_{i+1} are successive realizations of (3.1.1)-(3.1.2), then

$$U_{i+1} = [(amU_i + c) \bmod m] / m \Rightarrow$$

$$\cos(2\pi U_{i+1}) = \cos[2\pi(aU_i + c/m)], \quad \sin(2\pi U_{i+1}) = \sin[2\pi(aU_i + c/m)]$$

and hence

$$X_i = \cos[2\pi(aU_i + c/m)] \sqrt{-2 \log U_i}, \quad Y_i = \sin[2\pi(aU_i + c/m)] \sqrt{-2 \log U_i}.$$

All possible values of (X_i, Y_i) fall on a spiral. As an approximation to a *pair* of independent variables the distribution of (X_i, Y_i) could hardly be worse. However, if one discards Y_i , the sequence $\{X_i\}$ suffers from no known problems of this kind. This is one of the reasons that acceptance sampling and composition rather than the Box-Muller transformation is used in statistical libraries. It illustrates the risks involved in seemingly straightforward combinations of distribution theory with pseudorandom uniform variables.

Given a sequence of standard normal random variables $\{z_i\}$, a sequence from the general univariate normal distribution $N(\mu, \sigma^2)$ can be generated through the familiar transformation $x_i = \mu + \sigma z_i$.

Gamma distributions. The gamma distribution is important in its own right, for included special cases like the chi-squared, and as a building block for other distributions like the beta. The gamma distribution with scale parameter λ and shape parameter a has probability density

$$f(x) = \lambda \exp(-\lambda x) (\lambda x)^{a-1} / \Gamma(a), \quad x \geq 0.$$

In general, random variables from this distribution may be generated efficiently using composition algorithms and acceptance methods. Fast and accurate methods are complicated but readily available in statistical software libraries. For example, IMSL uses the composition-acceptance methods of Ahrens and Dieter (1974) and Schmeiser and Lal (1980). A few special cases are worth note.

- (a) If $a = 1$, then the distribution is exponential with parameter λ and the inverse c.d.f. method discussed above is much more efficient.
- (b) If $a = 0.5$, then $x = z^2/2$, $z \sim N(0, \lambda^2)$.
- (c) If $\lambda = 0.5$, then $x \sim \chi^2(v)$, $v = 2a$. If a is an integer, then x is the sum of a independent exponentially distributed random variables each with parameter $\lambda = 0.5$. If v is an odd integer, then x is the sum of $[v/2]$ independent exponentially distributed random variables plus the square of an independent standard normal. For integers up to $v = 17$, these representations provide the basis for more efficient generation from the chi-squared distribution, but for larger integers it is more efficient to use the more general composition-acceptance methods.

3.4 Selected multivariate distributions

Generation of random vectors typically builds upon the ability to generate univariate random variables. Just how this should be done is not always obvious, however, and

sometimes the obvious method is not the most efficient. The examples that follow are intended only to illustrate this fact. Statistical software libraries should be consulted for implementation of these methods.

Multinomial distribution. The multinomial distribution indicates the probability of k_j realizations of outcome j , from m possible outcomes, in n independent trials. If p_j is the probability of outcome j in any given trial, then

$$p(k_j) = \frac{n!}{\prod_{j=1}^m k_j!} \prod_{j=1}^m p_j^{k_j}, \quad k_j \geq 0 \text{ and } \sum_{j=1}^m k_j = n.$$

The decomposition of this distribution into its full conditionals, $p(k_1)$, $p(k_2|k_1)$, $p(k_3|k_1, k_2)$, $p(k_4|k_1, k_2, k_3)$, ..., $p(k_m|k_1, k_2, \dots, k_{m-1})$, may be used to generate the k_j . We have

$$p(k_1) = \binom{n}{k_1} p_1^{k_1} (1-p_1)^{(n-k_1)}, \quad 0 \leq k_1 \leq n,$$

$$p(k_j|k_1, \dots, k_{j-1}) = \binom{\tilde{n}_j}{k_j} \tilde{p}_j^{k_j} (1-\tilde{p}_j)^{(\tilde{n}_j-k_j)}, \quad 0 \leq k_j \leq \tilde{n}_j,$$

$$\text{where } \tilde{n}_j \equiv n - \sum_{i=1}^{j-1} k_i, \quad \tilde{p}_j \equiv p_j / \left(1 - \sum_{i=1}^{j-1} p_i\right).$$

These distributions are all binomial.

Multivariate normal distribution. The generation of a multivariate normal random vector \mathbf{x} from the distribution $N(\boldsymbol{\mu}, \boldsymbol{\Sigma})$ is based on the familiar decomposition

$$\mathbf{z} \sim N(\mathbf{0}, \mathbf{I}_m), \quad \mathbf{x} = \boldsymbol{\mu} + \mathbf{A}\mathbf{z} \text{ with } \mathbf{A}\mathbf{A}' = \boldsymbol{\Sigma}.$$

While any factorization \mathbf{A} of $\boldsymbol{\Sigma}$ will suffice, it is most efficient to make \mathbf{A} upper or lower triangular so that $m(m+1)/2$ rather than m^2 products are required in the transformation from z to x . The Cholesky decomposition, in which the diagonal elements of the upper or lower triangular \mathbf{A} are positive, is typically used.

Wishart distribution. If $\mathbf{x}_{i1} \stackrel{iid}{\sim} N(\mathbf{0}, \boldsymbol{\Sigma})$, the distribution of $\mathbf{A} = \sum_{i=1}^n (\mathbf{x}_i - \bar{\mathbf{x}})(\mathbf{x}_i - \bar{\mathbf{x}})'$ is Wishart, with p.d.f.

$$f(\mathbf{A}) = \frac{|\mathbf{A}|^{\frac{1}{2}(n-m)} \exp\left(-\frac{1}{2} \text{tr } \boldsymbol{\Sigma}^{-1} \mathbf{A}\right)}{2^{\frac{1}{2}(n-1)m} \pi^{m(m-1)/4} |\boldsymbol{\Sigma}|^{\frac{1}{2}(n-1)} \prod_{i=1}^m \Gamma\left[\frac{1}{2}(n-i)\right]};$$

for brevity, $\mathbf{A} \sim W(\boldsymbol{\Sigma}, n-1)$. (For obvious reasons this distribution arises frequently in simulations. It is also important in Bayesian inference, where the posterior distribution of the inverse of the variance matrix for a normal population often has this form.) Direct

construction of \mathbf{A} through generation of $\{\mathbf{x}_i\}_{i=1}^n$ becomes impractical for large n . A more efficient indirect method follows Anderson (1984). Let Σ have lower triangular Choleski decomposition $\Sigma = \mathbf{L}\mathbf{L}'$, and suppose $\mathbf{Q} \sim \mathbf{W}(\mathbf{I}_m, n-1)$. Then $\mathbf{L}\mathbf{Q}\mathbf{L}' \sim \mathbf{W}(\Sigma, n-1)$ (Anderson, 1984, pp. 254-255). Furthermore \mathbf{Q} has representation

$$\mathbf{Q} = \mathbf{U}\mathbf{U}' \quad u_{ij} = 0 \quad (i < j < m)$$

$$u_{ij} \sim \mathbf{N}(0,1) \quad u_{ii} \sim \chi^2(n-i)$$

($i = 1, \dots, m$), with the u_{ij} mutually independent for $i \geq j$ (Anderson, 1984, p. 247). Even if n is quite small, this indirect construction is much more efficient than the direct construction.

4. Independence Monte Carlo

Building on the ability to produce sequences of vectors that are well described as i.i.d. random variables, we return to the integration problem with particular attention to high dimensions. There are two distinct but closely related problems that arise in economics and econometrics.

Problem I is to evaluate

$$I = \int_D f(\mathbf{x}) d\mathbf{x}.$$

Problem E is to evaluate

$$E = E[g(\mathbf{x})],$$

where \mathbf{x} is a random vector with c.d.f. $P(\mathbf{x})$. To simplify notation, assume that P is absolutely continuous and that \mathbf{x} has a probability density function $p(\mathbf{x})$. It is implicit in Problem E that $\int_D g(\mathbf{x})p(\mathbf{x})d\mathbf{x}$ is absolutely convergent in its domain D . Detailed examples of Problems E and I are provided in Section 7.

If a random vector \mathbf{z} has p.d.f. $p(\mathbf{z})$, then any function $r(\mathbf{z}) = a \cdot p(\mathbf{z})$, $a > 0$, is said to be a kernel density function for \mathbf{z} . In order to express some key moments compactly, let $E_r[g(\mathbf{z})]$ denote the expectation of $g(\mathbf{z})$ if \mathbf{z} has kernel density function $r(\mathbf{z})$; similarly $\text{var}_r[g(\mathbf{z})]$ for variance.

Many of the procedures discussed in this section are straightforward applications of two results in basic mathematical statistics. Let $\{y_i\}$ be an i.i.d. sequence from a population, and let $\bar{y}_N = \frac{1}{N} \sum_{i=1}^N y_i$ and $s_N^2 = \frac{1}{N-1} \sum_{i=1}^N (y_i - \bar{y}_N)^2$. If the population has finite first moment, then $E(\bar{y}_N) = E(y)$ and the strong law of large numbers states that

$$\bar{y}_N \xrightarrow{a.s.} E(y);$$

i.e., $P[\lim_{N \rightarrow \infty} \bar{y}_N = E(y)] = 1$. If the same population also has a finite variance σ^2 , then the central limit theorem establishes that

$$\sqrt{N}[\bar{y}_N - E(y)] \xrightarrow{d} N(0, \sigma^2);$$

i.e., $\lim_{N \rightarrow \infty} P\{\sqrt{N}[\bar{y}_N - E(y)] \leq c\sigma\} = \Phi(c)$, where $\Phi(\cdot)$ is the c.d.f. of the $N(0,1)$ distribution. In this case $E(s_N^2) = \sigma^2$, and from the strong law of large numbers,

$$s_N^2 \xrightarrow{a.s.} \sigma^2.$$

4.1 Simple Monte Carlo

In the case of Problem I, suppose that

$$f(\mathbf{x}) = g(\mathbf{x})p(\mathbf{x}),$$

with $p(\mathbf{x}) \geq 0$ and $\int_D p(\mathbf{x}) d\mathbf{x} = p^*$, where p^* is a known positive constant. Then $p(\mathbf{x})$ is a kernel density function. Suppose further that it is possible to draw ps eudorandom vectors $\{\mathbf{x}_i\}$ from the distribution with probability density function $p(\mathbf{x})/p^*$, as described in Section 3. Since

$$I = \int_D f(\mathbf{x}) d\mathbf{x} = \int_D p^* g(\mathbf{x}) [p(\mathbf{x})/p^*] d\mathbf{x} = E_p[p^* g(\mathbf{x})],$$

it follows that

$$I_N = N^{-1} p^* \sum_{i=1}^N g(\mathbf{x}_i) \xrightarrow{a.s.} I. \quad (4.1.1)$$

The requirement that p^* is known may be weakened by replacing p^* with a sequence $p_N^* \xrightarrow{a.s.} p^*$ in the last expression. (Some practical methods of producing p_N^* at essentially no incremental cost are taken up in Section 4.2.) If p^* is known, then $E(I_N) = I$, but if p^* must be replaced by a consistent estimator, then in general $E(I_N) \neq I$ but (4.1.1) is still true.

If in addition $\int_D g^2(\mathbf{x}) p(\mathbf{x}) d\mathbf{x}$ is absolutely convergent, this result can be extended to provide a measure of the accuracy of I_N . Let

$$\sigma^2 = \text{var}_p[p^* g(\mathbf{x})] = p^{*-1} \int_D [p^* g(\mathbf{x}) - I]^2 p(\mathbf{x}) d\mathbf{x}.$$

Then

$$\sqrt{N}(I_N - I) \xrightarrow{d} N(0, \sigma^2), \quad N^{-1} \sum_{i=1}^N [p^* g(\mathbf{x}_i) - I_N]^2 \xrightarrow{a.s.} \sigma^2.$$

(The result may be extended to include cases in which p^* is approximated by a sequence of p_N^* , but some changes are required; see Section 4.2.) This result makes exact the intuitive notion that $p(\cdot)$ should be chosen to mimic the shape of $f(\cdot)$.

The solution of Problem E by simple Monte Carlo is even simpler, as long as it is possible to construct an i.i.d. sequence from the probability distribution of \mathbf{x} in $E[g(\mathbf{x})]$, for then $E_N = \frac{1}{N} \sum_{i=1}^N g(\mathbf{x}_i) \xrightarrow{a.s.} E$ and $E(E_N) = E \forall N$. It is not necessary to know the integrating constant of the kernel probability density for \mathbf{x} . If $\sigma^2 = \text{var}[g(\mathbf{x})]$ exists, then $\sqrt{N}(E_N - E) \xrightarrow{d} N(0, \sigma^2)$ as well.

As an example, consider the problem

$$I = \int_{\mathfrak{R}^k} f(\mathbf{x}) d\mathbf{x} = \int_{\mathfrak{R}^k} g(\mathbf{x}) p(\mathbf{x}) d\mathbf{x} = \int_{\mathfrak{R}^k} g(\mathbf{x}) \exp\left[-\frac{1}{2}(\mathbf{x} - \mu)' \mathbf{H}(\mathbf{x} - \mu)\right] d\mathbf{x},$$

where \mathbf{H} is positive definite. Since $p(\mathbf{x})$ is a multivariate normal kernel density function,

$$I_N = (2\pi)^{k/2} |\mathbf{H}|^{1/2} N^{-1} \sum_{i=1}^N g(\mathbf{x}_i), \quad \mathbf{x}_i \stackrel{i.i.d.}{\sim} N(\mu, \mathbf{H}^{-1}).$$

Because $p(\mathbf{x}) \geq 0 \forall \mathbf{x} \in \mathfrak{R}^k$, $I_N \xrightarrow{a.s.} I$ regardless of the form of $f(\mathbf{x})$. However, convergence will be impractically slow if $g(\mathbf{x})$ is ill conditioned or (equivalently) μ and \mathbf{H} are chosen so that $p(\cdot)$ poorly mimics $f(\cdot)$. If $\text{var}_p[g(\mathbf{x})]$ exists, then

$$\sigma^2 = (2\pi)^k |\mathbf{H}| \text{var}_p[\mathbf{g}(\mathbf{x})]$$

provides the pertinent measure of the adequacy of I_N as an approximation of I . Only this expression -- not the dimensionality k -- matters.

4.2 Acceptance methods

Acceptance methods may be used to evaluate integrals in much the same way as they are used to produce pseudorandom numbers. In Problem I, suppose that $0 \leq \mathbf{g}(\mathbf{x}) \leq a < \infty \forall \mathbf{x} \in D$. Suppose further that p^* is known or equivalently that $p(\mathbf{x})$ is a probability density function and not merely a kernel. Let $\{\mathbf{x}_i\}$ be an i.i.d. sequence drawn from a distribution function with p.d.f. $p(\mathbf{x})$, and let u_i be a corresponding Bernoulli random variable,

$$u_i = 0 \text{ or } 1, \quad \mathbf{P}(u_i = 1) = a^{-1} \mathbf{g}(\mathbf{x}_i).$$

Then

$$\begin{aligned} I_N &= N^{-1} a \sum_{i=1}^N u_i \xrightarrow{a.s.} a \mathbf{E}_p(u_i) = a \int_D a^{-1} \mathbf{g}(\mathbf{x}) p(\mathbf{x}) d\mathbf{x} = I, \\ \mathbf{E}(I_N) &= I \forall N, \quad \sqrt{N}(I_N - I) \xrightarrow{d} \mathbf{N}(0, \sigma^2), \\ \sigma^2 &= aI - I^2, \quad aI_N - I_N^2 \xrightarrow{a.s.} \sigma^2. \end{aligned} \quad (4.2.1)$$

This method may be extended to $\mathbf{g}(\mathbf{x})$ for which $-\infty < l \leq \mathbf{g}(\mathbf{x}) \leq u < \infty$, by defining $\mathbf{g}^+(\mathbf{x}) = \sup[0, \mathbf{g}(\mathbf{x})]$, $\mathbf{g}^-(\mathbf{x}) = -\inf[0, \mathbf{g}(\mathbf{x})]$, and approximating $\int_D \mathbf{g}^+(\mathbf{x}) d\mathbf{x}$ and $\int_D \mathbf{g}^-(\mathbf{x}) d\mathbf{x}$ separately. Observe that σ^2 is an increasing function of a and the unconditional probability of acceptance $\mathbf{P}(u_i = 1) = a^{-1}I$ is a decreasing function of a . If $p(\mathbf{x}) \propto \mathbf{g}(\mathbf{x})$, then $\mathbf{P}(u_i) = 1$ and $\sigma^2 = 0$, but this is tantamount to being able to integrate $f(\mathbf{x})$ analytically. In general one seeks to minimize a . If a is too large, then very few u_i will be accepted, and the method will be impractical.

In Problem E, acceptance methods may be applied to draw from the distribution with probability density $p(\mathbf{x})$. If $h(\mathbf{x})$ is a source density as described in Section 3.2, $0 \leq p(\mathbf{x})/h(\mathbf{x}) \leq a < \infty \forall \mathbf{x} \in D$, then a sequence of i.i.d. draws from the distribution with p.d.f. $p(\mathbf{x})$ may be constructed. If we take $\{\mathbf{x}_i\}_{i=1}^N$ to be the accepted draws, then

$$\begin{aligned} E_N &= N^{-1} \sum_{i=1}^N \mathbf{g}(\mathbf{x}_i) \xrightarrow{a.s.} E, \quad \mathbf{E}(E_N) = E \forall N, \quad \sqrt{N}(E_N - E) \xrightarrow{d} \mathbf{N}(0, \sigma^2), \\ \sigma^2 &= \text{var}_p[\mathbf{g}(\mathbf{x})], \quad s_N^2 = \sum_{i=1}^N [\mathbf{g}(\mathbf{x}_i) - E_N]^2 / N \xrightarrow{a.s.} \sigma^2. \end{aligned} \quad (4.2.2)$$

If we take $\{\mathbf{z}_i\}_{i=1}^N$ to be draws from the source density, and $u_i = 1$ if \mathbf{z}_i is accepted and $u_i = 0$ if not, then

$$E_N = \sum_{i=1}^N u_i g(\mathbf{z}_i) / \sum_{i=1}^N u_i \xrightarrow{a.s.} E, \quad \sqrt{N}(E_N - E) \xrightarrow{d} \mathcal{N}(0, \sigma^2),$$

$$\sigma^2 = a \int_D [g(\mathbf{z}) - E]^2 p(\mathbf{z}) d\mathbf{z}, \quad a \sum_{i=1}^N u_i [g(\mathbf{z}_i) - E]^2 / \sum_{i=1}^N u_i \xrightarrow{a.s.} \sigma^2. \quad (4.2.3)$$

(In this case one again seeks to choose $h(\mathbf{x})$ so as to minimize a .) Which expression is more relevant depends on the particulars of the problem. We shall return to this topic in Section 4.4.

The acceptance method just described assumes that the probability density is known, including its constant of integration -- i.e., $\int_D p(\mathbf{x}) d\mathbf{x} = 1$. This assumption may be strong in practice. In Problem I, one may recognize $p(\mathbf{x})$ as a probability density kernel, not knowing the constant of integration. Acceptance or adaptive methods might be applied to draw from the distribution with kernel density $p(\mathbf{x})$; these methods do not require that one know the constant of integration for $p(\mathbf{x})$. If $p(\mathbf{x})$ is the kernel and $p^* = \int_D p(\mathbf{x}) d\mathbf{x}$, it is then the case for acceptance methods in Problem I that

$$I_N = N^{-1} a p^* \sum_{i=1}^N u_i \xrightarrow{a.s.} I.$$

Whether or not consistent evaluation of p^* is possible depends on the method used to draw variables from the distribution with kernel $p(\mathbf{x})$. If the method is acceptance sampling or a variant on acceptance sampling (e.g., the adaptive method for log-concave densities described in Section 3.2), one can approximate p^* using the methods just described as long as the actual probability density (not just the kernel) of the source distribution for the target kernel $p(\mathbf{x})$ is known. This produces a sequence p_i^* with the property $\bar{p}_N^* \equiv N^{-1} \sum_{i=1}^N p_i^* \xrightarrow{a.s.} p^*$. In this case clearly

$$I_N = N^{-1} a p_N^* \sum_{i=1}^N u_i \xrightarrow{a.s.} I, \quad \sqrt{N}(I_N - I) \xrightarrow{d} \mathcal{N}(0, \sigma^2),$$

but σ^2 is affected by the substitution of \bar{p}_n^* for p^* .

One may work out expressions for σ^2 and a corresponding consistent (in N) approximation of σ^2 , as has been done already in several cases. Such expressions are quite useful in the analytical comparison of approximation methods. But if the goal is simply to assess approximation error, straightforward asymptotic expansion is much simpler. To illustrate the method, return to the case of simple Monte Carlo integration with p^* unknown, (4.1.1). Let M be the number of i.i.d. draws from source density $h(\mathbf{z})$ for target density $p(\mathbf{z})$, define $a = \sup_D [p(\mathbf{z})/h(\mathbf{z})]$, and let

$$y_i = p(\mathbf{z}_i)/h(\mathbf{z}_i)$$

$$u_i = \begin{cases} 1 & \text{with probability } p(\mathbf{z}_i)/ah(\mathbf{z}_i), \\ 0 & \text{otherwise} \end{cases}$$

$$w_i = u_i g(\mathbf{z}_i).$$

Defining $\bar{y}_M = M^{-1} \sum_{i=1}^M y_i$, $\bar{u}_M = M^{-1} \sum_{i=1}^M u_i$, $\bar{w}_M = M^{-1} \sum_{i=1}^M w_i$,

$$I_M = \bar{y}_M \bar{w}_M / \bar{u}_M \xrightarrow{a.s.} I.$$

As long as $\int_D g^2(x) p(x) dx$ is absolutely convergent, $\sqrt{M}(I_M - I) \xrightarrow{d} N(0, \sigma^2)$, and

$$I_M^2 \left[\frac{\hat{\text{var}}(y_i)}{\bar{y}_M^2} + \frac{\hat{\text{var}}(w_i)}{\bar{w}_M^2} + \frac{\hat{\text{var}}(u_i)}{\bar{u}_M^2} + \frac{2 \hat{\text{cov}}(y_i, w_i)}{\bar{y}_M \bar{w}_M} - \frac{2 \hat{\text{cov}}(y_i, u_i)}{\bar{y}_M \bar{u}_M} - \frac{2 \hat{\text{cov}}(w_i, u_i)}{\bar{w}_M \bar{u}_M} \right] \xrightarrow{a.s.} \sigma^2.$$

(This expression may be derived by the delta method, i.e., by linearizing I_M in \bar{y}_M , \bar{u}_M and \bar{w}_M . The terms $\hat{\text{var}}(y_i)$, $\hat{\text{cov}}(y_i, w_i)$, etc., are computed in the usual way from $\{y_i, w_i, u_i\}_{i=1}^M$.)

4.3 Importance sampling

The method of importance sampling may be used to solve Problem I or Problem E, under similar circumstances: one has available a probability distribution with p.d.f. somewhat similar to the integrand $f(\mathbf{x})$ in Problem I or the probability density function $p(\mathbf{x})$ in Problem E and wishes to use an independent, identically distributed sample from this distribution to approximate I or E. Rather than use acceptance to generate an i.i.d. sample from the distribution with p.d.f. $p(\mathbf{x})$, importance sampling uses all of the draws from the source probability distribution but weights that sample to obtain a convergent approximation. In this method the probability density function of the source distribution is called the importance sampling density, a term due to Hammersly and Handscomb (1964), who were among the first to propose the method. It appears to have been introduced to the economics literature by Kloek and van Dijk (1978). We shall denote the importance sampling density $j(\mathbf{x})$.

Suppose that for Problem I one can draw an i.i.d. sequence of random vectors $\{\mathbf{x}_i\}$ from the importance distribution and that the support of this distribution includes D . Then

$$E_j[f(\mathbf{x}_i)/j(\mathbf{x}_i)] = \int_D [f(\mathbf{x})/j(\mathbf{x})] j(\mathbf{x}) d\mathbf{x} = \int_D f(\mathbf{x}) d\mathbf{x} = I.$$

Since $f(\mathbf{x}_i)/j(\mathbf{x}_i)$ is also an i.i.d. sequence,

$$I_N \equiv N^{-1} \sum_{i=1}^N [f(\mathbf{x}_i)/j(\mathbf{x}_i)] \xrightarrow{a.s.} I$$

by the strong law of large numbers. Furthermore, $E(I_N) = I \forall N$. This result is remarkable for its weakness: no upper bound on $f(\mathbf{x})/j(\mathbf{x})$ is required as is the case for $f(\mathbf{x})/h(\mathbf{x})$ in acceptance sampling. The requirement that the support of $j(\mathbf{x})$ include D is necessary and usually trivial to verify.

In Problem E importance sampling may be attractive if there is no simple method of constructing pseudorandom numbers drawn from the distribution $P(\cdot)$ underlying the expectation operator. If the constant of integration for the probability density is known, then

$$E_N = N^{-1} \sum_{i=1}^N [g(\mathbf{x}_i) p(\mathbf{x}_i) / j(\mathbf{x}_i)] \xrightarrow{a.s.} E \text{ and } E(E_N) = E \forall N$$

as long as the support of the importance sampling distribution includes that of $P(\cdot)$. If the constant of integration is not known and $p(\mathbf{x})$ is merely the kernel of the probability density function, $\int_D p(\mathbf{x}) d\mathbf{x} = p^*$, then

$$N^{-1} \sum_{i=1}^N [g(\mathbf{x}_i) p(\mathbf{x}_i) / j(\mathbf{x}_i)] \xrightarrow{a.s.} p^* E, \quad N^{-1} \sum_{i=1}^N [p(\mathbf{x}_i) / j(\mathbf{x}_i)] \xrightarrow{a.s.} p^*,$$

and hence

$$E_N \equiv \frac{\sum_{i=1}^N [g(\mathbf{x}_i) p(\mathbf{x}_i) / j(\mathbf{x}_i)]}{\sum_{i=1}^N [p(\mathbf{x}_i) / j(\mathbf{x}_i)]} \xrightarrow{a.s.} E, \quad (4.3.1)$$

but of course $E(E_N) \neq E$ in general. In either case $w(\mathbf{x}) = p(\mathbf{x}) / j(\mathbf{x})$ may be regarded as a weight function, large weights being assigned to those $g(\mathbf{x}_i)$ for which the importance sampling distribution assigns smaller probability than does the probability distribution $P(\cdot)$.

To assess the accuracy of importance sampling approximations using a central limit theorem, more is required. In the case of Problem I, suppose that $\int_D [f^2(\mathbf{x}) / j(\mathbf{x})] d\mathbf{x}$ is absolutely convergent. Then $f(\mathbf{x}_i) / j(\mathbf{x}_i)$ is an i.i.d. sequence and

$$I_N \xrightarrow{a.s.} I, \quad \sqrt{N}(I_N - I) \xrightarrow{d} N(0, \sigma^2),$$

$$\sigma^2 = \int_D \left[\frac{f^2(\mathbf{x})}{j(\mathbf{x})} \right] d\mathbf{x} - I^2 = E_j \left[\frac{f(\mathbf{x})}{j(\mathbf{x})} - I \right]^2, \quad s_N^2 = N^{-1} \sum_{i=1}^N \frac{f^2(\mathbf{x}_i)}{j^2(\mathbf{x}_i)} - I_N^2 \xrightarrow{a.s.} \sigma^2. \quad (4.3.2)$$

It is therefore practical to assess the accuracy of I_N as an approximation of I . The convergence of $\int_D [f^2(\mathbf{x}) / j(\mathbf{x})] d\mathbf{x}$ must be established analytically, however. If $|f(\mathbf{x}) / j(\mathbf{x})|$ is bounded above on D or if D is compact and $f^2(\mathbf{x}) / j(\mathbf{x})$ is bounded above, then convergence obtains. If neither of these conditions is satisfied, then verifying convergence may be difficult. In choosing an importance sampling density, it is especially important to insure that the tails of $j(\mathbf{x})$ decline no faster than those of $f(\mathbf{x})$. If these conditions are not met, but one still proceeds with the approximation, then convergence is usually quite slow. Violation of the central limit theorem convergence condition then may be evidenced by values of s_N^2 that increase with N .

Assessing the accuracy of E_N as an approximation of E is complicated by the ratio of terms in (4.3.1). If both

$$E_p[w(\mathbf{x})] = \int_D [p^2(\mathbf{x}) / j(\mathbf{x})] d\mathbf{x} \text{ and } E_p[g^2(\mathbf{x}) w(\mathbf{x})] = \int_D [g^2(\mathbf{x}) p(\mathbf{x})] d\mathbf{x} \quad (4.3.3)$$

are absolutely convergent, then

$$E_N \xrightarrow{a.s.} E, \quad \sqrt{N}(E_N - E) \xrightarrow{d} N(0, \sigma^2),$$

$$\sigma^2 = E_p \left\{ [g(\mathbf{x}) - E]^2 w(\mathbf{x}) \right\} = p^{*-1} \int_D \left\{ [g(\mathbf{x}) - E]^2 w(\mathbf{x}) p(\mathbf{x}) \right\} d\mathbf{x}, \quad (4.3.4)$$

$$s_N^2 = \frac{N \sum_{i=1}^N [g(\mathbf{x}_i) - E]^2 w(\mathbf{x}_i)}{\left[\sum_{i=1}^N w(\mathbf{x}_i) \right]^2} \xrightarrow{a.s.} \sigma^2.$$

(Derivations are given in Geweke, 1989.) This result provides a practical way to assess approximation error and also indicates conditions in which the method of importance sampling will work well for Problem E. A small value of $E_p[w(\mathbf{x})]$, perhaps as reflected in a small upper bound on $w(\mathbf{x})$, combined with small $\text{var}_p[g(\mathbf{x})]$, will lead to small values of σ^2 . As in the case of Problem I, central limit theorem convergence conditions must be verified analytically.

There has been little practical work to date on the optimal choice of importance sampling distributions. Using a result of Rubinstein (1981, Theorem 4.3.1) one can show that the importance sampling density with kernel $|g(\mathbf{x}) - E|p(\mathbf{x})$ provides the smallest possible value of σ^2 . This is not very useful, since drawing pseudorandom vectors from this distribution is likely to be awkward at best. There has been some attention to optimization within families of importance sampling densities (Geweke, 1989), but optimization procedures themselves generally involve integrals that in turn require numerical approximation. Adaptive methods use previously drawn \mathbf{x}_i to identify large values of $f(\mathbf{x})/j(\mathbf{x})$, $w(\mathbf{x})$, or $g^2(\mathbf{x})w(\mathbf{x})$ and modify $j(\mathbf{x})$ accordingly (Evans, 1991). Such procedures can be convenient but are limited by the fact that \mathbf{x}_i is least likely to be drawn where $j(\mathbf{x})$ is small. Informal, deterministic methods for tailoring $j(\mathbf{x})$ have worked well in some problems in Bayesian econometrics (Geweke, 1989).

In Problem I the objective in choosing the importance sampling density j is to find $j(\mathbf{x})$ that mimics the shape of $f(\mathbf{x})$ as closely as possible; the relevant metric is (4.3.2). Finding $j(\mathbf{x}) \propto f(\mathbf{x})$ will drive σ^2 to zero, but this amounts to analytical solution of the problem since $\int_D j(\mathbf{x}) d\mathbf{x} = 1$. In Problem E the relevant metric (4.3.4) is more complicated, involving both the variance of $g(\mathbf{x})$ and the closeness of $j(\mathbf{x})$ to $p(\mathbf{x})$ as reflected in $w(\mathbf{x}) = p(\mathbf{x})/j(\mathbf{x})$. As long as $\text{var}_p[g(\mathbf{x})] > 0$, no choice of $j(\mathbf{x})$ will drive σ^2 to zero, and if $\text{var}_p[g(\mathbf{x})] = 0$, then Problem E reduces to Problem I. If $j(\mathbf{x}) \propto p(\mathbf{x})$, then $\sigma^2 = \text{var}_p[g(\mathbf{x})]$, which can serve as a benchmark in evaluating the adequacy of $j(\mathbf{x})$. The ratio $\sigma^2/\text{var}_p[g(\mathbf{x})]$ has been termed the relative numerical efficiency of $j(\mathbf{x})$ (Geweke, 1989): it indicates the ratio of iterations using $p(\mathbf{x})$ itself as the importance sampling density, to the number using $j(\mathbf{x})$, required to achieve the same accuracy of approximation of E . Relative numerical efficiency much less than 1.0 (less than 0.1, certainly less than 0.01) indicates poor imitation of $p(\mathbf{x})$ by $j(\mathbf{x})$ in the metric (4.3.4), possibly the existence

of a better importance sampling distribution or the failure of the underlying convergence conditions (4.3.3).

4.4 A note on the choice of method

There is considerable scope for combining the methods discussed in Sections 3 and 4. For example, the pseudorandom number generation in making draws from the population with probability density $h(\mathbf{x})$, in the case of acceptance sampling, or $j(\mathbf{x})$, in the case of importance sampling, generally will involve several of the methods discussed in Section 3.2. In even moderately complex problems, the investigator needs to tailor these methods, balancing computational efficiency against demands for the development and checking of reliable code.

Acceptance sampling and importance sampling are clearly similar. In fact, given a candidate source density, one has the choice of undertaking either acceptance or importance sampling. A straightforward comparison of approximation errors indicates the issues involved in the choice. In Problem I, the variance in acceptance sampling is

$$\sigma_1^2 = \int_D [g(\mathbf{x}) - I]^2 p(\mathbf{x}) d\mathbf{x} = \int_D g^2(\mathbf{x}) p(\mathbf{x}) d\mathbf{x} - I^2$$

if by *draw* we mean *accepted* draw. But if instead we mean every draw from the source distribution, the variance is

$$\sigma_2^2 = aI - I^2, \quad a = \sup_D g(\mathbf{x}),$$

from (4.2.1). In importance sampling, where all draws are used but differentially weighted, the variance is

$$\sigma_3^2 = \int_D g^2(\mathbf{x}) p(\mathbf{x}) d\mathbf{x} - I^2,$$

from (4.3.2). Hence given a choice between acceptance and importance sampling in Problem I, importance sampling is clearly preferred: it conserves information from all draws, whereas the rejected draws in acceptance sampling require execution time but do not further improve the accuracy of the approximation.

For Problem E the situation is different. The variance is

$$\sigma_4^2 = \int_D [g(\mathbf{x}) - E]^2 p(\mathbf{x}) d\mathbf{x}$$

for acceptance sampling (see (4.2.2)) if we count only accepted draws and

$$\sigma_5^2 = a \int_D [g(\mathbf{z}) - E]^2 p(\mathbf{z}) d\mathbf{z}, \quad a = \sup_D [p(\mathbf{z})/h(\mathbf{z})]$$

if we count all draws (see (4.2.3)). For importance sampling, expressing (4.3.4) in the notation of acceptance sampling, we have

$$\sigma_6^2 = \int_D [g(\mathbf{x}) - E]^2 w(\mathbf{x}) p(\mathbf{x}) d\mathbf{x}, \quad w(\mathbf{x}) = p(\mathbf{x})/h(\mathbf{x}).$$

Since $\sigma_4^2 \leq \sigma_6^2 \leq \sigma_5^2$, a choice between acceptance and importance sampling on grounds of computational efficiency rests on the particulars of the problem. If evaluation of $g(\mathbf{x})$ is sufficiently expensive relative to evaluation of $p(\mathbf{x})/h(\mathbf{x})$, acceptance sampling will be more efficient; otherwise, importance sampling will be the choice.

In fact one may combine acceptance and importance sampling. Let c be any positive constant, and define

$$w(\mathbf{z}_i) = \begin{cases} p(\mathbf{z}_i)/c h(\mathbf{z}_i) & \text{if } p(\mathbf{z}_i)/h(\mathbf{z}_i) \geq c \\ 1 & \text{with probability } p(\mathbf{z}_i)/c h(\mathbf{z}_i) \text{ if } p(\mathbf{z}_i)/h(\mathbf{z}_i) < c \\ 0 & \text{otherwise} \end{cases}$$

Then $\sum_{i=1}^n w(\mathbf{z}_i)g(\mathbf{z}_i) / \sum_{i=1}^n w(\mathbf{z}_i) \xrightarrow{a.s.} E$. For any given problem there will be a value of c that minimizes the variance of approximation error relative to required computing time. This may be found experimentally; or for some analytical methods, see Müller (1991, Chapter 2). The hybrid method can result in dramatic increases in efficiency when computation of $g(\mathbf{x})$ is relatively expensive (or there are many such functions to be evaluated) and the weight function $w(\mathbf{x})$ is small with high probability.

A more fundamental choice is that between the simulation methods discussed in this and the previous section and the deterministic algorithms outlined in Section 2. Many problems in economics require integration in very high dimensions. (Two examples are presented in Section 7.) For such problems the most practical deterministic procedures are the low discrepancy methods of Section 2.3. Tables 2 and 3 provide some specific comparisons for dimensions as high as $d = 100$. (Execution time for quadrature methods in these problems is approximately $8 \times 4^{d-10}$ seconds on a Sun 10/51 workstation: .01 seconds for $d = 5$, 8 seconds for $d = 10$, 3 months for $d = 20$, about 10^4 times the estimated age of the universe for $d = 40$, ...)

Table 2 extends the analysis of the same problem taken up in Section 2.3. As noted there, the bounds in (2.3.1) and (2.3.2) are useless for this problem and most others. The actual Halton errors presented in Table 2 were found by direct computation, using the first d primes as the bases. The Monte Carlo errors were found analytically. Two error bounds are presented, one based on a 95% confidence interval ($\pm 1.96\sigma$) and a second based on a $100(1-10^{-12})\%$ confidence interval ($\pm 7.13\sigma$). For lower dimensions the comparison is dominated by the convergence of the Halton sequence at rate $\log m/m$ compared with Monte Carlo at rate $m^{-1/2}$: the Halton sequence is much more accurate. But for any reasonable fixed value of m , the comparison in higher dimensions is dominated by an approximately exponential rate of error increase in d for the Halton sequence, contrasted with the rate $d^{1/2}$ for Monte Carlo. For $m = 1,000$ iterations, Monte Carlo is more efficient

for d exceeding about 25 if one applies the $p = .05$ standard and for d exceeding about 45 for the $p = 10^{-12}$ standard. For $m = 50,000$ the breakpoints occur around $d = 35$ and $d = 110$, respectively. (The Halton error is not monotone decreasing in m because of the systematic way in which points are selected.)

Table 3 provides a comparison of these methods for an example of Problem E. The Halton sequence is first mapped into the normal distribution applying the inverse-c.d.f. transformation in each dimension. Each of the five panels provides approximations to successively higher moments, p , of the multivariate normal distribution. Within each panel, the comparison is dominated by the same features noted for Table 2. Comparisons across panels are dominated by important characteristics of each method. Monte Carlo errors are proportional to $E[(z)^{2p}] = [(2p-1) \cdot (2p-3) \cdot \dots \cdot 3 \cdot 1]^{1/2}$, where $z \sim N(0, 1)$. Halton errors reflect an interaction between the ordering of the points and the characteristics of x_i^p . When p is odd, x_i^p is an odd monotone increasing function of x_i , whereas the standard normal probability density function is even. For any fixed m , the Halton points systematically exclude positive x_i values for which the corresponding $-x_i$ value has been included. Hence the error is always negative (as it was in Table 2 for the same reason). When p is even, this is not the case and the size of the error is smaller as well. The tendency of the Halton sequence to systematically exclude larger x_i has more severe consequences for evaluation of the integral the higher the value of odd p . Thus, for $p = 5$ independence Monte Carlo becomes dominant for values of d exceeding a fairly small threshold.

The largest problems worked for Table 3 ($d = 100, m = 50,000$) required about 75 seconds on a Sun 10/51 when solved using a Halton sequence. Independence Monte Carlo was about 15 times faster in every case. The difference reflects the inherent speed of linear congruential generators, contrasted with the floating point operations required to generate a Halton sequence. For more complex and realistic problems the relative speed of independence Monte Carlo is less important, since computation time typically will be dominated by subsequent computations involving the sequences produced by either method.

These comparisons illustrate the general rule that simulation methods are preferred for higher dimensional problems. If the dimension is very low, then quadrature methods are much faster and more accurate. For intermediate dimensions, quadrature is impractical and low discrepancy methods are more accurate than simulation methods. Just where the breakpoints occur is problem-specific, and the situation is complicated by the fact that there are no useful independent assessments of approximation error for low discrepancy methods. Simulation methods always provide an assessment of numerical error as a by-product, for square-integrable functions. Combined with the checks for robustness of

results with respect to alternative uniform random number generators and seed values, these methods are practical and reliable for a much wider range of problems than is any deterministic algorithm. As we shall see, their application in complex problems can be very natural.

5. Variance reduction

In any of the independence Monte Carlo methods a single draw can be replaced by the mean of M identically but not independently distributed draws. For example, in simple Monte Carlo for Problem I,

$$I_{N,M} = N^{-1} \sum_{i=1}^N \left\{ M^{-1} \left[\sum_{j=1}^M g(\mathbf{x}_{ij}) \right] \right\}.$$

For any $i \neq k$ \mathbf{x}_{ij} and \mathbf{x}_{k1} are independent, whereas \mathbf{x}_{ij} and \mathbf{x}_{i1} are dependent. Since all \mathbf{x}_{ij} are drawn from the distribution with probability density $p(\mathbf{x})$,

$$I_{N,M} \xrightarrow{a.s.} I, \quad \sqrt{N} (I_{N,M} - I) \xrightarrow{d} \mathbf{N}(0, \sigma^{*2}),$$

$$\sigma^{*2} = \text{var} \left[M^{-1} \sum_{j=1}^M g(\mathbf{x}_{ij}) \right], \quad s_N^{*2} = N^{-1} \sum_{i=1}^N \left[M^{-1} \sum_{j=1}^M g(\mathbf{x}_{ij}) - I_{N,M} \right]^2 \xrightarrow{a.s.} \sigma^{*2}.$$

The idea is to set up the relation among $\mathbf{x}_{i1}, \dots, \mathbf{x}_{iM}$ in such a way that $\sigma^{*2} < M^{-1} \text{var}_p [g(\mathbf{x}_{ij})]$. If in addition the cost of generating the M -tuple is insignificantly greater than the cost of generating M independent variables from $p(\mathbf{x})$, then $I_{N,M}$ provides a computationally more efficient approximation of I than does I_N .

There are numerous variants on this technique. This section takes up four that account for most use of the method: antithetic variables, systematic sampling, conditional expectations, and control variables. The scope for combining these variance reduction techniques with the methods of Section 4 or Section 6 is enormous. Rather than list all the possibilities, the purpose here is to provide some appreciation of the circumstances in which each variant may be practical and productive.

5.1 Antithetic Monte Carlo

This technique is due to Hammersly and Morton (1956) and has been widely used in statistics, experimental design, and simulation (e.g., Mikhail, 1972; Mitchell, 1973; Geweke, 1988). In antithetic simple Monte Carlo integration $M = 2$ correlated variables are drawn in each of N replications. Then,

$$\sigma^{*2} = \frac{1}{2} \left\{ \text{var} [g(\mathbf{x}_{i1})] + \text{cov} [g(\mathbf{x}_{i1}), g(\mathbf{x}_{i2})] \right\}.$$

As long as $\text{cov} [g(\mathbf{x}_{i1}), g(\mathbf{x}_{i2})] < 0$, antithetic simple Monte Carlo integration with $N/2$ replications will have smaller error variance than simple Monte Carlo integration with N replications, and the computational requirements will be about the same.

To focus on the main ideas, consider the situation in which $p(\mathbf{x})$ is symmetric about a point μ in Problem I set out in Section 4. In this case $\mathbf{x}_{i1} = \mu + \mathbf{w}_i$, $\mathbf{x}_{i2} = \mu - \mathbf{w}_i$ describes a pair of variables drawn from the distribution with p.d.f. $p(\mathbf{x})$ with correlation matrix $-\mathbf{I}$. If $g(\mathbf{x})$ were a linear function, then $\text{var} \left\{ \frac{1}{2} [g(\mathbf{x}_{i1}) + g(\mathbf{x}_{i2})] \right\} = 0$, and variance reduction

would be complete. (Clearly $I = g(\mu)$; this case is of interest only as a limit for numerical integration problems.) At the other extreme, if $g(\mathbf{x})$ is also symmetric about μ , then $\text{var}\left\{\frac{1}{2}[g(\mathbf{x}_{i1}) + g(\mathbf{x}_{i2})]\right\} = \text{var}[g(\mathbf{x})]$: N replications of antithetic simple Monte Carlo integration will yield as much information as N replications of simple Monte Carlo, but will usually require about double the number of computations. As an intermediate case, suppose that $d(y) = g(\mathbf{xy})$ is either monotone nondecreasing or monotone nonincreasing for all \mathbf{x} . Then $g(\mathbf{x}_{i1}) - I$ and $g(\mathbf{x}_{i2}) - I$ must be of opposite sign if they are nonzero. This implies $\text{cov}[g(\mathbf{x}_{i1}), g(\mathbf{x}_{i2})] < 0$, whence $\sigma^{*2} \leq \frac{1}{2} \text{var}[g(\mathbf{x})] = \sigma^2/2$, and so antithetic simple Monte Carlo integration produces gains in efficiency.

The use of antithetic Monte Carlo integration is especially powerful in an important class of Bayesian learning and inference problems. In these problems \mathbf{x} typically represents a vector of parameters unknown to an economic agent or an econometrician, and $p(\mathbf{x})$ is the probability density of that vector conditional on information available. The integral I could correspond to an expected utility or a posterior probability. If the available information is based on an i.i.d. sample of size T , then it is natural to write $p_T(\mathbf{x})$ for $p(\mathbf{x})$. As T increases, the distribution $p_T(\mathbf{x})$ generally becomes increasingly symmetric and concentrated about the true value of the vector of unknown parameters, reflecting the operation of a central limit theorem. In these circumstances $g(\mathbf{x})$ is increasingly well described by a linear approximation of itself over most of the support of $p_T(\mathbf{x})$, as T increases. Suppose that the agent or econometrician approximates I using simple Monte Carlo with accuracy indicated by σ_T^2 or by antithetic simple Monte Carlo with accuracy indicated by σ_T^{*2} . Given some side conditions, mainly continuous differentiability of $g(\mathbf{x})$ in a neighborhood of the true value of the parameter vector \mathbf{x} and a nonzero derivative of $g(\mathbf{x})$ at this point, it may be shown that $\sigma_T^{*2}/\sigma_T^2 \rightarrow 0$ (Geweke, 1988). Given additional side conditions, mainly twice continuous differentiability of $g(\mathbf{x})$ in a neighborhood of the true value of the parameter vector \mathbf{x} , it may be shown that $T\sigma_T^{*2}/\sigma_T^2$ converges to a constant. The constant is inversely related to the magnitude of $\partial g(\mathbf{x})/\partial \mathbf{x}$ and directly related to the magnitude of $\partial^2 g(\mathbf{x})/\partial \mathbf{x} \partial \mathbf{x}'$, each evaluated at the true value of the parameter vector \mathbf{x} (Geweke, 1988). This result is an example of acceleration, because it indicates an interesting sequence of conditions under which the relative advantage of a variance reduction method increases without bound.

Application of the method of antithetic variables with techniques more complicated than simple Monte Carlo is generally straightforward. In the case of importance sampling, \mathbf{x}_{i1} and \mathbf{x}_{i2} are drawn from the importance sampling density $j(\mathbf{x})$. In Problem I the term

$\frac{1}{2} [f(\mathbf{x}_{i1})/j(\mathbf{x}_{i1}) + f(\mathbf{x}_{i2})/j(\mathbf{x}_{i2})]$ replaces $f(\mathbf{x}_i)/j(\mathbf{x}_i)$. In Problem E, define $w(\mathbf{x}) = p(\mathbf{x})/j(\mathbf{x})$ as before. Then

$$E_N \equiv \frac{\sum_{i=1}^N [g(\mathbf{x}_{i1})w(\mathbf{x}_{i1}) + g(\mathbf{x}_{i2})w(\mathbf{x}_{i2})]}{\sum_{i=1}^N [w(\mathbf{x}_{i1}) + w(\mathbf{x}_{i2})]} \xrightarrow{a.s.} E, \quad \sqrt{N}(E_N - E) \xrightarrow{d} N(0, \sigma^{*2}),$$

$$\sigma^{*2} = E_p \left\{ \left[\frac{g(\mathbf{x}_{i1})w(\mathbf{x}_{i1}) + g(\mathbf{x}_{i2})w(\mathbf{x}_{i2})}{w(\mathbf{x}_{i1}) + w(\mathbf{x}_{i2})} - E \right]^2 \frac{w(\mathbf{x}_{i1}) + w(\mathbf{x}_{i2})}{2} \right\},$$

$$s_N^2 = \frac{N \sum_{i=1}^N \left[\frac{g(\mathbf{x}_{i1})w(\mathbf{x}_{i1}) + g(\mathbf{x}_{i2})w(\mathbf{x}_{i2})}{w(\mathbf{x}_{i1}) + w(\mathbf{x}_{i2})} - E_N \right]^2 [w(\mathbf{x}_{i1}) + w(\mathbf{x}_{i2})]}{4 \left\{ \sum_{i=1}^N [w(\mathbf{x}_{i1}) + w(\mathbf{x}_{i2})] \right\}^2} \xrightarrow{a.s.} \sigma^{*2}.$$

These results are valid for any antithetic variables algorithm, even if $j(\mathbf{x})$ is not symmetric and even if the variance of the approximation error σ^2 is increased rather than decreased in moving to the use of antithetic variables. The essential requirements are that the \mathbf{x}_{ij} 's be drawn from the importance sampling distribution and that \mathbf{x}_{ij} and \mathbf{x}_{kl} be independent for $i \neq k$.

In complex problems involving multivariate \mathbf{x} , pseudorandom variables often may be generated by use of successive conditionals for $\mathbf{x}' = (\mathbf{x}'_{(1)}, \dots, \mathbf{x}'_{(m)})$,

$$p(\mathbf{x}) = p(\mathbf{x}_{(1)}) p(\mathbf{x}_{(2)} | \mathbf{x}_{(1)}) \dots p(\mathbf{x}_{(m)} | \mathbf{x}_{(1)}, \dots, \mathbf{x}_{(m-1)}).$$

In such cases a pair of antithetic variables \mathbf{x}_{i1} and \mathbf{x}_{i2} may be created by constructing a pair for a single, convenient subvector $\mathbf{x}_{(j)}$. Especially if $g(\mathbf{x}) = g(\mathbf{x}_{(j)})$, the benefits of antithetic Monte Carlo will then be realized in both Problem I and Problem E. An example of this use of antithetic variables is taken up in Section 7.2.

5.2 Systematic sampling

Systematic sampling (McGrath, 1970) combines certain advantages of deterministic and Monte Carlo methods. The former achieve great efficiency by systematically choosing points for evaluation in specific low-dimensional problems; the latter produce indications of accuracy as a byproduct and are amenable to high-dimensional problems. Systematic sampling specifies an m -tuple of points as a deterministic function of a random vector \mathbf{u} , $\mathbf{x}_j = f_j(\mathbf{u})$ ($j = 1, \dots, m$), with the property that the induced distribution of every \mathbf{x}_j is that of the probability density function $p(\mathbf{x})$.

As a leading example consider the case of univariate x , with pseudorandom variables from the distribution of x constructed using the inverse c.d.f method (Section 3.2). Denote

$F(c) = P[x \leq c]$, suppose u_i ($i = 1, \dots, N$) are independently and uniformly distributed on the unit interval, and take

$$x_{ij} = F^{-1}([u_i + j/m]) \quad (j = 1, \dots, m),$$

where “[\cdot]” denotes greatest fractional part. Clearly the method need not be limited to evenly spaced grids; e.g., Richtmeyer’s method (Section 2.3) could just as easily be applied. Extension to higher dimensions is straightforward, but is subject to all of the problems of deterministic methods there. The advantage of systematic methods is that approximation error is generally $O(m^{-1})$ whereas that in Monte Carlo is $O_p(N^{-1/2})$.

In high-dimensional problems systematic sampling can be advantageous when confined to a subset of the vector \mathbf{x} that is especially troublesome for Monte Carlo and/or is an important source of variation in the function $g(\mathbf{x})$. As an example of the former condition, suppose it is difficult to find an importance sampling density that mimics $p(\mathbf{x})$, but $\mathbf{x}' = \begin{pmatrix} \mathbf{x}'_{(1)} & \mathbf{x}'_{(2)} \\ 1 \times m_1 & 1 \times m_2 \end{pmatrix}$, a good importance sampling density for the marginal p.d.f. $p(\mathbf{x}_{(1)})$ is available, and the inverse c.d.f. $F^{-1}(p|\mathbf{x}_{(1)})$ of the conditional distribution of $\mathbf{x}_{(2)}$ can be evaluated. One may generate $\mathbf{x}_{(1)i}$ together with corresponding importance sampling weight w_i ; draw (u_1, \dots, u_{m_2}) independently distributed on the unit interval; create the systematic sample

$$\mathbf{x}_{(2)ij_1 \dots j_{m_2}} = F^{-1}([u_1 + j_1/l_1], \dots, [u_{m_2} + j_{m_2}/l_{m_2}]) \quad (j_k = 1, \dots, l_k; k = 1, \dots, m_2).$$

Then record

$$g_i = \left[\prod_{k=1}^{m_2} l_k \right]^{-1} \sum_{j_1=1}^{l_1} \dots \sum_{j_{m_2}=1}^{l_{m_2}} g(\mathbf{x}_{(1)i}, \mathbf{x}_{(2)ij_1 \dots j_{m_2}})$$

along with each weight w_i . Previous expressions in Section 4.3 for I_N , σ^2 , and s_N^2 are then valid with g_i in place of $g(\mathbf{x}_i)$. In particular (4.3.2) is still true, and s_N^2 may be used to assess the increase in accuracy yielded by systematic sampling with higher values of the l_k .

5.3 The use of conditional expectations

Suppose there is a partition of \mathbf{x} , $\mathbf{x}' = (\mathbf{x}'_{(1)}, \mathbf{x}'_{(2)})$, such that

$$g(\mathbf{x}) = g(\mathbf{x}_{(1)}, \mathbf{x}_{(2)}) = g^*(\mathbf{x}_{(1)}) \mathfrak{l}(\mathbf{x}_{(2)}),$$

where $\mathfrak{l}(\cdot)$ is linear; $p(\mathbf{x}) = p(\mathbf{x}_{(1)}, \mathbf{x}_{(2)}) = p(\mathbf{x}_{(1)})p(\mathbf{x}_{(2)}|\mathbf{x}_{(1)})$; it is possible to draw pseudorandom vectors from the marginal distribution for $\mathbf{x}_{(1)}$ with p.d.f. $p(\mathbf{x}_{(1)})$; and $E[\mathbf{x}_{(2)}|\mathbf{x}_{(1)}]$ is known analytically. Then

$$\int_D g(\mathbf{x}) p(\mathbf{x}) d\mathbf{x} = \int_D g^*(\mathbf{x}_{(1)}) p(\mathbf{x}_{(1)}) \mathfrak{l}[E(\mathbf{x}_{(2)}|\mathbf{x}_{(1)})] d\mathbf{x}_{(1)}, \quad (5.3.1)$$

and

$$\text{var}_{p(\mathbf{x}_{(1)})} \left\{ g^*(\mathbf{x}_{(1)}) \left[E(\mathbf{x}_{(2)} | \mathbf{x}_{(1)}) \right] \right\} \leq \text{var}_{p(\mathbf{x})} [g(\mathbf{x})].$$

Consequently, application of Monte Carlo methods directly in (5.3.1) will produce an approximation error with smaller variance than would Monte Carlo in the general framework set forth in Section 4.

The use of conditional expectations in fact bears a close relationship to antithetic Monte Carlo integration. In particular, if one could draw antithetic variables $\mathbf{x}_{(2)i1}$ and $\mathbf{x}_{(2)i2}$ from the distribution with p.d.f. $p(\mathbf{x}_{(2)} | \mathbf{x}_{(1)})$ with perfectly negative correlation, then $\frac{1}{2}(\mathbf{x}_{(2)i1} + \mathbf{x}_{(2)i2}) = E(\mathbf{x}_{(2)} | \mathbf{x}_{(1)})$, and exactly the same result would be obtained.

More generally, whenever $g(\mathbf{x})$ is a function of $\mathbf{x}_{(1)}$ only, it is usually worth noting whether $E[g(\mathbf{x}_{(1)}) | \mathbf{x}_{(2)}]$ can be evaluated analytically. If so, then the variance of approximation error can be reduced by using the function of interest $E[g(\mathbf{x}_{(1)}) | \mathbf{x}_{(2)i}]$ rather than $g(x_{(1)i})$. Since $g(x_{(1)i}) = E[g(x_{(1)}) | x_{(2)}] + \eta$ with $\text{cov}\{\eta, E[g(x_{(1)}) | x_{(2)}]\} = 0$,

$$\text{var}_{p(\mathbf{x}_{(2)})} \left\{ E[g(\mathbf{x}_{(1)}) | \mathbf{x}_{(2)}] \right\} \leq \text{var}_{p(\mathbf{x}_{(1)})} [g(\mathbf{x}_{(1)})].$$

Against this improvement should be balanced the time required for the additional computations, which are generally of no further use in generation of the \mathbf{x}_i ; this time is usually small.

5.4 Control variables

It is often the case that one is able to solve approximations to Problem I or Problem E analytically. For example, if the mean μ of the distribution with p.d.f $p(\mathbf{x})$ is known and one has available a linear approximation $g^{(1)}(\mathbf{x})$ of the function $g(\mathbf{x})$, then the mean of $g^{(1)}(\mathbf{x})$ is $g^{(1)}(\mu)$. Moreover if $\{\mathbf{x}_i\}_{i=1}^N$ is a pseudorandom sample drawn from the distribution with p.d.f. $p(\mathbf{x})$, then $g(\mathbf{x}_i)$ and $g^{(1)}(\mathbf{x}_i)$ will be positively correlated if the linear approximation is good for most \mathbf{x}_i . In this situation the method of control variables, introduced by Kahn and Marshall (1953) and Hammersly and Handscomb (1964), can be used to reduce the variance of the approximation error in I_N or E_N .

We develop the specific method for simple Monte Carlo integration in Problem I; extension to more involved methods is straightforward. Let $J_N = N^{-1} \sum_{i=1}^N h(\mathbf{x}_i)$ have known mean J . (In the example given $h(\mathbf{x}_i) = g^{(1)}(\mathbf{x}_i)$, $J_N = N^{-1} \sum_{i=1}^N g^{(1)}(\mathbf{x}_i)$ and $J = g^{(1)}(\mu)$.) Consider approximations of the form

$$I'_N = I_N + \beta(J_N - J),$$

where I_N is computed as before. It is the case that $I'_N \xrightarrow{a.s.} I$, and as long as $\text{var}_p[\mathbf{h}(\mathbf{x}_i)]$ exists, a central limit theorem may still be used to evaluate numerical accuracy. One can easily verify that $\text{var}(I'_N)$ is minimized by $\beta = -\text{cov}(J_N, I_N)/\text{var}(J_N)$, and in this case

$$\text{var}(I'_N) = \text{var}(I_N) - \frac{\text{cov}^2(J_N, I_N)}{\text{var}(J_N)} = \text{var}(I_N)[1 - \text{corr}^2(J_N, I_N)].$$

Usually the parameter β is unknown. It may be estimated in the obvious way from the replications.

This method is easily extended to the case in which a vector of estimates $\mathbf{J}_N = (J_N^{(1)}, \dots, J_N^{(q)})'$ with known mean $\mathbf{J} = (J^{(1)}, \dots, J^{(q)})'$ is available. If we denote

$$\Sigma_{q \times q} = \text{var}(\mathbf{J}_N), \quad \mathbf{c}_{q \times 1} = \text{cov}(\mathbf{J}_N, I_N),$$

then the variance of the approximation

$$I'_N = I_N + \beta'(\mathbf{J}_N - \mathbf{J})$$

is minimized by $\beta = \Sigma^{-1}\mathbf{c}$, and in this case

$$\text{var}(I'_N) = \text{var}(I_N) - \mathbf{c}'\Sigma^{-1}\mathbf{c} = \text{var}(I_N) \left[1 - \frac{\mathbf{c}'\Sigma^{-1}\mathbf{c}}{\text{var}(I_N)} \right].$$

6. Markov chain Monte Carlo methods

All of the independence Monte Carlo methods for integration assume the ability to efficiently generate pseudorandom variables from a distribution with specified probability density function $p(\mathbf{x})$. But in many economic problems it is difficult or impossible to find a generation algorithm that is sufficiently efficient to be practical. An instructive limiting case is the one in which the constituents of \mathbf{x} are independently distributed,

$$p(\mathbf{x}) = \prod_{i=1}^m p_i(x_i).$$

One could construct an acceptance sampling algorithm with a source density $h_i(z_i)$ corresponding to each $p_i(z_i)$, and accept the draw with probability $p(\mathbf{z})/ah(\mathbf{z})$, where

$$a = \sup_z [p(z)/h(z)] = \prod_{i=1}^m a_i, \quad a_i = \sup_{z_i} [p(z_i)/h(z_i)] \quad (i = 1, \dots, m).$$

Since a is directly proportional to the time required to obtain an accepted draw (see Section 3.2) this expression makes clear that acceptance sampling can be subject to its own curse of dimensionality if the source density is constructed element-by-element. Essentially the same difficulty can arise in importance sampling, where it is manifested in only a few weights $w(\mathbf{x}_i)$ accounting for the sum.

This example is of interest only as a limiting case. If the x_i really were independent, one could employ acceptance sampling element-by-element, and computation time would then be proportional to $\sum_{i=1}^m a_i$. An obvious extension of this idea to the general case is to write

$$p(\mathbf{x}) = p(x_1) \prod_{i=1}^m p_{i|1,\dots,i-1}(x_i | x_1, \dots, x_{i-1})$$

and employ acceptance or importance sampling for each conditional. The difficulty here is that construction of probability density kernels for the marginal in x_1 and all but the last conditional require analytic integration. Notable simple cases aside, this is not possible, and it remains impossible for subvectors as well as individual components.

This section takes up a recently developed generalization of independence Monte Carlo that has become known as *Markov chain Monte Carlo*. The idea is to construct a Markov chain with state space D and invariant distribution with p.d.f. $p(\mathbf{x})$. Following an initial transient or *burn-in* phase, simulated values from the chain form a basis for approximating $E_p[g(\mathbf{x})]$, thus solving Problem E. If the p.d.f. $p(\mathbf{x})$ does not contain an unknown factor of proportionality p^* , then Problem I is solved as well. What is required is to construct an appropriate algorithm and verify that its invariant distribution is unique, with p.d.f. $p(\mathbf{x})$.

Markov chain methods have a history in mathematical physics dating back to the algorithm of Metropolis *et al.* (1953). This method, which is described in Hammersly and Handscomb (1964, Section 9.3) and Ripley (1987, Section 4.7), was generalized by

Hastings (1970), who focused on statistical problems, and was further explored by Peskun (1973). A version particularly suited to image reconstruction and problems in spatial statistics was introduced by Geman and Geman (1984). This was subsequently shown to have great potential for Bayesian computation by Gelfand and Smith (1990). Their work, combined with data augmentation methods (Tanner and Wong, 1987), has proven very successful in the treatment of latent variables and other unobservables in economic models. (An example is given in Section 7.1.) Since 1990 application of Markov chain Monte Carlo methods has grown rapidly; new refinements, extensions, and applications appear almost continuously.

This section concentrates on developing the methods, deferring serious examples to Section 7. We begin with a heuristic introduction to two widely used variants of these methods, the Gibbs sampler and the Metropolis-Hastings algorithm (Section 6.1). Some theory of continuous state Markov chains required to demonstrate convergence is given in Section 6.2. Easily verified sufficient conditions for convergence of the Gibbs sampler are set forth in Section 6.3 and for convergence of the Metropolis-Hastings algorithm in Section 6.4. Some practical issues in assessing the error of approximation are treated in Section 6.5. Much of the treatment here draws heavily on the work of Tierney (1991a, 1991b), who first used the theory of general state space Markov chains to demonstrate convergence, and Roberts and Smith (1992), who elucidated sufficient conditions for convergence that turn out to be applicable in a wide variety of problems in economics.

6.1 Two Markov chain Monte Carlo algorithms

Motivated by the role of $p(\mathbf{x})$ in Problem I or Problem E, discussion here proceeds assuming that \mathbf{x} is continuously distributed. However, there is no harm in regarding \mathbf{x} as discrete on a first reading. A full development covering both the continuous and discrete cases is given in Section 6.2.

The Gibbs sampler begins with a partition, or *blocking*, of \mathbf{x} , $\mathbf{x}' = (\mathbf{x}'_{(1)}, \mathbf{x}'_{(k)})$. For $i = 1, \dots, k$, $\mathbf{x}'_{(i)} = (\mathbf{x}_{i1}, \dots, \mathbf{x}_{im(i)})$ and $m(i) \geq 1$; $\sum_{i=1}^k m(i) = m$; and the x_{ij} are the components of \mathbf{x} . Let $p(\mathbf{x}_{(i)} | \mathbf{x}_{(-i)})$ denote the conditional p.d.f.'s induced by $p(\mathbf{x})$, where $\mathbf{x}_{(-i)} = \{\mathbf{x}_{(j)}, j \neq i\}$.

Suppose we were given a single drawing \mathbf{x}^0 , $\mathbf{x}'^0 = (\mathbf{x}'^0_{(1)}, \mathbf{x}'^0_{(k)})$, from the distribution with p.d.f. $p(\mathbf{x})$. Successively make drawings from the conditional distribution as follows:

$$\begin{aligned}
\mathbf{x}_{(1)}^1 &\sim p\left(\cdot \mid \mathbf{x}_{(-1)}^0\right) \\
\mathbf{x}_{(2)}^1 &\sim p\left(\cdot \mid \mathbf{x}_{(1)}^1, \mathbf{x}_{(3)}^0, \mathcal{K}, \mathbf{x}_{(k)}^0\right) \\
&\quad \text{M} \\
\mathbf{x}_{(j)}^1 &\sim p\left(\cdot \mid \mathbf{x}_{(1)}^1, \mathcal{K}, \mathbf{x}_{(j-1)}^1, \mathbf{x}_{(j+1)}^0, \mathcal{K}, \mathbf{x}_{(k)}^0\right) \\
&\quad \text{M} \\
\mathbf{x}_{(k)}^1 &\sim p\left(\cdot \mid \mathbf{x}_{(-k)}^1\right).
\end{aligned} \tag{6.1.1}$$

This defines a transition process from \mathbf{x}^0 to $\mathbf{x}^1 = (\mathbf{x}_{(1)}^1, \mathcal{K}, \mathbf{x}_{(k)}^1)$. The Gibbs sampler is defined by the choice of blocking and the forms of the conditional densities induced by $p(\mathbf{x})$ and the blocking. Since $\mathbf{x}^0 \sim p(\mathbf{x})$, $(\mathbf{x}_{(1)}^1, \mathcal{K}, \mathbf{x}_{(j-1)}^1, \mathbf{x}_{(j)}^1, \mathbf{x}_{(j+1)}^0, \mathcal{K}, \mathbf{x}_{(k)}^1) \sim p(\mathbf{x})$ at each step in (6.1.1) by definition of the conditional density. In particular, $\mathbf{x}^1 \sim p(\mathbf{x})$.

Iteration of the algorithm produces a sequence $\mathbf{x}^0, \mathbf{x}^1, \mathcal{K}, \mathbf{x}^t, \mathcal{K}$ which is a realization of a Markov chain with probability density function kernel for the transition from point \mathbf{x} to point \mathbf{y} given by

$$K_G(\mathbf{x}, \mathbf{y}) = \prod_{j=1}^k p\left[\mathbf{y}_{(j)} \mid \mathbf{x}_{(j)} (j > 1), \mathbf{y}_{(j)} (j < 1)\right].$$

Any single iterate \mathbf{x}^t retains the property that it is drawn from the distribution with p.d.f. $p(\mathbf{x})$.

For the Gibbs sampler to be practical, it is essential that the blocking be chosen in such a way that one can make the drawings (6.1.1) in an efficient manner. For many problems in economics, the blocking is natural and the conditional distributions are familiar; Section 7.1 provides an example. In making the drawings (6.1.1) all the methods of Sections 3 and 4 are at our disposal. Observe that in this context acceptance sampling is attractive relative to importance sampling, since the former produces independent, identically distributed, unweighted drawings from the conditional distribution.

Of course, it is generally difficult or impossible to make even one initial draw from the distribution with p.d.f. $p(\mathbf{x})$. The purpose of that assumption here is to marshal an informal argument that $p(\mathbf{x})$ is the p.d.f. of the invariant distribution of the Markov chain. A leading practical problem is to elucidate conditions in which the distribution of \mathbf{x}^t will converge to that corresponding to $p(\mathbf{x})$ for any choice of \mathbf{x}^0 in the domain D , and we turn to this in Section 6.3.

The Metropolis-Hastings algorithm begins with an arbitrary transition probability density function $q(\mathbf{x}, \mathbf{y})$ and a starting value \mathbf{x}^0 . If $\mathbf{x}^t = \mathbf{x}$, the random vector generated

from $q(\mathbf{x}, \mathbf{y})$ is considered as a candidate value for \mathbf{x}^{t+1} . The algorithm actually sets $\mathbf{x}^{t+1} = \mathbf{y}$ with probability

$$\alpha(\mathbf{x}, \mathbf{y}) = \min \left\{ \frac{p(\mathbf{y})q(\mathbf{y}, \mathbf{x})}{p(\mathbf{x})q(\mathbf{x}, \mathbf{y})}, 1 \right\};$$

otherwise, the algorithm sets $\mathbf{x}^{t+1} = \mathbf{x} = \mathbf{x}^t$. This defines a Markov chain with a generally mixed continuous-discrete transition probability from \mathbf{x} to \mathbf{y} given by

$$K(\mathbf{x}, \mathbf{y}) = \begin{cases} q(\mathbf{x}, \mathbf{y})\alpha(\mathbf{x}, \mathbf{y}) & \text{if } \mathbf{y} \neq \mathbf{x} \\ 1 - \int_D q(\mathbf{x}, \mathbf{z})\alpha(\mathbf{x}, \mathbf{z})d\mathbf{z} & \text{if } \mathbf{y} = \mathbf{x} \end{cases}.$$

This form of the algorithm is due to Hastings (1970). The Metropolis *et al.* (1953) form takes $q(\mathbf{x}, \mathbf{y}) = q(\mathbf{y}, \mathbf{x})$. A simple variant that is often useful is the independence chain (Tierney, 1991a, 1991b), $q(\mathbf{x}, \mathbf{y}) = j(\mathbf{y})$. Then

$$\alpha(\mathbf{x}, \mathbf{y}) = \min \left\{ \frac{p(\mathbf{y})j(\mathbf{x})}{p(\mathbf{x})j(\mathbf{y})}, 1 \right\} = \min \left\{ \frac{w(\mathbf{y})}{w(\mathbf{x})}, 1 \right\},$$

where $w(\mathbf{x}) = p(\mathbf{x})/j(\mathbf{x})$. The independence chain is closely related to acceptance sampling (Section 4.2) and importance sampling (Section 4.3). But rather than place a low (high) probability of acceptance or a low (high) weight on a draw that is too likely (unlikely) relative to $p(\mathbf{x})$, the independence chain assigns a high (low) probability of accepting the candidate for the next draw.

There is a simple two-step argument that motivates the convergence of the sequence $\{\mathbf{x}^t\}$ generated by the Metropolis-Hastings algorithm to $p(\cdot)$. (This approach is due to Chib and Greenberg, 1994.) First, observe that if any transition probability function $p(\mathbf{x}, \mathbf{y})$ satisfies the reversibility condition

$$p(\mathbf{x})p(\mathbf{x}, \mathbf{y}) = p(\mathbf{y})p(\mathbf{y}, \mathbf{x}),$$

then it has $p(\cdot)$ as its invariant distribution. To see this, note that

$$\int p(\mathbf{x})p(\mathbf{x}, \mathbf{y})d\mathbf{x} = \int p(\mathbf{y})p(\mathbf{y}, \mathbf{x})d\mathbf{x} = p(\mathbf{y}) \int p(\mathbf{y}, \mathbf{x})d\mathbf{x} = p(\mathbf{y}).$$

The second step is to consider the implications of the requirement that $K(\mathbf{x}, \mathbf{y})$ be reversible: $p(\mathbf{x})K(\mathbf{x}, \mathbf{y}) = p(\mathbf{y})K(\mathbf{y}, \mathbf{x})$. For $\mathbf{y} \neq \mathbf{x}$ it implies that

$$p(\mathbf{x})q(\mathbf{x}, \mathbf{y})\alpha(\mathbf{x}, \mathbf{y}) = p(\mathbf{y})q(\mathbf{y}, \mathbf{x})\alpha(\mathbf{y}, \mathbf{x}).$$

Suppose (without loss of generality) that $p(\mathbf{x})q(\mathbf{x}, \mathbf{y}) \geq p(\mathbf{y})q(\mathbf{y}, \mathbf{x})$. If we take $\alpha(\mathbf{y}, \mathbf{x}) = 1$ and $\alpha(\mathbf{x}, \mathbf{y}) = p(\mathbf{y})q(\mathbf{y}, \mathbf{x})/p(\mathbf{x})q(\mathbf{x}, \mathbf{y})$, this equality is satisfied.

In implementing the Metropolis-Hastings algorithm, the transition probability density function must share two important properties. First, it must be possible to generate \mathbf{y} efficiently from $q(\mathbf{x}, \mathbf{y})$. All the methods of Sections 3 and 4 are potential tools for these drawings. (Once again, acceptance sampling is attractive relative to importance sampling.) A second key characteristic of a satisfactory transition process is that the unconditional

acceptance rate not be so low that the time required to generate a sufficient number of distinct \mathbf{x}^t is too great.

6.2 Mathematical background

Let $\{\mathbf{x}^t\}_{t=0}^{\infty}$ be a Markov chain defined on $D \subseteq \mathfrak{R}^m$ with transition kernel $K: D \times D \rightarrow \mathfrak{R}^+$ such that, with respect to a σ -finite measure ν on the Borel σ -field of \mathfrak{R}^m , for ν -measurable A ,

$$P(\mathbf{x}^t \in A | \mathbf{x}^{t-1} = \mathbf{x}) = \int_A K(\mathbf{x}, \mathbf{y}) d\nu(\mathbf{y}) + r(\mathbf{x})\chi_A(\mathbf{x}),$$

$$\text{where } r(\mathbf{x}) = 1 - \int_D K(\mathbf{x}, \mathbf{y}) d\nu(\mathbf{y}) \text{ and } \chi_A(\mathbf{x}) = \begin{cases} 1 & \text{if } \mathbf{x} \in A \\ 0 & \text{if } \mathbf{x} \notin A \end{cases}.$$

The measure ν will be Lebesgue for continuous distributions and discrete for discrete distributions.

The transition kernel K is substochastic: it defines only the distribution of accepted candidates. Assume that K has no absorbing states, so that $r(\mathbf{x}) < 1 \forall \mathbf{x} \in D$. The corresponding substochastic kernel over t steps is then defined iteratively,

$$\mathbf{K}^{(t)}(\mathbf{x}, \mathbf{y}) = \int \mathbf{K}^{(t-1)}(\mathbf{x}, \mathbf{z}) K(\mathbf{z}, \mathbf{y}) d\nu(\mathbf{z}) + \mathbf{K}^{(t-1)}(\mathbf{x}, \mathbf{y}) r(\mathbf{y}) + [r(\mathbf{x})]^{t-1} K(\mathbf{x}, \mathbf{y}).$$

This describes all t -step transitions that involve at least one accepted move. As a function of \mathbf{y} it is the p.d.f. with respect to ν of \mathbf{x}^t , given $\mathbf{x}^0 = \mathbf{x}$, excluding realizations with $\mathbf{x}^j = \mathbf{x} \forall j = 1, 2, \dots, t$.

An invariant distribution for the Markov chain is a function $p(\mathbf{x})$ that satisfies

$$\begin{aligned} P(A) &= \int_A p(\mathbf{x}) d\nu(\mathbf{x}) = \int_D \left\{ \int_A K(\mathbf{x}, \mathbf{y}) d\nu(\mathbf{y}) + r(\mathbf{x})\chi_A(\mathbf{x}) \right\} p(\mathbf{x}) d\nu(\mathbf{x}) \\ &= \int_D P(\mathbf{x}^t \in A | \mathbf{x}^{t-1} = \mathbf{x}) p(\mathbf{x}) d\nu(\mathbf{x}) \end{aligned}$$

for all ν -measurable A . Let $D^* = \{\mathbf{x} \in D: p(\mathbf{x}) > 0\}$. The kernel K is *p-irreducible* if for all $\mathbf{x} \in D^*$, $P(A) > 0$ implies that $P(\mathbf{x}^t \in A | \mathbf{x}^0 = \mathbf{x}) > 0$ for some $t \geq 1$. It is *aperiodic* if there exists no ν -measurable partition $D = \bigcup_{s=0}^{r-1} B_s$ ($r \geq 2$) such that

$$P(\mathbf{x}^t \in B_{t \bmod(r)} | \mathbf{x}^0 = \mathbf{x} \in B_0) = 1 \quad \forall t.$$

Define $|f| = \int_D |f(\mathbf{x})| d\nu(\mathbf{x})$ for all ν -measurable functions f defined on D . If K is *p-irreducible* and *aperiodic*, then

$$(A) \quad \text{For all } \mathbf{x}^0 \in D, \lim_{t \rightarrow \infty} |\mathbf{K}^{(t)} - p| = 0;$$

$$(B) \quad \text{If } g \text{ is } p\text{-integrable, then for all } \mathbf{x}^0 \in D,$$

$$N^{-1} \sum_{t=1}^N g(\mathbf{x}^t) \xrightarrow{a.s.} \int_D g(\mathbf{x}) p(\mathbf{x}) d\nu(\mathbf{x})$$

(Tierney, 1991b, based on Numelin, 1984).

The kernel K is Harris recurrent if $P[\mathbf{x}^t \in B \text{ i.o.}] = 1$ for all ν -measurable B with $\int_B p(\mathbf{x}) d\nu(\mathbf{x}) > 0$ and all $\mathbf{x}^0 \in D$. (A general discussion of recurrence is provided by Numelin (1984, Chapter 3).) If K is p -irreducible and Harris recurrent, then

(C) The invariant probability distribution $p(\mathbf{x})$ is unique.

(Numelin, 1984, Corollary 5.2; Tierney, 1991b, Section 3.1). Harris recurrence eliminates situations like the one shown in Figure 4, where the support is disconnected and the Markov chain is the Gibbs sampler. Note that if $\mathbf{x}^0 \in D_i$, it is impossible that $\mathbf{x}^t \in D_j$ ($j \neq i$, any $t > 0$). In the situation portrayed in Figure 4, there are two invariant distributions, one for D_1 (reached if $\mathbf{x}^0 \in D_1$) and one for D_2 (reached if $\mathbf{x}^0 \in D_2$).

6.3 Convergence of the Gibbs sampler

The Gibbs sampler requires that the conditional probability density functions

$$p[\mathbf{x}_{(i)} | \mathbf{x}_{(-i)}] = p(\mathbf{x}) / \int_{x_{(i)}} p(\mathbf{x}) d\nu_i(x_{(i)}) \quad (i = 1, \dots, k)$$

be well-defined on their supports. In this case the transition kernel density is

$$K_G(\mathbf{x}, \mathbf{y}) = \prod_{j=1}^k p[\mathbf{y}_{(j)} | \mathbf{x}_{(j)} (j > 1), \mathbf{y}_{(j)} (j < 1)].$$

If $\mathbf{x}^0 \in D$, then $p(\mathbf{x})$ is the density of an invariant distribution of the chain defined by K_G :

$$\begin{aligned} & \int_D K_G(\mathbf{x}, \mathbf{y}) p(\mathbf{x}) d\nu(\mathbf{x}) \\ &= p(\mathbf{y}_{(k)} | \mathbf{y}_{(-k)}) \int p[\mathbf{y}_{(k-1)} | \mathbf{x}_{(k)}, \mathbf{y}_{(j)} (j < k-1)] \int p[\mathbf{y}_{(k-2)} | \mathbf{x}_{(k)}, \mathbf{x}_{(k-1)}, \mathbf{y}_{(j)} (j < k-2)] \\ & \quad \int p[\mathbf{y}_{(2)} | \mathbf{y}_{(1)}, \mathbf{x}_{(j)} (j > 2)] \int p[\mathbf{y}_{(1)} | \mathbf{x}_{(j)} (j > 1)] \int p[\mathbf{x}_{(1)} | \mathbf{x}_{(j)} (j > 1)] d\nu_1(\mathbf{x}_{(1)}) \\ & \quad p[\mathbf{x}_{(2)} | \mathbf{x}_{(j)} (j > 2)] d\nu_2(\mathbf{x}_{(2)}) p[\mathbf{x}_{(3)} | \mathbf{x}_{(j)} (j > 3)] d\nu_3(\mathbf{x}_{(3)}) \\ & \quad \int p[\mathbf{x}_{(k-1)} | \mathbf{x}_{(k)}] d\nu_{k-1}(\mathbf{x}_{(k-1)}) p[\mathbf{x}_{(k)}] d\nu_k(\mathbf{x}_{(k)}) \\ &= p(\mathbf{y}_{(k)} | \mathbf{y}_{(-k)}) \int p[\mathbf{y}_{(k-1)} | \mathbf{x}_{(k)}, \mathbf{y}_{(j)} (j < k-1)] \int p[\mathbf{y}_{(k-2)} | \mathbf{x}_{(k)}, \mathbf{x}_{(k-1)}, \mathbf{y}_{(j)} (j < k-2)] \\ & \quad \int p[\mathbf{y}_{(2)} | \mathbf{y}_{(1)}, \mathbf{x}_{(j)} (j > 2)] \int p[\mathbf{y}_{(1)} | \mathbf{x}_{(j)} (j > 2)] p[\mathbf{x}_{(3)} | \mathbf{x}_{(j)} (j > 3)] d\nu_3(\mathbf{x}_{(3)}) \\ & \quad \int p[\mathbf{x}_{(k-1)} | \mathbf{x}_{(k)}] d\nu_{k-1}(\mathbf{x}_{(k-1)}) p[\mathbf{x}_{(k)}] d\nu_k(\mathbf{x}_{(k)}) \end{aligned}$$

$$\begin{aligned}
&= \mathbb{P}(\mathbf{y}_{(k)} | \mathbf{y}_{(-k)}) \int \mathbb{P}[\mathbf{y}_{(k-1)} | \mathbf{x}_{(k)}, \mathbf{y}_{(j)} (j < k-1)] \int \mathbb{P}[\mathbf{y}_{(k-2)} | \mathbf{x}_{(k)}, \mathbf{x}_{(k-1)}, \mathbf{y}_{(j)} (j < k-2)] \\
&\stackrel{\text{L}}{=} \int \mathbb{P}[\mathbf{y}_{(1)}, \mathbf{y}_{(2)} | \mathbf{x}_{(j)} (j > 3)] \stackrel{\text{L}}{=} \int \mathbb{P}[\mathbf{x}_{(k-1)} | \mathbf{x}_{(k)}] d\nu_{k-1}(\mathbf{x}_{(k-1)}) \mathbb{P}[\mathbf{x}_{(k)}] d\nu_k(\mathbf{x}_{(k)}) \\
&= \stackrel{\text{L}}{=} \\
&= \mathbb{P}(\mathbf{y}_{(k)} | \mathbf{y}_{(-k)}) \int \mathbb{P}[\mathbf{y}_{(k-1)} | \mathbf{x}_{(k)}, \mathbf{y}_{(j)} (j < k-1)] \int \mathbb{P}[\mathbf{y}_{(k-2)} | \mathbf{x}_{(k)}, \mathbf{x}_{(k-1)}, \mathbf{y}_{(j)} (j < k-2)] \\
&\quad \cdot \mathbb{P}[\mathbf{y}_{(1)}, \mathbf{y}_{(2)}, \mathbb{K}, \mathbf{y}_{(k-3)} | \mathbf{x}_{(k-1)}, \mathbf{x}_{(k)}] \mathbb{P}[\mathbf{x}_{(k-1)} | \mathbf{x}_{(k)}] d\nu_{k-1}(\mathbf{x}_{(k-1)}) \mathbb{P}[\mathbf{x}_{(k)}] d\nu_k(\mathbf{x}_{(k)}) \\
&= \mathbb{P}(\mathbf{y}_{(k)} | \mathbf{y}_{(-k)}) \int \mathbb{P}[\mathbf{y}_{(k-1)} | \mathbf{x}_{(k)}, \mathbf{y}_{(j)} (j < k-1)] \mathbb{P}[\mathbf{y}_{(1)}, \mathbf{y}_{(2)}, \mathbb{K}, \mathbf{y}_{(k-2)} | \mathbf{x}_{(k)}] \mathbb{P}[\mathbf{x}_{(k)}] d\nu_k(\mathbf{x}_{(k)}) \\
&= \mathbb{P}(\mathbf{y}_{(k)} | \mathbf{y}_{(-k)}) \mathbb{P}[\mathbf{y}_{(1)}, \mathbf{y}_{(2)}, \mathbb{K}, \mathbf{y}_{(k-1)}] = \mathbb{P}(\mathbf{y}).
\end{aligned}$$

If ν is discrete, p -irreducibility of K_G is sufficient for results (A), (B), and (C) in Section 6.2 (Tierney, 1991b). The continuous (Lebesgue measure) case is technically more difficult, but it may be shown that three simple conditions are jointly sufficient for results (A), (B), and (C) (Roberts and Smith, 1992):

- (1) $p(\mathbf{x})$ is lower semicontinuous at 0;
- (2) $\int p(\mathbf{x}) dx_i$ is locally bounded ($i = 1, \mathbb{K}, k$);
- (3) D^* is connected.

A function $h(\mathbf{x})$ is lower semicontinuous at 0 if, for all \mathbf{x} with $h(\mathbf{x}) > 0$, there exists an open neighborhood $N_{\mathbf{x}} \supset \mathbf{x}$ and $\varepsilon > 0$ such that for all $\mathbf{y} \in N_{\mathbf{x}}$, $h(\mathbf{y}) \geq \varepsilon > 0$. This condition rules out situations like the one shown in Figure 5, where the probability density is uniform on a closed set. For any point \mathbf{x} on the boundary there is no open neighborhood $N_{\mathbf{x}} \supset \mathbf{x}$ such that for all $\mathbf{y} \in N_{\mathbf{x}}$, $h(\mathbf{y})$ is bounded away from 0. The point A is absorbing.

The local boundedness condition, together with lower semicontinuity at 0, ensures that the Markov chain is aperiodic. It does so by guaranteeing that for the sequence of support sets $B^t(\mathbf{x}) = \{\mathbf{y} \in D^* : K_G^{(t)}(\mathbf{x}, \mathbf{y}) > 0\}$, $B^t(\mathbf{x}) \subseteq B^{t+1}(\mathbf{x})$ for all $t \geq 1$ and all $\mathbf{x} \in D^*$ (Roberts and Smith, 1992, Lemma 3).

Connectedness of D^* , together with conditions (1) and (2), implies that the Gibbs sampler is p -irreducible (Roberts and Smith, 1992, Theorem 2). Conditions (2) and (3) further imply that the probability measure P corresponding to $p(\mathbf{x})$ is absolutely

continuous, and consequently (Tierney, 1991b, Corollary 1) the Gibbs sampler is Harris recurrent. Therefore $p(\mathbf{x})$ is the unique invariant probability density of the Gibbs sampler.

These conditions are by no means necessary for convergence of the Gibbs sampler; Tierney (1991b) provides substantially weaker conditions. However, the conditions stated here are satisfied for a very wide range of problems in economics and are much easier to verify than the weaker conditions.

6.4 Convergence of the Metropolis-Hastings algorithm

Take the transition probability density function $q(\mathbf{x}, \mathbf{y})$ of Section 6.1 to be a Markov chain kernel with respect to ν , $q: D^* \times D^* \rightarrow \mathfrak{R}^+$. Defining $\alpha: D^* \times D^* \rightarrow [0,1]$ as before, define $K_H: D^* \times D^* \rightarrow \mathfrak{R}^+$ by

$$K_H(\mathbf{x}, \mathbf{y}) = q(\mathbf{x}, \mathbf{y})\alpha(\mathbf{x}, \mathbf{y}).$$

This is the substochastic kernel governing transitions of the chain from \mathbf{x} to \mathbf{y} that are accepted according to the probability $\alpha(\mathbf{x}, \mathbf{y})$. The distribution $p(\mathbf{x})d\nu(\mathbf{x})$ is invariant if for all ν -measurable sets A ,

$$P(A) = \int_A p(\mathbf{x})d\nu(\mathbf{x}) = \int_D P[\mathbf{y} \in A|\mathbf{x}]p(\mathbf{x})d\nu(\mathbf{x}).$$

Recalling that

$$P[\mathbf{y} \in A|\mathbf{x}] = \int_A K_H(\mathbf{x}, \mathbf{y})d\nu(\mathbf{y}) + \left[1 - \int_D K_H(\mathbf{x}, \mathbf{z})d\nu(\mathbf{z})\right]\chi_A(\mathbf{x}),$$

$$\begin{aligned} & \int_D P[\mathbf{y} \in A|\mathbf{x}]p(\mathbf{x})d\nu(\mathbf{x}) \\ &= \int_D \int_A K_H(\mathbf{x}, \mathbf{y})d\nu(\mathbf{y})p(\mathbf{x})d\nu(\mathbf{x}) \\ & \quad + \int_D \chi_A(\mathbf{x})p(\mathbf{x})d\nu(\mathbf{x}) - \int_D \int_D K_H(\mathbf{x}, \mathbf{y})d\nu(\mathbf{y})\chi_A(\mathbf{x})p(\mathbf{x})d\nu(\mathbf{x}) \\ &= \int_D \int_A K_H(\mathbf{x}, \mathbf{y})d\nu(\mathbf{y})p(\mathbf{x})d\nu(\mathbf{x}) \\ & \quad + \int_A p(\mathbf{x})d\nu(\mathbf{x}) - \int_A \int_D K_H(\mathbf{x}, \mathbf{y})d\nu(\mathbf{y})p(\mathbf{x})d\nu(\mathbf{x}). \end{aligned}$$

Since $p(\mathbf{x})K_H(\mathbf{x}, \mathbf{y}) = \min[p(\mathbf{y})q(\mathbf{y}, \mathbf{x}), p(\mathbf{x})q(\mathbf{x}, \mathbf{y})]$ is symmetric in \mathbf{x} and \mathbf{y} , the last expression reduces to $\int_A p(\mathbf{x})d\nu(\mathbf{x}) = P(\mathbf{x} \in A)$.

From this derivation it is clear that invariance is unaffected by an arbitrary scaling of $K_H(\mathbf{x}, \mathbf{y})$ by a constant c . The choice of c affects the properties of the Metropolis-Hastings algorithm in important practical ways. Larger values of c result in fewer rejected

draws but slower convergence to $p(\mathbf{x})$, whereas smaller values of c increase the proportion of rejected candidates but accelerate the rate of convergence to $p(\mathbf{x})$.

Roberts and Smith (1992) show that the convergence properties of the Hastings-Metropolis algorithm are inherited from those of $q(\mathbf{x}, \mathbf{y})$: if q is aperiodic and p -irreducible, then so is the Hastings-Metropolis algorithm. If $q(\mathbf{x}, \mathbf{y})$ is constructed as a Gibbs sampler (as is often the case), then the conditions set forth in Section 6.3 may be used to verify aperiodicity and p -irreducibility. A Hastings-Metropolis chain is always Harris recurrent, and therefore the invariant distribution p is unique.

6.5 Assessing convergence and numerical accuracy

In any practical application one is concerned with the discrepancy between $E[g(\mathbf{x})] = \int_D g(\mathbf{x})p(\mathbf{x})d\mathbf{x}$ and its numerical approximation $N^{-1} \sum_{i=1}^N g(\mathbf{x}_i)$. Consider the decomposition

$$N^{-1} \sum_{t=1}^N g(\mathbf{x}^t) - E[g(\mathbf{x})] = \left\{ E \left[N^{-1} \sum_{t=1}^N g(\mathbf{x}^t) \middle| \mathbf{x}^0 \right] - E[g(\mathbf{x})] \right\} + \left\{ N^{-1} \sum_{t=1}^N g(\mathbf{x}^t) - E \left[N^{-1} \sum_{t=1}^N g(\mathbf{x}^t) \middle| \mathbf{x}^0 \right] \right\} = A_N(\mathbf{x}^0) + B_N(\mathbf{x}^0).$$

The term $A_N(\mathbf{x}^0)$ is nonstochastic and in general nonzero, but $\lim_{N \rightarrow \infty} A_N(\mathbf{x}^0) = 0$ if conditions set forth earlier in this section are satisfied. The purpose of a transient or burn-in phase is to reduce $A_N(\mathbf{x}^0)$, but for any finite transient period it will still be the case in general that $A_N(\mathbf{x}^0) \neq 0$. This difficulty is termed the convergence or sensitivity to initial conditions problem. The term $B_N(\mathbf{x}^0)$ is stochastic and is the analog of $E_N - E$ or $I_N - I$ for acceptance or importance sampling. This term vanishes as $N \rightarrow \infty$, but assessing its size is complicated by the fact that $\{\mathbf{x}^t\}$ is neither independently nor identically distributed. This difficulty may be termed the numerical accuracy problem.

A leading cause of slow convergence is multimodality of the probability distribution, for example, as shown in Figure 6 for a Gibbs sampler. In the limit multimodality approaches disconnectedness of the support, and increasingly large values of N are required for $A_N(\mathbf{x}^0)$ to be close to 0. This difficulty is essentially undetectable given a single Markov chain: for a chain of any fixed length, one can imagine multimodal distributions for which the probability of leaving the neighborhood of a single mode is arbitrarily small. This sort of convergence problem is precisely the same as the multimodality problem in optimization, where iteration from a single starting value can by itself never guarantee the determination of a global optimum. Multimodal disturbances are difficult to manage by any method, including those discussed in Section 4. In the context of the Markov chain Monte Carlo algorithms, the question may be recast as one of sensitivity

to initial conditions: \mathbf{x}_A^0 , \mathbf{x}_B^0 , and \mathbf{x}_C^0 will lead to quite different chains, in Figure 6, unless the simulations are sufficiently long.

A Markov chain Monte Carlo algorithm can be made fully robust against sensitivity to initial conditions by constructing many very long chains. Just how one should trade off the number of chains against their length for a given budget of computation time is problem specific and as a practical matter not yet fully understood. Many of the issues involved are discussed by Gelman and Rubin (1992), Geyer (1992), and their disciples and cited works. In an extreme variant of the multiple chains approach, the chain is restarted many times, with initial values chosen independently and identically distributed from an appropriate distribution. But finding an appropriate distribution may be difficult: one that is too concentrated reintroduces the difficulties exemplified by Figure 6; one that is too diffuse may require excessively long chains for convergence. These problems aside, proper use of the output of Markov chain Monte Carlo in a situation of multimodality requires specialized diagnostics; Zellner and Min (1992) have obtained some interesting results of this kind. At the other extreme a single starting value is used. This approach provides the largest number of iterations toward convergence, but diagnostics of the type of problem illustrated in Figure 6 will not be as clear.

In specific circumstances a central limit theorem applies to $B_N(\mathbf{x}^0)$, which may therefore be used to assess the numerical accuracy problem. To develop one set of such circumstances, suppose that the Markov chain is stationary. This could be guaranteed by drawing \mathbf{x}^0 from the stationary distribution. Such a drawing would be time consuming (if not, i.i.d. sampling from $p(\mathbf{x})$ is possible), but only one is required. Alternatively, one could iterate the chain many times beginning from an arbitrary initial value, discard all but the last iteration, and take this value as drawn from the stationary distribution to begin a new chain. Suppose $\text{var}_p[g(\mathbf{x})]$ is finite and denote $\gamma_i = \text{cov}_K[g(\mathbf{x}^i), g(\mathbf{x}^{i+1})]$. A Markov chain with kernel K is *reversible* if $K(\mathbf{x}, \mathbf{y}) = K(\mathbf{y}, \mathbf{x})$ for all $\mathbf{x}, \mathbf{y} \in D$. Hastings-Metropolis chains are always reversible; Gibbs sampling chains are not (Geyer, 1992, Section 2). If the Markov chain is stationary, p -irreducible, and reversible, then

$$N \text{var}(g_N) \xrightarrow{a.s.} \sigma^2 = \sum_{i=-\infty}^{\infty} \gamma_i,$$

and if $\sigma^2 < \infty$, then

$$\sqrt{N}(g_N - G) \xrightarrow{d} N(0, \sigma^2)$$

(Kipnis and Varadhan, 1986).

In the absence of reversibility, known sufficient conditions for central limit theorems are strong and difficult to verify. For example, if for some $m < \infty$ $P(\mathbf{x}^{t+m} \in A | \mathbf{x}^t = \mathbf{x}) / \int_A p(\mathbf{x}) d\nu(\mathbf{x})$ is bounded below uniformly in \mathbf{x} , then D is a small state

space and $\{\mathbf{x}^t\}$ is uniformly ergodic (Tierney, 1991b, Proposition 2). Then if $\text{var}_p[g(\mathbf{x})]$ is finite, there exists $\sigma^2 < \infty$ such that $\sqrt{N}(g_N - G) \xrightarrow{d} N(0, \sigma^2)$. The boundedness condition, however, is generally difficult to establish.

In neither circumstance is there a known sufficient condition for approximation of the variance term σ^2 of the central limit theorem. The problem is formally quite similar to estimating the variance of the sample mean $\bar{z}_N = N^{-1} \sum_{t=1}^N z_t$ of a stationary time series $\{z_t\}$. In the time series problem, well-established mixing conditions (rates of decay for $\text{cov}(z_t, z_{t+i})$) are sufficient for consistent estimation of $\text{var}(\bar{z}_N)$ (e.g., Hannan, 1970, pp. 207-210). In time series applications these conditions remain assumptions. The difficulty in applying these conditions to Markov chain Monte Carlo is that they cannot be established from verifiable fundamentals.

Nevertheless, applications of the time series procedures as if sufficient mixing conditions obtain appear to give quite reliable results for real problems in economics. That is, applying a central limit theorem as if the output of the Markov chain Monte Carlo algorithm were a stationary process satisfying the mixing conditions yields accurate probability statements about the output of the same algorithm applied to the same problem with a new starting value and initial seed for the random number generator (Geweke, 1992a; Geyer, 1992). This leads to a conservative but practical procedure for assessing the accuracy and reliability of Markov chain Monte Carlo. First, execute several short runs -- a burn-in of 50 to 100 iterations followed by a chain of length $N = 500$ or $N = 1000$ is sufficient for many problems. Examine the g_N and their standard errors as assessed by conventional time series procedures for a single time series to see whether the scatter of each g_N across the short runs is consistent with these standard errors. If necessary, increase the length of the short runs until this consistency is achieved. Second, choose the last value of one of the short runs, and use it as the starting value of a long run of from $N = 10^4$ to $N = 10^6$ iterations. As a final check, compare the g_N from the single long run with the confidence intervals implied by the short runs. Report the final value of g_N , together with its numerical standard error as computed by time series methods for a single series.

7. Some examples

The usefulness of all of these methods lies as much in their appropriate combination as in the application of any one individually. We turn now to some examples that illustrate some useful combinations and in the process treat a few topics closely related to integration and simulation.

7.1 Stochastic volatility

Models in which the volatility of asset returns varies smoothly over time have received considerable attention in recent years. (For a survey of several approaches, see Bollerslev, Chou, and Kroner, 1992.) Persistent but changing volatility is an evident characteristic of returns data. Since the conditional distribution of returns is relevant in the theory of portfolio allocation, proper treatment of volatility is important. Time-varying volatility also affects the properties of real growth and business cycle models.

A simple model of time-varying volatility is the stochastic volatility model, the descriptive properties of which have been examined by a series of investigators beginning with Taylor (1986). The approach here closely follows that of Jacquier, Polson, and Rossi (1994). Let r_t denote the one-period return of a single asset, and let \mathbf{x}_t be a vector of deterministic time series such as indicators for day of the week, holidays, etc. A simple stochastic volatility model is

$$r_t = \beta' \mathbf{x}_t + \varepsilon_t, \quad \varepsilon_t = h_t^{1/2} u_t \quad (7.1.1)$$

$$\log h_t = \alpha + \delta \log h_{t-1} + \sigma_v v_t \quad (7.1.2)$$

$$\begin{pmatrix} u_t \\ v_t \end{pmatrix} \stackrel{iid}{\sim} \mathbf{N}(0, \mathbf{I}_2). \quad (7.1.3)$$

At time T an economic agent is concerned with future returns r_{T+1}, \dots, r_{T+q} through an expected utility function

$$\mathbb{E}[\mathbf{V}(r_{T+1}, \dots, r_{T+q}; \mathbf{z}) | \Phi_T] = \mathbb{E}[\mathbf{V}(\mathbf{r}_q; \mathbf{K}) | \Phi_T], \quad (7.1.4)$$

where \mathbf{z} is a generic vector of other arguments which may be known or unknown at time T .

Evaluation of this expected utility function requires the solution of an integration problem. We will consider this problem for three different specifications of the information set Φ_T in turn. Denoting $\mathbf{r}_T = (r_{1,K}, r_T)'$, $\mathbf{x}_{T+q} = (x_{1,K}, \mathbf{x}_{T+q})'$, $\theta' = (\beta', \alpha, \delta, \sigma_v)$, and $\mathbf{h}_T = (h_{1,K}, h_T)'$, these are

$$\Phi_T^{(1)} = \{\mathbf{r}_T, \mathbf{x}_{T+q}, \theta, \mathbf{h}_T\}; \quad \Phi_T^{(2)} = \{\mathbf{r}_T, \mathbf{x}_{T+q}, \theta\}; \quad \Phi_T^{(3)} = \{\mathbf{r}_T, \mathbf{x}_{T+q}\}.$$

As one may readily verify, deterministic approximations of the type discussed in Section 2 are inconvenient for this problem. Even explicit expressions of the integrals in closed form

are awkward and unrevealing. Simulation methods are much more direct and have the added advantage that one set of simulations can suffice for several alternative values of the other arguments \mathbf{z} in (7.1.4). These arguments might include taste parameters or the values of decision variables which themselves do not affect \mathbf{r}_q . (Section 7.2 provides an example involving explicit optimization.)

The solution for the problem for $\Phi_T^{(1)}$ is simple. In the notation of Section 4, repeated period-by-period simulation of $\mathbf{x} = \mathbf{r}_q$ provides an independent identically distributed sample $\{\tilde{\mathbf{r}}_q^{(i)}\}_{i=1}^N$ with a probability density $p(\mathbf{x}) = p(\mathbf{r}_q | \Phi_T^{(1)})$ that we have not even expressed. Then

$$\mathbb{E}\left[V(r_{T+1}, \mathbb{K}, r_{T+q}; \mathbb{K}) | \Phi_T^{(1)}\right] = \int g(\mathbf{x}) p(\mathbf{x}) d\mathbf{x},$$

where $g(\mathbf{x}) = V(\mathbf{x}; \mathbb{K}) = V(r_{T+1}, \mathbb{K}, r_{T+q}; \mathbb{K})$. Consequently,

$$\mathbb{E}\left[V(r_{T+1}, \mathbb{K}, r_{T+q}; \mathbb{K}) | \Phi_T^{(1)}\right] \approx N^{-1} \sum_{i=1}^N V(\tilde{\mathbf{r}}_q^{(i)}).$$

The problem for $\Phi_T^{(2)}$ is more difficult. Rather than \mathbf{h}_T itself, the agent has available only

$$\begin{aligned} p(\mathbf{h}_T | \mathbf{r}_T, \mathbf{x}_T, \theta) &= p(\mathbf{h}_T, \mathbf{r}_T | \mathbf{x}_T, \theta) / p(\mathbf{r}_T) \\ &= p(\mathbf{r}_T | \mathbf{h}_T, \mathbf{x}_T, \theta) p(\mathbf{h}_T | \mathbf{x}_T, \theta) / p(\mathbf{r}_T) \propto p(\mathbf{r}_T | \mathbf{h}_T, \mathbf{x}_T, \theta) p(\mathbf{h}_T | \mathbf{x}_T, \theta) \\ &= (2\pi)^{-T/2} \prod_{t=1}^T h_t^{-1/2} \exp\left[-\sum_{t=1}^T \varepsilon_t^2 / 2h_t\right] \\ &\quad \cdot (2\pi)^{-T/2} \sigma_v^{-T} h_t^{-1} \exp\left[-\sum_{t=1}^T (\log h_t - \alpha - \delta \log h_{t-1})^2 / 2\sigma_v^2\right] \\ &\propto \prod_{t=1}^T h_t^{-3/2} \exp\left[-\sum_{t=1}^T \varepsilon_t^2 / 2h_t\right] \exp\left[-\sum_{t=1}^T (\log h_t - \alpha - \delta \log h_{t-1})^2 / 2\sigma_v^2\right], \end{aligned} \quad (7.1.5)$$

where $\varepsilon_t = r_t - \beta' \mathbf{x}_t$. The simple Monte Carlo solution of the previous problem could be extended to this one if one could draw an i.i.d. sample $\{\tilde{\mathbf{h}}_T^{(i)}\}_{i=1}^N$ from the distribution implied by the last kernel. This is clearly not possible, nor are there obvious source or importance sampling distributions for the methods of Sections 4.2 or 4.3.

This problem can be solved in a number of ways, and a comparison of three alternatives is instructive. All begin with the kernels of the conditional probability densities for individual h_t implied by (7.1.5). For $t = 2, \dots, T-1$ the kernel is

$$p[h_t | h_s (t \neq s), \theta, \varepsilon_t] \propto h_t^{-3/2} \exp(-\varepsilon_t^2 / 2h_t) \exp\left[-(\log h_t - \mu_t)^2 / 2\sigma^2\right], \quad (7.1.6)$$

where

$$\mu_t = \frac{\alpha(1 - \delta) + \delta(\log h_{t-1} + \log h_{t+1})}{1 + \delta^2}, \quad \sigma^2 = \frac{\sigma_v^2}{1 + \delta^2}.$$

(Similar expressions for h_1 and h_T may be constructed.)

The first two approaches construct a Gibbs sampler for the h_t , drawing and successively replacing h_1, h_2, \dots, h_T . Each cycle of drawing and replacement produces the next realization of $\tilde{\mathbf{h}}_T^{(i)}$ in the Markov chain. Note from (7.1.5) that $\lim_{h_t \rightarrow 0} p(h_t | \mathbf{r}_t, \mathbf{x}_t, \theta) = 0$ for any $t = 1, \dots, T$, and since the support of \mathbf{h}_T is the positive orthant of \mathfrak{R}^T the probability density function of \mathbf{h}_T is lower semicontinuous at 0. The remaining sufficient conditions for convergence of the Gibbs sampler are clearly satisfied. Conditional on each $\tilde{\mathbf{h}}_T^{(i)}$ in the chain, draw a single $\tilde{\mathbf{r}}_q^{(i)}$ as in the problem for $\Phi_T^{(1)}$. Since $\left| p(\tilde{\mathbf{h}}_T^{(i)}) - p(\mathbf{h}_T | \mathbf{r}_T, \mathbf{x}_T, \theta) \right| \rightarrow 0$, it follows that $\left| p(\tilde{\mathbf{r}}_q^{(i)}) - p(\mathbf{r}_q | \mathbf{r}_T, \mathbf{x}_{T+q}, \theta) \right| \rightarrow 0$. Both approaches work directly with the conditional distribution of $H_t = \log h_t$, which from (7.1.6) is given by

$$\log p(H_t | H_s (s \neq t), \theta, \varepsilon_t) = -\exp(-\varepsilon_t^2/2) \exp(-H_t) - (H_t - \mu_t^*)^2 / 2\sigma^2 \quad (7.1.7)$$

(up to an additive constant), where $\mu_t^* = \mu_t - .5\sigma^2$, but differ in the method for obtaining H_t .

The first approach is to use acceptance sampling. A reasonable source distribution is $N(\mu_t^*, \sigma^2)$, for which the acceptance probability is

$$\exp[-(\varepsilon_t^2/2) \exp(-H_t)] = \exp(-\varepsilon_t^2/2h_t).$$

The acceptance probability falls below .01 if and only if ε_t^2/h_t exceeds 9.2, which is highly unlikely if the model reasonably well describes the distribution of the returns r_t . The acceptance probability could be improved somewhat using the optimizing procedures set out in Section 3.2, but given the favorable acceptance probabilities for the $N(\mu_t^*, \sigma^2)$ source distribution, the additional overhead might not be warranted.

The second approach is to note that the log-conditional kernel densities (7.1.7) are strictly concave and apply the adaptive method of Gilks and Wild (1992). Their algorithm (described in Section 3.2) may be initialized by noting that $H_t = \mu_t^*$ lies to the left of the mode of the log-conditional and a solution of $(1 - H_t + H_t^2/2) \exp(-\varepsilon_t^2/2) - (H_t - \mu_t^*)/\sigma^2$ lies to the right of the mode. Except for the method of drawing H_t , the solution of the problem proceeds as in the first approach.

The third approach is to construct a Metropolis-Hastings independence chain. This is done by forming a Metropolis step M_t for each h_t and then combining all T steps into a single transition $M = M_1 M_2 \dots M_T$. At each M_t either a candidate new value is accepted or the old value of h_t is retained. Thus, when M operates on the old \mathbf{h}_T it generally produces a mixture of old and new h_t in the new \mathbf{h}_T . The transition kernel M is p-irreducible and aperiodic, and an argument like the one in Section 6.4 shows that $p(\mathbf{h}_T | \mathbf{r}_T, \mathbf{x}_T, \theta)$ is the invariant distribution of M (Jacquier, Polson, and Rossi, 1994, Section 2). A useful distribution for the Metropolis-Hastings independence chain is the gamma distribution for h_t^{-1} with shape parameter $a = [1 - 2 \exp(\sigma^2)] / [1 - \exp(\sigma^2)] + .5$ and scale parameter

$\lambda = (a - 1)\exp(\mu_t + .5\sigma^2) + .5\varepsilon_t^2$. Combined with an appropriate scaling of the transition kernel, as discussed in Section 6.4, this chain produces convergence at a practical rate (see Jacquier, Polson, and Rossi, 1994, Section 2.4, for details).

The solution of the problem for $\Phi_T^{(2)}$ is directly usable in the solution of the problem for $\Phi_T^{(3)}$, in the context of the Gibbs sampler. From the form of (7.1.1)-(7.1.3) the probability density kernel for θ and \mathbf{h}_T underlying the expectations operator in (7.1.4) is

$$\prod_{t=1}^T h_t^{-1/2} \exp\left[-\sum_{t=1}^T (r_t - \beta' \mathbf{x}_t)^2 / 2h_t\right] \cdot \sigma_v^{-T} \exp\left[-\sum_{t=1}^T (\log h_t - \alpha - \delta \log h_{t-1})^2 / 2\sigma_v^2\right] p(\beta, \alpha, \delta, \sigma_v), \quad (7.1.8)$$

where $p(\beta, \alpha, \delta, \sigma_v)$ is the prior probability density function of $\theta' = (\beta', \alpha, \delta, \sigma_v)$. A Gibbs sampler with blocking (\mathbf{h}_T, θ) will alternate drawing and substitution for $\mathbf{h}_T | \mathbf{r}_T, \mathbf{x}_T, \theta$ and $\theta | \mathbf{r}_T, \mathbf{x}_T, \mathbf{h}_T$. The drawing for \mathbf{h}_T is the same one constructed to solve the problem for $\Phi_T^{(2)}$.

The second drawing is facilitated by noting that the kernel of (7.1.8) in θ may be expressed

$$\propto \prod_{t=1}^T \exp\left[-\sum_{t=1}^T (r_t - \beta' \mathbf{x}_t)^2 / 2h_t\right] \cdot \sigma_v^{-(T+1)} \exp\left[-\sum_{t=1}^T (\log h_t - \alpha - \delta \log h_{t-1})^2 / 2\sigma_v^2\right]$$

if the prior probability distribution has the conventional improper kernel $p(\beta, \alpha, \delta, \sigma_v) \propto \sigma_v^{-1}$. Thus, β and $(\alpha, \delta, \sigma_v)$ are conditionally independent. In each case the distribution follows from standard treatments of Bayesian learning about a linear model (e.g., Poirier, 1995, Section 9.9):

$$\beta \sim N(\mathbf{b}, \mathbf{Q}^{-1}), \text{ where } \mathbf{Q} = \sum_{t=1}^T h_t^{-1} \mathbf{x}_t \mathbf{x}_t' \text{ and } \mathbf{b} = \mathbf{Q}^{-1} \sum_{t=1}^T h_t^{-1} \mathbf{x}_t r_t,$$

for β and

$$S^2 / \sigma_v^2 \sim \chi^2(T - 2), \quad (\alpha, \delta) | \sigma_v \sim N(\mathbf{c}, \sigma_v^2 \mathbf{P}^{-1}), \text{ where}$$

$$\mathbf{P} = \begin{bmatrix} \mathbf{T} & \sum_{t=1}^T \log h_{t-1} \\ \sum_{t=1}^T \log h_{t-1} & \sum_{t=1}^T \log^2 h_{t-1} \end{bmatrix}, \quad \mathbf{c} = \mathbf{P}^{-1} \begin{bmatrix} \sum_{t=1}^T \log h_{t-1} \\ \sum_{t=1}^T \log h_t \log h_{t-1} \end{bmatrix},$$

$$\text{and } S^2 = \sum_{t=1}^T (\log h_t - c_1 - c_2 \log h_{t-1}),$$

for $(\alpha, \delta, \sigma_v)$.

7.2 Integration and optimization

The solution of all but the simplest dynamic optimization problems cannot be expressed in closed form. Since the objective function in these problems is expected utility, integration is required to evaluate a candidate solution. Finding a good numerical approximation to the solution therefore requires optimizations of a function which can be evaluated only inexactly. Moreover this evaluation must in general be repeated many times

in the process of approximating the solution. Several approaches to this important problem have been proposed: a good introduction is provided by Taylor and Uhlig (1990) and the papers following that article; more recent work includes McGrattan (1993). Here we discuss a widely applicable procedure that uses Monte Carlo integration to solve dynamic optimization problems subject to an imposed parameterization of the decision rule and then loosens the parametric restrictions so as to approach the optimum. The description here closely follows Smith (1991) who invented the method. The notation and assumptions are largely those of Stokey and Lucas (1989, Chapter 9).

The problem. Many dynamic optimization problems can be expressed

$$\max_{\{\mathbf{x}_t\}_{t=1}^{\infty}} E_0 \sum_{t=0}^{\infty} \beta^t r(\mathbf{x}_t, \mathbf{x}_{t+1}, \mathbf{z}_t) \quad (7.2.1)$$

given $\mathbf{x}_0, \mathbf{z}_0$ and subject to $\mathbf{x}_{t+1} \in \Gamma(\mathbf{x}_t, \mathbf{z}_t) \forall t$.

The sequence of state vectors $\{\mathbf{z}_t\}_{t=1}^{\infty}$ is a Markov process with transition density

$$v(\mathbf{z}_{t+1} | \mathbf{z}_t); \mathbf{z}_t \in Z \subseteq \mathfrak{R}^1 \forall t; \quad (7.2.2)$$

and Z is either compact or countable. The decision vector $\mathbf{x}_t \in X \subseteq \mathfrak{R}^p$; X is closed and convex. The agent observes the state vector $\mathbf{s}'_t = (\mathbf{x}'_t, \mathbf{z}'_t) \in S = X \times Z$ prior to choosing \mathbf{x}_{t+1} . The operator E_0 denotes expectations conditional on the period 0 information set \mathbf{s}_0 . The return function r is bounded, continuous in $(\mathbf{x}_t, \mathbf{x}_{t+1}, \mathbf{z}_t)$, and concave in $(\mathbf{x}_t, \mathbf{x}_{t+1}) \forall \mathbf{z}_t \in Z$. The correspondence Γ is nonempty, compact- and convex-valued, and continuous. The convexity of Γ precludes problems with discrete choice sets; for a treatment of discrete choice similar to the one here for continuous choice, see Geweke, Slonim, and Zarkin (1992).

These assumptions imply the existence of a unique, time-invariant continuous decision rule $w: S \rightarrow X$ that expresses optimal $\mathbf{x}_{t+1} = w(\mathbf{x}_t, \mathbf{z}_t)$ (Stokey and Lucas, 1989, Chapter 9).

The optimization problem is to determine the decision rule. The approach taken here is to replace w with a rule of thumb characterized by a vector of parameters $\boldsymbol{\psi}: \mathfrak{R}^k$:

$$\mathbf{x}_{t+1} = h(\mathbf{x}_t, \mathbf{z}_t; \boldsymbol{\psi}), \boldsymbol{\psi} \in C \subseteq \mathfrak{R}^k, C \text{ compact.} \quad (7.2.3)$$

This rule closes the model. Given $\mathbf{s}_0, \mathbf{z} = \{\mathbf{z}_t\}_{t=1}^T$, and $\boldsymbol{\psi}$, (7.2.2)-(7.2.3) determines $\mathbf{x} = \{\mathbf{x}_t\}_{t=1}^{T+1} = q(\mathbf{z}; \boldsymbol{\psi}, \mathbf{s}_0)$ through the obvious iterations.

Let $b(\mathbf{x}, \mathbf{z}; \mathbf{s}_0) = \sum_{t=0}^T \beta^t r(\mathbf{x}_t, \mathbf{x}_{t+1}, \mathbf{z}_t)$ denote the utility delivered by the sequences \mathbf{x} and \mathbf{z} given \mathbf{s}_0 for the dynamic optimization problem with horizon truncated at T . Repressing \mathbf{s}_0 to maintain notational simplicity, $g(\mathbf{z}, \boldsymbol{\psi}) = b[q(\mathbf{z}; \boldsymbol{\psi}, \mathbf{s}_0), \mathbf{z}; \mathbf{s}_0]$ is delivered utility for decision rule h with parameterization $\boldsymbol{\psi}$. Given h the agent chooses the best possible $\boldsymbol{\psi}$, which we shall denote

$$\psi_0 = \arg \max_{\psi} E_0[g(\mathbf{z}, \psi)]. \quad (7.2.4)$$

Problem (7.2.4) is a simplification of Problem (7.2.1), but it still cannot be solved analytically. The chief complication is the evaluation of the integral associated with E_0 in (7.2.4). The key idea in the solution described here is to simulate the behavior of \mathbf{z} for different values of ψ , thereby providing approximations to $E_0[g(\mathbf{z}, \psi)]$. As we shall see, arbitrarily good approximations to ψ_0 may be obtained in this way. By increasing T and employing a sequence of functions h that are increasingly flexible through a longer parameter vector ψ , the solution of (7.2.4) may be made to approximate that of (7.2.1) (Smith, 1991).

The algorithm. Generate n i.i.d. sequences $\tilde{\mathbf{z}}^{(i)} = \{\tilde{\mathbf{z}}_t^{(i)}\}_{t=1}^T$ according to (7.2.2), and take $\Theta = \{\tilde{\mathbf{z}}^{(i)}\}_{i=1}^n$ to be the collection of these sequences. If we let $Q_n(\Theta, \psi) = \sum_{i=1}^n g(\tilde{\mathbf{z}}^{(i)}, \psi)$, then $n^{-1} Q_n(\Theta, \psi) \xrightarrow{a.s.} E_0[g(\mathbf{z}, \psi)]$. Since the set of sequences Θ is fixed,

$$\hat{\psi}_n = \arg \max_{\psi} n^{-1} Q_n(\Theta, \psi)$$

is a well-defined, deterministic optimization problem that can be solved using standard hill climbing methods. These methods will be more efficient to the extent that $\partial r/\partial h$ and $\partial h/\partial \psi$ (better yet, $\partial^2 r/\partial h^2$ and $\partial^2 h/\partial \psi \partial \psi'$ in addition) can be evaluated analytically.

Asymptotic properties. Given four further assumptions, $\hat{\psi}_n \xrightarrow{a.s.} \psi_0$ and central limit theorems may be used to assess the accuracy of the approximation of ψ_0 by $\hat{\psi}_n$ and of $E_0[g(\mathbf{z}, \psi)]$ by $n^{-1} Q_n(\Theta, \psi)$.

- (1) $g(\mathbf{z}, \psi)$ is twice continuously differentiable in ψ for all \mathbf{z} .
- (2) The following functions are regular:
 - (a) $g(\mathbf{z}, \psi)$, $\partial g(\mathbf{z}, \psi)/\partial \psi$, $\partial^2 g(\mathbf{z}, \psi)/\partial \psi \partial \psi'$;
 - (b) $[\partial g(\mathbf{z}, \psi)/\partial \psi][\partial g(\mathbf{z}, \psi)/\partial \psi']$;
 - (c) $g^2(\mathbf{z}, \psi)$.

Regular is used in the sense of Tauchen (1985). Denoting the probability density function of \mathbf{z} by $f(\mathbf{z})$, $d(\mathbf{z}, \psi)$ is regular if

- (i) $d(\mathbf{z}, \psi)$ is measurable in $\mathbf{z} \forall \psi \in C$;
- (ii) d is separable (Huber, 1967);
- (iii) d is dominated -- i.e., $\exists b \ni \int b(\mathbf{z}) d\mathbf{z} < \infty$ and $|d(\mathbf{z}, \psi)| < b(\mathbf{z}) \forall \psi \in C$;
- (iv) $d(\mathbf{z}, \psi)$ is continuous in $\psi \forall \mathbf{z}$.

- (3) $E[g(\mathbf{z}, \psi)]$ (the existence of which is guaranteed by Assumption 2(a)) is unique ly maximized at ψ_0 , an interior point of C .
- (4) $E[\partial^2 g(\mathbf{z}, \psi_0)/\partial\psi\partial\psi']$ (the existence of which is also guaranteed by Assumption 2(a)) is nonsingular.

Given these four further assumptions, one can usefully approximate ψ_0 :

$$\hat{\psi}_n \xrightarrow{p} \psi_0, \quad n^{1/2}(\hat{\psi}_n - \psi_0) \xrightarrow{d} N(\mathbf{0}, \mathbf{V});$$

$$\mathbf{V} = \mathbf{A}^{-1} \mathbf{B} \mathbf{A}^{-1}, \quad \text{with } \mathbf{A} = E\left[\frac{\partial^2 g(\mathbf{z}, \psi_0)}{\partial\psi\partial\psi'}\right], \quad \mathbf{B} = E\left[\frac{\partial g(\mathbf{z}, \psi_0)}{\partial\psi} \frac{\partial g(\mathbf{z}, \psi_0)}{\partial\psi'}\right];$$

$$\hat{\mathbf{A}}_n = n^{-1} \frac{\partial^2 Q_n(\Theta, \hat{\psi}_n)}{\partial\psi\partial\psi'} \xrightarrow{p} \mathbf{A}, \quad \hat{\mathbf{B}}_n = n^{-1} \sum_{i=1}^n \frac{\partial g(\tilde{\mathbf{z}}^{(i)}, \hat{\psi}_n)}{\partial\psi} \frac{\partial g(\tilde{\mathbf{z}}^{(i)}, \hat{\psi}_n)}{\partial\psi'} \xrightarrow{p} \mathbf{B}.$$

Under exactly the same conditions, one can also usefully approximate $E[g(\mathbf{z}, \psi)]$:

$$n^{-1} Q_n(\Theta, \hat{\psi}_n) \xrightarrow{a.s.} E[g(\mathbf{z}, \psi)], \quad n^{1/2}\{n^{-1} Q_n(\Theta, \hat{\psi}_n) - E[g(\mathbf{z}, \psi)]\} \xrightarrow{d} N(0, \sigma^2);$$

$$\hat{\sigma}_n^2 = n^{-1} \sum_{i=1}^n g^2(\tilde{\mathbf{z}}^{(i)}, \hat{\psi}_n) - [n^{-1} Q_n(\Theta, \hat{\psi}_n)]^2 \xrightarrow{p} \sigma^2 = \text{var}[g(\mathbf{z}, \psi)].$$

Proofs are given by Smith (1991) who uses asymptotic theory developed by Amemiya (1985) and Tauchen (1985). The second result is especially useful in valuing the approximation error: see Smith (1991, Section 5).

Antithetic variables. In many applications the conditional distribution of the exogenous state vector \mathbf{z}_t , with probability density function $v(\mathbf{z}_t|\mathbf{z}_{t-1})$, is smooth and symmetric or nearly symmetric. The return function r is commonly monotone increasing or decreasing in each element of \mathbf{z}_t and may be nearly linear over most of the support of the distribution of \mathbf{z}_t . In such circumstances there are substantial gains in the use of antithetic variables as described in Section 5.1. Let $\tilde{\mathbf{z}}^{(i1)}$ and $\tilde{\mathbf{z}}^{(i2)}$ denote such an antithetic pair. (Exactly how the pair is drawn will depend on the particulars of the problem. What is essential, as discussed in Section 5.1, is that $\tilde{\mathbf{z}}^{(i1)}$ and $\tilde{\mathbf{z}}^{(i2)}$ be identically distributed.) Consider $n/2$ replications of $\tilde{\mathbf{z}}^{(i1)}$ and $\tilde{\mathbf{z}}^{(i2)}$ in lieu of n replications of $\tilde{\mathbf{z}}^{(i)}$. Redefine

$$Q_n(\Theta, \psi) = \sum_{i=1}^{n/2} [g(\tilde{\mathbf{z}}^{(i1)}, \psi) + g(\tilde{\mathbf{z}}^{(i2)}, \psi)]$$

with $\Theta = \{\tilde{\mathbf{z}}^{(i1)}, \tilde{\mathbf{z}}^{(i2)}\}_{i=1}^{n/2}$ and take $\hat{\psi}_n = \arg \max_{\psi} n^{-1} Q_n(\Theta, \psi)$. Then $\hat{\psi}_n$ and $n^{-1} Q_n(\theta, \psi)$ are consistent for ψ and $E[g(\mathbf{z}, \psi)]$ as before. There are again central limit theorems, but now

$$\mathbf{V} = \mathbf{A}^{-1} \mathbf{B}^* \mathbf{A}^{-1}, \quad \text{with } \mathbf{B}^* = \mathbf{B} + \frac{1}{2}(\mathbf{C} + \mathbf{C}'), \quad \text{where}$$

$$\mathbf{C} = E\left[\frac{\partial g(\tilde{\mathbf{z}}^{(i1)}, \psi_0)}{\partial\psi} \frac{\partial g(\tilde{\mathbf{z}}^{(i2)}, \psi_0)}{\partial\psi'}\right], \quad \hat{\mathbf{C}}_n = \left(\frac{n}{2}\right)^{-1} \sum_{i=1}^{n/2} \frac{\partial g(\tilde{\mathbf{z}}^{(i1)}, \psi_0)}{\partial\psi} \frac{\partial g(\tilde{\mathbf{z}}^{(i2)}, \psi_0)}{\partial\psi'} \xrightarrow{p} \mathbf{C}$$

and

$$\begin{aligned} \sigma^2 &= \text{var}[g(\mathbf{z}, \psi_0)] + \text{cov}[g(\tilde{\mathbf{z}}^{(i1)}, \psi_0), g(\tilde{\mathbf{z}}^{(i2)}, \psi_0)], \\ n^{-1} \sum_{i=1}^{n/2} [g^2(\tilde{\mathbf{z}}^{(i1)}, \hat{\psi}_n) + g^2(\tilde{\mathbf{z}}^{(i2)}, \hat{\psi}_n)] &+ \left(\frac{n}{2}\right)^{-1} \sum_{i=1}^n g(\tilde{\mathbf{z}}^{(i1)}, \hat{\psi}_n) g(\tilde{\mathbf{z}}^{(i2)}, \hat{\psi}_n) \\ &- 2[n^{-1} Q_n(\Theta, \hat{\psi}_n)]^2 \xrightarrow{p} \sigma^2. \end{aligned}$$

Smith (1991) applies this method to a variant of the Brock and Mirman (1972) growth model. The characteristic of the model that is important for the success of the use of antithetic variables is that the exogenous state variables move smoothly over time and the return function is only modestly nonlinear over most of the support of \mathbf{z} . Using only 100 antithetic pairs and $T = 800$, Smith determines ψ up to four significant figures. The suboptimality of the resulting decision rules turns out to be equivalent to a per-period decrease in consumption of $2 \times 10^{-5}\%$.

References

- Ahrens, J.H., and U. Dieter, 1974, "Computer Methods for Sampling from Gamma, Beta, Poisson, and Binomial Distributions," *Computing* 12:223-246.
- Ahrens, J.H., and U. Dieter, 1980, "Sampling from Binomial and Poisson Distributions: A Method with Bounded Computation Times," *Computing* 25: 193-208.
- Amemiya, T., 1985, *Advanced Econometrics*. Cambridge: Harvard University Press.
- Anderson, T.W., 1984, *An Introduction to Multivariate Statistical Analysis*. New York: Wiley. (Second edition)
- Bollerslev, T., R. Chou, and K.F. Kroner, 1992, "ARCH Modelling in Finance," *Journal of Econometrics* 52: 5-59.
- Box, G.E.P., and M.E. Muller, 1958, "A Note on the Generation of Random Normal Deviates," *Annals of Mathematical Statistics* 29: 610-611.
- Bratley, P., B.L. Fox, and L.E. Schrage, 1987, *A Guide to Simulation*. New York: Springer-Verlag. (Second edition)
- Brock, W.A., and L.J. Mirman, 1972, "Optimal Economic Growth and Uncertainty: the Discounted Case," *Journal of Economic Theory* 4: 497-513.
- Chib, S., and E. Greenberg, 1994, "Understanding the Metropolis-Hastings Algorithm," Washington University John M. Olin School of Business working paper.
- Coveyou, R.R., and R.D. MacPherson, 1967, "Fourier Analysis of Uniform Random Number Generators," *Journal of the ACM* 14: 100-119.
- Davis, P.J., and P. Rabinowitz, 1984, *Methods of Numerical Integration*. Orlando: Academic Press. (Second edition)
- Devroye, L., 1986, *Non-uniform Random Variate Generation*. New York: Springer-Verlag.
- Evans, M., 1991, "Adaptive Importance Sampling and Chaining," *Contemporary Mathematics* 115 (Statistical Multiple Integration): 137-142. Providence: American Mathematical Society.
- Fishman, G.F., and L.R. Moore, III, 1982, "A Statistical Evaluation of Multiplicative Random Number Generators with Modulus $2^{31}-1$," *Journal of the American Statistical Association* 77: 129-136.
- Fishman, G.F., and L.R. Moore, III, 1986, "An Exhaustive Analysis of Multiplicative Congruential Random Number Generators with Modulus $2^{31}-1$," *SIAM Journal on Scientific and Statistical Computing* 7: 24-45.
- Forsythe, G.E., 1972, "Von Neumann's Comparison Method for Random Sampling from the Normal and Other Distributions," *Mathematical Computation* 26: 817-826.
- Gelfand, A.E., and A.F.M. Smith, 1990, "Sampling Based Approaches to Calculating Marginal Densities," *Journal of the American Statistical Association* 85: 398-409.

- Gelman, A., and D.B. Rubin, 1992, "Inference from Iterative Simulation Using Multiple Sequences," *Statistical Science* 7: 457-472.
- 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.
- Genz, A., 1991, "Subregion Adaptive Algorithms for Multiple Integrals," in N. Flournoy and R.K. Tsutakawa (eds.), *Contemporary Mathematics* 115 (Statistical Multiple Integration): 23-31. Providence: American Mathematical Society.
- Genz, A., 1993, "Subregion Adaptive Integration of Functions Having a Dominant Peak," University of Washington working paper.
- Genz, A., and A. Malik, 1980, "An Adaptive Algorithm for Numerical Integration over an N-Dimensional Rectangular Region," *Journal of Computational and Applied Mathematics* 6: 295-302.
- Genz, A., and A.A. Malik, 1983, "An Imbedded Family of Fully Symmetric Numerical Integration Rules," *SIAM Journal of Numerical Analysis* 20: 580-588.
- Geweke, J., 1986, "Exact Inference in the Inequality Constrained Normal Linear Regression Model," *Journal of Applied Econometrics* 1: 127-141.
- Geweke, J., 1988, "Antithetic Acceleration of Monte Carlo Integration in Bayesian Inference," *Journal of Econometrics* 38: 73-89.
- Geweke, J., 1989, "Bayesian Inference in Econometric Models Using Monte Carlo Integration," *Econometrica* 57: 1317-1340.
- Geweke, J., 1991, "Efficient Simulation from the Multivariate Normal and Student- t Distributions Subject to Linear Constraints," in E.M. Keramidas (ed.), *Computing Science and Statistics: Proceedings of the 23rd Symposium on the Interface*, 571-578.
- Geweke, J., 1992a, "Evaluating the Accuracy of Sampling-Based Approaches to the Calculation of Posterior Moments," in J.M. Bernardo *et al.* (eds.), *Bayesian Statistics 4: Proceedings of the Fourth Valencia International Meeting*. Oxford: Clarendon Press.
- Geweke, J., 1992b, "Priors for Macroeconomic Time Series," Federal Reserve Bank of Minneapolis Institute for Empirical Macroeconomics Discussion Paper No. 64.
- Geweke, J., R. Slonim, and G. Zarkin, 1992, "Econometric Solution Methods for Dynamic Discrete Choice Problems," University of Minnesota Department of Economics working paper.
- Geyer, C.J., 1992, "Practical Markov Chain Monte Carlo," *Statistical Science* 7: 473-481.
- Gilks, W.R., and P. Wild, 1992, "Adaptive Rejection Sampling for Gibbs Sampling," *Applied Statistics (JRSS Series C)* 41: 337-348.
- Golub, G.H., and J.H. Welsch, 1969, "Calculation of Gaussian Quadrature Rules," *Mathematics of Computation* 23: 221-230.

- Gradshteyn, I.S., and I.M. Ryzhik, 1965, *Tables of Integrals, Series, and Products*. New York: Academic Press, 1965.
- Greenberger, M., 1961, "Notes on a New Pseudo-Random Number Generator," *Journal of the ACM* 8: 163-167.
- Halton, J.M., 1960, "On the Efficiency of Evaluating Certain Quasi-random Sequences of Points in Evaluating Multi-dimensional Integrals," *Numerische Mathematik* 2: 84-90.
- Hammersley, J.M., 1960, "Monte Carlo Methods for Solving Multivariable Problems," *Annals of the New York Academy of Sciences* 86: 844-874.
- Hammersly, J.M., and D.C. Handscomb, 1964, *Monte Carlo Methods*. London: Methuen and Company.
- Hammersly, J.M., and K.W. Morton, 1956, "A New Monte Carlo Technique: Antithetic Variates," *Proceedings of the Cambridge Philosophical Society* 52: 449-474.
- Hannan, E.J., 1970, *Multiple Time Series*. New York: Wiley.
- Hart, H.F., E.W. Cheney, C.L. Lawson, H.J. Maehly, C.K. Mesztenyi, J.R. Rice, H.G. Thacher, Jr., and C. Witzgall, 1968, *Computer Approximations*. New York: Wiley.
- Hastings, W.K., 1970, "Monte Carlo Sampling Methods Using Markov Chains and Their Applications," *Biometrika* 57: 97-109.
- Hlawka, E., 1961, "Funktionen von Beschränkter Variation in der Theorie der Gleichverteilung," *Annali di Matematica Pura Ed Applicata* 54: 325-333.
- Huber, P.J., 1967, "The Behavior of Maximum Likelihood Estimates under Nonstandard Conditions," in L.M. LeCam and J. Neyman (eds.), *Proceedings of the Fifth Berkeley Symposium on Mathematical Statistics and Probability*, Vol. 1, 221-234. Berkeley: University of California Press.
- IMSL, 1994, *IMSL Stat/Library*. Houston: Visual Numerics, Inc.
- Jacquier, E., N.G. Polson, and P.E. Rossi, 1994, "Bayesian Analysis of Stochastic Volatility Models," *Journal of Business and Economic Statistics*, forthcoming.
- Judd, K.L., 1991, "Numerical Methods in Economics," Hoover Institution Stanford University manuscript.
- Kachitvichyanukul, V., 1982, "Computer Generation of Poisson, Binomial, and Hypergeometric Random Variates." Unpublished Ph.D. dissertation, Purdue University.
- Kahn, M., and A.W. Marshall, 1953, "Methods of Reducing Sample Size in Monte Carlo Computations," *Operations Research* 1: 263-278.
- Kinderman, A.J., and J.G. Ramage, 1976, "Computer Generation of Normal Random Variables," *Journal of the American Statistical Association* 71: 893-896.

- Kipnis, C., and S.R.S. Varadhan, 1986, "Central Limit Theorem for Additive Functionals of Reversible Markov Processes and Applications to Simple Exclusions," *Communications in Mathematical Physics* **104**: 1-19.
- Kloek, T., and H.K. van Dijk, 1978, "Bayesian Estimates of Equation System Parameters: An Application of Integration by Monte Carlo," *Econometrica* **46**: 1-20.
- Knuth, D.E., 1981, *The Art of Computer Programming, Volume 2: Seminumerical Algorithms*. Reading: Addison-Wesley. (Second edition)
- Kronmal, R.A., and A.V. Peterson, 1979, "On the Alias Method for Generating Random Variables from a Discrete Distribution," *American Statistician* **33**: 214-218.
- L'Ecuyer, P., 1986, "Efficient and Portable Combined Pseudorandom Number Generators," *Communications of the ACM* **29**: 304-313.
- Marsaglia, G., 1961, "Expressing a Random Variable in Terms of Uniform Random Variables," *Annals of Mathematical Statistics* **32**: 894-899.
- Marsaglia, G., 1964, "Generating a Variable from the Tail of a Normal Distribution," *Technometrics* **6**: 101-102.
- Marsaglia, G., 1968, "Random Numbers Fall Mainly in the Planes," *Proceedings of the National Academy of Sciences* **60**: 25-28.
- Marsaglia, G., 1972, "The Structure of Linear Congruential Sequences," in S.K. Zarema (ed.), *Applications of Number Theory to Numerical Analysis*. New York: Academic Press.
- Marsaglia, G., and T.A. Bray, 1964, "A Convenient Method for Generating Normal Variables," *SIAM Review* **6**: 260-264.
- Marsaglia, G., and T.A. Bray, 1968, "On-line Random Number Generators and Their Use in Combinations," *Communications of the ACM* **11**: 757-759.
- Marsaglia, G., M.D. MacLaren, and T.A. Bray, 1964, "A Fast Procedure for Generating Normal Random Variables," *Communications of the ACM* **7**: 4-10.
- Marsaglia, G. and A. Zaman, 1991, "A New Class of Random Number Generators," *The Annals of Applied Probability* **1**: 462-480.
- McGrath, E.I., 1970, *Fundamentals of Operations Research*. San Francisco: West Coast University Press.
- McGrattan, E., 1993, "Solving the Stochastic Growth Model with a Finite Element Method," Federal Reserve Bank of Minneapolis working paper.
- McNamee, J., and F. Stenger, 1967, "Construction of Fully Symmetric Numerical Integration Formulas," *Numerical Mathematics* **10**: 327-344.
- Metropolis, N., A.W. Rosenbluth, M.N. Rosenbluth, A.H. Teller, and E. Teller, 1953, "Equation of State Calculations by Fast Computing Machines," *The Journal of Chemical Physics* **21**: 1087-1092.

- Mikhail, W.M., 1972, "Simulating the Small-sample Properties of Econometric Estimators," *Journal of the American Statistical Association* **67**: 620-624.
- Mitchell, B., 1973, "Variance Reduction by Antithetic Variates in G1/G1 Queueing Simulation," *Operations Research* **21**: 988-997.
- Müller, P., 1991, "Numerical Integration in Bayesian Analysis," Purdue University unpublished Ph.D. dissertation.
- Niederreiter, H., 1992, *Random Number Generation and Quasi-Monte Carlo Methods*. Philadelphia: Society for Industrial and Applied Mathematics.
- Numelin, E., 1984, *General Irreducible Markov Chains and Non-negative Operators*. Cambridge: Cambridge University Press.
- Peskun, P.H., 1973, "Optimum Monte-Carlo Sampling using Markov Chains," *Biometrika* **60**: 607-612
- Piessens, R., E. deDoncker-Kapenga, C.W. Überhuber, and D.K. Kahander, 1983, *QUADPACK*. New York: Springer-Verlag.
- Poirier, D., 1995, *Intermediate Statistics and Econometrics: A Comparative Approach*. Cambridge: MIT Press.
- Press, W.H., B.P. Flannery, S.A. Teukolsky, and W.T. Vetterling, 1986, *Numerical Recipes: The Art of Scientific Computing*. Cambridge: Cambridge University Press.
- Richtmeyer, R.D., 1952, "On the Evaluation of Definite Integrals and a quasi-Monte Carlo Method Based on Properties of Algebraic Numbers," Report LA-1342. Los Alamos: Los Alamos Scientific Laboratories.
- Richtmeyer, R.D., 1958, "A Non-random Sampling Method, Based on Congruences for Monte Carlo Problems, Report NYO-8674. New York: Institute of Mathematical Sciences, New York University.
- Ripley, R.D., 1987, *Stochastic Simulation*. New York: Wiley.
- Roberts, G.O., and A.F.M. Smith, 1992, "Simple Conditions for the Convergence of the Gibbs Sampler and Metropolis-Hastings Algorithms," University of Cambridge Statistical Laboratory Research Report No. 92-30.
- Rubinstein, R.Y., 1981, *Simulation and the Monte Carlo Method*. New York: Wiley.
- Schmeiser, B.W., and R. Lal, 1980, "Squeeze Methods for Generating Gamma Variates," *Journal of the American Statistical Association* **75**: 679-682.
- Smith, A.A., 1991, "Solving Stochastic Dynamic Programming Problems using Rules of Thumb," Queen's University Department of Economics Discussion Paper No. 816.
- Stokey, N.L., and R.E. Lucas, Jr., 1989, *Recursive Methods in Economic Dynamics*. Cambridge: Harvard University Press.
- Strecok, A.J., 1968, "On the Calculation of the Inverse of the Error Function," *Mathematics of Computation* **22**: 144-158.

- Tanner, M.A., and W.H. Wong, 1987, "The Calculation of Posterior Distributions by Data Augmentation," *Journal of the American Statistical Association* **82**: 528-550.
- Tauchen, G., 1985, "Diagnostic Testing and Evaluation of Maximum Likelihood Models," *Journal of Econometrics* **30**: 415-443.
- Taylor, S., 1986, *Modelling Financial Time Series*. New York: Wiley.
- Taylor, J., and H. Uhlig, 1990, "Solving Nonlinear Stochastic Growth Models: A Comparison of Alternative Solution Methods," *Journal of Business and Economic Statistics* **8**: 1-17.
- Tezuka, S., P. L'Ecuyer, and R. Couture, 1993, "On the Lattice Structure of the Add-with-carry and Subtract-with-borrow Random Number Generators," *ACM Transactions on Modeling and Computer Simulation* **3**: 315-331.
- Tierney, L., 1991a, "Exploring Posterior Distributions Using Markov Chains," in E.M. Keramidas (ed.), *Computing Science and Statistics: Proceedings of the 23rd Symposium on the Interface*, 563-570. Fairfax: Interface Foundation of North America, Inc.
- Tierney, L., 1991b, "Markov Chains for Exploring Posterior Distributions," Technical Report No. 560, University of Minnesota School of Statistics. Forthcoming, *Annals of Statistics*.
- Tierney, L., and J.B. Kadane, 1986, "Accurate Approximations for Posterior Moments and Marginal Densities," *Journal of the American Statistical Association* **81**: 82-86.
- Tierney, L., R.E. Kass, and J.B. Kadane, 1989, "Fully Exponential Laplace Approximations to Expectations and Variances of Nonpositive Functions," *Journal of the American Statistical Association* **84**: 710-716.
- van Dooren, P., and L. de Ridder, 1976, "An Adaptive Algorithm for Numerical Integration over an N-Dimensional Rectangular Region," *Journal of Computational and Applied Mathematics* **2**: 207-217.
- von Neumann, J., 1951, "Various Techniques Used in Connection with Random Digits," *National Bureau of Standards Applied Mathematics*, Series 12, pp. 36-38.
- Walker, A.J., 1974, "New Fast Method for Generating Discrete Random Numbers with Arbitrary Frequency Distributions," *Electronics Letters* **10**: 127-128.
- Walker, A.J., 1977, "An Efficient Method for Generating Discrete Random Variables with General Distributions," *ACM Transactions on Mathematical Software* **3**: 253-256.
- Wichmann, B.A., and I.D. Hill, 1982, "An Efficient and Portable Pseudo-random Number Generator," *Applied Statistics* **31**: 188-190.
- Wild, P., and W.R. Gilks, 1993, "Adaptive Rejection Sampling from Log-concave Density Functions," *Applied Statistics (JRSS Series C)* **42**: 701-708.
- Zellner, A., and C. Min, 1992, "Gibbs Sampler Convergence Criteria," University of Chicago Graduate School of Business working paper.

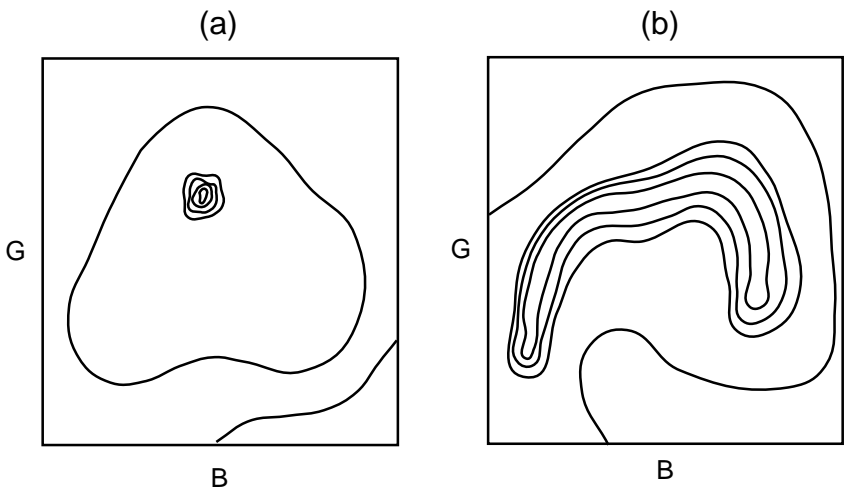


Figure 1. Contours of the function to be integrated are shown.

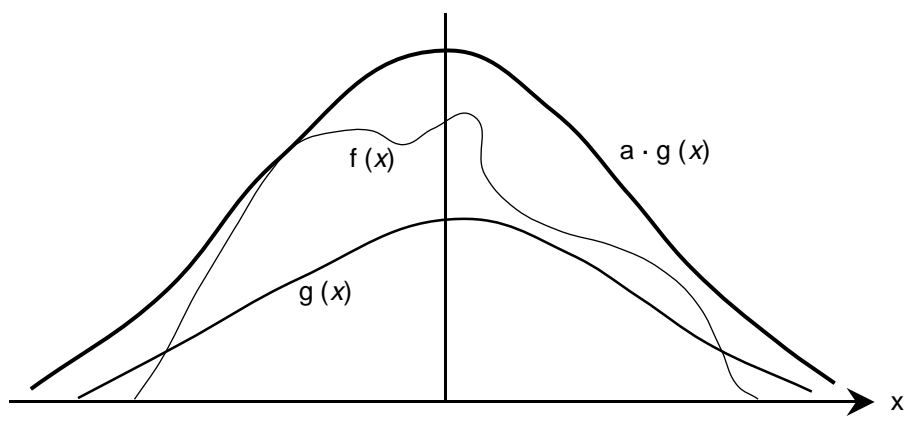


Figure 2. The target density is $f(x)$, the source density is $g(x)$, and $a = \sup[f(x)/g(x)]$.

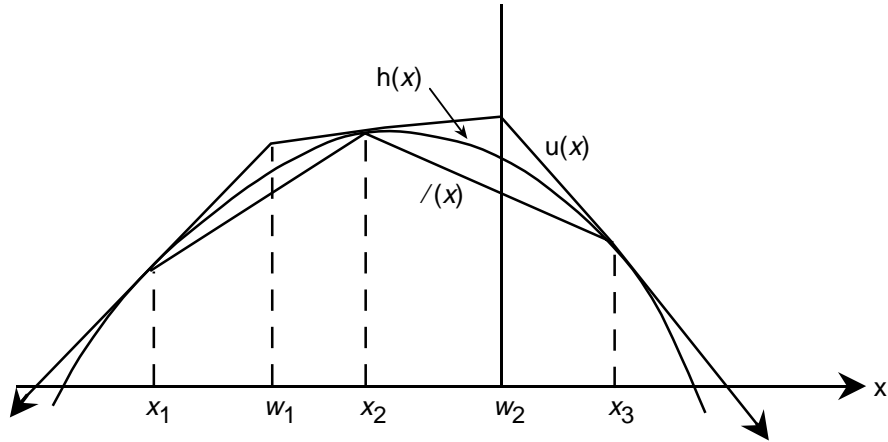


Figure 3. The function $h(x) = \log f(x)$, where $f(x)$ is a log-concave p.d.f. The lower hull $l(x)$ is formed by the chords joined at the x_j , and the upper hull $u(x)$ is formed by the tangents at the x_j which are joined at the w_j .

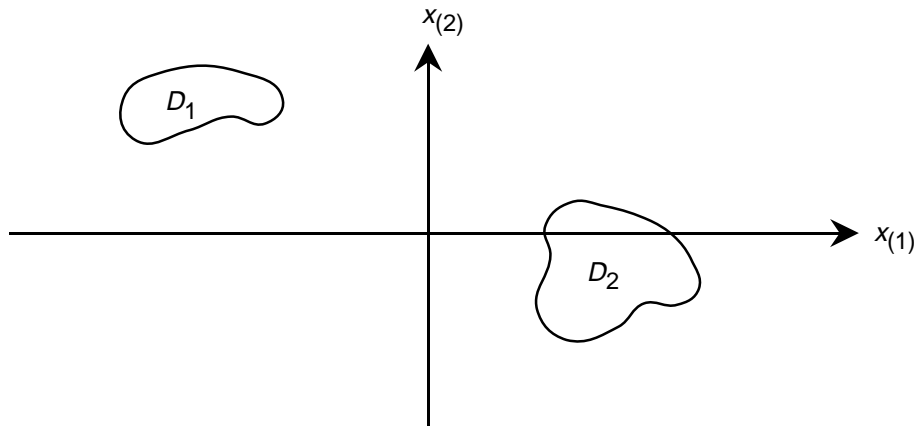


Figure 4. The disconnected support $D = D_1 \cup D_2$ for the probability distribution implies that a Gibbs sampler with blocking $(x_{(1)}, x_{(2)})$ will not be Harris recurrent. In the example shown it cannot converge from any starting value.

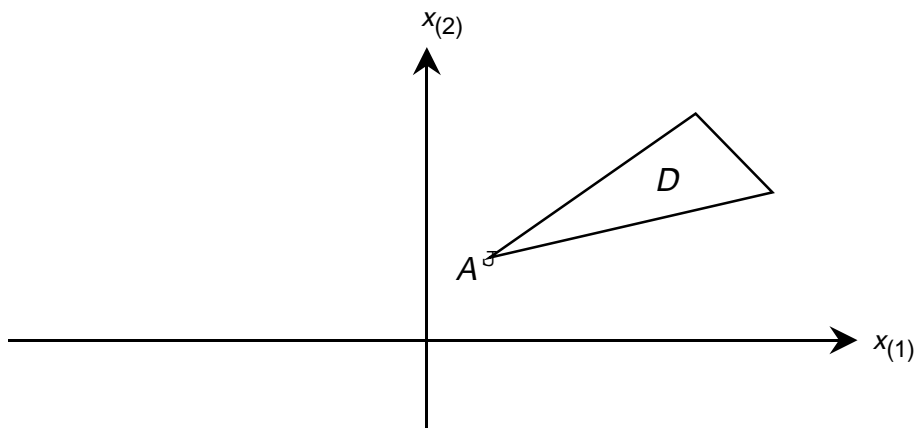


Figure 5. The probability density $p(\mathbf{x})$ is uniform on the closed set D and consequently is not lower semicontinuous at 0. The point A is absorbing for the Gibbs sampler with blocking $(x_{(1)}, x_{(2)})$, so if $\mathbf{x}^0 = A$ convergence will not occur.

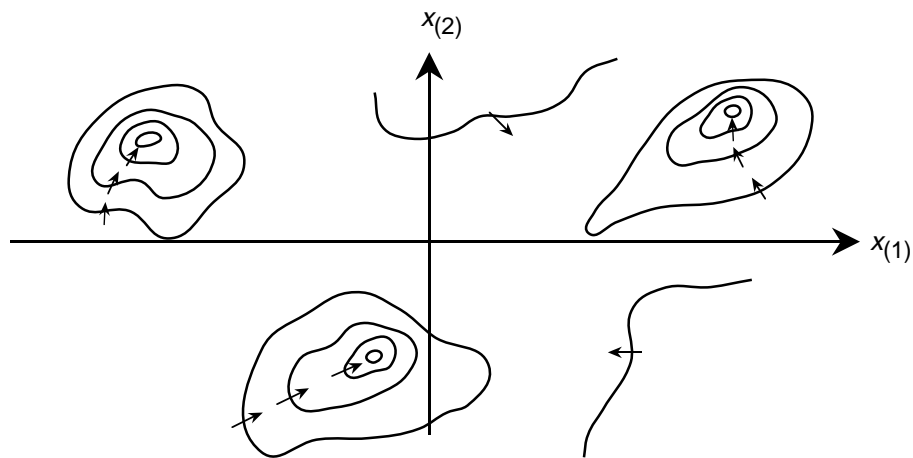


Figure 6. Iso-probability density contours of a multimodal bivariate distribution are shown. (Arrows indicate directions of increased density.) Given sufficiently steep gradients the Gibbs sampler will converge very slowly.

Table 1

Evaluations required to approximate $\int_{\mathbb{I}^d} f(\mathbf{x}) d\mathbf{x}$, $f(\mathbf{x}) = \sum_{j=1}^d f(x_j)$,
 with maximum error c : Actual number and upper bound

d	2	3	4	5
$c = .01$:				
Actual	228	442	661	1060
Bound	19,335	1,014,825	9.154×10^7	1.522×10^{10}
$c = 10^{-5}$:				
Actual	640,426	1,039,188	1,523,433	2,379,162
Bound	52,477,915	3.469×10^9	3.513×10^{11}	6.114×10^{13}

Table 2

Error comparison for Halton sequence and independence Monte Carlo

$$\int_{\mathbb{T}^d} f(\mathbf{x}) d\mathbf{x}, f(\mathbf{x}) = \sum_{i=1}^d x_i$$

 $m = 1,000$

d	Halton error	Halton bound	MC error($p = .05$)	MC error($p = 10^{-12}$)
5	-7.526×10^{-3}	9.302×10^2	.04000	.1455
10	-.02807	6.053×10^{19}	.05658	.2058
20	-.1097	2.616×10^{29}	.08002	.2911
40	-.3824	8.225×10^{72}	.1132	.4117
60	-.8202	2.467×10^{121}	.1386	.5042
80	-1.476	1.250×10^{173}	.1600	.5822
100	-2.062	1.447×10^{227}	.1789	.6509

 $m = 50,000$

d	Halton error	Halton bound	MC error($p = .05$)	MC error($p = 10^{-12}$)
5	-2.786×10^{-4}	1.071×10^2	5.658×10^{-3}	.02058
10	-8.861×10^{-4}	3.533×10^{10}	8.002×10^{-3}	.02911
20	-3.537×10^{-3}	3.225×10^{30}	.01132	.04117
40	-.02216	3.356×10^{75}	.01600	.05822
60	-.02768	2.2990×10^{127}	.01960	.07131
80	-.05681	2.4186×10^{181}	.02263	.08234
100	-.08779	5.235×10^{237}	.02530	.09205

Table 3

Error comparison for Halton sequence and Monte Carlo

$$f(\mathbf{x}) = \sum_{i=1}^d x_i, \quad \mathbf{x} \sim N(\mathbf{0}, \mathbf{I}_d); \text{ evaluate } E[f(\mathbf{x})]$$

d	Halton error	$m = 1,000$		$m = 50,000$		
		Monte Carlo error		Monte Carlo error		
		($p = .05$)	($p = 10^{-12}$)	($p = .05$)	($p = 10^{-12}$)	
5	-.04190	.1386	.5042	-1.808×10^{-3}	.01960	.07131
10	-.1411	.1960	.7131	-5.552×10^{-3}	.02772	.1008
20	-.5497	.2772	1.008	-.02076	.03920	.1426
40	-1.7306	.3920	1.426	-.06548	.05544	.2017
60	-3.3617	.4801	1.747	-.1461	.06790	.2470
80	-5.6578	.5544	2.017	-.2573	.07840	.2852
100	-7.8073	.6198	2.255	-.2336	.08765	.3189

$$f(\mathbf{x}) = \sum_{i=1}^d x_i^2, \quad \mathbf{x} \sim N(\mathbf{0}, \mathbf{I}_d); \text{ evaluate } E[f(\mathbf{x})]$$

d	Halton error	$m = 1,000$		$m = 50,000$		
		Monte Carlo error		Monte Carlo error		
		($p = .05$)	($p = 10^{-12}$)	($p = .05$)	($p = 10^{-12}$)	
5	-.0496	.2400	.8733	-1.664×10^{-3}	.03395	.1235
10	-.0941	.3395	1.2350	-2.418×10^{-3}	.04801	.1747
20	-.0864	.4801	1.746	-4.611×10^{-3}	.06790	.2470
40	.2436	.6790	2.470	-6.367×10^{-3}	.0962	.3493
60	.5680	.8316	3.0252	-3.662×10^{-3}	.1176	.4278
80	.4982	.9602	3.4932	.0243	.1358	.4940
100	1.449	1.074	3.906	-.04932	.1518	.5523

Table 3 (continued)

$$f(\mathbf{x}) = \sum_{i=1}^d x_i^3, \quad \mathbf{x} \sim N(\mathbf{0}, \mathbf{I}_d); \text{ evaluate } E[f(\mathbf{x})]$$

		$m = 1,000$		$m = 50,000$		
		Monte Carlo error		Monte Carlo error		
d	Halton error	($p = .05$)	($p = 10^{-12}$)	Halton error	($p = .05$)	($p = 10^{-12}$)
5	-.3500	.5368	1.953	-.02286	.07591	.2761
10	-1.083	.7591	2.762	-.06800	.1073	.3906
20	-4.072	1.074	3.906	-.2386	.1518	.5523
40	-11.865	1.518	5.523	-.6821	.2174	.7811
60	-19.564	1.859	6.765	-1.411	.2630	.9567
80	-27.78	2.147	7.811	-2.641	.3036	1.104
100	-36.18	2.400	8.733	-2.218	.3395	1.235

$$f(\mathbf{x}) = \sum_{i=1}^d x_i^4, \quad \mathbf{x} \sim N(\mathbf{0}, \mathbf{I}_d); \text{ evaluate } E[f(\mathbf{x})]$$

		$m = 1,000$		$m = 50,000$		
		Monte Carlo error		Monte Carlo error		
d	Halton error	($p = .05$)	($p = 10^{-12}$)	Halton error	($p = .05$)	($p = 10^{-12}$)
5	-.7442	1.420	5.167	-.03612	.2008	.7307
10	-1.046	2.008	7.307	-.04667	.2840	1.0333
20	-.8494	2.840	10.33	-.07076	.4017	1.461
40	7.504	4.016	14.61	.03523	.5681	2.067
60	16.88	4.919	17.90	.1150	.0957	2.521
80	23.48	5.681	20.66	-.1898	.8034	2.923
100	32.94	6.351	23.105	-.7909	.8982	3.268

Table 3 (continued)

$$f(\mathbf{x}) = \sum_{i=1}^d x_i^5, \quad \mathbf{x} \sim N(\mathbf{0}, \mathbf{I}_d); \text{ evaluate } E[f(\mathbf{x})]$$

d	Halton error	$m = 1,000$		$m = 50,000$		
		Monte Carlo error ($p = .05$)	Monte Carlo error ($p = 10^{-12}$)	Halton error	Monte Carlo error ($p = .05$)	Monte Carlo error ($p = 10^{-12}$)
5	-3.216	4.260	15.50	-.3365	.6026	2.192
10	-1.043	6.025	21.92	-1.006	.8521	3.100
20	-36.44	8.521	31.00	-3.433	1.205	4.384
40	-118.9	12.05	43.84	-9.549	1.704	6.200
60	-202.8	14.76	53.69	-14.50	2.087	7.593
80	-281.6	17.04	62.00	-13.11	2.410	8.760
100	-359.7	19.05	69.32	-23.97	2.695	9.803