



## Examples of Adaptive MCMC

Gareth O. Roberts & Jeffrey S. Rosenthal

To cite this article: Gareth O. Roberts & Jeffrey S. Rosenthal (2009) Examples of Adaptive MCMC, Journal of Computational and Graphical Statistics, 18:2, 349-367, DOI: [10.1198/jcgs.2009.06134](https://doi.org/10.1198/jcgs.2009.06134)

To link to this article: <https://doi.org/10.1198/jcgs.2009.06134>



Published online: 01 Jan 2012.



Submit your article to this journal [↗](#)



Article views: 2105



View related articles [↗](#)



Citing articles: 363 View citing articles [↗](#)

# Examples of Adaptive MCMC

Gareth O. ROBERTS and Jeffrey S. ROSENTHAL

We investigate the use of adaptive MCMC algorithms to automatically tune the Markov chain parameters during a run. Examples include the Adaptive Metropolis (AM) multivariate algorithm of Haario, Saksman, and Tamminen (2001), Metropolis-within-Gibbs algorithms for nonconjugate hierarchical models, regionally adjusted Metropolis algorithms, and logarithmic scalings. Computer simulations indicate that the algorithms perform very well compared to nonadaptive algorithms, even in high dimension.

**Key Words:** Adaption; Convergence; Hierarchical models; Markov chain Monte Carlo; Metropolis algorithm; Metropolis-within-Gibbs; Nonconjugate priors; Non-Markovian.

## 1. INTRODUCTION

MCMC algorithms such as the Metropolis–Hastings algorithm (Metropolis et al. 1953; Hastings 1970) are extremely widely used in statistical inference, to sample from complicated high-dimensional distributions. Tuning of associated parameters such as proposal variances is crucial to achieve efficient mixing, but can also be very difficult.

Adaptive MCMC algorithms attempt to deal with this problem by automatically “learning” better parameter values of Markov chain Monte Carlo algorithms while they run. In this article, we consider a number of examples of such algorithms, including some in high dimensions. We shall see that adaptive MCMC can be very successful at finding *good* parameter values with little user intervention. In our context, good will be defined in terms of some appropriate measure of Markov chain mixing, such as the integrated autocorrelation of a functional of interest.

It is known that adaptive MCMC algorithms will not always preserve stationarity of  $\pi(\cdot)$ ; see, for example, Rosenthal (2004) and proposition 3 of Roberts and Rosenthal (2005). However, they will converge if the adaptations are done at regeneration times (Gilks, Roberts, and Sahu 1998; Brockwell and Kadane 2005), or under various technical condi-

---

Gareth O. Roberts is Professor, Department of Mathematics and Statistics, Fylde College, Lancaster University, Lancaster, LA1 4YF, England (E-mail: [g.o.roberts@lancaster.ac.uk](mailto:g.o.roberts@lancaster.ac.uk)). Jeffrey S. Rosenthal is Professor, Department of Statistics, University of Toronto, Toronto, Ontario M5S 3G3, Canada (E-mail: [jeff@math.toronto.edu](mailto:jeff@math.toronto.edu)).

© 2009 American Statistical Association, Institute of Mathematical Statistics,  
and Interface Foundation of North America

*Journal of Computational and Graphical Statistics*, Volume 18, Number 2, Pages 349–367  
DOI: 10.1198/jcgs.2009.06134

tions about the adaption procedure (Haario, Saksman, and Tamminen 2001; Andrieu and Moulines 2003; Atchadé and Rosenthal 2005; Andrieu and Atchadé 2007).

Roberts and Rosenthal (2005) proved ergodicity of adaptive MCMC under conditions which we find simpler to apply, and which do not require that the adaptive parameters converge. To state their result precisely, suppose the algorithm updates  $X_n$  to  $X_{n+1}$  using the kernel  $P_{\Gamma_n}$ , where each fixed kernel  $P_\gamma$  has stationary distribution  $\pi(\cdot)$ , but where the  $\Gamma_n$  are random indices, chosen iteratively from some collection  $\mathcal{Y}$  based on past algorithm output. Write  $\|\cdot\|$  for total variation distance,  $\mathcal{X}$  for the state space, and  $M_\epsilon(x, \gamma) = \inf\{n \geq 1 : \|P_\gamma^n(x, \cdot) - \pi(\cdot)\| \leq \epsilon\}$  for the convergence time of the kernel  $P_\gamma$  when beginning in state  $x \in \mathcal{X}$ . Then theorem 13 of Roberts and Rosenthal (2005), combined slightly with their corollaries 8 and 9 and theorem 23, guaranteed that  $\lim_{n \rightarrow \infty} \|\mathcal{L}(X_n) - \pi(\cdot)\| = 0$  (asymptotic convergence), and also  $\lim_{n \rightarrow \infty} \frac{1}{n} \sum_{i=1}^n g(X_i) = \pi(g)$  for all bounded  $g : \mathcal{X} \rightarrow \mathbf{R}$  (WLLN), assuming only the *Diminishing Adaptation* condition

$$\lim_{n \rightarrow \infty} \sup_{x \in \mathcal{X}} \|P_{\Gamma_{n+1}}(x, \cdot) - P_{\Gamma_n}(x, \cdot)\| = 0 \quad \text{in probability,} \quad (1.1)$$

and the *Bounded Convergence* condition

$$\{M_\epsilon(X_n, \Gamma_n)\}_{n=0}^\infty \text{ is bounded in probability, } \quad \epsilon > 0. \quad (1.2)$$

Furthermore, they proved that (1.2) is satisfied whenever  $\mathcal{X} \times \mathcal{Y}$  is finite, or is compact in some topology in which either the transition kernels  $P_\gamma$ , or the Metropolis–Hastings proposal kernels  $Q_\gamma$ , have jointly continuous densities. (Condition (1.1) can be ensured directly, by appropriate design of the adaptive algorithm.) A SLLN is precluded because the convergence statements above are only stated “in probability,” whereas CLTs do not necessarily hold because  $\Gamma_n$  does not necessarily converge at all.

Such results provide a “hunting license” to look for useful adaptive MCMC algorithms. In this article, we shall consider a variety of such algorithms. We shall see that they do indeed converge correctly, and often have significantly better mixing properties than comparable nonadaptive algorithms.

We present a collection of examples. For each one, our adaptive strategy steers the algorithm toward a desired operational “optimal” according to some prescribed criterion. Crucially, our approach differs from that of Andrieu and Moulines (2003) and Andrieu and Atchadé (2007) in that, unlike our method, convergence of the adaptive strategy is specifically sought in their approach. Our regularity conditions are thus weaker and easier to verify, though as a result, the results we can demonstrate are necessarily weaker also.

## 2. ADAPTIVE METROPOLIS (AM)

In this section, we consider a version of the Adaptive Metropolis (AM) algorithm of Haario, Saksman, and Tamminen (2001). We begin with a  $d$ -dimensional target distribution  $\pi(\cdot)$ . We perform a Metropolis algorithm with proposal distribution given at iteration  $n$  by  $Q_n(x, \cdot) = N(x, (0.1)^2 I_d/d)$  for  $n \leq 2d$ , whereas for  $n > 2d$ ,

$$Q_n(x, \cdot) = (1 - \beta)N(x, (2.38)^2 \Sigma_n/d) + \beta N(x, (0.1)^2 I_d/d), \quad (2.1)$$

where  $\Sigma_n$  is the current empirical estimate of the covariance structure of the target distribution based on the run so far, and where  $\beta$  is a small positive constant (we take  $\beta = 0.05$ ).

The motivation for (2.1) is as follows. It is known from Roberts, Gelman, and Gilks (1997) and Roberts and Rosenthal (2001) that the proposal  $N(x, (2.38)^2 \Sigma_n/d)$  is optimal in a particular large-dimensional context. Thus, the  $N(x, (2.38)^2 \Sigma_n/d)$  proposal is an effort to approximate this. The mixture of a little bit of nonrandom normal,  $N(x, (0.1)^2 I_d/d)$ , is a “safety measure” to avoid the algorithm getting stuck at problematic (e.g., singular) values of  $\Sigma_n$ .

Because empirical estimates change at the  $n$ th iteration by only  $O(1/n)$ , it follows that (1.1) will be satisfied. Restricting  $\beta > 0$  in (2.1) ensures that (1.2) is satisfied, at least for a large family of target densities which includes all those which are log-concave outside some arbitrary bounded region (see Sec. 8). Hence, this algorithm will indeed converge to  $\pi(\cdot)$  and satisfy the WLLN. (Haario et al. instead let  $Q_n(x, \cdot) = N(x, \Sigma_n + \epsilon I_d)$  for small  $\epsilon$ , to force  $c_1 I_d \leq \Sigma_n \leq c_2 I_d$  for some  $c_1, c_2 > 0$ , which also ensures (1.1) and (1.2) for target distributions with bounded support, but we prefer to avoid this strong assumption.)

To test this algorithm, we let  $\pi(\cdot) = N(0, MM^t)$ , where  $M$  is a  $d \times d$  matrix generated randomly by letting  $\{M_{ij}\}_{i,j=1}^d$  be i.i.d.  $\sim N(0, 1)$ . This ensures that the target covariance matrix  $\Sigma = MM^t$  will be highly erratic, so that sampling from  $\pi(\cdot)$  presents a significant challenge for sampling if the dimension is at all high.

The resulting trace plot of the first coordinate of the Markov chain is presented in Figure 1 for dimensions  $d = 100$  (top left) and  $d = 200$  (top right). In both cases, the Markov chain takes a long time to adapt properly and settle down to rapid mixing. In the early stages, the algorithm vastly underestimates the true stationary variance, thus illustrating the pitfalls of premature diagnoses of MCMC convergence. In the later stages, by contrast, the algorithm has “learned” how to sample from  $\pi(\cdot)$ , and does so much more successfully.

Another way of monitoring the success of this algorithm’s adapting is as follows. Consider a multidimensional random-walk Metropolis algorithm with proposal covariance ma-

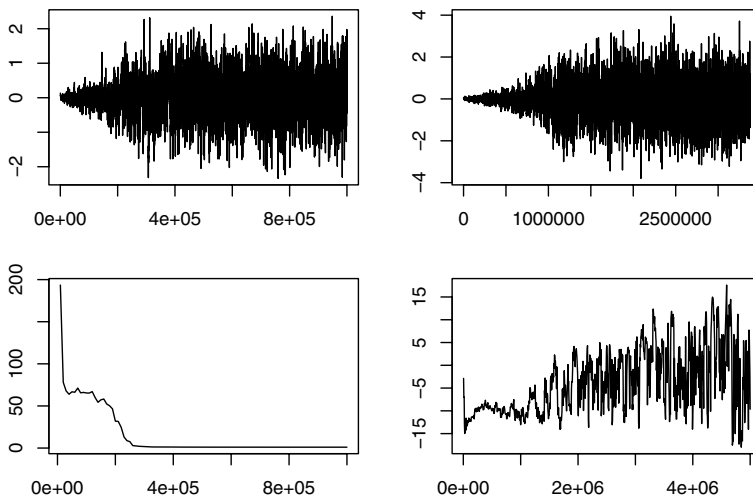


Figure 1. The first coordinate of the AM Markov chain in dimension 100 (top left) and dimension 200 (top right), and the suboptimality factor  $b$  in dimension 100 (bottom left), and the first coordinate of the banana-shaped example (bottom right), all plotted against iteration number.

trix  $(2.38)^2 \Sigma_p/d$ , acting on a normal target distribution with true covariance matrix  $\Sigma$ . Theorem 5 of Roberts and Rosenthal (2001) proved that it is optimal to take  $\Sigma_p = \Sigma$ , and for other  $\Sigma_p$  the mixing rate will be slower than this by a suboptimality factor of

$$b \equiv d \frac{\sum_{i=1}^d \lambda_i^{-2}}{(\sum_{i=1}^d \lambda_i^{-1})^2},$$

where  $\{\lambda_i\}$  are the eigenvalues of the matrix  $\Sigma_p^{1/2} \Sigma^{-1/2}$ . Usually we will have  $b > 1$ , and the closer  $b$  is to 1, the better. The criterion being optimized in AM is therefore  $b^{-1}$ .

So how does the AM algorithm perform by this measure? For the run in dimension 100, the value of this suboptimality coefficient  $b$  begins at the huge value of 193.53, and then eventually decreases toward 1, reaching 1.086 after 500,000 iterations, and 1.024 after 1,000,000 iterations (Figure 1, bottom left). In dimension 200, the value of  $b$  is even more erratic, starting around 183,000 and oscillating wildly before decreasing to about 1.04 after 800,000 iterations.

We conclude from this that the AM algorithm does indeed “learn” about the true target covariance matrix, and converges to an algorithm which samples very (almost optimally) efficiently from  $\pi(\cdot)$ . It is true that it takes many iterations for the algorithm to learn this information (nearly 400,000 iterations in dimension 100, and nearly 2,000,000 in dimension 200). On the other hand, what the algorithm is learning is a  $d \times d$  covariance matrix with many parameters (5,050 parameters in dimension 100, and 20,100 in dimension 200). We feel that this indicates very impressive performance of the AM algorithm in high dimensions.

## 2.1 AN IRREGULARLY SHAPED EXAMPLE

AM can be expected to work well on target densities in which the density contours form roughly elliptical contours. In such examples the global covariance gives a good measure of dependence valid in all parts of the state space. However, it is interesting to see how the approach performs on a more challenging problem with more irregularly shaped contours.

We also applied our full Adaptive Metropolis algorithm to a “banana-shaped” distribution, as proposed by Haario, Saksman, and Tamminen (1999, 2001), with density

$$f_B = f_d \circ \phi_B,$$

where  $f_d$  is the  $d$ -dimensional density of a  $N(\mathbf{0}, \text{diag}(100, 1, 1, \dots, 1))$  distribution, and where  $\phi_B(x_1, \dots, x_d) = (x_1, x_2 + Bx_1^2 - 100B, x_3, \dots, x_d)$  with  $B > 0$  the “bananicity” constant. So,

$$f_B(x_1, \dots, x_d) \propto \exp[-x_1^2/200 - \frac{1}{2}(x_2 + Bx_1^2 - 100B)^2 - \frac{1}{2}(x_3^2 + x_4^2 + \dots + x_d^2)].$$

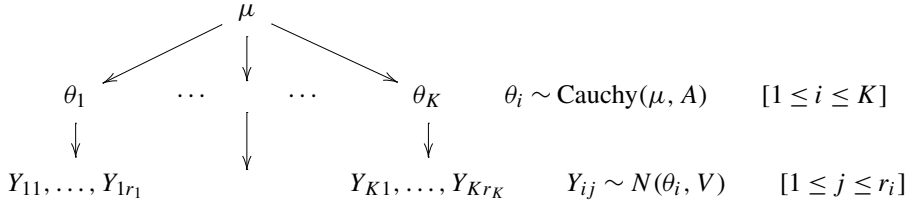
Specifically, we take dimension  $d = 20$ , and take  $B = 0.1$ , and run the algorithm for 5,000,000 iterations. A trace plot of the first coordinate is given in Figure 1 (bottom right).

It is clear that the adaptation has improved mixing here. However, mixing is still very poor after 5,000,000 iterations, which is to be expected given that the Metropolis method struggles to traverse this distributions support. Although the AM algorithm attempts to

move around as effectively as it can, classes of algorithms which can adjust the covariance of the proposal distribution according to the current state of the algorithm should be required. This in part motivates some of the methods we shall introduce in later sections.

### 3. ADAPTIVE METROPOLIS-WITHIN-GIBBS

Consider the following statistical model:



with priors  $N(0, 1)$  on  $\mu$ , and  $IG(1, 1)$  on  $A$  and  $V$ . Here  $\{Y_{ij}\}$  are observed data,  $IG(a, b)$  is the inverse gamma distribution with density proportional to  $e^{-b/x} x^{-(a+1)}$ , and  $\text{Cauchy}(m, s)$  is a translated and scaled Cauchy distribution with density proportional to  $[1 + ((x - m)/s)^2]^{-1}$ . This model gives rise to a posterior distribution  $\pi(\cdot)$  on the  $(K + 3)$ -dimensional vector  $(A, V, \mu, \theta_1, \dots, \theta_K)$ , conditional on the observed data  $\{Y_{ij}\}$ .

We take  $K = 500$ , and let the  $r_i$  vary between 5 and 500. The resulting model is too complicated for analytic computation, and far too high-dimensional for numerical integration. Furthermore, the presence of the Cauchy (as opposed to Normal) distribution destroys conjugacy, and thus makes a classical Gibbs sampler (as in Gelfand and Smith 1990) infeasible. Instead, a Metropolis-within-Gibbs algorithm (Metropolis et al. 1953; Tierney 1994) seems appropriate.

Such an algorithm might proceed as follows. We consider each of the 503 variables in turn. For each, we propose updating its value by adding a  $N(0, \sigma^2)$  increment. That proposal is then accepted or rejected according to the usual Metropolis ratio. This process is repeated many times, allowing the variables to hopefully converge in distribution to  $\pi(\cdot)$ . But how should  $\sigma^2$  be chosen? Should it be different for different variables? How can we feasibly determine appropriate scalings in such high dimension?

To answer these questions, an adaptive algorithm can be used. We proceed as follows. For each of the variables  $i$  [ $1 \leq i \leq K + 3$ ], we create an associated variable  $ls_i$  giving the logarithm of the standard deviation to be used when proposing a normal increment to variable  $i$ . We begin with  $ls_i = 0$  for all  $i$  (corresponding to unit proposal variance). After the  $n$ th “batch” of 50 iterations, we update each  $ls_i$  by adding or subtracting an adaption amount  $\delta(n)$ . The adapting attempts to make the acceptance rate of proposals for variable  $i$  as close as possible to 0.44 (which is optimal for one-dimensional proposals in certain settings; cf. Roberts, Gelman, and Gilks 1997; Roberts and Rosenthal 2001). Specifically, we increase  $ls_i$  by  $\delta(n)$  if the fraction of acceptances of variable  $i$  was more than 0.44 on the  $n$ th batch, or decrease  $ls_i$  by  $\delta(n)$  if it was less.

Condition (1.1) is satisfied provided  $\delta(n) \rightarrow 0$ ; we take  $\delta(n) = \min(0.01, n^{-1/2})$ . Our approach is to specify a global maximal parameter value  $M < \infty$ , and restrict each  $ls_i$  to

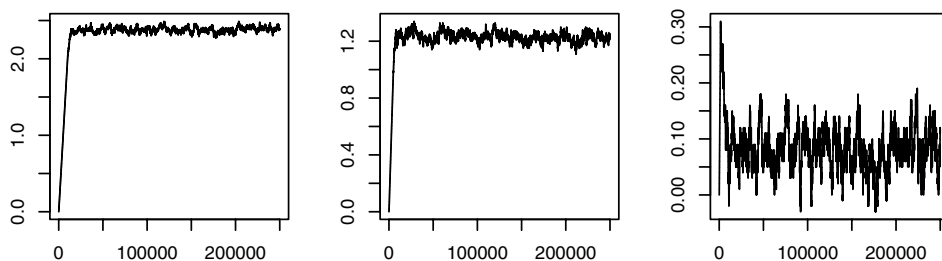


Figure 2. The log proposal standard deviations corresponding to the Metropolis-within-Gibbs variables  $\theta_1$  (left),  $\theta_2$  (middle), and  $\theta_3$  (right), plotted against batch number.

the interval  $[-M, M]$ . For a large class of target densities (which includes all those which are log-concave outside an arbitrary bounded region) this ensures (1.2) holds. In practice, the  $ls_i$  stabilize nicely so the bound on  $M$  is not actually needed.

To test this adaptive algorithm, we generate independent test data  $Y_{ij} \sim N(i - 1, 10^2)$ , for  $1 \leq i \leq 500$  and  $1 \leq j \leq r_i$ . For such data, our simulations show that the scaling variables quickly settle down near “good” values where acceptance rates are roughly 0.44. Indeed, for the location variables  $\theta_1$ ,  $\theta_2$ , and  $\theta_3$ , the corresponding  $ls$  variables converge to values near 2.4, 1.2, and 0.1, respectively (Figure 2). So the algorithm appears to be converging well.

Just how good are the values chosen? The following table presents the integrated autocorrelation times (ACT) and average squared jumping distances (after discarding the first fifth of the run as burn-in), for both the adaptive algorithm and the corresponding “fixed” algorithm where each  $ls_i$  is fixed at 0:

Variable	$r_i$	Algorithm	ACT	Avr sq dist
$\theta_1$	5	Adaptive	2.59	14.932
$\theta_1$	5	Fixed	31.69	0.863
$\theta_2$	50	Adaptive	2.72	1.508
$\theta_2$	50	Fixed	7.33	0.581
$\theta_3$	500	Adaptive	2.72	0.150
$\theta_3$	500	Fixed	2.67	0.147

This table shows that, when comparing adaptive to fixed algorithms, for variables  $\theta_1$  and  $\theta_2$ , the autocorrelation times are significantly smaller (better) and the average squared jumping distances are significantly larger (better). Thus, adapting has significantly improved the MCMC algorithm, by automatically choosing appropriate proposal scalings separately for each coordinate. For variable  $\theta_3$  the performance of the two algorithms is virtually identical, which is not surprising because (Figure 2 (right)) the optimal log proposal standard deviation happens to be very close to 0 in that case.

In summary, this adaptive algorithm appears to correctly scale the proposal standard deviations, leading to a Metropolis-within-Gibbs algorithm which mixes much faster than

a naive one with unit proposal scalings. Coordinates are improved wherever possible, and are left about the same when they happen to already be optimal. This works even in high dimensions, and does not require any direct user intervention or high-dimensional insight. This algorithm has recently been applied to a statistical genetics problem (Turro et al. 2007).

### 3.1 A COMPARISON WITH SCAM

A different component-wise adaptive scaling method, the Single Component Adaptive Metropolis (SCAM) algorithm, was presented in Haario, Saksman, and Tamminen (2005). That algorithm, which resembles the Adaptive Metropolis algorithm of Haario, Saksman, and Tamminen (2001), is very interesting and promising, but differs significantly from ours because the SCAM adapting is done based on the empirical variance of each component based on the run so far.

For comparative purposes, we also ran the SCAM algorithm of Haario, Saksman, and Tamminen (2005) on the same example as that above for adaptive Metropolis-within-Gibbs. The SCAM algorithm uses the proposal distribution  $Y_n^i \sim N(X_{n-1}^i, v_n^i)$  for the  $i$ th coordinate, where

$$v_n^i = \begin{cases} 5^2, & n \leq 10 \\ (2.4)^2(g_n^i + 0.05), & n \geq 11. \end{cases}$$

Here  $g_n^i$  is the sample variance of  $X_0^{(i)}, X_1^{(i)}, \dots, X_{n-1}^{(i)}$ . (Intuitively, for  $n \geq 11$ ,  $v_n^i$  attempts to mimic an “optimal” one-dimensional variance  $(2.38)^2 \text{var}_\pi(X_i)$  similar to what was discussed earlier; the published version of SCAM omits the square in “2.4” but we assume the above is what was intended.) Writing  $\bar{x}_n^i = \frac{1}{n} \sum_{j=0}^{n-1} x_j^{(i)}$ , we see (cf. Haario, Saksman, and Tamminen 2005) that we can use the recursive equations

$$\bar{x}_n^i = \frac{n-1}{n} \bar{x}_{n-1}^i + \frac{1}{n} x_{n-1}^i$$

and

$$g_n^i = \frac{n-2}{n-1} g_{n-1}^i + (\bar{x}_{n-1}^i)^2 + \frac{1}{n-1} (x_{n-1}^i)^2 - \frac{n}{n-1} (\bar{x}_n^i)^2.$$

We again consider the first three coordinates, as above. The graphs of their proposal variances (again on a log scale, for consistency with the above) are presented in Figure 3.

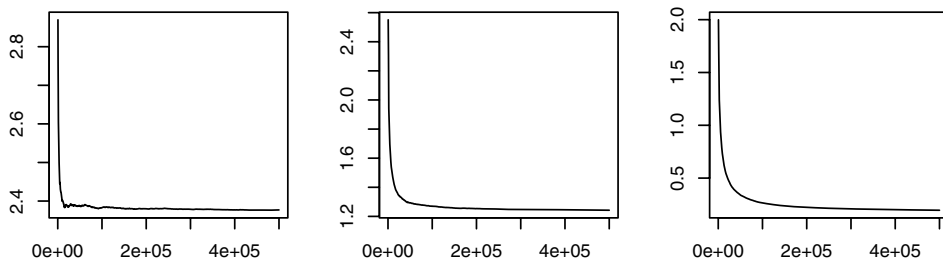


Figure 3. The log proposal standard deviations corresponding to the SCAM variables  $\theta_1$  (left),  $\theta_2$  (middle), and  $\theta_3$  (right), each plotted against iteration number.



We also compute the mean  $\log \sigma$ , ACT, and average squared jumping distance:

Variable	$r_i$	Algorithm	$\log(\sigma)$	ACT	Avr sq dist
$\theta_1$	5	Adaptive	2.35	2.59	14.932
$\theta_1$	5	Fixed	0	31.69	0.863
$\theta_1$	5	SCAM	2.38	2.77	14.951
$\theta_2$	50	Adaptive	1.21	2.72	1.508
$\theta_2$	50	Fixed	0	7.33	0.581
$\theta_2$	50	SCAM	1.27	2.77	1.486
$\theta_3$	500	Adaptive	0.08	2.72	0.150
$\theta_3$	500	Fixed	0	2.67	0.147
$\theta_3$	500	SCAM	0.26	2.77	0.145

The table shows that the results of SCAM are comparable to those of our adaptive Metropolis-within-Gibbs algorithm. In this case, they were virtually identical for  $\theta_1$ , and just *slightly* worse for  $\theta_2$  and  $\theta_3$ . As for choice of proposal variance  $\sigma^2$ , there are some differences, with the SCAM choices generally larger than those for our algorithm. Overall, we feel that both of these algorithms are useful approaches to high-dimensional adaptive MCMC, and both should be kept in the applied user's arsenal.

#### 4. STATE-DEPENDENT SCALING

We next consider examples of full-dimensional Metropolis–Hastings algorithms, where the proposal distribution is given by  $Q(x, \cdot) = N(x, \sigma_x^2)$ , that is, such that the proposal variance depends on the current state  $x \in \mathcal{X}$ . For such an algorithm, according to the usual Metropolis–Hastings formula (Hastings 1970), a proposal from  $x$  to  $y$  is accepted with probability

$$\alpha(x, y) = \min \left[ 1, \frac{\pi(y)}{\pi(x)} (\sigma_x / \sigma_y)^d \exp \left( -\frac{1}{2} (x - y)^2 (\sigma_y^{-2} - \sigma_x^{-2}) \right) \right]. \quad (4.1)$$

As a first case, we let  $\mathcal{X} = \mathbf{R}$ , and  $\pi(\cdot) = N(0, 1)$ . We consider proposal kernels of the form

$$Q_{a,b}(x, \cdot) = N \left( x, e^a \left( \frac{1 + |x|}{\exp(\hat{\pi})} \right)^b \right),$$

where  $\hat{\pi}$  is our current empirical estimate of  $\pi(g)$  where  $g(x) = \log(1 + |x|)$ . (We divide by  $\exp(\hat{\pi})$  to make the choices of  $a$  and  $b$  “orthogonal” in some sense.) After the  $n$ th batch of 100 iterations, we update  $a$  by adding or subtracting  $\delta(n)$  in an effort to, again, make the acceptance rate as close as possible to 0.44. We also add or subtract  $\delta(n)$  to  $b$  to make the acceptance rates,  $acc_-$  and  $acc_+$  respectively, in the regions  $A_- = \{x \in \mathcal{X} : \log(1 + |x|) > \hat{\pi}\}$  and  $A_+ = \{x \in \mathcal{X} : \log(1 + |x|) \leq \hat{\pi}\}$  as equal as possible. This then increases the proposal variance in the region where acceptance rates are highest (thus lowering the acceptance rate) and correspondingly increases the acceptance rate where the acceptance rate is lowest. The criterion being minimized here is therefore  $acc_+^2 + acc_-^2$ .

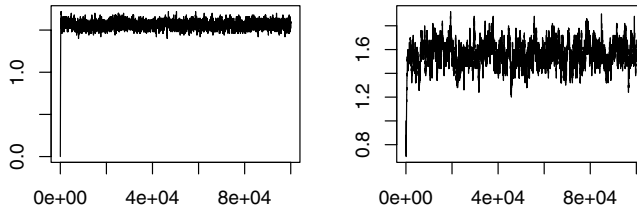


Figure 4. The tuning parameters  $a$  (left) and  $b$  (right) in the State-Dependent Scaling example, plotted against batch number, showing quick approach to “good” values near 1.5 and 1.6, respectively, but with continuing oscillation.

As in previous examples, condition (1.1) is automatically satisfied, at least if we insist on  $\delta(n) \rightarrow 0$ . For us to be able to demonstrate (1.2), however (at least for a particular family of target densities), we shall impose an extra condition, requiring that  $a$  and  $b$  be constrained within  $[-M, M]$  for some global parameter  $M < \infty$ .

So how does this algorithm perform in practice? Empirical expected values quickly converge to their true values, showing excellent mixing. Furthermore, the tuning parameters  $a$  and  $b$  quickly find their “good” values (Figure 4), though they do continue to oscillate due to the extremely slow rate at which  $\delta(n) \rightarrow 0$ .

To determine how well the adaptive algorithm is performing, we compare its integrated autocorrelation time and average squared jumping distance to corresponding nonadaptive algorithms, having either fixed constant variance  $\sigma^2$  [including the optimal constant value,  $(2.38)^2$ ], and to the corresponding variable-variance algorithm. The results are as follows:

Algorithm	Acceptance rate	ACT	Avr sq dist
Adaptive (as above)	0.456	2.63	0.769
$\sigma^2 = \exp(-5)$	0.973	49.92	0.006
$\sigma^2 = \exp(-1)$	0.813	8.95	0.234
$\sigma^2 = 1$	0.704	4.67	0.450
$\sigma^2 = (2.38)^2$	0.445	2.68	0.748
$\sigma^2 = \exp(5)$	0.237	7.22	0.305
$\sigma_x^2 = e^{1.5 \left( \frac{1+ x }{0.534822} \right)^{1.6}}$	0.456	2.58	0.778

We see that our adaptive scheme is much better than arbitrarily chosen fixed-variance algorithms, slightly better than the optimally chosen fixed-variance algorithm (chosen by an ad hoc search for maximizing Average Square Jumping Distance, and given on the fifth line), and nearly as good as an ideally chosen variable- $\sigma^2$  scheme chosen using a similar maximization of average squared jumping distance on a grid of possible  $(a, b)$  values (bottom line). The results are quite impressive, because we did not do any manual tuning of our algorithm at all other than telling the computer to seek a 0.44 acceptance rate.

Although these functional forms of  $\sigma_x^2$  seem promising, it is not clear how to generalize them to higher dimensional problems. Instead, we next consider a different algorithm in which the  $\sigma_x^2$  are piecewise constant over various regions of the state space.

## 5. REGIONAL ADAPTIVE METROPOLIS ALGORITHM (RAMA)

The Regional Adaptive Metropolis Algorithm (RAMA) begins by partitioning the state space  $\mathcal{X}$  into a finite number of disjoint regions:  $\mathcal{X} = \mathcal{X}_1 \dot{\cup} \cdots \dot{\cup} \mathcal{X}_m$ . The algorithm then proceeds by running a Metropolis algorithm with proposal  $Q(x, \cdot) = N(x, \exp(2a_i))$  whenever  $x \in \mathcal{X}_i$ . Thus, if  $x \in \mathcal{X}_i$  and  $y \in \mathcal{X}_j$ , then  $\sigma_x^2 = e^{2a_i}$  and  $\sigma_y^2 = e^{2a_j}$ , and it follows from (4.1) that a proposal from  $x$  to  $y$  is accepted with probability

$$\alpha(x, y) = \min \left[ 1, \frac{\pi(y)}{\pi(x)} \exp \left( d(a_i - a_j) - \frac{1}{2}(x - y)^2 [\exp(-2a_j) - \exp(-2a_i)] \right) \right].$$

The adaptations proceed as follows, in an effort to make the acceptance probability close to 0.234 in each region. (Such an acceptance rate is optimal in certain high-dimensional settings—see Roberts, Gelman, and Gilks 1997; Roberts and Rosenthal 1998a, 2001; Bédard 2006a, 2006b—and we envisage that typically  $\mathcal{X}_i$  would be a space of the same dimension as  $\mathcal{X}$ .) For  $1 \leq i \leq d$ , the parameter  $a_i$  is updated by, after the  $n$ th batch of 100 iterations, considering the fraction of acceptances of those proposals which originated from  $\mathcal{X}_i$ . If that fraction is less than 0.234, then  $a_i$  is decreased by  $\delta(n)$ , whereas if it is more, then  $a_i$  is increased by  $\delta(n)$ . Then, if  $a_i > M$  we set  $a_i = M$ , whereas if  $a_i < -M$  we set  $a_i = -M$ . Finally, if there were no proposals from  $\mathcal{X}_i$  during the entire batch, then  $a_i$  is left unchanged. Thus the algorithm attempts to minimize  $\sum_{i=1}^m acc_i^2$  where  $acc_i$  represents the acceptance rate for moves starting in the region  $\mathcal{X}_i$ .

Provided that  $\delta(n) \rightarrow 0$ , condition (1.1) will trivially hold. Moreover, if we assume that  $M < \infty$ , then it is natural to demonstrate (1.2) again by using a simultaneous drift condition. Such an argument will require some conditions on the target density, but is easy to demonstrate for log-concave densities such as the example below. See Section 8 for further discussion.

For a first example, we let  $\mathcal{X} = \mathbf{R}^d$ , and  $\pi(\cdot) = N(0, I_d)$ . We consider proposal kernels of the form

$$Q_{a,b}(x, \cdot) = N(x, e^{2a} \mathbf{1}_{\|x\|^2 \leq d} + e^{2b} \mathbf{1}_{\|x\|^2 > d}).$$

Once every 100 iterations, we update  $a$  by adding or subtracting  $\delta(n)$  to make the acceptance rate in the region  $\{\|x\|^2 \leq d\}$  as close as possible to 0.234. We also add or subtract  $\delta(n)$  to  $b$  to make the acceptance rate in the region  $\{\|x\|^2 > d\}$  as close as possible to 0.234. We again restrict  $a$  and  $b$  to some  $[-M, M]$ . (We take  $\delta(n) = \min(0.01, n^{-1/2}) \equiv 0.01$  and  $M = 100$ .) We choose dimension  $d = 10$ , and begin with  $a = b = 0$ .

How well does it work? The tuning parameters  $a$  and  $b$  quickly migrate toward their “good” values of  $-0.3$  and  $-0.13$ , respectively, but they continue to oscillate somewhat around these values (Figure 5).

How good are the values of  $a$  and  $b$  found by the computer? The following table gives comparisons of the integrated autocorrelation time and average squared jumping distance

for various choices of  $a$  and  $b$ :

$a, b$	ACT	Avr sq dist
Adaptive (as above)	15.54	0.1246
$-0.3, -0.13$	15.07	0.1258
$-0.3, 0.0$	15.44	0.1213
$0.0, -0.13$	17.04	0.1118
$0.0, 0.0$	17.037	0.1100

The table indicates that the adaptive algorithm (top line) is quite competitive with the corresponding fixed-parameter choice (second line), which in turn has smaller integrated autocorrelation time, and larger average squared jumping distance, than any of the other choices of  $a$  and  $b$ . This indicates that the computer has again succeeded in finding good values for the tuning parameters.

Next, we consider the following statistical model related to James–Stein estimators, as studied in, for example, Rosenthal (1996):

$$\begin{array}{ccccc}
 & & \mu & & \\
 & \swarrow & \downarrow & \searrow & \\
 \theta_1 & \cdots & & \cdots & \theta_K \\
 \downarrow & & & & \downarrow \\
 Y_1 & \cdots & & \cdots & Y_K
 \end{array}
 \quad
 \begin{array}{l}
 \theta_i \sim N(\mu, A) \quad [1 \leq i \leq K] \\
 Y_i \sim N(\theta_i, A) \quad [1 \leq i \leq K]
 \end{array}$$

Here the  $\{Y_i\}$  are observed data. We use the prior distributions  $\mu \sim N(\mu_0, \sigma_0^2)$  and  $A \sim IG(a_1, b_1)$ , and replace  $V$  by its (fixed) empirical Bayes estimate. We let  $\pi(\cdot)$  be the resulting posterior distribution for  $(A, \mu, \theta_1, \dots, \theta_K)$ , on the  $(K + 2)$ -dimensional state space  $\mathcal{X} = [0, \infty) \times \mathbf{R}^{K+1}$ . The density of  $\pi(\cdot)$ , with respect to Lebesgue measure, is then given by

$$\begin{aligned}
 f(A, \mu, \theta_1, \dots, \theta_K) &= N(\mu_0, \sigma_0^2; \mu) IG(a_1, b_1; A) \times \prod_{i=1}^K [N(\mu, A; \theta_i) N(\theta_i, V; Y_i)] \\
 &\propto \exp(-(\mu - \mu_0)^2 / 2\sigma_0^2) \exp(-b_1 / A) / A^{a_1+1} \\
 &\quad \times \prod_{i=1}^K [A^{-1/2} \exp(-(\theta_i - \mu)^2 / 2A) V^{-1/2} \exp(-(Y_i - \theta_i)^2 / 2V)].
 \end{aligned}$$

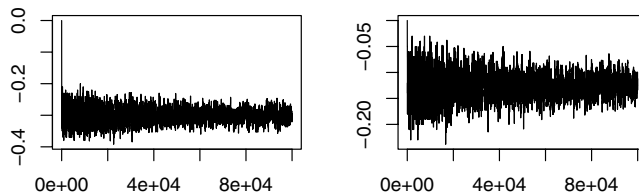


Figure 5. The tuning parameters  $a$  (left) and  $b$  (right) in the Normal RAMA example, each plotted against batch number.

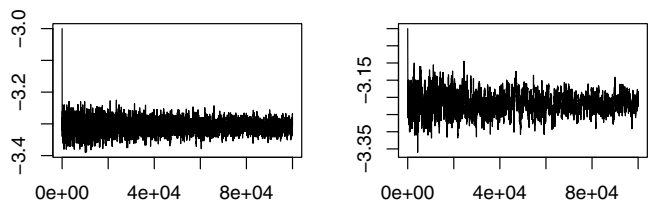


Figure 6. The tuning parameters  $a$  (left) and  $b$  (right) in the James–Stein RAMA example, plotted against batch number.

For a numerical example, we let  $K = 18$ , and let  $Y_1, \dots, Y_{18}$  be the (real) baseball data of table 1 of Morris (1983) (see also Efron and Morris 1975). Thus,  $\mathcal{X} \subseteq \mathbf{R}^{20}$ . We choose the prior parameters as  $\mu_0 = 0$ ,  $\sigma_0^2 = 1$ ,  $a_1 = -1$ , and  $b_1 = 2$ .

We again perform the RAMA algorithm. Specifically, after the  $n$ th batch of 100 iterations, we update  $a$  by adding or subtracting  $\delta(n)$  to make the acceptance rate in the region  $\{\sum_i (\theta_1 - \mu_0)^2 \leq 0.15\}$  as close as possible to 0.234. We also add or subtract  $\delta(n)$  to  $b$  to make the acceptance rate in the region  $\{\sum_i (\theta_1 - \mu_0)^2 > 0.15\}$  as close as possible to 0.234.

The simulations again show good mixing, and rapid convergence of functional averages to their true posterior means. Furthermore, the adaptive parameters  $a$  and  $b$  quickly settle down to near  $-3.3$  and  $-3.2$ , respectively (Figure 6).

How good are the values of the tuning parameters chosen? We again compare integrated autocorrelation times and average squared jumping distances, as follows (acceptance rates are also shown):

$a, b$	Acc rate	ACT	Avr sq dist $\times 10^4$
Adaptive (as above)	0.228	31.60	2.756
$-3.3, -3.2$	0.194	25.75	2.793
$-2.3, -2.3$	0.003	50.67	0.192
$-4.3, -4.3$	0.655	38.92	1.168
$-3.3, -4.3$	0.647	36.91	1.153
$-4.3, -3.3$	0.281	38.04	2.407
$-0.6, -0.6$	$2.5 \times 10^{-5}$	53.97	0.010

We again see that the adaptive algorithm (top line) is quite competitive with the corresponding fixed-parameter choice (second line), which in turn is better than any of the other choices of  $a$  and  $b$ . This shows that, once again, the adaptive algorithm has automatically chosen good values of the MCMC tuning parameters, without requiring user intervention.

**REMARKS**

1. In our simulations, the condition  $M < \infty$  has never been necessary, because RAMA has never tried to push any of the  $\{a_j\}$  toward unbounded values. Indeed, we conjecture that under appropriate regularity assumptions (e.g., if the densities are jointly

continuous), condition (1.2) will be satisfied automatically due to drifting of the parameters  $a_i$  back to reasonable values due to the adaptive process (cf. Roberts and Rosenthal 2005, corollary 14).

2. If some value  $a_j$  is much too large, then  $\alpha(x, y)$  may be very small for all  $y \in \mathcal{X}_j$  and  $x \notin \mathcal{X}_j$ . This means that the region  $\mathcal{X}_j$  may virtually never be entered, so that  $a_j$  will remain virtually constant, leading to isolation of  $\mathcal{X}_j$  and thus very poor convergence. Hence, it is important with RAMA to begin with sufficiently small values of the  $\{a_i\}$ . Alternatively, it might be wise to decrease each  $a_i$  slightly (rather than leaving it unchanged) after each batch in which there were no proposals from  $\mathcal{X}_i$ .
3. The version of RAMA presented here requires that the user specify the regions  $\{\mathcal{X}_i\}_{i=1}^m$  by hand. However, it may also be possible to have the computer automatically select appropriate regions, by, for example, doing a preliminary run with fixed proposal variance, and then grouping together state-space subsets which appear to have similar acceptance rates.
4. Roberts, Gelman, and Gilks (1997) showed that in certain situations the optimal scaling for Metropolis algorithms can be characterized as that which has acceptance probability 0.234. One can ask whether these results carry over to RAMA, and whether equal acceptance rates on different regions (as sought by RAMA) truly lead to optimality. We believe this to be true quite generally, but can only prove it for very specific settings (e.g., birth-death processes). The method of proof of Roberts et al. (see also Bédard 2006a, 2006b) appears to carry over away from the region boundaries, but the behavior at the region boundaries is more complicated.
5. If we set  $\delta(n)$  to a constant, as opposed to having  $\delta(n) \rightarrow 0$ , then condition (1.1) might fail, so the chain might not converge to  $\pi(\cdot)$ . On the other hand, the chain together with the parameter values  $\{a_j\}$  is jointly Markovian, and under appropriate scaling may have its own joint diffusion limit. It would be interesting (Stewart 2006) to study that diffusion limit, to, for example, see how much asymptotic error results from failing to satisfy (1.1).

## 6. TO LOG OR NOT TO LOG

Suppose  $\pi$  is the density function for a real-valued random variable  $W$ . Then if  $\pi$  is heavy-tailed, it may be advantageous to take logarithms, that is, to instead consider the density function for  $\tilde{W} \equiv \log W$ . This leads to the question, when is it advantageous to consider  $\tilde{W}$  in place of  $W$ ? Once again, adaptive algorithms can provide insights into this question.

To avoid problems of negative or near-negative values, we modify the logarithm function and instead consider the function

$$\ell(w) \equiv \text{sgn}(w) \log(1 + |w|),$$

where  $\text{sgn}(w) = 1$  for  $w > 0$ , and  $\text{sgn}(w) = -1$  for  $w < 0$ . The function  $\ell$  is an increasing, continuously differentiable mapping from  $\mathbf{R}$  onto  $\mathbf{R}$ , with inverse  $\ell^{-1}(w) = \text{sgn}(w)(e^{|w|} - 1)$ . If  $\pi$  is the density for  $W$ , and  $\tilde{W} = \log(W)$ , then taking Jacobians shows that the density for  $\tilde{W}$  is given by  $\tilde{\pi}(w) = e^{|w|}\pi(e^{|w|} - 1)$ .

A result of Mengersen and Tweedie (1996) says, essentially, that a random-walk Metropolis (RWM) algorithm for a density  $\pi$  will be geometrically ergodic if and only if  $\pi$  has exponential or subexponential tails, that is, satisfies

$$\log \pi(x) - \log \pi(y) \geq \eta(y - x), \quad y > x \geq x_1, \quad (6.1)$$

for some  $x_1 > 0$  and  $\eta > 0$  (and similarly for  $y < x \leq -x_1$ ). (A similar result holds in multidimensional contexts; cf. Roberts and Tweedie 1996.) But if  $\pi$  on  $\mathbf{R}$  satisfies (6.1), then so does  $\tilde{\pi}$ , because if  $y > x \geq -\log(\eta) + \beta \geq x_1 > 0$ , then

$$\begin{aligned} \log \tilde{\pi}(x) - \log \tilde{\pi}(y) &= (x - y) + \log \pi(e^x - 1) - \log \pi(e^y - 1) \\ &\geq (x - y) + \eta((e^y - 1) - (e^x - 1)) \\ &= -(y - x) + \eta e^x (e^{y-x} - 1) \\ &\geq -(y - x) + \eta^x (y - x) \\ &= (y - x)(\eta e^x - 1) \geq (y - x)(e^\beta - 1). \end{aligned}$$

Hence, (6.1) is satisfied for  $\tilde{\pi}$  with  $\tilde{\eta} = e^\beta - 1$ . In fact, by making  $\beta$  arbitrarily large, we can make  $\tilde{\eta}$  as large as we like, showing that the tails of  $\tilde{\pi}$  are in fact subexponential.

This suggests that, at least as far as geometric ergodicity is concerned, it is essentially always better to work with  $\tilde{\pi}$  than with  $\pi$ . As a specific example, if  $\pi$  is the standard Cauchy distribution, then RWM on  $\pi$  is *not* geometrically ergodic, but RWM on  $\tilde{\pi}$  is.

Despite this evidence in favor of log transforms for RWM, it is not clear that taking logarithms (or applying  $\ell$ ) necessarily helps with the *quantitative* convergence of RWM. To investigate this, we use an adaptive algorithm.

Specifically, given  $\pi$ , we consider two different algorithms: one a RWM on  $\pi$ , and the other a RWM on  $\tilde{\pi}$ , each using proposal distributions of the form  $Q(x, \cdot) = N(x, \sigma^2)$ . After the  $n$ th batch of 100 iterations, we allow each version to adapt its own scaling parameter  $\sigma$  by adding or subtracting  $\delta(n)$  to  $\log(\sigma)$ , in an effort to achieve acceptance rate near 0.44 for each version. Then, once every 100 batches, we consider whether to switch versions (i.e., to apply  $\ell$  if we currently have not, or to undo  $\ell$  if we currently have), based on whether the current average squared jumping distance is smaller than that from the last time we used the other version. (We force the switch to the other version if it fails 100 times in succession, to avoid getting stuck forever with just one version.)

How does this adaptive algorithm work in practice? In the following table we considered three different one-dimensional symmetric target distributions: a standard Normal, a standard Cauchy, and a Uniform $[-100, 100]$ . For each target, we report the percentage of the time that the adaptive algorithm spent on the logged density  $\tilde{\pi}$  (as opposed to the regular density  $\pi$ ). We also report the mean value of the log proposal standard deviation for

both the regular and the logged RWM versions:

Target	Log %	$ls_{\text{reg}}$	$ls_{\text{log}}$
Normal	3.62%	2.52	2.08
Cauchy	99.0%	3.49	2.66
Uniform	4.95%	6.66	2.65

We see from this table that, for the Normal and Uniform distributions, the adaptive algorithm saw no particular advantage to taking logarithms, and indeed stayed in the regular (unlogged)  $\pi$  version the vast majority of the time. On the other hand, for the Cauchy target, the algorithm uses the logged  $\tilde{\pi}$  essentially as much as possible. This shows that this adaptive algorithm is able to distinguish between when taking logs is helpful (e.g., for the heavy-tailed Cauchy target) and when it is not (e.g., for the light-tailed Normal and Uniform targets).

For multidimensional target distributions, it is possible to take logs (or apply the function  $\ell$ ) separately to each coordinate. Because lighter tails are still advantageous in multidimensional settings (Roberts and Tweedie 1996), it seems likely to be advantageous to apply  $\ell$  to precisely those coordinates which correspond to heavy tails in the target distribution. In high dimensions, this cannot feasibly be done by hand, but an adaptive algorithm could still do it automatically. Such multidimensional versions of this adaptive logarithm algorithm appear worthy of further investigation.

## 7. HOW TO ADAPT

Whereas the theory of adaptation has progressed significantly in recent years, practical implementation raises many important and largely unstudied problems. One issue is that we still know very little about how to optimize MCMC algorithms. The use of monitored acceptance probabilities has the appeal of simplicity and the support of some theory. In one dimension the use of 0.44 and in higher dimensional problems the adoption of 0.234 together with the scaling rule  $\sigma = 2.38/d^{1/2}$  are based on results in Gelman, Roberts, and Gilks (1996). However, the use of these simple rules, although often effective, are based on approximations and are not rigorously proved for complex nonhomogeneous models used in statistical analysis. In our two heterogeneous scaling examples, we are guided by established theoretical properties of MCMC—particularly that Metropolis algorithms are not geometrically ergodic for heavy-tailed target densities. We believe that there is considerable further scope for algorithm development based on MCMC theory.

One important question asks whether an effective adaptive scheme should require that  $\Gamma$  converges. It is intuitively appealing to think of the adaptive scheme searching for “the best” algorithm from a collection of candidates. Our approach here is, however, not to require convergence of  $\Gamma$  because we are eager to have adaptive procedures which work in as general a context as possible. It may well be (and we suspect so) that all the examples in this article involve situations where  $\Gamma_n$  does converge, but we have not attempted to



demonstrate this. A complementary approach to our work in this respect is that adopted by Andrieu and Moulines (2003). This has the appeal of generality, but it may be that an algorithm in which we do not have convergence of  $\Gamma_n$  converges less rapidly than one in which  $\Gamma_n$  does converge. More experience with practical examples is necessary to resolve these issues.

AM and RAMA are set up naturally in a multidimensional context. Multivariate generalizations of the state-dependent strategy used in Section 4 are clearly possible. The simplest idea is to apply the same strategy to each of the  $d$  components, obtaining a collection of parameters  $(a_1, b_1), \dots, (a_d, b_d)$  defining a proposal of independent components with variance in the  $i$ th direction given by  $e^a(1 + |x_i|)^b$ . More complex proposals which try to respect the dependence in the target density (as in AM) are also possible.

## 8. CHECKING THE BOUNDED CONVERGENCE CONDITION

The Diminishing Adaptation condition is relatively easy to check, and in fact adaptive procedures are generally constructed with this condition directly in mind. On the other hand, the Bounded Convergence condition is typically more difficult to check.

One way of showing this is to show that all MCMC kernels satisfy the same Lyapunov drift condition. For instance, Roberts and Rosenthal (2005) showed that an adaptive MCMC algorithm satisfying Diminishing Adaptation satisfies Bounded Convergence (and is hence ergodic) if the family  $\{P_\gamma\}_{\gamma \in \mathcal{Y}}$  is *simultaneously strongly aperiodically geometrically ergodic*, that is, there is  $C \in \mathcal{F}$ ,  $V : \mathcal{X} \rightarrow [1, \infty)$ ,  $\delta > 0$ ,  $\lambda < 1$ , and  $b < \infty$ , such that  $\sup_C V = v < \infty$ , and

(i) for each  $\gamma \in \mathcal{Y}$ , there exists a probability measure  $\nu_\gamma(\cdot)$  on  $C$  with  $P_\gamma(x, \cdot) \geq \delta \nu_\gamma(\cdot)$  for all  $x \in C$ ; and

(ii)  $(P_\gamma)V \leq \lambda V + b\mathbf{1}_C$ .

A natural approach to the establishment of simultaneously strongly aperiodic geometrical ergodicity is to use the drift function  $\pi^{-1/2}$  as in Roberts and Tweedie (1996) for AM and Roberts and Rosenthal (1998b) for Adaptive Metropolis-within-Gibbs. As an example of a precise result which can be shown in this way, for the AM algorithm, the condition will hold by this argument for all target densities which are log-concave (except perhaps on some bounded region).

The state-dependent proposal variance case can also be analyzed to give a drift condition with Lyapunov function  $\pi^{-1/2}$ , at least for  $b < 2$ . This is essentially because asymptotically (as  $|x| \rightarrow \infty$ ) the accept/reject ratio is dominated by the ratio  $\pi(y)/\pi(x)$  and thus all moves to smaller  $|x|$  values are accepted with all moves to larger  $\pi(x)$  values possibly rejected. Then standard calculations as those in Roberts and Tweedie (1996), together with continuity and compactness arguments for the parameters  $a$  and  $b$ , are sufficient to demonstrate the simultaneous drift condition.

It seems that these results will be easily generalizable to the other examples in this paper, essentially because all methods are essentially constructed from random-walk Metropolis. These extensions are subject to further work (Bai, Roberts, and Rosenthal 2008).

## 9. CONCLUSION

This article has considered automated tuning of MCMC algorithms, especially Metropolis–Hastings algorithms, with quite positive results.

For example, for Metropolis-within-Gibbs algorithms, the following (generally well-known) statements are all reinforced through our simulation. (1) The choice of proposal variance  $\sigma^2$  is crucial to the success of the algorithm. (2) Good values of  $\sigma^2$  can vary greatly from one coordinate to the next. (3) There are far too many coordinates to be able to select good values of  $\sigma^2$  for each coordinate by hand. (4) Adaptive methods can be used to get the computer to find good values of  $\sigma^2$  automatically. (5) If done carefully, the adaptive methods can be provably ergodic, and quite effective in practice, thus allowing for good tuning and rapid convergence of MCMC algorithms that would otherwise be impractical.

The practical experience of this article, though very promising, also raises important questions. In particular, how robust are the strategies suggested in the various methods here? For example, in the adaptive Metropolis-within-Gibbs example, how crucial is the choice of  $\delta(n)$  in the success of the method, and how does this vary from problem to problem?

We still know comparatively little about what adaptive strategies to use in any particular context. Our feeling is that the choice of adaptive strategy should be guided by theoretical knowledge about MCMC. For instance, when using RWM, particular problems are observed with heavy-tailed target distributions (such as lack of geometric ergodicity, breakdown of CLTs, etc.). In this case it makes sense to use a strategy which attempts to stabilize the algorithm excursions, and this points to the use of heterogeneous scaling and/or a strategy which lightens the tails (such as that introduced in Section 6).

One important issue for adaptive scaling concerns the practical issue that scaling the proposal correlation structure to match that of the target will be a very poor strategy when some of the target distribution variances are infinite. Typically in practical MCMC situations, this may not be very easy to check analytically by inspection of the target density. For this reason, perhaps it makes more sense to scale according to acceptance rate criteria rather than variances. However, it is impossible to use this to match correlation structure, and further work is required to introduce robust versions of the AM and other methods.

Adaptive strategies are generally simple to implement. However, it is very important that such a strategy is constructed in such a way that the conditions for ergodicity are satisfied. Further work is clearly required to give sufficiently simple conditions to enable routine adaptation to take place in applied problems. In terms of adaptive strategies, there is now extensive MCMC theory to help guide the construction of suitable adaptive algorithms. One potential problem evolves from an adaptive strategy which is too “greedy” in that it tries to adapt too closely to initial information from the output. Such algorithms can take considerable time to “recover” from misleading initial information.

Overall, we feel that these results indicate the widespread applicability of adaptive MCMC algorithms to many different MCMC settings, including complicated high-dimensional distributions. We hope that this article will inspire users of MCMC to experiment with adaptive algorithms in their future applications (e.g., Turro et al. 2007).

All of the software used to run the algorithms described herein is freely available at [probability.ca/adapt](http://probability.ca/adapt).

## ACKNOWLEDGMENTS

We thank Sylvia Richardson for a very helpful suggestion, and thank the editors and referees for many constructive comments that improved the article. Research supported in part by NSERC of Canada.

[Received September 2006. Revised October 2008.]

## REFERENCES

- Andrieu, C., and Atchadé, Y. F. (2007), "On the Efficiency of Adaptive MCMC Algorithms," *Electronic Communications in Probability*, 12, 336–349.
- Andrieu, C., and Moulines, E. (2003), "On the Ergodicity Properties of Some Adaptive Markov Chain Monte Carlo Algorithms," preprint.
- Atchadé, Y. F., and Rosenthal, J. S. (2005), "On Adaptive Markov Chain Monte Carlo Algorithms," *Bernoulli*, 11, 815–828.
- Bai, Y., Roberts, G. O., and Rosenthal, J. S. (2008), "The Bounded Convergence Condition for Adaptive MCMC," preprint, available at <http://probability.ca/jeff/research.html>.
- Bédard, M. (2006a), "Weak Convergence of Metropolis Algorithms for Non-iid Target Distributions," preprint.
- (2006b), "Optimal Acceptance Rates for Metropolis Algorithms: Moving Beyond 0.234," preprint.
- Brockwell, A. E., and Kadane, J. B. (2005), "Identification of Regeneration Times in MCMC Simulation, With Application to Adaptive Schemes," *Journal of Computational and Graphical Statistics*, 14, 436–458.
- Efron, B., and Morris, C. (1975), "Data Analysis Using Stein's Estimator and Its Generalizations," *Journal of the American Statistical Association*, 70 (350), 311–319.
- Gelfand, A. E., and Smith, A. F. M. (1990), "Sampling Based Approaches to Calculating Marginal Densities," *Journal of the American Statistical Association*, 85, 398–409.
- Gelman, A., Roberts, G. O., and Gilks, W. R. (1996), "Efficient Metropolis Jumping Rules," in *Bayesian Statistics 5*, eds. J. M. Bernardo, J. O. Berger, A. P. Dawid, and A. F. M. Smith, Oxford, U.K.: Oxford University Press, pp. 599–608.
- Gilks, W. R., Roberts, G. O., and Sahu, S. K. (1998), "Adaptive Markov Chain Monte Carlo," *Journal of the American Statistical Association*, 93, 1045–1054.
- Haario, H., Saksman, E., and Tamminen, J. (1999), "Adaptive Proposal Distribution for Random Walk Metropolis Algorithm," *Computational Statistics*, 14, 375–395.
- (2001), "An Adaptive Metropolis Algorithm," *Bernoulli*, 7, 223–242.
- (2005), "Componentwise Adaptation for High Dimensional MCMC," *Computational Statistics*, 20, 265–274.
- Hastings, W. K. (1970), "Monte Carlo Sampling Methods Using Markov Chains and Their Applications," *Biometrika*, 57, 97–109.
- Mengersen, K. L., and Tweedie, R. L. (1996), "Rates of Convergence of the Hastings and Metropolis Algorithms," *Annals of Statistics*, 24, 101–121.
- Metropolis, N., Rosenbluth, A., Rosenbluth, M., Teller, A., and Teller, E. (1953), "Equations of State Calculations by Fast Computing Machines," *Journal of Chemical Physics*, 21, 1087–1091.
- Morris, C. (1983), "Parametric Empirical Bayes Confidence Intervals," in *Scientific Inference, Data Analysis, and Robustness*, eds. G. E. P. Box, T. Leonard, and C. F. Wu, New York: Academic Press, pp. 25–50.

- Roberts, G. O., and Rosenthal, J. S. (1998a), "Optimal Scaling of Discrete Approximations to Langevin Diffusions," *Journal of the Royal Statistical Society, Ser. B*, 60, 255–268.
- (1998b), "Two Convergence Properties of Hybrid Samplers," *The Annals of Applied Probability*, 8, 397–407.
- (2001), "Optimal Scaling for Various Metropolis–Hastings Algorithms," *Statistical Science*, 16, 351–367.
- (2005), "Coupling and Ergodicity of Adaptive MCMC," preprint.
- Roberts, G. O., and Tweedie, R. L. (1996), "Geometric Convergence and Central Limit Theorems for Multidimensional Hastings and Metropolis Algorithms," *Biometrika*, 83, 95–110.
- Roberts, G. O., Gelman, A., and Gilks, W. R. (1997), "Weak Convergence and Optimal Scaling of Random Walk Metropolis Algorithms," *Annals of Applied Probability*, 7, 110–120.
- Rosenthal, J. S. (1996), "Analysis of the Gibbs Sampler for a Model Related to James–Stein Estimators," *Statistics and Computing*, 6, 269–275.
- (2004), "Adaptive MCMC Java Applet," available at <http://probability.ca/jeff/java/adapt.html>.
- Stewart, A. (2006), personal communication.
- Tierney, L. (1994), "Markov Chains for Exploring Posterior Distributions" (with discussion), *The Annals of Statistics*, 22, 1701–1762.
- Turro, E., Bochkina, N., Hein, A. M. K., and Richardson, S. (2007), "BGX: A Bioconductor Package for the Bayesian Integrated Analysis of Affymetrix GeneChips," *BMC Bioinformatics*, 8, 439–448. Available at <http://www.biomedcentral.com/1471-2105/8/439>.