

Introducing Markov chain Monte Carlo

Walter R Gilks

Sylvia Richardson

David J Spiegelhalter

1.1 Introduction

Markov chain Monte Carlo (MCMC) methodology provides enormous scope for realistic statistical modelling. Until recently, acknowledging the full complexity and structure in many applications was difficult and required the development of specific methodology and purpose-built software. The alternative was to coerce the problem into the over-simple framework of an available method. Now, MCMC methods provide a unifying framework within which many complex problems can be analysed using generic software.

MCMC is essentially Monte Carlo integration using Markov chains. Bayesians, and sometimes also frequentists, need to integrate over possibly high-dimensional probability distributions to make inference about model parameters or to make predictions. Bayesians need to integrate over the posterior distribution of model parameters given the data, and frequentists may need to integrate over the distribution of observables given parameter values. As described below, *Monte Carlo* integration draws samples from the the required distribution, and then forms sample averages to approximate expectations. *Markov chain* Monte Carlo draws these samples by running a cleverly constructed Markov chain for a long time. There are many ways of constructing these chains, but all of them, including the Gibbs sampler (Geman and Geman, 1984), are special cases of the general framework of Metropolis *et al.* (1953) and Hastings (1970).

It took nearly 40 years for MCMC to penetrate mainstream statistical practice. It originated in the statistical physics literature, and has been used for a decade in spatial statistics and image analysis. In the last few years, MCMC has had a profound effect on Bayesian statistics, and has also found applications in classical statistics. Recent research has added considerably to its breadth of application, the richness of its methodology, and its theoretical underpinnings.

The purpose of this book is to introduce MCMC methods and their applications, and to provide pointers to the literature for further details. Having in mind principally an applied readership, our role as editors has been to keep the technical content of the book to a minimum and to concentrate on methods which have been shown to help in real applications. However, some theoretical background is also provided. The applications featured in this volume draw from a wide range of statistical practice, but to some extent reflect our own biostatistical bias. The chapters have been written by researchers who have made key contributions in the recent development of MCMC methodology and its application. Regrettably, we were not able to include all leading researchers in our list of contributors, nor were we able to cover all areas of theory, methods and application in the depth they deserve.

Our aim has been to keep each chapter self-contained, including notation and references, although chapters may assume knowledge of the basics described in this chapter. This chapter contains enough information to allow the reader to start applying MCMC in a basic way. In it we describe the Metropolis–Hastings algorithm, the Gibbs sampler, and the main issues arising in implementing MCMC methods. We also give a brief introduction to Bayesian inference, since many of the following chapters assume a basic knowledge. Chapter 2 illustrates many of the main issues in a worked example. Chapters 3 and 4 give an introduction to important concepts and results in discrete and general state-space Markov chain theory. Chapters 5 through 8 give more information on techniques for implementing MCMC or improving its performance. Chapters 9 through 13 describe methods for assessing model adequacy and choosing between models, using MCMC. Chapters 14 and 15 describe MCMC methods for non-Bayesian inference, and Chapters 16 through 25 describe applications or summarize application domains.

1.2 The problem

1.2.1 Bayesian inference

Most applications of MCMC to date, including the majority of those described in the following chapters, are oriented towards Bayesian inference. From a Bayesian perspective, there is no fundamental distinction between

observables and parameters of a statistical model: all are considered random quantities. Let D denote the observed data, and θ denote model parameters and missing data. Formal inference then requires setting up a joint probability distribution $P(D, \theta)$ over all random quantities. This joint distribution comprises two parts: a *prior* distribution $P(\theta)$ and a *likelihood* $P(D|\theta)$. Specifying $P(\theta)$ and $P(D|\theta)$ gives a *full probability model*, in which

$$P(D, \theta) = P(D|\theta) P(\theta).$$

Having observed D , Bayes theorem is used to determine the distribution of θ conditional on D :

$$P(\theta|D) = \frac{P(\theta)P(D|\theta)}{\int P(\theta)P(D|\theta)d\theta}.$$

This is called the *posterior* distribution of θ , and is the object of all Bayesian inference.

Any features of the posterior distribution are legitimate for Bayesian inference: moments, quantiles, highest posterior density regions, etc. All these quantities can be expressed in terms of posterior expectations of functions of θ . The posterior expectation of a function $f(\theta)$ is

$$E[f(\theta)|D] = \frac{\int f(\theta)P(\theta)P(D|\theta)d\theta}{\int P(\theta)P(D|\theta)d\theta}.$$

The integrations in this expression have until recently been the source of most of the practical difficulties in Bayesian inference, especially in high dimensions. In most applications, analytic evaluation of $E[f(\theta)|D]$ is impossible. Alternative approaches include numerical evaluation, which is difficult and inaccurate in greater than about 20 dimensions; analytic approximation such as the Laplace approximation (Kass *et al.*, 1988), which is sometimes appropriate; and Monte Carlo integration, including MCMC.

1.2.2 Calculating expectations

The problem of calculating expectations in high-dimensional distributions also occurs in some areas of frequentist inference; see Geyer (1995) and Diebolt and Ip (1995) in this volume. To avoid an unnecessarily Bayesian flavour in the following discussion, we restate the problem in more general terms. Let X be a vector of k random variables, with distribution $\pi(\cdot)$. In Bayesian applications, X will comprise model parameters and missing data; in frequentist applications, it may comprise data or random effects. For Bayesians, $\pi(\cdot)$ will be a posterior distribution, and for frequentists it will be a likelihood. Either way, the task is to evaluate the expectation

$$E[f(X)] = \frac{\int f(x)\pi(x)dx}{\int \pi(x)dx} \quad (1.1)$$

for some function of interest $f(\cdot)$. Here we allow for the possibility that the distribution of X is known only up to a constant of normalization. That is, $\int \pi(x)dx$ is unknown. This is a common situation in practice, for example in Bayesian inference we know $P(\theta|D) \propto P(\theta)P(D|\theta)$, but we cannot easily evaluate the normalization constant $\int P(\theta)P(D|\theta)d\theta$. For simplicity, we assume that X takes values in k -dimensional Euclidean space, i.e. that X comprises k continuous random variables. However, the methods described here are quite general. For example, X could consist of discrete random variables, so then the integrals in (1.1) would be replaced by summations. Alternatively, X could be a mixture of discrete and continuous random variables, or indeed a collection of random variables on any probability space. Indeed, k can itself be variable: see Section 1.3.3. Measure theoretic notation in (1.1) would of course concisely accommodate all these possibilities, but the essential message can be expressed without it. We use the terms *distribution* and *density* interchangeably.

1.3 Markov chain Monte Carlo

In this section, we introduce MCMC as a method for evaluating expressions of the form of (1.1). We begin by describing its constituent parts: Monte Carlo integration and Markov chains. We then describe the general form of MCMC given by the Metropolis–Hastings algorithm, and a special case: the Gibbs sampler.

1.3.1 Monte Carlo integration

Monte Carlo integration evaluates $E[f(X)]$ by drawing samples $\{X_t, t = 1, \dots, n\}$ from $\pi(\cdot)$ and then approximating

$$E[f(X)] \approx \frac{1}{n} \sum_{t=1}^n f(X_t).$$

So the population mean of $f(X)$ is estimated by a sample mean. When the samples $\{X_t\}$ are independent, laws of large numbers ensure that the approximation can be made as accurate as desired by increasing the sample size n . Note that here n is under the control of the analyst: it is not the size of a fixed data sample.

In general, drawing samples $\{X_t\}$ independently from $\pi(\cdot)$ is not feasible, since $\pi(\cdot)$ can be quite non-standard. However the $\{X_t\}$ need not necessarily be independent. The $\{X_t\}$ can be generated by any process which, loosely speaking, draws samples throughout the support of $\pi(\cdot)$ in the correct proportions. One way of doing this is through a Markov chain having $\pi(\cdot)$ as its stationary distribution. This is then *Markov chain Monte Carlo*.

1.3.2 Markov chains

Suppose we generate a sequence of random variables, $\{X_0, X_1, X_2, \dots\}$, such that at each time $t \geq 0$, the next state X_{t+1} is sampled from a distribution $P(X_{t+1}|X_t)$ which depends only on the current state of the chain, X_t . That is, *given* X_t , the next state X_{t+1} does not depend further on the history of the chain $\{X_0, X_1, \dots, X_{t-1}\}$. This sequence is called a *Markov chain*, and $P(\cdot|\cdot)$ is called the *transition kernel* of the chain. We will assume that the chain is time-homogenous: that is, $P(\cdot|\cdot)$ does not depend on t .

How does the starting state X_0 affect X_t ? This question concerns the distribution of X_t given X_0 , which we denote $P^{(t)}(X_t|X_0)$. Here we are not given the intervening variables $\{X_1, X_2, \dots, X_{t-1}\}$, so X_t depends directly on X_0 . Subject to regularity conditions, the chain will gradually ‘forget’ its initial state and $P^{(t)}(\cdot|X_0)$ will eventually converge to a unique *stationary* (or *invariant*) distribution, which does not depend on t or X_0 . For the moment, we denote the stationary distribution by $\phi(\cdot)$. Thus as t increases, the sampled points $\{X_t\}$ will look increasingly like dependent samples from $\phi(\cdot)$. This is illustrated in Figure 1.1, where $\phi(\cdot)$ is univariate standard normal. Note that convergence is much quicker in Figure 1.1(a) than in Figures 1.1(b) or 1.1(c).

Thus, after a sufficiently long *burn-in* of say m iterations, points $\{X_t; t = m+1, \dots, n\}$ will be dependent samples approximately from $\phi(\cdot)$. We discuss methods for determining m in Section 1.4.6. We can now use the output from the Markov chain to estimate the expectation $E[f(X)]$, where X has distribution $\phi(\cdot)$. Burn-in samples are usually discarded for this calculation, giving an estimator

$$\bar{f} = \frac{1}{n-m} \sum_{t=m+1}^n f(X_t). \quad (1.2)$$

This is called an *ergodic average*. Convergence to the required expectation is ensured by the ergodic theorem.

See Roberts (1995) and Tierney (1995) in this volume for more technical discussion of several of the issues raised here.

1.3.3 The Metropolis–Hastings algorithm

Equation (1.2) shows how a Markov chain can be used to estimate $E[f(X)]$, where the expectation is taken over its stationary distribution $\phi(\cdot)$. This would seem to provide the solution to our problem, but first we need to discover how to construct a Markov chain such that its stationary distribution $\phi(\cdot)$ is precisely our distribution of interest $\pi(\cdot)$.

Constructing such a Markov chain is surprisingly easy. We describe the form due to Hastings (1970), which is a generalization of the method

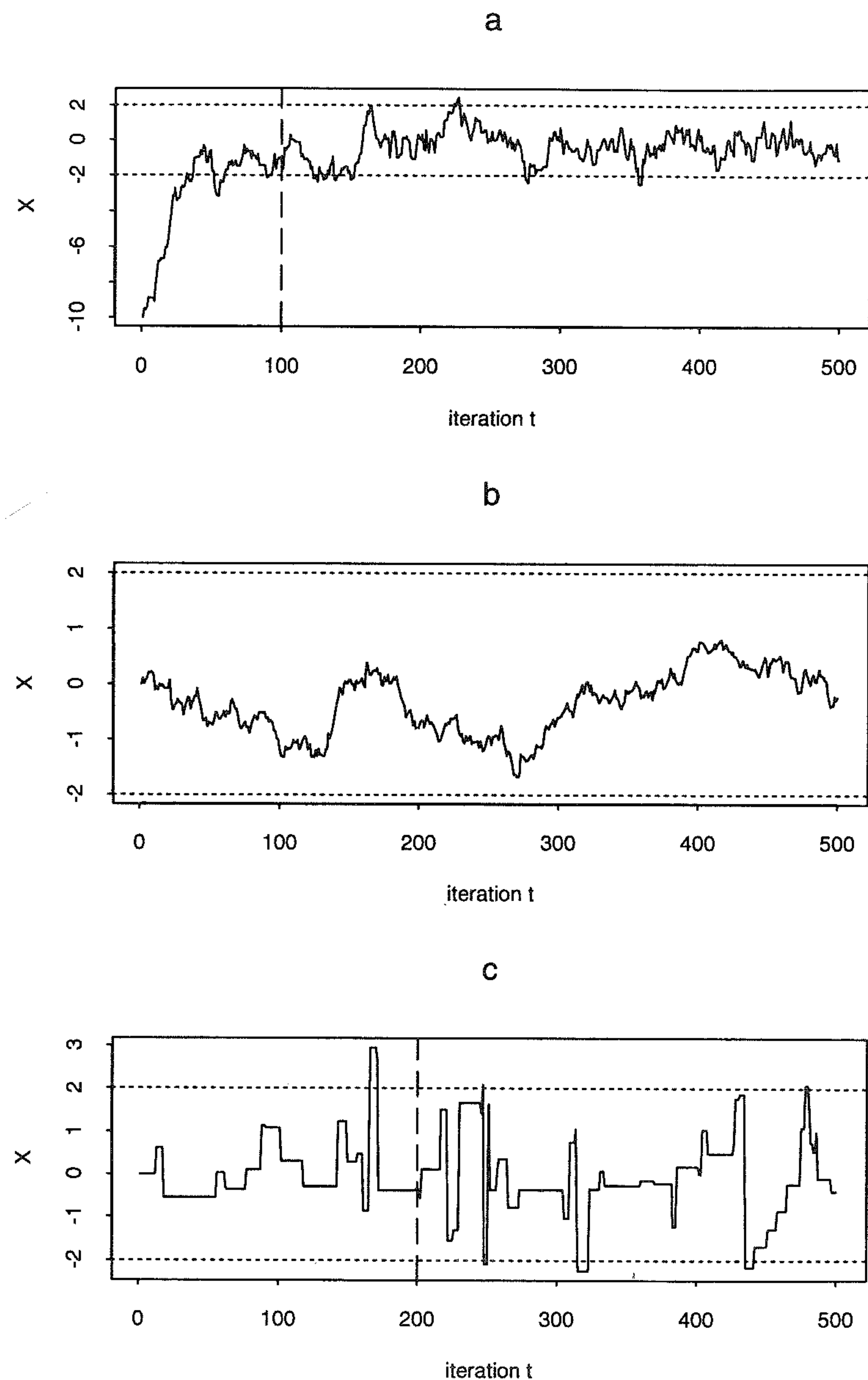


Figure 1.1 500 iterations from Metropolis algorithms with stationary distribution $N(0, 1)$ and proposal distributions (a) $q(.|X) = N(X, 0.5)$; (b) $q(.|X) = N(X, 0.1)$; and (c) $q(.|X) = N(X, 10.0)$. The burn-in is taken to be to the left of the vertical broken line.

first proposed by Metropolis *et al.* (1953). For the *Metropolis-Hastings* (or *Hastings-Metropolis*) algorithm, at each time t , the next state X_{t+1} is chosen by first sampling a *candidate* point Y from a *proposal* distribution $q(.|X_t)$. Note that the proposal distribution may depend on the current point X_t . For example, $q(.|X)$ might be a multivariate normal distribution with mean X and a fixed covariance matrix. The candidate point Y is then *accepted* with probability $\alpha(X_t, Y)$ where

$$\alpha(X, Y) = \min \left(1, \frac{\pi(Y)q(X|Y)}{\pi(X)q(Y|X)} \right) \quad (1.3)$$

If the candidate point is accepted, the next state becomes $X_{t+1} = Y$. If the candidate is rejected, the chain does not move, i.e. $X_{t+1} = X_t$. Figure 1.1 illustrates this for univariate normal proposal and target distributions; Figure 1.1(c) showing many instances where the chain did not move for several iterations.

Thus the Metropolis-Hastings algorithm is extremely simple:

```
Initialize  $X_0$ ; set  $t = 0$ .
Repeat {
  Sample a point  $Y$  from  $q(.|X_t)$ 
  Sample a Uniform(0,1) random variable  $U$ 
  If  $U \leq \alpha(X_t, Y)$  set  $X_{t+1} = Y$ 
  otherwise set  $X_{t+1} = X_t$ 
  Increment  $t$ 
}
```

Remarkably, the proposal distribution $q(.|.)$ can have any form and the stationary distribution of the chain will be $\pi(.)$. (For regularity conditions see Roberts, 1995: this volume.) This can be seen from the following argument. The transition kernel for the Metropolis-Hastings algorithm is

$$P(X_{t+1}|X_t) = q(X_{t+1}|X_t)\alpha(X_t, X_{t+1}) + I(X_{t+1} = X_t)[1 - \int q(Y|X_t)\alpha(X_t, Y)dY], \quad (1.4)$$

where $I(.)$ denotes the indicator function (taking the value 1 when its argument is true, and 0 otherwise). The first term in (1.4) arises from acceptance of a candidate $Y = X_{t+1}$, and the second term arises from rejection, for all possible candidates Y . Using the fact that

$$\pi(X_t)q(X_{t+1}|X_t)\alpha(X_t, X_{t+1}) = \pi(X_{t+1})q(X_t|X_{t+1})\alpha(X_{t+1}, X_t)$$

which follows from (1.3), we obtain the *detailed balance* equation:

$$\pi(X_t)P(X_{t+1}|X_t) = \pi(X_{t+1})P(X_t|X_{t+1}). \quad (1.5)$$

Integrating both sides of (1.5) with respect to X_t gives:

$$\int \pi(X_t)P(X_{t+1}|X_t)dX_t = \pi(X_{t+1}). \quad (1.6)$$

The left-hand side of equation (1.6) gives the marginal distribution of X_{t+1} under the assumption that X_t is from $\pi(\cdot)$. Therefore (1.6) says that if X_t is from $\pi(\cdot)$, then X_{t+1} will be also. Thus, once a sample from the stationary distribution has been obtained, all subsequent samples will be from that distribution. This only proves that the stationary distribution is $\pi(\cdot)$, and is not a complete justification for the Metropolis–Hastings algorithm. A full justification requires a proof that $P^{(t)}(X_t|X_0)$ will converge to the stationary distribution. See Roberts (1995) and Tierney (1995) in this volume for further details.

So far we have assumed that X is a fixed-length vector of k continuous random variables. As noted in Section 1.2, there are many other possibilities, in particular X can be of *variable dimension*. For example, in a Bayesian mixture model, the number of mixture components may be variable: each component possessing its own scale and location parameters. In this situation, $\pi(\cdot)$ must specify the joint distribution of k and X , and $q(Y|X)$ must be able to propose moves between spaces of differing dimension. Then Metropolis–Hastings is as described above, with formally the same expression (1.3) for the acceptance probability, but where dimension-matching conditions for moves between spaces of differing dimension must be carefully considered (Green, 1994a,b). See also Geyer and Møller (1993), Grenander and Miller (1994), and Phillips and Smith (1995: this volume) for MCMC methodology in variably dimensioned problems.

1.4 Implementation

There are several issues which arise when implementing MCMC. We discuss these briefly here. Further details can be found throughout this volume, and in particular in Chapters 5–8. The most immediate issue is the choice of proposal distribution $q(\cdot|\cdot)$.

1.4.1 Canonical forms of proposal distribution

As already noted, any proposal distribution will ultimately deliver samples from the target distribution $\pi(\cdot)$. However, the rate of convergence to the stationary distribution will depend crucially on the relationship between $q(\cdot|\cdot)$ and $\pi(\cdot)$. Moreover, having ‘converged’, the chain may still *mix* slowly (i.e. move slowly around the support of $\pi(\cdot)$). These phenomena are illustrated in Figure 1.1. Figure 1.1(a) shows rapid convergence from a somewhat extreme starting value: thereafter the chain mixes rapidly. Figure 1.1(b),(c) shows slow mixing chains: these would have to be run much longer to obtain reliable estimates from (1.2), despite having been started at the mode of $\pi(\cdot)$.

In high-dimensional problems with little symmetry, it is often necessary to perform exploratory analyses to determine roughly the shape and ori-

entation of $\pi(\cdot)$. This will help in constructing a proposal $q(\cdot|\cdot)$ which leads to rapid mixing. Progress in practice often depends on experimentation and craftsmanship, although untuned canonical forms for $q(\cdot|\cdot)$ often work surprisingly well. For computational efficiency, $q(\cdot|\cdot)$ should be chosen so that it can be easily sampled and evaluated.

Here we describe some canonical forms for $q(\cdot|\cdot)$. Roberts (1995), Tierney (1995) and Gilks and Roberts (1995) in this volume discuss rates of convergence and strategies for choosing $q(\cdot|\cdot)$ in more detail.

The Metropolis Algorithm

The *Metropolis algorithm* (Metropolis *et al.*, 1953) considers only symmetric proposals, having the form $q(Y|X) = q(X|Y)$ for all X and Y . For example, when X is continuous, $q(\cdot|X)$ might be a multivariate normal distribution with mean X and constant covariance matrix Σ . Often it is convenient to choose a proposal which generates each component of Y conditionally independently, given X_t . For the Metropolis algorithm, the acceptance probability (1.3) reduces to

$$\alpha(X, Y) = \min \left(1, \frac{\pi(Y)}{\pi(X)} \right) \quad (1.7)$$

A special case of the Metropolis algorithm is *random-walk Metropolis*, for which $q(Y|X) = q(|X - Y|)$. The data in Figure 1.1 were generated by random-walk Metropolis algorithms.

When choosing a proposal distribution, its scale (for example Σ) may need to be chosen carefully. A cautious proposal distribution generating small steps $Y - X_t$ will generally have a high acceptance rate (1.7), but will nevertheless mix slowly. This is illustrated in Figure 1.1(b). A bold proposal distribution generating large steps will often propose moves from the body to the tails of the distribution, giving small values of $\pi(Y)/\pi(X_t)$ and a low probability of acceptance. Such a chain will frequently not move, again resulting in slow mixing as illustrated in Figure 1.1(c). Ideally, the proposal distribution should be scaled to avoid both these extremes.

The independence sampler

The *independence sampler* (Tierney, 1994) is a Metropolis–Hastings algorithm whose proposal $q(Y|X) = q(Y)$ does not depend on X . For this, the acceptance probability (1.3) can be written in the form

$$\alpha(X, Y) = \min \left(1, \frac{w(Y)}{w(X)} \right), \quad (1.8)$$

where $w(X) = \pi(X)/q(X)$.

In general, the independence sampler can work very well or very badly (see Roberts, 1995: this volume). For the independence sampler to work

well, $q(\cdot)$ should be a good approximation to $\pi(\cdot)$, but it is safest if $q(\cdot)$ is heavier-tailed than $\pi(\cdot)$. To see this, suppose $q(\cdot)$ is lighter-tailed than $\pi(\cdot)$, and that X_t is currently in the tails of $\pi(\cdot)$. Most candidates will not be in the tails, so $w(X_t)$ will be much larger than $w(Y)$ giving a low acceptance probability (1.8). Thus heavy-tailed independence proposals help to avoid long periods stuck in the tails, at the expense of an increased overall rate of candidate rejection.

In some situations, in particular where it is thought that large-sample theory might be operating, a multivariate normal proposal might be tried, with mean at the mode of $\pi(\cdot)$ and covariance matrix somewhat greater than the inverse Hessian matrix

$$\left[-\frac{d^2 \log \pi(x)}{dx^T dx} \right]^{-1}$$

evaluated at the mode.

Single-component Metropolis-Hastings

Instead of updating the whole of X *en bloc*, it is often more convenient and computationally efficient to divide X into components $\{X_{.1}, X_{.2}, \dots, X_{.h}\}$ of possibly differing dimension, and then update these components one by one. This was the framework for MCMC originally proposed by Metropolis *et al.* (1953), and we refer to it as *single-component Metropolis-Hastings*. Let $X_{.-i} = \{X_{.1}, \dots, X_{.i-1}, X_{.i+1}, \dots, X_{.h}\}$, so $X_{.-i}$ comprises all of X except $X_{.i}$.

An iteration of the single-component Metropolis-Hastings algorithm comprises h updating steps, as follows. Let $X_{t,i}$ denote the state of $X_{.i}$ at the end of iteration t . For step i of iteration $t+1$, $X_{.i}$ is updated using Metropolis-Hastings. The candidate $Y_{.i}$ is generated from a proposal distribution $q_i(Y_{.i}|X_{t,i}, X_{t,-i})$, where $X_{t,-i}$ denotes the value of $X_{.-i}$ after completing step $i-1$ of iteration $t+1$:

$$X_{t,-i} = \{X_{t+1,1}, \dots, X_{t+1,i-1}, X_{t,i+1}, \dots, X_{t,h}\},$$

where components $1, 2, \dots, i-1$ have already been updated. Thus the i^{th} proposal distribution $q_i(\cdot|\cdot, \cdot)$ generates a candidate only for the i^{th} component of X , and may depend on the *current* values of any of the components of X . The candidate is accepted with probability $\alpha(X_{t,-i}, X_{t,i}, Y_{.i})$ where

$$\alpha(X_{t,-i}, X_{t,i}, Y_{.i}) = \min \left(1, \frac{\pi(Y_{.i}|X_{t,-i})q_i(X_{t,i}|Y_{.i}, X_{t,-i})}{\pi(X_{t,i}|X_{t,-i})q_i(Y_{.i}|X_{t,i}, X_{t,-i})} \right). \quad (1.9)$$

Here $\pi(X_{.i}|X_{.-i})$ is the *full conditional* distribution for $X_{.i}$ under $\pi(\cdot)$ (see below). If $Y_{.i}$ is accepted, we set $X_{t+1,i} = Y_{.i}$; otherwise, we set $X_{t+1,i} = X_{t,i}$. The remaining components are not changed at step i .

Thus each updating step produces a move in the direction of a coordinate axis (if the candidate is accepted), as illustrated in Figure 1.2. The proposal

distribution $q_i(\cdot|\cdot, \cdot)$ can be chosen in any of the ways discussed earlier in this section.

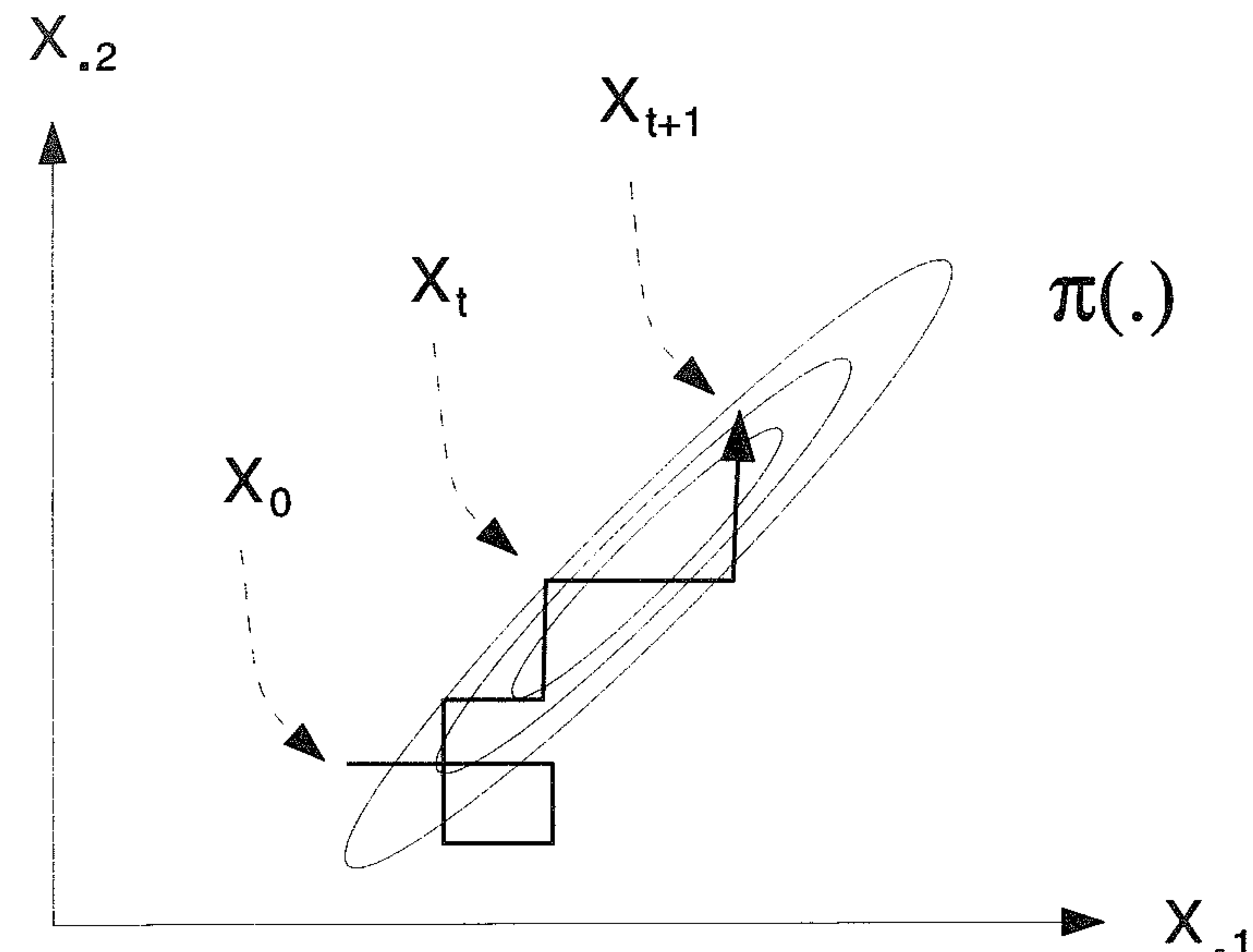


Figure 1.2 Illustrating a single-component Metropolis-Hastings algorithm for a bivariate target distribution $\pi(\cdot)$. Components 1 and 2 are updated alternately, producing alternate moves in horizontal and vertical directions.

The full conditional distribution $\pi(X_{.i}|X_{.-i})$ is the distribution of the i^{th} component of X conditioning on all the remaining components, where X has distribution $\pi(\cdot)$:

$$\pi(X_{.i}|X_{.-i}) = \frac{\pi(X)}{\int \pi(X) dX_{.i}}. \quad (1.10)$$

Full conditional distributions play a prominent role in many of the applications in this volume, and are considered in detail by Gilks (1995: this volume). That the single-component Metropolis-Hastings algorithm with acceptance probability given by (1.9) does indeed generate samples from the target distribution $\pi(\cdot)$ results from the fact that $\pi(\cdot)$ is uniquely determined by the set of its full conditional distributions (Besag, 1974).

In applications, (1.9) often simplifies considerably, particularly when $\pi(\cdot)$ derives from a conditional independence model: see Spiegelhalter *et al.* (1995) and Gilks (1995) in this volume. This provides an important computational advantage. Another important advantage of single-component updating occurs when the target distribution $\pi(\cdot)$ is naturally specified in terms of its full conditional distributions, as commonly occurs in spatial

models; see Besag (1974), Besag *et al.* (1995) and Green (1995: this volume).

Gibbs sampling

A special case of single-component Metropolis–Hastings is the *Gibbs sampler*. The Gibbs sampler was given its name by Geman and Geman (1984), who used it for analysing Gibbs distributions on lattices. However, its applicability is not limited to Gibbs distributions, so ‘Gibbs sampling’ is really a misnomer. Moreover, the same method was already in use in statistical physics, and was known there as the *heat bath algorithm*. Nevertheless, the work of Geman and Geman (1984) led to the introduction of MCMC into mainstream statistics via the articles by Gelfand and Smith (1990) and Gelfand *et al.* (1990). To date, most statistical applications of MCMC have used Gibbs sampling.

For the Gibbs sampler, the proposal distribution for updating the i^{th} component of X is

$$q_i(Y_i|X_{-i}, X_{-i}) = \pi(Y_i|X_{-i}) \quad (1.11)$$

where $\pi(Y_i|X_{-i})$ is the full conditional distribution (1.10). Substituting (1.11) into (1.9) gives an acceptance probability of 1; that is, Gibbs sampler candidates are always accepted. Thus Gibbs sampling consists purely in sampling from full conditional distributions. Methods for sampling from full conditional distributions are described in Gilks (1995: this volume).

1.4.2 Blocking

Our description of single-component samplers in Section 1.4.1 said nothing about how the components should be chosen. Typically, low-dimensional or scalar components are used. In some situations, multivariate components are natural. For example, in a Bayesian random-effects model, an entire precision matrix would usually comprise a single component. When components are highly correlated in the stationary distribution $\pi(\cdot)$, mixing can be slow; see Gilks and Roberts (1995: this volume). Blocking highly correlated components into a higher-dimensional component may improve mixing, but this depends on the choice of proposal.

1.4.3 Updating order

In the above description of the single-component Metropolis–Hastings algorithm and Gibbs sampling, we assumed a fixed updating order for the components of X_t . Although this is usual, a fixed order is not necessary: random permutations of the updating order are quite acceptable. Moreover, not all components need be updated in each iteration. For example,

we could instead update only one component per iteration, selecting component i with some fixed probability $s(i)$. A natural choice would be to set $s(i) = \frac{1}{h}$. Zeger and Karim (1991) suggest updating highly correlated components more frequently than other components, to improve mixing. Note that if $s(i)$ is allowed to depend on X_t then the acceptance probability (1.9) should be modified, otherwise the stationary distribution of the chain may no longer be the target distribution $\pi(\cdot)$. Specifically, the acceptance probability becomes

$$\min \left(1, \frac{\pi(Y_i|X_{-i})s(i|Y_i, X_{-i})q_i(X_i|Y_i, X_{-i})}{\pi(X_i|X_{-i})s(i|X_i, X_{-i})q_i(Y_i|X_i, X_{-i})} \right).$$

1.4.4 Number of chains

So far we have considered running only one chain, but multiple chains are permissible. Recommendations in the literature have been conflicting, ranging from many short chains (Gelfand and Smith, 1990), to several long ones (Gelman and Rubin, 1992a,b), to one very long one (Geyer, 1992). It is now generally agreed that running many short chains, motivated by a desire to obtain independent samples from $\pi(\cdot)$, is misguided unless there is some special reason for needing independent samples. Certainly, independent samples are not required for ergodic averaging in (1.2). The debate between the several-long-runs school and the one-very-long-run school seems set to continue. The latter maintains that one very long run has the best chance of finding new modes, and comparison between chains can never prove convergence, whilst the former maintains that comparing several seemingly converged chains might reveal genuine differences if the chains have not yet approached stationarity; see Gelman (1995: this volume). If several processors are available, running one chain on each will generally be worthwhile.

1.4.5 Starting values

Not much has been written on this topic. If the chain is irreducible, the choice of starting values X_0 will not affect the stationary distribution. A rapidly mixing chain, such as in Figure 1.1(a), will quickly find its way from extreme starting values. Starting values may need to be chosen more carefully for slow-mixing chains, to avoid a lengthy burn-in. However, it is seldom necessary to expend much effort in choosing starting values. Gelman and Rubin (1992a,b) suggest using ‘over-dispersed’ starting values in multiple chains, to assist in assessing convergence; see below and Gelman (1995: this volume).

1.4.6 Determining burn-in

The length of burn-in m depends on X_0 , on the rate of convergence of $P^{(t)}(X_t|X_0)$ to $\pi(X_t)$ and on how similar $P^{(t)}(.|.)$ and $\pi(.)$ are required to be. Theoretically, having specified a criterion of 'similar enough', m can be determined analytically. However, this calculation is far from computationally feasible in most situations (see Roberts, 1995: this volume). Visual inspection of plots of (functions of) the Monte-Carlo output $\{X_t, t = 1, \dots, n\}$ is the most obvious and commonly used method for determining burn-in, as in Figure 1.1. Starting the chain close to the mode of $\pi(.)$ does not remove the need for a burn-in, as the chain should still be run long enough for it to 'forget' its starting position. For example, in Figure 1.1(b) the chain has not wandered far from its starting position in 500 iterations. In this case, m should be set greater than 500.

More formal tools for determining m , called *convergence diagnostics*, have been proposed. Convergence diagnostics use a variety of theoretical methods and approximations, but all make use of the Monte Carlo output in some way. By now, at least 10 convergence diagnostics have been proposed; for a recent review, see Cowles and Carlin (1994). Some of these diagnostics are also suited to determining run length n (see below).

Convergence diagnostics can be classified by whether or not they are based on an arbitrary function $f(X)$ of the Monte Carlo output; whether they use output from a single chain or from multiple chains; and whether they can be based purely on the Monte Carlo output.

Methods which rely on monitoring $\{f(X_t), t = 1, \dots, n\}$ (e.g. Gelman and Rubin, 1992b; Raftery and Lewis, 1992; Geweke, 1992) are easy to apply, but may be misleading since $f(X_t)$ may appear to have converged in distribution by iteration m , whilst another unmonitored function $g(X_t)$ may not have. Whatever functions $f(.)$ are monitored, there may be others which behave differently.

From a theoretical perspective, it is better to compare globally the full joint distribution $P^{(t)}(.)$ with $\pi(.)$. To avoid having to deal with $P^{(t)}(.)$ directly, several methods obtain samples from it by running multiple parallel chains (Ritter and Tanner, 1992; Roberts, 1992; Liu and Liu, 1993), and make use of the transition kernel $P(.|.)$. However, for stability in the procedures, it may be necessary to run many parallel chains. When convergence is slow, this is a serious practical limitation.

Running parallel chains obviously increases the computational burden, but can be useful, even informally, to diagnose slow convergence. For example, several parallel chains might individually appear to have converged, but comparisons between them may reveal marked differences in the apparent stationary distributions (Gelman and Rubin, 1992a).

From a practical perspective, methods which are based purely on the Monte Carlo output are particularly convenient, allowing assessment of

convergence without recourse to the transition kernel $P(.|.)$, and hence without model-specific coding.

This volume does not contain a review of convergence diagnostics. This is still an active area of research, and much remains to be learnt about the behaviour of existing methods in real applications, particularly in high dimensions and when convergence is slow. Instead, the chapters by Raftery and Lewis (1995) and Gelman (1995) in this volume contain descriptions of two of the most popular methods. Both methods monitor an arbitrary function $f(.)$, and are based purely on the Monte Carlo output. The former uses a single chain and the latter multiple chains.

Geyer (1992) suggests that calculation of the length of burn-in is unnecessary, as it is likely to be less than 1% of the total length of a run sufficiently long to obtain adequate precision in the estimator \bar{f} in (1.2), (see below). If extreme starting values are avoided, Geyer suggests setting m to between 1% and 2% of the run length n .

1.4.7 Determining stopping time

Deciding when to stop the chain is an important practical matter. The aim is to run the chain long enough to obtain adequate precision in the estimator \bar{f} in (1.2). Estimation of the variance of \bar{f} (called the *Monte Carlo variance*) is complicated by lack of independence in the iterates $\{X_t\}$.

The most obvious informal method for determining run length n is to run several chains in parallel, with different starting values, and compare the estimates \bar{f} from (1.2). If they do not agree adequately, n must be increased. More formal methods which aim to estimate the variance of \bar{f} have been proposed: see Roberts (1995) and Raftery and Lewis (1995) in this volume for further details.

1.4.8 Output analysis

In Bayesian inference, it is usual to summarize the posterior distribution $\pi(.)$ in terms of means, standard deviations, correlations, credible intervals and marginal distributions for components X_i of interest. Means, standard deviations and correlations can all be estimated by their sample equivalents in the Monte Carlo output $\{X_{t,i}, t = m+1, \dots, n\}$, according to (1.2). For example, the marginal mean and variance of X_i are estimated by

$$\bar{X}_{.i} = \frac{1}{n-m} \sum_{t=m+1}^n X_{t,i}$$

and

$$S_{.i}^2 = \frac{1}{n-m-1} \sum_{t=m+1}^n (X_{t,i} - \bar{X}_{.i})^2.$$

Note that these estimates simply ignore other components in the Monte Carlo output.

A $100(1 - 2p)\%$ credible interval $[c_p, c_{1-p}]$ for a scalar component X_i can be estimated by setting c_p equal to the p^{th} quantile of $\{X_{t,i}, t = m + 1, \dots, n\}$, and c_{1-p} equal to the $(1 - p)^{th}$ quantile. Besag *et al.* (1995) give a procedure for calculating rectangular credible regions in two or more dimensions.

Marginal distributions can be estimated by kernel density estimation. For the marginal distribution of X_i , this is

$$\pi(X_i) \approx \frac{1}{n - m} \sum_{t=m+1}^n K(X_i | X_t),$$

where $K(\cdot | X_t)$ is a density concentrated around $X_{t,i}$. A natural choice for $K(X_i | X_t)$ is the full conditional distribution $\pi(X_i | X_{t,-i})$. Gelfand and Smith (1990) use this construction to estimate expectations under $\pi(\cdot)$. Thus their *Rao-Blackwellized* estimator of $E[f(X_i)]$ is

$$\bar{f}_{RB} = \frac{1}{n - m} \sum_{t=m+1}^n E[f(X_i) | X_{t,-i}], \quad (1.12)$$

where the expectation is with respect to the full conditional $\pi(X_i | X_{t,-i})$. With reasonably long runs, the improvement from using (1.12) instead of (1.2) is usually slight, and in any case (1.12) requires a closed form for the full conditional expectation.

1.5 Discussion

This chapter provides a brief introduction to MCMC. We hope we have convinced readers that MCMC is a simple idea with enormous potential. The following chapters fill out many of the ideas sketched here, and in particular give some indication of where the methods work well and where they need some tuning or further development.

MCMC methodology and Bayesian estimation go together naturally, as many of the chapters in this volume testify. However, Bayesian model validation is still a difficult area. Some techniques for Bayesian model validation using MCMC are described in Chapters 9–13.

The philosophical debate between Bayesians and non-Bayesians has continued for decades and has largely been sterile from a practical perspective. For many applied statisticians, the most persuasive argument is the availability of robust methods and software. For many years, Bayesians had difficulty solving problems which were straightforward for non-Bayesians, so it is not surprising that most applied statisticians today are non-Bayesian. With the arrival of MCMC and related software, notably the Gibbs sampling program BUGS (see Spiegelhalter *et al.*, 1995: this volume), we hope

more applied statisticians will become familiar and comfortable with Bayesian ideas, and apply them.

References

- Besag, J. (1974) Spatial interaction and the statistical analysis of lattice systems (with discussion). *J. R. Statist. Soc. B*, **36**, 192–236.
- Besag, J., Green, P., Higdon, D. and Mengersen, K. (1995) Bayesian computation and stochastic systems. *Statist. Sci.* (in press).
- Cowles, M. K. and Carlin, B. P. (1994) Markov chain Monte Carlo convergence diagnostics: a comparative review. *Technical Report 94-008*, Division of Biostatistics, School of Public Health, University of Minnesota.
- Diebolt, J. and Ip, E. H. S. (1995) Stochastic EM: methods and application. In *Markov Chain Monte Carlo in Practice* (eds W. R. Gilks, S. Richardson and D. J. Spiegelhalter), pp. 259–273. London: Chapman & Hall.
- Gelfand, A. E. and Smith, A. F. M. (1990) Sampling-based approaches to calculating marginal densities. *J. Am. Statist. Ass.*, **85**, 398–409.
- Gelfand, A. E., Hills, S. E., Racine-Poon, A. and Smith, A. F. M. (1990) Illustration of Bayesian inference in normal data models using Gibbs sampling. *J. Am. Statist. Ass.*, **85**, 972–985.
- Gelman, A. (1995) Inference and monitoring convergence. In *Markov Chain Monte Carlo in Practice* (eds W. R. Gilks, S. Richardson and D. J. Spiegelhalter), pp. 131–143. London: Chapman & Hall.
- Gelman, A. and Rubin, D. B. (1992a) A single series from the Gibbs sampler provides a false sense of security. In *Bayesian Statistics 4* (eds J. M. Bernardo, J. Berger, A. P. Dawid and A. F. M. Smith), pp. 625–631. Oxford: Oxford University Press.
- Gelman, A. and Rubin, D. B. (1992b) Inference from iterative simulation using multiple sequences. *Statist. Sci.*, **7**, 457–472.
- Geman, S. and Geman, D. (1984) Stochastic relaxation, Gibbs distributions and the Bayesian restoration of images. *IEEE Trans. Pattern Anal. Mach. Intel.*, **6**, 721–741.
- Geweke, J. (1992) Evaluating the accuracy of sampling-based approaches to the calculation of posterior moments. In *Bayesian Statistics 4* (eds J. M. Bernardo, J. Berger, A. P. Dawid and A. F. M. Smith), pp. 169–193. Oxford: Oxford University Press.
- Geyer, C. J. (1992) Practical Markov chain Monte Carlo. *Statist. Sci.*, **7**, 473–511.
- Geyer, C. J. (1995) Estimation and optimization of functions. In *Markov Chain Monte Carlo in Practice* (eds W. R. Gilks, S. Richardson and D. J. Spiegelhalter), pp. 241–258. London: Chapman & Hall.
- Geyer, C. J. and Møller, J. (1993) Simulation procedures and likelihood inference for spatial point processes. *Technical Report*, University of Aarhus.
- Gilks, W. R. (1995) Full conditional distributions. In *Markov Chain Monte Carlo in Practice* (eds W. R. Gilks, S. Richardson and D. J. Spiegelhalter), pp. 75–88. London: Chapman & Hall.