

Markov chain Monte Carlo

- Markov chain Monte Carlo (MCMC) methods introduced below can be used to generate a draw from a distribution that approximates some target distribution $f(\cdot)$
- The basic idea is that we know that $\mathbf{X} \sim f(\mathbf{x})$, but we can't get $f(\mathbf{x})$ in closed form or we can't simulate easily from $f(\mathbf{x})$.
- The MCMC methods introduced below can be used to generate a draw from a distribution that approximates f , but they are more properly viewed as methods for generating a sample from which expectations of functions of \mathbf{X} can reliably be estimated.

Brief review of Markov chains: Definitions

1. A discrete-time, discrete-space Markov chain is $X^{(0)}, X^{(1)}, \dots$ where $X^{(t)}$ obeys the *Markov property* that

$$P \left[X^{(t)} \mid x^{(0)}, \dots, x^{(t-1)} \right] = P \left[X^{(t)} \mid x^{(t-1)} \right]$$

2. A Markov chain is governed by a *transition probability matrix*, \mathbf{P} , with the ij th element equaling $P \left[X^{(t+1)} = j \mid X^{(t)} = i \right]$.
3. Let π denote a vector of probabilities with i th element π_i denoting the marginal probability that $X^{(t)} = i$. Then the marginal distribution of $X^{(t+1)}$ must be $\pi^T \mathbf{P}$.

If $\pi^T \mathbf{P} = \pi^T$ then π is a *stationary distribution* for \mathbf{P} .

4. A Markov chain is *irreducible* if any state j can be reached from any state i in a finite number of steps for all i and j .
5. A Markov chain is *periodic* if it can visit certain portions of the state space only at regularly spaced intervals.

Brief review of Markov chains: Stationary distribution

If $X^{(1)}, X^{(2)}, \dots$ are realizations from an irreducible and aperiodic Markov chain with stationary distribution π , then

1. $X^{(n)}$ converges in distribution to the distribution given by π
2. For any function h

$$\frac{1}{n} \sum_{t=1}^n h(X^{(t)}) \rightarrow E_{\pi}\{h(X)\}$$

almost surely as $n \rightarrow \infty$, provided $E_{\pi}\{|h(X)|\}$ exists.

This is one form of the *ergodic theorem*, which is a generalization of the strong law of large numbers.

Markov chain Monte Carlo

- Let $\{\mathbf{X}^{(t)}\}$ denote a Markov chain for $t = 0, 1, 2, \dots$, where $\mathbf{X}^{(t)} = (X_1^{(t)}, \dots, X_p^{(t)})$ and the state space is either continuous or discrete.
- The MCMC sampling strategy is to construct an irreducible, aperiodic Markov chain for which the stationary distribution equals the target distribution, f .
- For sufficiently large t , a realization, $\mathbf{X}^{(t)}$, from this chain will have approximate marginal distribution f .
- There are various strategies for constructing a MCMC algorithm including the Metropolis-Hastings algorithm, the Gibbs sampler, etc.
- A popular application of MCMC methods is to facilitate Bayesian inference where f is a Bayesian posterior distribution for parameters \mathbf{X}

Metropolis-Hastings algorithm

The method begins at $t = 0$ with the selection of $\mathbf{X}^{(0)} = \mathbf{x}^{(0)}$ drawn at random from some starting distribution, g , with the requirement that $f(\mathbf{x}^{(0)}) > 0$.

Given $\mathbf{X}^{(t)} = \mathbf{x}^{(t)}$, the algorithm generates $\mathbf{X}^{(t+1)}$ as follows.

1. Sample a candidate value, \mathbf{X}^* , from a *proposal distribution* $g(\cdot | \mathbf{x}^{(t)})$.
2. Compute the *Metropolis-Hastings ratio*, $R(\mathbf{x}^{(t)}, \mathbf{X}^*)$, where

$$R(\mathbf{u}, \mathbf{v}) = \frac{f(\mathbf{v}) g(\mathbf{u} | \mathbf{v})}{f(\mathbf{u}) g(\mathbf{v} | \mathbf{u})}. \quad (92)$$

3. Sample a value for $\mathbf{X}^{(t+1)}$ according to the following:

$$\mathbf{X}^{(t+1)} = \begin{cases} \mathbf{X}^* & \text{with probability } \min\{R(\mathbf{x}^{(t)}, \mathbf{X}^*), 1\} \\ \mathbf{x}^{(t)} & \text{otherwise.} \end{cases} \quad (93)$$

4. Increment t and return to step 1.

Metropolis-Hastings algorithm

A few details:

- If the proposal distribution is chosen so that the chain is irreducible and aperiodic, then the unique limiting stationary distribution of the chain generated by the Metropolis-Hastings algorithm is the target distribution.
- Sequence $\mathbf{x}^{(0)}, \mathbf{x}^{(1)}, \dots$ will likely include multiple copies of some points in the state space.

It is important to include these copies in the chain and in any estimates based on the output, since the frequencies of sampled points are used to correct for the fact that the proposal density differs from the target density.

Metropolis-Hastings algorithm

The distribution of realizations from the Metropolis-Hastings chain approximates the stationary distribution of the chain as t progresses, therefore

$$\mathbf{E} \{h(\mathbf{X})\} \approx \frac{1}{n} \sum_{i=1}^n h(\mathbf{x}^{(i)})$$

Some useful quantities that can be estimated include

- means: $\mathbf{E} \{h(\mathbf{X})\}$
- variances: $\mathbf{E} \{ (h(\mathbf{X}) - \mathbf{E} \{h(\mathbf{X})\})^2 \}$
- tail probabilities: $\mathbf{E} \{ 1_{\{h(\mathbf{X}) \leq q\}} \}$ for constant q , where $1_{\{A\}} = 1$ if A is true and 0 otherwise
- Virtually any other statistic that may be of interest!

Can also estimate f itself using density estimation methods

Choosing a proposal distribution

A well chosen proposal distribution produces

- candidate values that cover the support of the stationary distribution in a reasonable number of iterations
- candidate values that are not accepted or rejected too frequently

Both of these factors are related to the spread of the proposal distribution:

- **Too diffuse:** candidate values will be rejected frequently and the chain will require many iterations to adequately explore the space of the target distribution
- **Too focused:** the chain will remain in one small region of the target distribution for many iterations while other regions of the target distribution will not be adequately explored

Independence chains

Suppose that the proposal distribution for the Metropolis-Hastings algorithm is chosen such that for some fixed density g

$$g(\mathbf{x}^* | \mathbf{x}^{(t)}) = g(\mathbf{x}^*)$$

This yields an independence chain, where each candidate value is drawn independently of the past.

In this case, the Metropolis-Hastings ratio is

$$R(\mathbf{x}^{(t)}, \mathbf{X}^*) = \frac{f(\mathbf{X}^*) g(\mathbf{x}^{(t)})}{f(\mathbf{x}^{(t)}) g(\mathbf{X}^*)}$$

The resulting Markov chain is irreducible and aperiodic if $g(\mathbf{x}) > 0$ whenever $f(\mathbf{x}) > 0$.

Bayesian inference and MCMC

In the Bayesian paradigm, parameters are assumed to follow some probability distribution. Suppose that data \mathbf{Y} has a distribution parameterized by θ .

Bayes Theorem:

$$\begin{aligned} p(\theta|\mathbf{y}) &= \frac{p(\theta)p(\mathbf{y}|\theta)}{\int p(\theta)p(\mathbf{y}|\theta) d\theta} \\ &= \frac{p(\theta)p(\mathbf{y}|\theta)}{p(\mathbf{y})} \\ &\propto p(\theta)p(\mathbf{y}|\theta) \end{aligned}$$

where

- $p(\theta|\mathbf{y})$ is the *posterior density* of θ
- $p(\mathbf{y}|\theta)$ is the *likelihood* or *sampling distribution*
- $p(\theta)$ is the *prior density*

In Bayesian applications we use MCMC to simulate from the posterior distribution $p(\theta|\mathbf{y})$, so the target distribution is typically the posterior distribution.

Bayesian inference and MCMC

Simple strategy: use the prior as a proposal distribution in an independence chain.

In our Metropolis-Hastings notation,

$$f(\boldsymbol{\theta}) = p(\boldsymbol{\theta}|\mathbf{y})$$

and

$$g(\boldsymbol{\theta}^*) = p(\boldsymbol{\theta}^*).$$

Then the Metropolis-Hastings ratio is given by:

$$R(\boldsymbol{\theta}^{(t)}, \boldsymbol{\theta}^*) = \frac{p(\mathbf{y}|\boldsymbol{\theta}^*)}{p(\mathbf{y}|\boldsymbol{\theta}^{(t)})}$$

Comments:

- If proposal distribution = prior distribution, the Metropolis-Hastings ratio equals the likelihood ratio.
- By definition, the support of the prior covers the support of the target posterior, so the stationary distribution of this chain is the desired posterior.
- More sophisticated algorithms often have better performance.

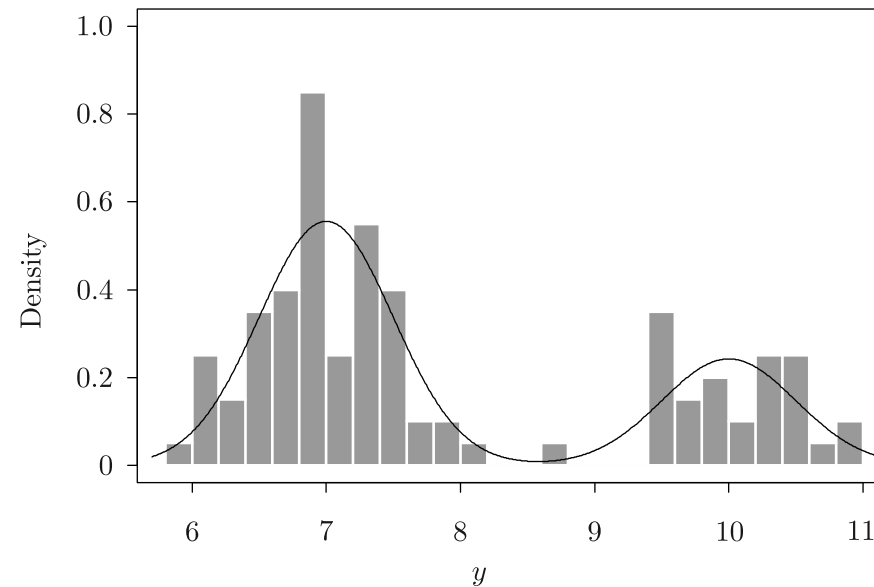
MCMC Example 1: Estimating a mixture parameter

Suppose we have observed data y_1, y_2, \dots, y_{100} sampled independently and identically distributed from the mixture distribution

$$\delta N(7, 0.5^2) + (1 - \delta)N(10, 0.5^2) \quad (94)$$

- Mixture densities are common in real-life applications where, for example, the data may come from more than one population
- We will use MCMC techniques to construct a chain whose stationary distribution equals the posterior density of δ assuming a $\text{Unif}(0,1)$ prior distribution for δ

MCMC Example 1: Estimating a mixture parameter



Histogram of 100 observations simulated from the mixture distribution (94). The data were generated with $\delta = 0.7$, so we should find that the posterior density is concentrated in this area.

MCMC Example 1: Estimating a mixture parameter

Try two different MCMC proposal densities:

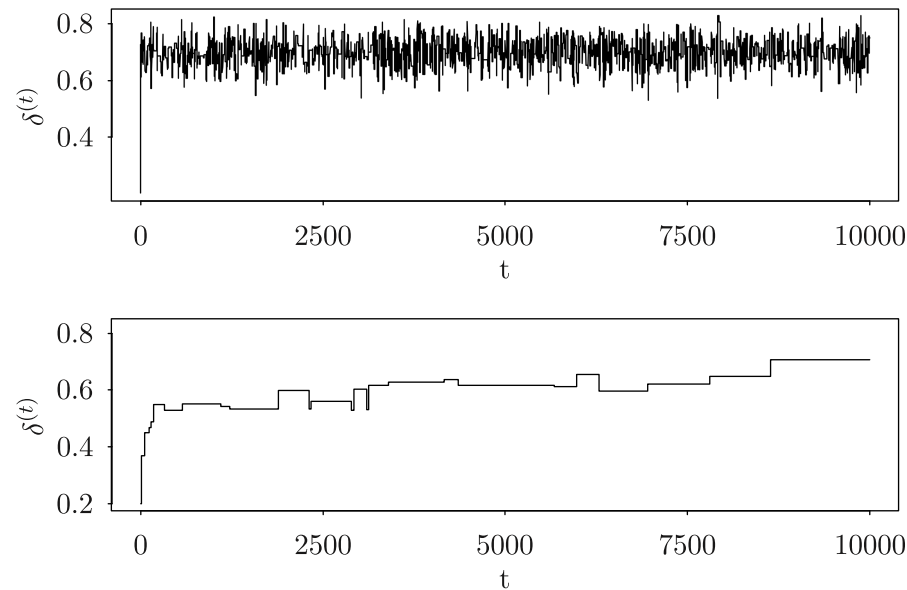
1. Beta(1,1): equivalent to a Unif(0,1) distribution
2. Beta(2,10):
 - skewed right with mean of 0.167
 - values of δ near 0.7 are unlikely to be generated from the proposal distribution

Generated 10,000 iterations of Metropolis-Hastings algorithm for each proposal.

Some output:

- Sample paths:
 - a plot of the chain realizations, $\delta^{(t)}$, against the iteration number, t
 - useful for investigating the behavior of the Markov chain
- Histogram of realizations

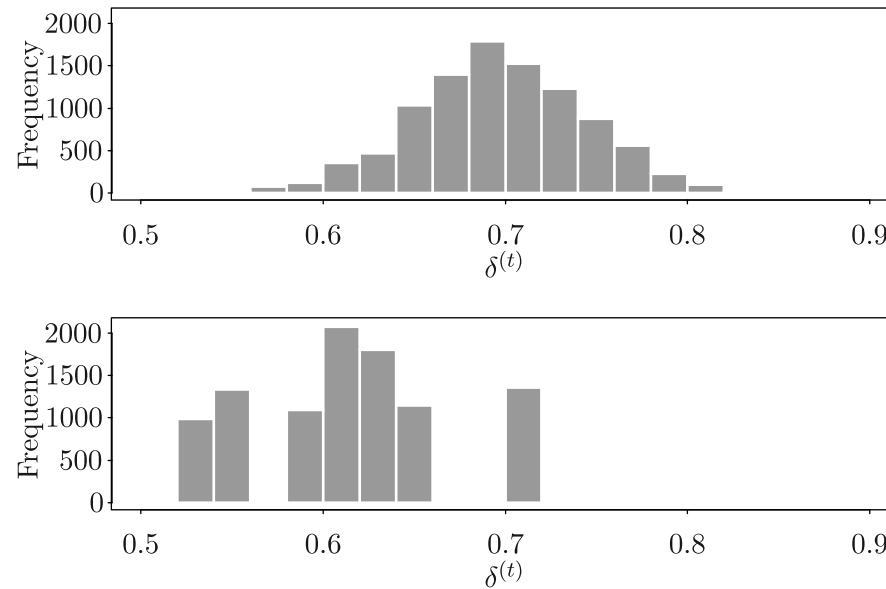
MCMC Example 1: Estimating a mixture parameter



Sample paths for δ from independence chains with proposal densities Beta(1,1) (top) and Beta(2,10) (bottom)

Burn-in period: Since initial iterates retain strong dependence on the starting value, eliminate them in estimation

MCMC Example 1: Estimating a mixture parameter



Histograms of $\delta^{(t)}$ for iterations 201-10,000 of independence chains with proposal densities Beta(1,1) (top) and Beta(2,10) (bottom)

Random walk chains

Random walk chain: Let \mathbf{X}^* be generated by drawing

$$\boldsymbol{\epsilon} \sim h(\boldsymbol{\epsilon})$$

for some density h and then setting

$$\mathbf{X}^* = \mathbf{x}^{(t)} + \boldsymbol{\epsilon}.$$

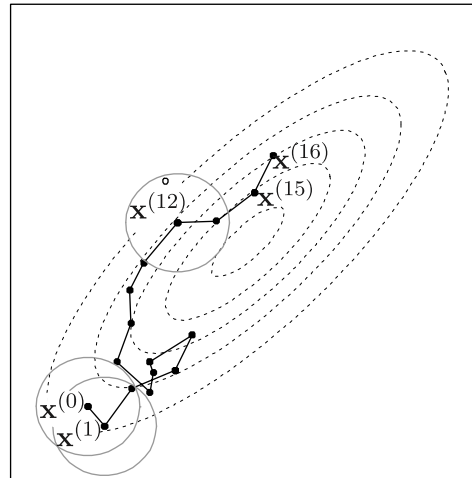
Note that the proposal is then given by

$$g\left(\mathbf{x}^* \mid \mathbf{x}^{(t)}\right) = h\left(\mathbf{x}^* - \mathbf{x}^{(t)}\right)$$

Common choices for h :

- uniform distribution over a ball centered at the origin
- scaled standard normal distribution
- scaled Student's t distribution

Random walk chains



Hypothetical random walk chain for sampling a two dimensional target distribution (dotted contours) using proposed increments sampled uniformly from a disk centered at the current value.

Gibbs Sampling

- Specifically adapted for multidimensional target distributions.
- Goal: still to construct a Markov chain whose stationary distribution—or some marginalization thereof—equals the target distribution, f
- Sequentially samples from univariate conditional distributions, which are often available in closed form

Basic Gibbs sampler

Suppose it is easy to sample from the univariate conditional distributions:

$$X_i | X_1, \dots, X_{i-1}, X_{i+1}, \dots, X_p \sim f(x_i | x_1, \dots, x_{i-1}, x_{i+1}, \dots, x_p)$$

A basic Gibbs sampler:

1. Select starting values $\mathbf{x}^{(0)}$ and set $t = 0$.
2. Generate, in turn,

$$\begin{aligned} X_1^{(t+1)} | \cdot &\sim f\left(x_1 | x_2^{(t)}, \dots, x_p^{(t)}\right) \\ X_2^{(t+1)} | \cdot &\sim f\left(x_2 | x_1^{(t+1)}, x_3^{(t)}, \dots, x_p^{(t)}\right) \\ &\vdots \\ X_{p-1}^{(t+1)} | \cdot &\sim f\left(x_{p-1} | x_1^{(t+1)}, x_2^{(t+1)}, \dots, x_{p-2}^{(t+1)}, x_p^{(t)}\right) \\ X_p^{(t+1)} | \cdot &\sim f\left(x_p | x_1^{(t+1)}, x_2^{(t+1)}, \dots, x_{p-1}^{(t+1)}\right) \end{aligned}$$

where $|\cdot$ denotes conditioning on the most recent updates to all other elements of \mathbf{X} .

3. Increment t and go to step 2.

Notes and enhancements to the Gibbs sampler

- The completion of step 2 for all components of \mathbf{X} is called a *cycle*
- Can reorder the conditional distribution updates with each cycle
- Can block or group elements of \mathbf{X} and update them jointly.

For example when $p = 4$, perhaps it is useful to consider

$$X_2^{(t+1)}, X_3^{(t+1)} \mid \cdot \sim f \left(x_2, x_3 \mid x_1^{(t+1)}, x_4^{(t)} \right)$$

- Can show Gibbs sampling is a Metropolis-Hastings algorithm with time inhomogeneous proposal distributions and with Metropolis-Hasting ratio R always equal to 1.

Hybrid Gibbs sampling

A *hybrid MCMC* algorithm might proceed with the following sequence of updates with $p = 6$:

1. Update $X_1^{(t+1)} \mid \left(x_2^{(t)}, x_3^{(t)}, x_4^{(t)}, x_5^{(t)}, x_6^{(t)} \right)$ with a Gibbs step.
2. Update $\left(X_2^{(t+1)}, X_3^{(t+1)} \right) \mid \left(x_1^{(t+1)}, x_4^{(t)}, x_5^{(t)}, x_6^{(t)} \right)$ with a Metropolis step.
3. Update $X_4^{(t+1)} \mid \left(x_1^{(t+1)}, x_2^{(t+1)}, x_3^{(t+1)}, x_5^{(t)}, x_6^{(t)} \right)$ with a step from a random walk chain.
4. Update $\left(X_5^{(t+1)}, X_6^{(t+1)} \right) \mid \left(x_1^{(t+1)}, x_2^{(t+1)}, x_3^{(t+1)}, x_4^{(t+1)} \right)$ with a Gibbs step.

The Metropolis-Hastings steps within a Gibbs algorithm are typically useful when the univariate conditional density for one or more elements of \mathbf{X} is not available in closed form.

MCMC algorithm implementation

All MCMC algorithms described above have the correct limiting stationary distribution.

Mixing:

- How quickly does the chain forget its starting value?
- How quickly does the chain fully explore the support of the target distribution?
- How far apart in a sequence do observations need to be before they can be considered to be approximately independent?

Stationarity: has the chain has run sufficiently long so that

- it is reasonable to believe that the output adequately represents the target distribution?
- the output can be used reliably for estimation?

Ensuring good mixing and convergence

Strategies and diagnostics

1. Number of chains
2. Simple graphs to assess mixing and convergence
3. Reparameterization
4. Burn-in and run length

1. Number of chains

A difficult problem to diagnose: chain has become stuck in one or more modes of the target distribution.

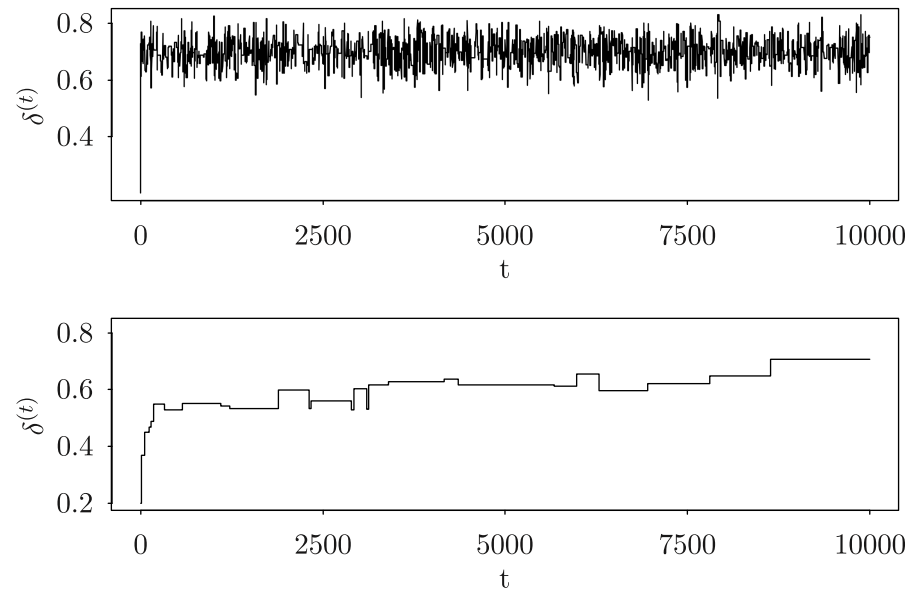
Partial solution: run multiple chains from diverse starting values and then compare the within- and between-chain behavior.

Diagnostics: sample path plots and Gelman-Rubin statistic

2. Simple graphs to assess mixing and convergence

Sample paths: plot of the iteration number versus the realizations of $X^{(t)}$ for $t = 0, 1, \dots$

- also called trace or history plots
- Want to see a very wiggly line that (perhaps) moves away from the starting value.



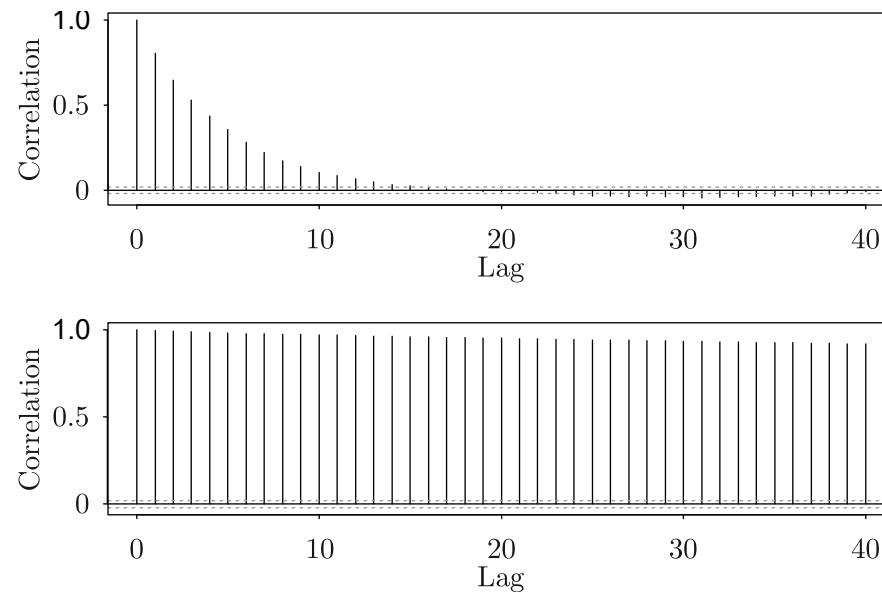
Sample paths for δ from independence chains with proposal densities Beta(1,1) (top) and Beta(2,10) (bottom)

2. Simple graphs to assess mixing and convergence

autocorrelation function (ACF) plot:

- A plot of the lag versus the correlation between iterates at that lag.
- Slow decay in the acf suggests poor mixing.
- Multi-parameter problems:
 - examine cross-correlations between parameters that might be related
 - High cross-correlations may indicate poor mixing of the chain

2. Simple graphs to assess mixing and convergence



Autocorrelation plots for example on estimating a mixture parameter with with proposal densities $\text{Beta}(1,1)$ (top) and $\text{Beta}(2,10)$ (bottom).

3. Reparameterization

- The mixing rates of both the Gibbs sampler and Metropolis-Hastings algorithms can be improved via reparameterization of the model.
- High correlation between the elements of X can lead to slow convergence. Reparameterization of the model can reduce the correlation and thus speed convergence.

Reparameterization approaches are typically adapted for specific models:

- If there are continuous covariates in a linear model, center and scale the covariates to reduce correlations between the parameters in the model.
- Hierarchical centering: parameters in a hierarchical Bayesian model are moved down the hierarchy. Particularly useful in models with random effects.

4. Burn-in and run length

Burn-in

- Initial MCMC iterates will not have exactly the correct marginal distribution, and the dependency on the initial point (or distribution) from which the chain was started may remain strong.
- To reduce the severity of this problem, the first several thousand values from the chain are typically discarded.

Gelman and Rubin statistic, R

- Compares the output of several chains
- If within and between chain variances are similar in magnitude, suggests that the chains are stationary
- Roughly: want $\sqrt{R} < 1.2$
- If not: increase run-length, reparameterize the model, etc.

4. Burn-in and run length: Gelman and Rubin statistic

Run the MCMC algorithm J separate times to produce equal-length chains ($J \geq 2$) with starting values dispersed over the support of the target density.

Let

- L denote the length of each chain after discarding D burn-in iterates.
- The variable of interest is X , and its value at the t th iteration of the j th chain is $x_j^{(t)}$.
- For the j th chain, the D values $x_j^{(0)}, \dots, x_j^{(D-1)}$ are discarded and the L values $x_j^{(D)}, \dots, x_j^{(D+L-1)}$ are retained.

4. Burn-in and run length: Gelman and Rubin statistic

Let

$$\bar{x}_j = \frac{1}{L} \sum_{t=D}^{D+L-1} x_j^{(t)} \text{ and } \bar{x}_{\cdot} = \frac{1}{J} \sum_{j=1}^J \bar{x}_j$$

and define the between-chain variance as

$$B = \frac{L}{J-1} \sum_{j=1}^J (\bar{x}_j - \bar{x}_{\cdot})^2$$

Now if the within-chain variance for the j th chain is

$$s_j^2 = \frac{1}{L-1} \sum_{t=D}^{D+L-1} \left(x_j^{(t)} - \bar{x}_j \right)^2, \text{ then let}$$

$$W = \frac{1}{J} \sum_{j=1}^J s_j^2$$

represent the mean of the J within-chain estimated variances.

Finally, let

$$R = \frac{\frac{L-1}{L}W + \frac{1}{L}B}{W}.$$

4. Burn-in and run length: Gelman and Rubin statistic

- If all the chains are stationary, then both the numerator and the denominator should estimate the marginal variance of X .
- If there are notable differences between the chains, then the numerator will exceed the denominator.
- As $L \rightarrow \infty$, $\sqrt{R} \rightarrow 1$.
- In practice, some authors suggest that $\sqrt{R} < 1.2$ is acceptable.

Remedies:

- If the chosen burn-in period did not yield an acceptable result, then D should be increased, L should be increased, or preferably both.
- If the $x_j^{(t)}$ iterates are transformed so that their distribution is approximately normal, the performance of this diagnostic is improved.
- Reparameterize the model and rerun the chain.

MCMC software

WinBUGS:

- WinBUGS is free software to carry out Gibbs sampling (www.mrc-bsu.cam.ac.uk/bugs/)
- Programming your own MCMC algorithm is the best way to learn what is going on in MCMC, but WinBUGS makes life a lot easier.
- Now available: code to run WinBUGS from within R

Other software available in R and other languages, but WinBUGS is the current popular choice

Output of MCMC procedures

Marginalization: If $\{\mathbf{X}^{(t)}\}$ represents a p -dimensional Markov chain, then the limiting distribution of $\{X_i^{(t)}\}$ is the i th marginal of f . If you are focused only on a property of this marginal, discard the rest of the simulation and analyze the realizations of $X_i^{(t)}$.

Mean: The most commonly used estimator is based on an empirical average. Discard the burn-in; then calculate the desired statistic by taking

$$\frac{1}{L} \sum_{t=D}^{D+L-1} h(\mathbf{X}^{(t)})$$

Standard deviation: sample standard deviation of realizations (after burn-in)

Probability estimates: The probability of any event can be estimated by the frequency of that event in the chain.

Output of MCMC procedures

Monte Carlo standard error: There are several ways to compute this.

Batch Method:

1. Separate the realizations into batches with, say, 50 consecutive iterations in each batch.
2. Compute the mean of each batch.
3. The estimated standard error is the standard deviation of these means divided by the square root of the number of batches.

Density/histogram:

- Histogram of the raw iterates.
- A kernel density estimate of the marginal posterior distributions.

Quantiles: Compute empirical quantiles

Other summary statistics: Be creative!

Practical implementation advice

Some rules of thumb:

- **Mixing:** For a Metropolis algorithm with normal target and proposal distributions, it has been suggested that a user should aim for acceptance rates of approximately
 - 45% for one- or two-dimensional problems
 - 23% for higher-dimensional problems
- **Burn-in lengths:** 0-50,000
- **Chain lengths:** 5,000-5,000,000
 - More samples means more accurate estimates of the posterior distribution.
 - One rule of thumb: the simulation should be run until the Monte Carlo error for each parameter of interest is less than 5% of the sample standard deviation.
 - Chain lengths will continue to increase as computing power grows.

MCMC Example 2: Fur seal pup capture-recapture analysis

- After centuries of severe population reductions due to commercial and subsistence hunting, the abundance of fur seals in New Zealand has been increasing in recent years.
- Our goal is to estimate the number of pups in a fur seal colony using a capture-recapture approach.
- Estimation: can estimate population size and survival rates, etc.
- Consider, for example, high recapture rates suggest that the true population size does not greatly exceed the total number of unique individuals ever captured.

MCMC Example 2: Fur seal pups

Let

- N be the unknown population size to be estimated
- I be the number of census attempts
- $\mathbf{c} = (c_1, \dots, c_I)$ be the total number of captures (including recaptures) at each attempt
- r be the total number of distinct animals captured during the study

We assume that the population is closed during the period of the sampling which means that deaths, births and migrations are inconsequential during this period.

MCMC Example 2: Fur seal pups

- We consider a model with separate, unknown capture probabilities for each census effort, $\alpha = (\alpha_1, \dots, \alpha_I)$.
- This model assumes that all animals are equally catchable on any one capture occasion, but capture probabilities may change over time.

The likelihood for this model is

$$L(N, \boldsymbol{\alpha} | \mathbf{c}, r) \propto \frac{N!}{(N-r)!} \prod_{i=1}^I \alpha_i^{c_i} (1 - \alpha_i)^{N-c_i} \quad (95)$$

This model is sometimes called the $M(t)$ model

MCMC Example 2: Fur seal pups

- In a study conducted on the Otago Peninsula on the South Island of New Zealand, fur seal pups were marked and released during $I = 7$ census attempts during one season.
- It is reasonable to assume the population of pups was closed during the study period.
- $r = \sum_{i=1}^7 m_i = 84$ unique fur seals were observed during the sampling period.

Number of pups captured and recaptured during 7 census efforts in one season

		Census attempt, i						
		1	2	3	4	5	6	7
Number captured	c_i	30	22	29	26	31	32	35
Number newly caught	m_i	30	8	17	7	9	8	5

MCMC Example 2: Fur seal pups, priors

For estimation, one might adopt a hierarchical Bayesian framework where N and α are assumed to be a priori independent with the following priors:

- N : noninformative Jeffreys prior $f(N) \propto 1/N$
- capture probabilities: $f(\alpha_i|\theta_1, \theta_2) = \text{Beta}(\theta_1, \theta_2)$ for $i = 1, \dots, 7$, assumed to be a priori exchangeable

Previous analyses with the $M(t)$ model have indicated sensitivity to the prior distribution for the capture probabilities. To mitigate this sensitivity, we introduce a hyperprior for (θ_1, θ_2) :

$$f(\theta_1, \theta_2) \propto \exp\{-(\theta_1 + \theta_2)/1000\},$$

with (θ_1, θ_2) assumed to be a priori independent of the remaining parameters.

MCMC Example 2: Fur seal pups, Gibbs

Goal:

$$p(N, \alpha_1, \dots, \alpha_7, \theta_1, \theta_2 | \mathbf{c}, r) \propto p(\mathbf{c}, \mathbf{m} | N, \alpha_1, \dots, \alpha_7, \theta_1, \theta_2) \\ p(N) \left[\prod_{i=1}^7 p(\alpha_i | \theta_1, \theta_2) \right] p(\theta_1, \theta_2)$$

A Gibbs sampler can then be constructed by simulating from the conditional posterior distributions

$$N - 84 | \cdot \sim \mathbf{NegBin} \left(84, 1 - \prod_{i=1}^7 (1 - \alpha_i) \right) \\ \alpha_i | \cdot \sim \mathbf{Beta}(c_i + \theta_1, N - c_i + \theta_2) \quad \text{for } i = 1, \dots, 7 \\ \theta_1, \theta_2 | \cdot \sim k \left[\frac{\Gamma(\theta_1 + \theta_2)}{\Gamma(\theta_1)\Gamma(\theta_2)} \right]^7 \prod_{i=1}^7 \alpha_i^{\theta_1} (1 - \alpha_i)^{\theta_2} \exp\{-(\theta_1 + \theta_2)/1000\}$$

where $|\cdot$ denotes conditioning on the remaining parameters and the data, and k is an unknown constant.

MCMC Example 2: Fur seal pups, Gibbs

Problem: it is difficult to produce a chain for (θ_1, θ_2) with adequate mixing and convergence behavior.

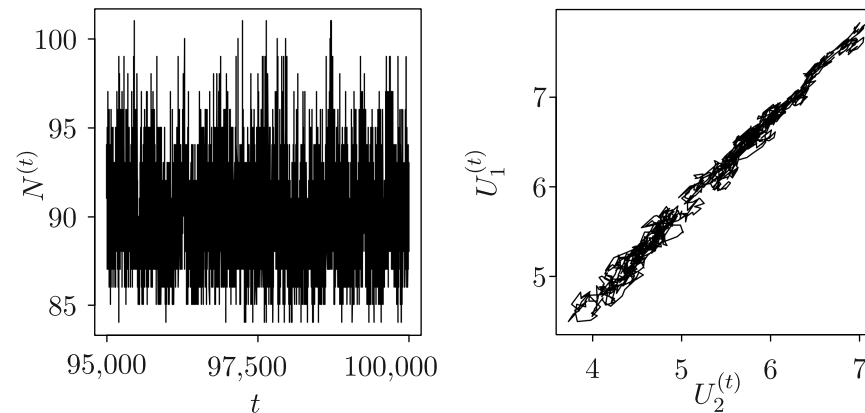
To improve performance, transform (θ_1, θ_2) to

$$\mathbf{U} = (U_1, U_2) = (\log \theta_1, \log \theta_2).$$

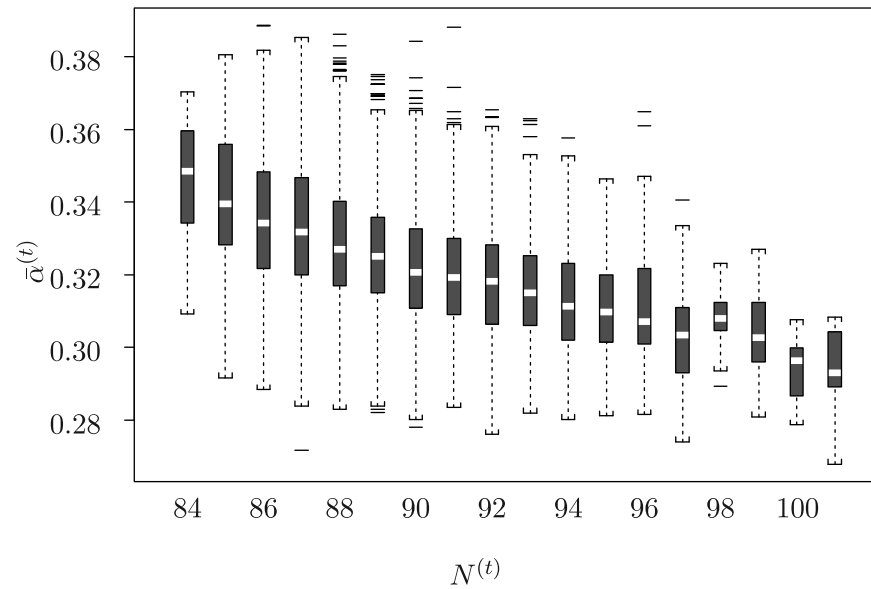
This requires transformation of the conditional densities.

MCMC Example 2: Fur seal pups, Gibbs

- The results below are based on a chain of 100,000 iterations with the first 1,000 iterations discarded for burn-in.
- Based on five runs of 100,000 iterations each, the Gelman and Rubin statistic for N is equal to 1.00047 which suggests the $N^{(t)}$ chain is roughly stationary.
- Sample paths for N (left panel) and U (right panel) for final 5,000 iterations in the seal pup example.

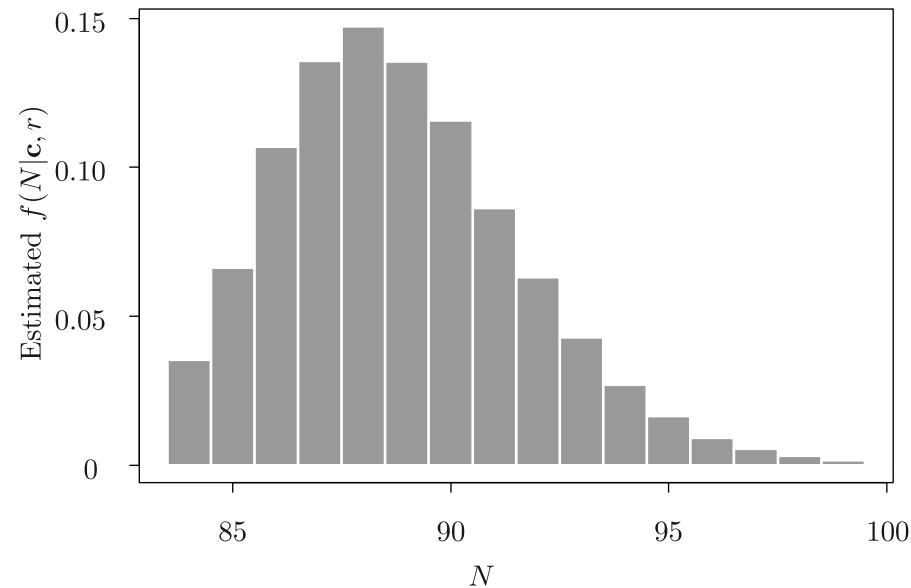


MCMC Example 2: Fur seal pups, Gibbs



Split boxplot of $\bar{\alpha}^{(t)}$ against $N^{(t)}$ for the seal pup example.

MCMC Example 2: Fur seal pups, Gibbs



Estimated marginal posterior probabilities for N for the seal pup example.

The posterior mean of N is 90 with a 95% highest posterior density interval of (84, 95).

Part II: Monte Carlo Integration Methods

6. Introduction to integration

7. Monte Carlo Integration

8. Markov Chain Monte Carlo Methods

9. Advanced MCMC Methods

Slice sampling and other auxiliary variable methods, Reversible jump MCMC, Perfect sampling, Langevin Metropolis-Hastings Algorithm, Hit-and-run algorithm, Multiple-try Metropolis-Hastings Algorithm

10. Integration: Concluding remarks

Slice sampling and other auxiliary variable methods

- In some cases standard MCMC methods can take too long to mix properly to be of practical use
- One remedy: augment the state space of the variable of interest
- Auxiliary variable methods can lead to chains which mix faster and require less tuning than standard MCMC methods

Consider a target distribution f which can be evaluated but not easily sampled. Construct an auxiliary variable algorithm as follows:

1. Augment the state space of \mathbf{X} by the state space of a vector of auxiliary variables, \mathbf{U} .
2. Construct a Markov chain over the joint state space of (\mathbf{X}, \mathbf{U}) having stationary distribution $(\mathbf{X}, \mathbf{U}) \sim f(\mathbf{x}, \mathbf{u})$ which marginalizes to the target $f(\mathbf{x})$.
3. When simulation has been completed, discard \mathbf{U} realizations and base inferences only on the marginal distribution of \mathbf{X} .

Slice sampling

Consider MCMC for a univariate variable $X \sim f(x)$, and suppose that it is impossible to sample directly from f .

- Introducing any univariate auxiliary variable U would allow us to consider a target density for $(X, U) \sim f(x, u)$.
- Writing $f(x, u) = f(x)f(u|x)$ suggests an auxiliary variable Gibbs sampling strategy that alternates between updates for X and U .
- The trick is to choose a U variable that speeds MCMC mixing for X .

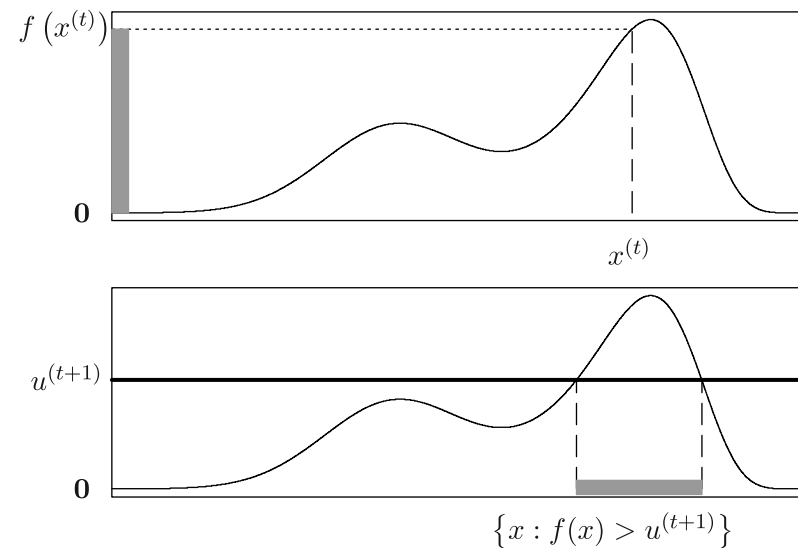
At iteration $t + 1$ of the slice sampler we alternately generate $X^{(t+1)}$ and $U^{(t+1)}$ according to

$$U^{(t+1)}|x^{(t)} \sim \mathbf{Unif}\left(0, f\left(x^{(t)}\right)\right), \quad (96)$$

$$X^{(t+1)}|u^{(t+1)} \sim \mathbf{Unif}\left\{x : f(x) \geq u^{(t+1)}\right\}. \quad (97)$$

The sampling in (97) can be challenging if f is not invertible. One approach to implementing equation (97) is to use rejection sampling.

Illustration of slice sampling

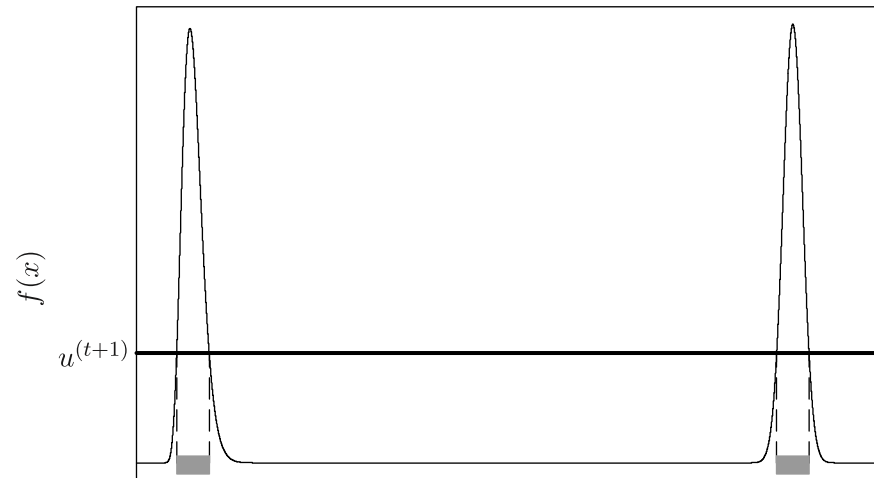


At iteration $t + 1$, the algorithm starts at $x^{(t)}$.

Upper panel: Draw $U^{(t+1)} \sim \text{Unif}(0, f(x^{(t)}))$. In the top panel this corresponds to sampling along the vertical shaded strip.

Lower panel: Draw $X^{(t+1)} | (U^{(t+1)} = u^{(t+1)})$ uniformly from the set of x values for which $f(x) \geq u^{(t+1)}$. In the lower panel this corresponds to sampling along the horizontal shaded strip.

Example: univariate multimodal target distribution



Generate samples from this multimodal target distribution using a

- **Standard Metropolis–Hastings algorithm:** algorithm may find one mode of the distribution, but it may take many iterations to find the other mode.
- **Slice sampler:** The horizontal shaded areas indicate the set defined in (97) from which $X^{(t+1)} | u^{(t+1)}$ is uniformly drawn. Hence the slice sampler will have about a 50% chance of switching modes each iteration. Therefore the slice sampler will mix much better with many fewer iterations required.

Reversible jump MCMC (RJMCMC)

- Standard MCMC algorithms require that the dimensionality of $\mathbf{X}^{(t)}$ and the interpretation of the elements of $\mathbf{X}^{(t)}$ do not change with t
- In many applications, it may be of interest to develop a chain that allows for changes in the dimension of the parameter space from one iteration to the next.
- Examples where RJMCMC methods are useful:
 - model selection
 - selection of the number of components in a mixture distribution
 - knot selection in nonparametric regression

RJMCMC: Notation

Consider constructing a Markov chain to explore a space of candidate models, each of which might be used to fit observed data y . Let

- $\mathcal{M}_1, \dots, \mathcal{M}_K =$ a collection of models under consideration.
- $\boldsymbol{\theta}_m =$ parameters in the m th model.
- $p_m =$ number of parameters in the m th model.
- In the Bayesian paradigm, we may envision random variables $\mathbf{X} = (M, \boldsymbol{\theta}_M)$ which together index the model and parameterize inference for that model.

We may assign prior distributions to these parameters, then seek to simulate from their posterior distribution using a MCMC method for which the t th random draw is $\mathbf{X}^{(t)} = \left(M^{(t)}, \boldsymbol{\theta}_{M^{(t)}}^{(t)} \right)$.

- $\boldsymbol{\theta}_{M^{(t)}}^{(t)} =$ parameters drawn for the model indexed by $M^{(t)}$
- $p_{M^{(t)}} =$ dimension $\boldsymbol{\theta}_{M^{(t)}}^{(t)}$

RJMCMC: Overview

The goal of RJMCMC is to generate samples with joint posterior density

$$f(m, \boldsymbol{\theta}_m | \mathbf{y})$$

This posterior arises from Bayes' theorem via

$$f(m, \boldsymbol{\theta}_m | \mathbf{y}) \propto f(\mathbf{y} | m, \boldsymbol{\theta}_m) f(\boldsymbol{\theta}_m | m) f(m),$$

where

- $f(\mathbf{y} | m, \boldsymbol{\theta}_m)$ denotes the density of the observed data under the m th model and its parameters
- $f(\boldsymbol{\theta}_m | m)$ denotes the prior density for the parameters in the m th model
- A prior weight of $f(m)$ is assigned to the m th model with $\sum_{m=1}^K f(m) = 1$.

RJMCMC: Overview

The posterior factorization

$$f(m, \boldsymbol{\theta}_m | \mathbf{y}) = f(m | \mathbf{y}) f(\boldsymbol{\theta}_m | m, \mathbf{y})$$

suggests two important types of inference:

1. $f(m | \mathbf{y})$ is the posterior probability for the m th model, normalized over all models under consideration
2. $f(\boldsymbol{\theta}_m | m, \mathbf{y})$ is the posterior density of the parameters in the m th model.

RJMCMC enables the construction of an appropriate Markov chain for $\mathbf{X} = (M, \boldsymbol{\theta}_M)$ that jumps between models with parameter spaces of different dimensions.

The key to the RJMCMC algorithm is the introduction of auxiliary random variables at times t and $t + 1$ with dimensions chosen so that the augmented variables (namely \mathbf{X} and the auxiliary variables) at times t and $t + 1$ have equal dimensions.

RJMCMC: Dimension matching

Consider a problem with two models, 1 and 2, such that:

- Model \mathcal{M}_1 has p_1 parameters
- Model \mathcal{M}_2 has p_2 parameters
- $p_2 > p_1$

Moving between models:

- Proposing a move from \mathcal{M}_1 to \mathcal{M}_2 :

Generate θ_2 from an invertible deterministic function of both θ_1 and an independent random component U_1 . We can write

$$\theta_2 = \mathbf{q}_{1,2}(\theta_1, U_1)$$

- Proposing a move from \mathcal{M}_2 to \mathcal{M}_1 :

Carry out via the inverse transformation, $(\theta_1, U_1) = \mathbf{q}_{1,2}^{-1}(\theta_2) = \mathbf{q}_{2,1}(\theta_2)$. Note that $\mathbf{q}_{2,1}$ is an entirely deterministic way to propose θ_1 from a given θ_2 .

RJMCMC: The algorithm

Assume that the chain is currently visiting model $m^{(t)}$, so the chain is in the state $\mathbf{x}^{(t)} = \left(m^{(t)}, \boldsymbol{\theta}_{m^{(t)}}^{(t)}\right)$. The next iteration of the RJMCMC algorithm can be summarized as follows:

1. Sample a candidate model $M^*|m^{(t)}$ from a proposal density with conditional density $g(\cdot|m^{(t)})$. The candidate model requires parameters $\boldsymbol{\theta}_{M^*}$ of dimension p_{M^*} .
2. Given $M^* = m^*$, generate an augmenting variable $\mathbf{U}|(m^{(t)}, \boldsymbol{\theta}_{m^{(t)}}^{(t)}, m^*)$ from a proposal distribution with density $h(\cdot|m^{(t)}, \boldsymbol{\theta}_{m^{(t)}}^{(t)}, m^*)$. Let

$$(\boldsymbol{\theta}_{m^*}^*, \mathbf{U}^*) = \mathbf{q}_{t,*} \left(\boldsymbol{\theta}_{m^{(t)}}^{(t)}, \mathbf{U} \right),$$

where $\mathbf{q}_{t,*}$ is an invertible mapping from $\left(\boldsymbol{\theta}_{m^{(t)}}^{(t)}, \mathbf{U}\right)$ to $(\boldsymbol{\theta}_{M^*}^*, \mathbf{U}^*)$ and the auxiliary variables have dimensions satisfying $p_{m^{(t)}} + p_{\mathbf{U}} = p_{m^*} + p_{\mathbf{U}^*}$

RJCMC: The algorithm

3. For a proposed model, $M^* = m^*$, and the corresponding proposed parameter values $\boldsymbol{\theta}_{m^*}^*$, compute the Metropolis–Hastings ratio given by

$$\frac{f(m^*, \boldsymbol{\theta}_{m^*}^* | \mathbf{y}) g(m^{(t)} | m^*) h(\mathbf{u}^* | m^*, \boldsymbol{\theta}_{m^*}^*, m^{(t)})}{f(m^{(t)}, \boldsymbol{\theta}_{m^{(t)}}^{(t)} | \mathbf{Y}) g(m^* | m^{(t)}) h(\mathbf{u} | m^{(t)}, \boldsymbol{\theta}_{m^{(t)}}^{(t)}, m^*)} |\mathbf{J}(t)|, \quad (98)$$

where

$$\mathbf{J}(t) = \left. \frac{d\mathbf{q}_{t,*}(\boldsymbol{\theta}, \mathbf{u})}{d(\boldsymbol{\theta}, \mathbf{u})} \right|_{(\boldsymbol{\theta}, \mathbf{u}) = (\boldsymbol{\theta}_{m^{(t)}}^{(t)}, \mathbf{U})}. \quad (99)$$

Accept the move to the model M^* with probability equal to the minimum of 1 and the expression in (98). If the proposal is accepted, set $\mathbf{X}^{(t+1)} = (M^*, \boldsymbol{\theta}_{M^*}^*)$. Otherwise, reject the candidate draw and set $\mathbf{X}^{(t+1)} = \mathbf{x}^{(t)}$.

The last term in (98) is the absolute value of the determinant of the Jacobian matrix arising from the change of variables from $(\boldsymbol{\theta}_{m^{(t)}}^{(t)}, \mathbf{U})$ to $(\boldsymbol{\theta}_{m^*}^*, \mathbf{U}^*)$. If $p_{M^{(t)}} = p_{M^*}$, then (98) simplifies to the standard Metropolis–Hastings ratio (92).

4. Discard \mathbf{U} and \mathbf{U}^* . Return to step 1.

RJMCMC: Simple example with two models

Consider a problem with $K = 2$ possible models:

- Model \mathcal{M}_1 has a one-dimensional parameter space $\theta_1 = \alpha$
- Model \mathcal{M}_2 has a two-dimensional parameter space $\theta_2 = (\beta, \gamma)$. Thus $p_1 = 1$ and $p_2 = 2$.

Let $m_1 = 1$ and $m_2 = 2$ (model indexes).

If the chain is currently in state $(1, \theta_1)$ and \mathcal{M}_2 is proposed:

A random variable $U \sim h(u|1, \theta_1, 2)$ is generated from some proposal density h . Let $\beta = \alpha - U$ and $\gamma = \alpha + U$, so $\mathbf{q}_{1,2}(\alpha, u) = (\alpha - u, \alpha + u)$ and $\left| \frac{d\mathbf{q}_{1,2}(\alpha, u)}{d(\alpha, u)} \right| = 2$.

For a proposed move from \mathcal{M}_1 to \mathcal{M}_2 , the Metropolis–Hasting ratio (98) is equal to

$$\frac{f(2, \beta, \gamma | \mathbf{Y}) g(1|2)}{f(1, \alpha | \mathbf{Y}) g(2|1) h(u|1, \theta_1, 2)} \times 2. \quad (100)$$

RJMCMC: Simple example with two models

If the chain is currently in state $(2, \theta_2)$ and \mathcal{M}_1 is proposed:

The inverse mapping is given by $(\alpha, u) = \mathbf{q}_{2,1}(\beta, \gamma) = \left(\frac{\beta+\gamma}{2}, \frac{\beta-\gamma}{2}\right)$. Therefore

$\left|\frac{d\mathbf{q}_{2,1}(\beta, \gamma)}{d(\beta, \gamma)}\right| = \frac{1}{2}$, and U^* is not required to match dimensions. This transition is entirely deterministic, so we replace $h(u^*|2, \theta_2, 1)$ in (98) with 1.

For a proposed move from \mathcal{M}_2 to \mathcal{M}_1 , the Metropolis–Hastings ratio equals the reciprocal of (100)

Challenges to implementing RJMCMC

- The number of dimensions can be enormous, so it can be critical to select an appropriate proposal distribution h and to construct efficient moves between the different dimensions of the model space.
- Diagnosis of convergence is difficult

Other advanced MCMC algorithms

Perfect Sampling

- Generates a chain that has exactly reached the stationary distribution
- Method has promise, but challenges in implementation

Langevin Metropolis-Hastings algorithm

- A random walk with drift. The drift favors proposals that move towards modes of the target distribution
- A random walk with drift can be generated using the proposal

$$\mathbf{X}^* = \mathbf{x}^{(t)} + \mathbf{d}^{(t)} + \sigma \boldsymbol{\epsilon}^{(t)}$$

where

$$\mathbf{d}^{(t)} = \left(\frac{\sigma^2}{2} \right) \frac{\partial \log f(\mathbf{x})}{\partial \mathbf{x}} \Big|_{\mathbf{x}=\mathbf{x}^{(t)}}$$

and $\boldsymbol{\epsilon}^{(t)}$ is a p -dimensional standard normal random variable.

- The scalar σ is a tuning parameter whose fixed value is chosen by the user to control the magnitude of proposed steps.

Other advanced MCMC algorithms

Hit-and-run algorithm

- Proposed move away from $\mathbf{x}^{(t)}$ is generated in two stages:
 1. Choose a direction to move
 2. Choose a distance to move in the chosen direction
- Resulting Markov chain isn't time invariant
- Can be useful when the state space of \mathbf{X} is sharply constrained

Multiple-try Metropolis-Hastings algorithm

Use when chain is slow to converge or trapped in a local mode of f

- Generate a large number of candidates and
- Choose among the candidates in a manner that ensures that the chain retains the correct limiting stationary distribution