

Markov chain Monte Carlo without likelihoods

Paul Marjoram*, John Molitor*, Vincent Plagnol†, and Simon Tavaré†*

*Biostatistics Division, Department of Preventive Medicine, Keck School of Medicine, and †Molecular and Computational Biology, Department of Biological Sciences, University of Southern California, Los Angeles, CA 90089

Communicated by Michael S. Waterman, University of Southern California, Los Angeles, CA, October 24, 2003 (received for review June 20, 2003)

Many stochastic simulation approaches for generating observations from a posterior distribution depend on knowing a likelihood function. However, for many complex probability models, such likelihoods are either impossible or computationally prohibitive to obtain. Here we present a Markov chain Monte Carlo method for generating observations from a posterior distribution without the use of likelihoods. It can also be used in frequentist applications, in particular for maximum-likelihood estimation. The approach is illustrated by an example of ancestral inference in population genetics. A number of open problems are highlighted in the discussion.

One of the basic problems in Bayesian statistics is the computation of posterior distributions. We imagine data \mathcal{D} generated from a model \mathcal{M} determined by parameters θ , the prior density of which is denoted by $\pi(\theta)$. We assume unless otherwise stated that the data are discrete. The posterior distribution of interest is $f(\theta|\mathcal{D})$, which is given by

$$f(\theta|\mathcal{D}) = \mathbb{P}(\mathcal{D}|\theta)\pi(\theta)/\mathbb{P}(\mathcal{D}), \quad [1]$$

where $\mathbb{P}(\mathcal{D}) = \int \mathbb{P}(\mathcal{D}|\theta)\pi(\theta)d\theta$ is the normalizing constant.

In most scientific contexts, explicit formulae for such posterior densities are few and far between, and we usually resort to stochastic simulation to generate observations from f . Perhaps the simplest approach for this is the **rejection method**:

- A1. Generate θ from $\pi(\cdot)$.
- A2. Accept θ with probability $h = \mathbb{P}(\mathcal{D}|\theta)$; return to A1.

Accepted observations have distribution $f(\theta|\mathcal{D})$ (cf. ref. 1). The computations can often be accelerated if an upper bound c for $\mathbb{P}(\mathcal{D}|\theta)$ is known; h then is replaced by h/c . If $\hat{\theta}$ denotes the maximum-likelihood estimator of θ , we could take $c = \mathbb{P}(\mathcal{D}|\hat{\theta})$.

There are many variations on this theme. Of particular relevance here is the case in which **the likelihood $\mathbb{P}(\mathcal{D}|\theta)$ cannot be computed explicitly**. One obvious approach then is:

- B1. Generate θ from $\pi(\cdot)$.
- B2. Simulate \mathcal{D}' from the model \mathcal{M} with parameter θ .
- B3. Accept θ if $\mathcal{D}' = \mathcal{D}$; return to B1.

The success of this approach depends on the fact that the underlying stochastic model \mathcal{M} is easy to simulate. This approach can be useful when computation of the likelihood is possible but time-consuming.

The practicality of algorithms such as these depends crucially on the size of $\mathbb{P}(\mathcal{D})$, because the probability of accepting an observation is proportional to $\mathbb{P}(\mathcal{D})$. In cases where the acceptance rate is too small, one might resort to approximate methods such as:

- C1. Generate θ from $\pi(\cdot)$.
- C2. Simulate \mathcal{D}' from the model \mathcal{M} with parameter θ .
- C3. Calculate the distance $\rho(\mathcal{D}, \mathcal{D}')$ between \mathcal{D}' and \mathcal{D} .
- C4. Accept θ if $\rho \leq \varepsilon$; return to C1.

This approach requires selection of a suitable metric ρ as well as a choice of ε . As $\varepsilon \rightarrow \infty$ it generates observations from the prior. If $\varepsilon = 0$, an observation \mathcal{D}' is accepted only if $\mathcal{D}' = \mathcal{D}$, and then accepted observations come from the density $f(\theta|\mathcal{D})$. The choice

of ε therefore reflects a tension between computability and accuracy. The method is still honest in that, for a given ρ and ε , we are generating independent and identically distributed observations from $f(\theta|\rho(\mathcal{D}, \mathcal{D}') \leq \varepsilon)$.

When \mathcal{D} is high-dimensional or continuous, this approach can be impractical as well, and then the comparison of \mathcal{D}' with \mathcal{D} can be made by using lower-dimensional summaries of the data. The motivation for this approach is that if the set of statistics $S = (S_1, \dots, S_p)$ is sufficient for θ , in that $\mathbb{P}(\mathcal{D}|S, \theta)$ is independent of θ , then $f(\theta|\mathcal{D}) = f(\theta|S)$. The normalizing constant $\mathbb{P}(S)$ is typically larger than $\mathbb{P}(\mathcal{D})$, resulting in more acceptances. In practice it will be hard, if not impossible, to identify a suitable set of sufficient statistics, and we then might resort to a more heuristic approach. Thus we seek to use knowledge of the particular problem at hand to suggest summary statistics that capture information about θ . With these statistics in hand, we have the following approximate Bayesian computation scheme for data \mathcal{D} summarized by S :

- D1. Generate θ from $\pi(\cdot)$.
- D2. Simulate \mathcal{D}' from stochastic model \mathcal{M} with parameter θ , and compute the corresponding statistics S' .
- D3. Calculate the distance $\rho(S, S')$ between S and S' .
- D4. Accept θ if $\rho \leq \varepsilon$, and return to D1.

There are several advantages to these rejection methods, among them the fact that they are usually easy to code, they generate independent observations (and thus can use embarrassingly parallel computation), and they readily provide estimates of Bayes factors that can be used for model comparison. On the other hand, sampling from the prior in complex probability models is unlikely to be sensible when the posterior is a long way from the prior. Later we discuss Markov chain Monte Carlo (MCMC) algorithms and provide an alternative MCMC approach that does not require the evaluation of likelihoods.

Examples from Evolutionary Biology

Examples of these algorithms have appeared in the evolutionary genetics literature. For example, inference problems in molecular population genetics can be described as follows. We sample the molecular variation present at several loci in a population, obtaining a discrete variation data set \mathcal{D} (DNA sequence data, for example). Inference and estimation for population parameters of interest such as mutation rates, recombination rates, migration rates, and demographic parameters are then based on a stochastic model \mathcal{M} for \mathcal{D} .

The coalescent (2) provides a commonly used modeling framework in this setting. The coalescent is a stochastic model for the ancestral relationships between the sampled sequences. In the absence of recombination, these ancestral relationships form a binary branching tree. Because the tree is not observed, inference for parameters of interest can be thought of as a

Abbreviations: MCMC, Markov chain Monte Carlo; MRCA, most recent common ancestor.

*To whom correspondence should be addressed at: Program in Molecular and Computational Biology, Department of Biological Sciences, SHS 172, University of Southern California, 835 West 37th Street, Los Angeles, CA 90089-1340. E-mail: stavare@usc.edu.

© 2003 by The National Academy of Sciences of the USA

missing data problem (for reviews see, for example, refs. 3 and 4).

Examples of algorithm A are given by Tavaré *et al.* (5), of algorithm C by Plagnol and Tavaré (6), and of algorithm D by Fu and Li (7), Weiss and von Haeseler (8) and Pritchard *et al.* (9), among others. Beaumont *et al.* (10) describe an interesting generalization of the rejection method in which all observations (θ, S') generated by the first two steps of algorithm D are used in a local-linear regression framework to generate observations that follow more closely the required distribution $f(\theta|\mathcal{D})$. This reference also contains a number of other examples of these approaches.

MCMC Methods

We begin by recalling the **Metropolis–Hastings algorithm** (11, 12) for generating observations from $f(\theta|\mathcal{D})$ using output from a Markov chain.

- E1. If now at θ , propose a move to θ' according to a transition kernel $q(\theta \rightarrow \theta')$.
- E2. Calculate

$$h = \min\left(1, \frac{\mathbb{P}(\mathcal{D}|\theta')\pi(\theta')q(\theta' \rightarrow \theta)}{\mathbb{P}(\mathcal{D}|\theta)\pi(\theta)q(\theta \rightarrow \theta')}\right). \quad [2]$$

- E3. Move to θ' with probability h , else remain at θ ; go to E1.

Under suitable regularity conditions, f is the stationary and limiting distribution of the chain. The practical complexities of implementing MCMC are described by Gilks *et al.* (13) for example. In concert with dramatically increased computing power, this approach has revolutionized Bayesian statistics over the last 15 years (see, for example, refs. 14 and 15).

One comparison that can be made between algorithms A and E is the way in which they use the likelihood $\mathbb{P}(\mathcal{D}|\theta)$. In the rejection method, the comparison is with $c = \mathbb{P}(\mathcal{D}|\theta)$ (a global comparison), whereas in the Metropolis–Hastings algorithm $\mathbb{P}(\mathcal{D}|\theta)$ is compared to $\mathbb{P}(\mathcal{D}|\theta')$ (a local comparison). One therefore expects that MCMC approaches accept observations more frequently, but the price paid for higher acceptance rates is dependent outcomes.

Approximating the Likelihood Ratio. The theme of this note is simulation of observations from a posterior when likelihoods are either hard or impossible to calculate. The first such approach is to approximate the likelihood ratio $\mathbb{P}(\mathcal{D}|\theta')/\mathbb{P}(\mathcal{D}|\theta)$ appearing in the acceptance probability in E3. This can be done by estimating each term in the ratio separately. For a given value of θ , estimate $\mathbb{P}(\mathcal{D}|\theta)$ by simulation of B data sets $\mathcal{D}_1, \dots, \mathcal{D}_B$ from the model \mathcal{M} with parameter θ , and form the point estimate

$$\hat{\mathbb{P}}(\mathcal{D}|\theta) = \frac{1}{B} \sum_{j=1}^B \mathbb{I}(\mathcal{D}_j = \mathcal{D}),$$

where $\mathbb{I}(A)$ is 1 if A is true and 0 otherwise. More sophisticated estimates might also be used depending on the details of the specific application. For example, an estimate of $\mathbb{P}(\mathcal{D}|\theta)$ might be precomputed and stored over a grid of θ values.

This method also applies when the underlying data are continuous, in which case the likelihood ratio is a ratio of densities. In this case the B simulated observations can be used in a kernel density-estimation routine, and the density at the point \mathcal{D} is returned. This approach can also be made dynamic, in that B need not be fixed ahead of time. See Diggle and Gratton (16) and the references contained therein for applications of this approach in frequentist settings. Of course, the same methods

can be applied for the approaches described in C and D above. An example appears later.

MCMC Without Likelihoods. In this section we describe an MCMC approach that is the natural analog of algorithm B in that no likelihoods are used or estimated in its implementation. It is based on the following steps:

- F1. If now at θ propose a move to θ' according to a transition kernel $q(\theta \rightarrow \theta')$.
- F2. Generate \mathcal{D}' using model \mathcal{M} with parameters θ' .
- F3. If $\mathcal{D}' = \mathcal{D}$, go to F4, and otherwise stay at θ and return to F1.
- F4. Calculate

$$h = h(\theta, \theta') = \min\left(1, \frac{\pi(\theta')q(\theta' \rightarrow \theta)}{\pi(\theta)q(\theta \rightarrow \theta')}\right).$$

- F5. Accept θ' with probability h and otherwise stay at θ , then return to F1.

The stationary distribution of the chain is indeed $f(\theta|\mathcal{D})$, as is demonstrated below.

Theorem. $f(\theta|\mathcal{D})$ is the stationary distribution of the chain.

Proof: Denote the transition mechanism of the chain by $r(\theta \rightarrow \theta')$, and (without loss of generality) choose $\theta' \neq \theta$ satisfying

$$\frac{\pi(\theta')q(\theta' \rightarrow \theta)}{\pi(\theta)q(\theta \rightarrow \theta')} \leq 1. \quad [3]$$

Then

$$\begin{aligned} f(\theta|\mathcal{D})r(\theta \rightarrow \theta') &= f(\theta|\mathcal{D})q(\theta \rightarrow \theta')\mathbb{P}(\mathcal{D}|\theta')h(\theta, \theta') \\ &= \frac{\mathbb{P}(\mathcal{D}|\theta)\pi(\theta)}{\mathbb{P}(\mathcal{D})} \left\{ q(\theta \rightarrow \theta')\mathbb{P}(\mathcal{D}|\theta') \right. \\ &\quad \times \left. \frac{\pi(\theta')q(\theta' \rightarrow \theta)}{\pi(\theta)q(\theta \rightarrow \theta')} \right\} \\ &= \frac{\mathbb{P}(\mathcal{D}|\theta')\pi(\theta')}{\mathbb{P}(\mathcal{D})} \{q(\theta' \rightarrow \theta)\mathbb{P}(\mathcal{D}|\theta)\} \\ &= f(\theta'|\mathcal{D})q(\theta' \rightarrow \theta)\mathbb{P}(\mathcal{D}|\theta)h(\theta', \theta) \\ &= f(\theta'|\mathcal{D})r(\theta' \rightarrow \theta). \end{aligned}$$

The argument when the ratio on the left of Eq. 3 is >1 is analogous. Thus $f(\theta|\mathcal{D})$ satisfies the detailed balance equations, which implies that indeed $f(\theta|\mathcal{D})$ is the stationary distribution of the chain, and the proof is complete.

Assuming that the chain is ergodic (which occurs under the same conditions that make the chain in algorithm E ergodic), we can now simulate observations having approximately the distribution $f(\theta|\mathcal{D})$. We also mention two special cases:

1. If $q(\theta' \rightarrow \theta) = q(\theta \rightarrow \theta')$ then h depends only on the prior.
2. If q is reversible with respect to π [so that $\pi(\theta)q(\theta \rightarrow \theta') = \pi(\theta')q(\theta' \rightarrow \theta)$ for all $\theta \neq \theta'$], then $h = 1$ and the algorithm reduces to a rejection method with correlated outputs.

For the reasons discussed earlier this approach also may be impractical, in which case we can resort to the equivalent of algorithms C and D by replacing step F3 above with:

- F3'. If $\rho(\mathcal{D}', \mathcal{D}) \leq \varepsilon$, go to F4, and otherwise stay at θ and return to F1,

in which case the stationary distribution is $f(\theta|\rho(\mathcal{D}', \mathcal{D}) \leq \varepsilon)$, or

