

# Efficient Sampling Using Metropolis Algorithms: Applications of Optimal Scaling Results

Mylène Bédard \*

July 17, 2006

## Abstract

We recently considered the optimal scaling problem of Metropolis algorithms for multidimensional target distributions with non-IID components. The results that were proven have wide applications and the aim of this paper is to show how practitioners can take advantage of them. In particular, we illustrate with several examples the case where the asymptotically optimal acceptance rate is the usual 0.234, and also the latest developments where smaller acceptance rates should be adopted for optimal sampling from the target distributions involved. We study the impact of the proposal scaling on the performance of the algorithm, and finally perform simulation studies exploring the efficiency of the algorithm when sampling from some popular statistical models.

**Keywords:** Asymptotically optimal acceptance rate, diffusion, hierarchical model, nonidentically distributed components, target distribution, speed measure

## 1 Introduction

Metropolis-Hastings algorithms are a popular class of MCMC methods that are used to generate data from basically any given probability distribution (Metropolis et al. 1953, Hastings 1970). The underlying principle is to build a Markov chain  $X(0), X(1), \dots$  having the probability distribution of interest as its unique stationary distribution. This distribution of interest is referred to as the target distribution and its density is denoted by  $\pi(\cdot)$ , which we assume to be defined on a continuous state space. Necessary to the construction of the Markov chain is the selection of a proposal density  $q(\cdot, \cdot)$ , which is used to generate potential moves for the chain. The algorithm is then performed by applying the following steps. Given that the time- $t$  position of the Markov chain is  $X(t)$ , we generate a new state  $Y(t+1)$  from the proposal distribution, which is accepted with probability

---

\*Department of Statistics, University of Toronto, Toronto, Ontario, Canada, M5S 3G3.  
Email: mylene@utstat.utoronto.ca      Original date: April 2006, revised July 2006

$\alpha(X(t), Y(t+1))$ , where  $\alpha(x, y) = 1 \wedge (\pi(y)q(y, x)/\pi(x)q(x, y))$ . If the suggested move is accepted, then the chain moves from  $X(t)$  to  $Y(t+1)$ ; otherwise,  $X(t+1) = X(t)$ .

The acceptance probability  $\alpha(x, y)$  is chosen such that the Markov chain is reversible with respect to the target density  $\pi(\cdot)$ , ensuring that  $\pi(\cdot)$  is stationary for the chain. To compel the convergence of the Markov chain to its stationary distribution, one might choose virtually any proposal density engendering a Markov chain that is irreducible and aperiodic. A popular choice for the proposal distribution is the normal family centered at the current state  $X(t)$  and scaled according to  $\sigma^2$ . For an algorithm with such a proposal distribution to be efficient, one must however carefully select the variance term. Indeed, small values for  $\sigma^2$  cause the chain to jump to nearby states, resulting in a lengthy exploration of the state space. On the other hand, large scaling values often suggest states located in low target density regions, encouraging the chain to refuse these moves and stand still for long periods of time.

In this paper, we consider the optimal scaling problem of Metropolis algorithms (Roberts et al. 1997, Breyers & Roberts 2000, Roberts & Rosenthal 2001, Neal & Roberts 2004) for targets with independent but not identically distributed components, and study how well existing asymptotic results serve the practical side. In particular, we consider the target model introduced recently in Bédard (2006a,b) and present various examples illustrating how the optimal value for the proposal variance and acceptance rate can be determined using the theorems proved in these papers. This demonstrates that, although of asymptotic nature, the results can be used to facilitate the tuning of algorithms in finite-dimensional problems.

In Section 2 we shall introduce the target and proposal distributions considered, and briefly discuss measures of efficiency for the algorithm. The optimal scaling results that shall be presented in Section 3 are divided in three cases: the asymptotically optimal acceptance rate (AOAR) is the usual 0.234; the AOAR is smaller than 0.234 and its exact value is solved on a case per case basis; finally the AOAR is 0, in which case the efficiency of the algorithm cannot be optimized. In that occurrence, we must then resort to inhomogeneous proposal distributions. In each of the three parts, the theoretical results are first discussed and examples illustrating their application are then presented. Section 4 aims to present simulation studies to investigate the efficiency of the algorithm applied to more complicated and widely used statistical models; specifically, we consider the normal hierarchical model, the variance components model and the gamma-gamma hierarchical model.

## 2 The Model

To illustrate the purpose of the following sections, suppose we wish to sample from

$$\pi(x_1, \dots, x_d) \propto x_1^4 \exp(-x_1) x_2^4 \exp(-x_2) \prod_{i=3}^d x_i^4 \exp(-x_i/5\sqrt{d}). \quad (1)$$

For a relatively low-dimensional target density of this sort, say  $d = 10$ , the density of the last 8 components is spread out over  $(0, \infty)$  while that of the first two have their mass concentrated within a much narrower interval of the state space. Choosing a proper variance for the proposal distribution is thus not an easy task: the last 8 components require a large proposal variance for appropriate exploration of their state space, but the selection of too large a value would result in

frequent rejection of the proposed moves by the variables  $X_1$  and  $X_2$ . A compromise between these requirements then becomes necessary.

For this example, the optimal proposal variance is close to  $\sigma^2 = 61$  for any dimension  $d$ , resulting in an AOAR lying around 0.098, as shall be seen in Example 5. In fact, tuning the algorithm to accept a proportion of 23.4% of the proposed moves would reduce the efficiency of the algorithm by about 20%, from where the importance of determining the right proposal variance. Before discussing the optimal scaling results, we however start by introducing the general model for the target density.

## 2.1 The Target Distribution

Suppose we want to generate a sample from the  $d$ -dimensional target density

$$\pi(d, \mathbf{x}^{(d)}) = \prod_{j=1}^d \theta_j(d) f(\theta_j(d) x_j), \quad (2)$$

where  $f$  is a smooth density and  $\theta_j^{-2}(d)$ ,  $j = 1, \dots, d$  are referred to as the scaling terms of the target distribution, which can be any function of  $d$  for which the limit exists. This model constitutes a natural extension of the  $d$ -dimensional density formed of independent and identically distributed (IID) components considered in the literature and used to investigate optimal acceptance rates.

The optimal scaling results that shall be presented for sampling from this target are originally valid for infinite-dimensional distributions. As  $d$  increases, some of the scaling terms will thus be repeated an infinite number of times; we assume that there are  $0 < m < \infty$  such different terms. Other terms will not be replicated as  $d \rightarrow \infty$ ; we suppose that there are  $n < \infty$  of them. We let the first  $n + m$  components of the vector  $\Theta^{-2}(d)$  consist in the  $n$  non-replicated scaling terms, followed by the  $m$  different scaling terms that shall be replicated. We further assume that the first  $n$  and next  $m$  components are respectively arranged according to an asymptotically increasing order.

Although the components  $\theta_{n+1}(d), \dots, \theta_{n+m}(d)$  appear infinitely often as  $d \rightarrow \infty$ , they might however not be replicated the same number of times. To determine the proportion of  $\Theta^{-2}(d)$  occupied by the  $i$ -th group ( $i \in \{1, \dots, m\}$ ), we define the cardinality functions

$$c(i, d) = \#\{j \in \{1, \dots, d\}; \theta_j(d) = \theta_{n+i}(d)\}, \quad i = 1, \dots, m. \quad (3)$$

The following example should help clarifying the notation just introduced.

**Example 1.** Consider a  $d$ -dimensional target density as in (2) with scaling terms  $1/\sqrt{d}$ ,  $4/\sqrt{d}$ , 10 and the other ones equally divided among  $2\sqrt{d}$  and  $(d+1)/2$ . As  $d \rightarrow \infty$ , the last two scaling terms are replicated so  $n = 3$ ,  $m = 2$  and  $\Theta^{-2}(d) = (1/\sqrt{d}, 4/\sqrt{d}, 10, 2\sqrt{d}, (d+1)/2, 2\sqrt{d}, (d+1)/2, \dots)$ . The cardinality functions for the scaling terms appearing infinitely often in the limit are

$$c(1, d) = \#\left\{j \in \{1, \dots, d\}; \theta_j(d) = (2\sqrt{d})^{-1/2}\right\} = \left\lceil \frac{d-3}{2} \right\rceil$$

and  $c(2, d) = \lfloor (d-3)/2 \rfloor$ , where  $\lceil \cdot \rceil$  and  $\lfloor \cdot \rfloor$  denote the ceiling and integer part functions respectively. Note however that such rigorousness is superfluous for applying the results and it is enough to affirm that both cardinality functions grow according to  $d/2$ .

For simplicity's sake, the model just presented for the target is not the most general form under which the conclusions of the theorems are satisfied. Indeed, it would be sufficient to assume that scaling terms belonging to a common group  $i \in \{1, \dots, m\}$  are of the same order. For more details, we refer the reader to Bédard (2006a,b).

## 2.2 The Proposal Distribution and its Scaling

The proposal distribution we consider for sampling from the target density  $\pi(\cdot)$  is such that  $\mathbf{Y}^{(d)}(t+1) \sim N(\mathbf{X}^{(d)}(t), \sigma^2(d) I_d)$ , where  $I_d$  is the  $d$ -dimensional identity matrix.

There exist two factors determining the form of the proposal variance as a function of  $d$ : the asymptotically smallest scaling term and the fact that some scaling terms appear infinitely often in the limit. If the first factor were ignored, the proposed moves would possibly be too large for the corresponding component, resulting in high rejection rates and slow convergence of the algorithm. The effect of the second factor is that as  $d \rightarrow \infty$ , the algorithm proposes more independent moves in a single step, increasing the odds of proposing an improbable move. In this case, a drop in the acceptance rate can be overturned by letting  $\sigma^2(d)$  be a decreasing function of the dimension. Combining these two constraints, the optimal form for the proposal variance as a function of  $d$  can be shown to be  $\sigma^2(d) = \ell^2 \sigma_\alpha^2(d)$ , with  $\ell$  some positive constant and  $\sigma_\alpha^2(d)$  the largest order function such that

$$\lim_{d \rightarrow \infty} \theta_1^2(d) \sigma_\alpha^2(d) < \infty \quad \text{and} \quad \lim_{d \rightarrow \infty} c(i, d) \theta_{n+i}^2(d) \sigma_\alpha^2(d) < \infty \quad \text{for } i = 1, \dots, m. \quad (4)$$

Our goal has thus evolved in optimizing the choice of the constant  $\ell$  in  $\sigma^2(d)$ .

**Example 2.** We now determine the optimal form for the proposal variance of the Metropolis algorithm in Example 1. According to (4), we have three limits to verify: the first one involves  $\theta_1^{-2}(d)$ , which is also the asymptotically smallest scaling term in the present case; the largest  $\sigma_\alpha^2(d)$  satisfying the finite property for  $\lim_{d \rightarrow \infty} \sqrt{d} \sigma_\alpha^2(d)$  is  $1/\sqrt{d}$ . For the second and third limits, we have  $\lim_{d \rightarrow \infty} (d-3) \sigma_\alpha^2(d) / 4\sqrt{d}$  and  $\lim_{d \rightarrow \infty} (d-3) \sigma_\alpha^2(d) / (d+1)$ ; the largest  $\sigma_\alpha^2(d)$  satisfying the finite property is then  $1/\sqrt{d}$  and 1 respectively. The function of largest order satisfying the constraint that all three limits be finite is  $\sigma_\alpha^2(d) = 1/\sqrt{d}$ .

## 2.3 Efficiency of the Algorithm

In order to optimize the mixing of our Metropolis algorithm, it would be convenient to determine criteria for measuring efficiency. The natural way to estimate a function of interest  $g(\cdot)$  is to compute the expectation of  $g(\cdot)$  with respect to  $\pi(\cdot)$ . Minimizing the asymptotic variance of  $g(\cdot)$  would then be a good way to optimize efficiency; however, an important drawback of this measure consists in its dependence on  $g(\cdot)$ . Since we do not want to lose generality by specifying such a quantity of interest, we instead choose the first order efficiency criterion, as used by Roberts & Rosenthal (1998) and Pasarica & Gelman (2003). This measures the average squared jumping distance of the  $(n+1)$ -st component of the algorithm and is defined by  $FOE = \mathbb{E} \left[ \left( X_{n+1}^{(d)}(t+1) - X_{n+1}^{(d)}(t) \right)^2 \right]$ .

The fact that *FOE* is based on the path of the  $(n + 1)$ -st component of the Markov chain is an important detail, given that the  $d$  components are not all identically distributed (although we could have chosen any of the last  $d - n$  components). Indeed, as  $d \rightarrow \infty$ , it can be shown that the path followed by any of the last  $d - n$  components of an appropriately rescaled version of the Metropolis algorithm converges to a diffusion process with speed measure  $v(\ell)$ . For any function of interest  $g(\cdot)$ , optimal efficiency is thus obtained by maximizing the speed measure, meaning that the effect of choosing a particular efficiency criterion vanishes as  $d$  gets larger. Consequently, the optimization problem is the same no matter which function  $g(\cdot)$  is under investigation and any efficiency measure considered in finite dimensions will be asymptotically equivalent (i.e. as  $d \rightarrow \infty$ ), including the first order efficiency introduced previously.

Even though the last  $d - n$  terms always converge to some diffusion limit, it might not be the case for the first  $n$  components, whose limit could remain discrete as  $d \rightarrow \infty$ . Trying to optimize the proposal variance by relying on these components would then result in conclusions that are specific to our choice of efficiency measure.

### 3 Optimal Scaling Results

#### 3.1 The Familiar Asymptotic Behavior

Consider a Metropolis algorithm with proposal distribution  $\mathbf{Y}^{(d)} \sim N(\mathbf{x}^{(d)}, \ell^2 \sigma_\alpha^2(d) I_d)$ , where  $\sigma_\alpha^2(d)$  satisfies (4). Suppose this algorithm is used to sample from a target density as in (2) and supporting the model described in Section 2.1.

It is now an established fact that 0.234 is the AOAR for target distributions with IID components (Roberts et al. 1997). Roberts & Rosenthal (2001) also showed that the same conclusion applies for  $\pi(\cdot)$  as in (2) but with  $\theta_j(d)$ 's independent of  $d$ . It is thus natural to wonder how big a discrepancy between the scaling terms is tolerated in order not to violate this established asymptotic behavior. It turns out that if

$$\lim_{d \rightarrow \infty} \frac{\theta_1^2(d)}{\sum_{j=1}^d \theta_j^2(d)} = 0, \quad (5)$$

then the optimal acceptance rate can be shown to converge towards 0.234. The optimal value  $\hat{\ell}$  maximizes the equation

$$v(\ell) = 2\ell^2 \Phi\left(-\frac{\ell\sqrt{E_R}}{2}\right) = \ell^2 a(\ell), \quad (6)$$

where  $a(\ell) = 2\Phi(-\ell\sqrt{E_R}/2)$  is the asymptotic acceptance rate of the algorithm and

$$E_R = \lim_{d \rightarrow \infty} \sum_{i=1}^m c(i, d) \theta_{n+i}^2(d) \sigma_\alpha^2(d) \mathbb{E} \left[ \left( \frac{f'(X)}{f(X)} \right)^2 \right], \quad (7)$$

with  $c(i, d)$  as in (3). We then obtain  $\hat{\ell} = 2.38/\sqrt{E_R}$ , from where  $a(\hat{\ell}) = 0.234$ . For a formal version of the theorem, see Bédard (2006a).

This result provides valuable guidelines for practitioners. It reveals that when the target distribution has no scaling term that is significantly smaller than the others, the asymptotic acceptance rate optimizing the efficiency of the chain is 0.234. Alternatively, setting  $\ell$  equal to  $2.38/\sqrt{E_R}$  leads to greatest efficiency of the algorithm. In some situations, finding  $\hat{\ell}$  will be easier while in others, tuning the algorithm according to the AOAR will reveal more convenient. In the present case, since the AOAR does not depend on the target density, it is simpler in practice to monitor the acceptance rate and to tune it to be about 0.234.

While the AOAR is independent of the target distribution,  $\hat{\ell}$  varies inversely proportionally to  $E_R$ . Two different factors influence this quantity: the function  $f(\cdot)$  in (2) and the  $\theta_j^{-2}(d)$ 's. The latter can have an effect through their size or the proportion of the vector  $\Theta^{-2}(d)$  they occupy. Specifically, suppose that  $c(i, d)\theta_{n+i}^2(d)$  is  $O(\sigma_\alpha^{-2}(d))$  for some  $i \in \{1, \dots, m\}$ , implying that the  $i$ -th group contributes to augment the value of  $E_R$ . The amount by which  $E_R$  increases is then proportional to the size of the  $\theta_j(d)$ 's, but inversely proportional to the quantity of scaling terms included in the group. The following examples shall clarify these concepts.

**Example 3.** Consider a  $d$ -dimensional target with independent, normally distributed components; we suppose that half of them have a variance of 1 and the other half a variance of  $2d$ . Applying (4), the proposal variance takes the form  $\sigma^2(d) = \ell^2/d$  and Condition (5) is verified by computing  $\lim_{d \rightarrow \infty} (d/2 + 1/4)^{-1} = 0$  (in fact this is trivial since  $n = 0$ ). We can thus optimize the efficiency of the algorithm by setting the acceptance rate to be close to 0.234; equivalently, since  $E_R = \lim_{d \rightarrow \infty} (1/2 + 1/4d) = 1/2$ , we find  $\hat{\ell} = 3.366$ . What is causing an increase of  $\hat{\ell}$  with respect to the baseline 2.38 for the case where all components are IID standard normal is the fact that only half of the components affect the accept/reject ratio in the limit.

The left graph in Figure 1 presents the relation between first order efficiency and  $\ell^2$ . The dotted curve has been obtained by performing 100,000 iterations of the Metropolis algorithm with  $d = 100$ , and as expected the maximum is located very close to  $\hat{\ell}^2 = 11.33$ . The simulations also agree with the theoretical curve (solid line) of  $v(\ell)$  in (6) versus  $\ell^2$ . For the second graph, we run the algorithm with various values of  $\ell$  and plot  $FOE$  as a function of the proportion of accepted moves for the different proposal variances. That is, each point in a given curve is the result of a simulation with a particular value for  $\ell$ . We again performed 100,000 iterations, but this time we repeated the simulations for different dimensions ( $d = 10, 20, 50, 100$ ), outlining the fact that the optimal acceptance rate converges very rapidly to its asymptotic counterpart. The theoretical curve of  $v(\ell)$  versus  $a(\ell)$  is represented by the solid line.

We note that efficiency is a relative measure in our case. Consequently, choosing an acceptance rate around 0.05 or 0.5, would necessitate to run the chain 1.5 times as long to obtain the same precision for a particular estimate.

Although MCMC methods are not necessarily required to sample from normal distributions, this type of target is widely used in the literature to investigate the optimal scaling problem, and thus allows us to see how our results compare to others. Note however that we could also have used any smooth density  $f(\cdot)$  with any  $\Theta^{-2}(d)$  satisfying Condition (5). While the convergence might get slightly slower as  $f(\cdot)$  gets further from normality, we generally observe curves similar to those of Figure 1, as well as a rapid convergence to the asymptotic curve. The following example presents a particular situation where the convergence of some components towards the AOAR is extremely slow, a phenomenon due to the form of  $\Theta^{-2}(d)$ .

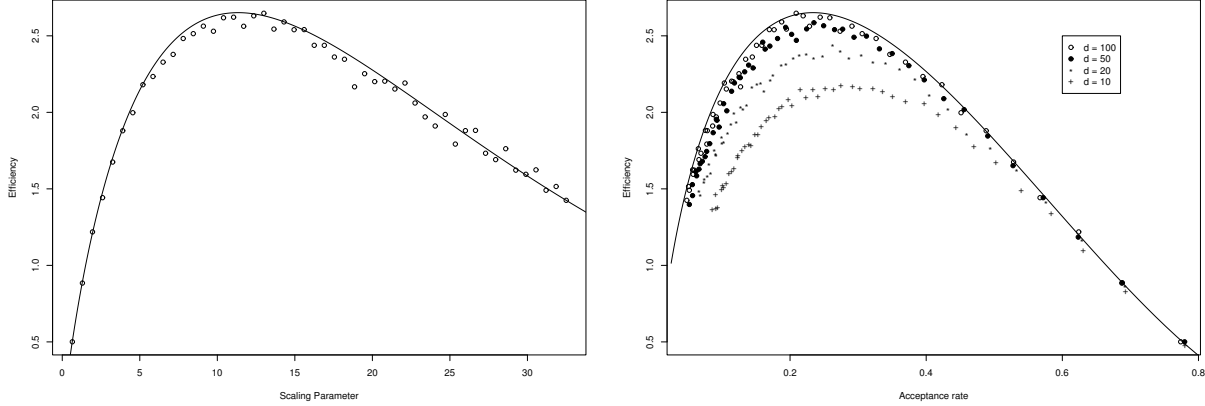


Figure 1: Left graph: efficiency of  $X_1$  versus  $\ell^2$ ; the dotted line is the result of simulations with  $d = 100$ . Right graph: efficiency of  $X_1$  versus the acceptance rate; the dotted lines come from simulations in dimensions 10, 20, 50 and 100. In both graphs, the solid line represents the theoretical curve  $v(\ell)$ .

**Example 4.** Consider a multivariate normal target with independent components and variances  $\Theta^{-2}(d) = (d^{-0.75}, 1, 1, 1, \dots)$ . We find  $\sigma^2(d) = \ell^2/d$  and  $\lim_{d \rightarrow \infty} d^{0.75} (d^{0.75} + (d-1))^{-1} = 0$ , implying that Condition (5) is verified. The quantity  $E_R$  being equal to 1, the optimal value for  $\ell$  is then the baseline 2.38.

The particularity of this case resides in the size of  $\theta_1^{-2}(d)$ , which is somewhat smaller than the other terms but not enough to remain significant as  $d \rightarrow \infty$ . As a consequence, the dimension of the target distribution must be quite large before the asymptotics kick in. In small dimensions, the optimal acceptance rate is thus closer to the case where there exist significantly small scaling terms, which shall be studied in Section 3.2.

Figure 2 demonstrates that even in small dimensions, the first order efficiency criterion based on any of the last  $d-1$  components is very close to 0.234 so as not to make much difference in practice. When first order efficiency is based on  $X_1$  however, setting  $d = 100$  yields an optimal acceptance rate around 0.3 and the dimensions must be raised as high as 100,000 to get an optimal acceptance rate reasonably close to the asymptotic one. Relying on  $X_1$  would then falsely suggest a higher optimal acceptance rate, as explained in Section 2.3.

### 3.2 A Reduction of the Asymptotically Optimal Acceptance Rate

In the presence of a finite number of significantly small scaling terms, choosing a correct proposal variance is a slightly more delicate task. We can think for instance of the density in (1), which seems to promote contradictory characteristics when it comes to the selection of an efficient proposal variance. In that example, the components  $X_1$  and  $X_2$  are said to rule the algorithm since they ensure that the proposal variance is not too big as a function of  $d$ . When dealing with such target

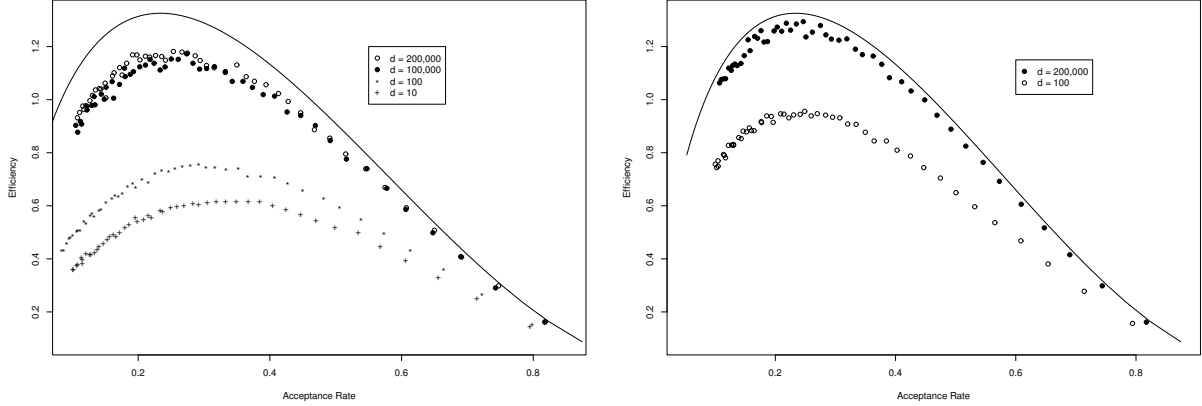


Figure 2: Left graph: efficiency of  $X_1$  versus the acceptance rate; the dotted curves are the results of simulations in dimensions 10, 100, 100,000 and 200,000. Right graph: efficiency of  $X_2$  versus the acceptance rate; the dotted curves come from simulations in dimensions 100 and 200,000. In both graphs, the solid line represents the theoretical curve  $v(\ell)$ .

densities, we realize that Condition (5) is violated, and thus

$$\lim_{d \rightarrow \infty} \frac{\theta_1^2(d)}{\sum_{j=1}^d \theta_j^2(d)} > 0. \quad (8)$$

The existence of scaling terms ruling the algorithm leads to one of two situations: if they are of extremely small order compared to the other scaling terms, this results in the inefficiency of the algorithm; on the other hand, a reasonable difference of size among them can be handled. In particular, an AOAR exists if there is at least one  $i \in \{1, \dots, m\}$  satisfying

$$\lim_{d \rightarrow \infty} c(i, d) \theta_{n+i}^2(d) \sigma_\alpha^2(d) > 0. \quad (9)$$

In words this condition requires that if we were ignoring  $X_1, \dots, X_n$ , the form selected for  $\sigma_\alpha^2(d)$  based on the last  $d - n$  components only would remain intact. This therefore ensures that the first  $n$  components are not solely affecting the selection of  $\sigma_\alpha^2(d)$ .

In this occurrence, the optimal value for  $\ell$  maximizes

$$v(\ell) = 2\ell^2 \mathbb{E}_{\mathbf{Y}^{(b)}, \mathbf{X}^{(b)}} \left[ \Phi \left( \frac{\sum_{j=1}^b \log(f(\theta_j Y_j) / f(\theta_j X_j)) - \ell^2 E_R / 2}{\sqrt{\ell^2 E_R}} \right) \right] = \ell^2 a(\ell), \quad (10)$$

where  $b = \max(j \in \{1, \dots, n\}; 0 < \lim_{d \rightarrow \infty} \theta_j(d) \theta_1^{-1}(d) < \infty)$  is the number of non-replicated components having a  $O(\theta_1^{-2}(d))$  scaling term ( $E_R$  and  $c(i, d)$  are as in (7) and (3) respectively).

The equation  $v(\ell)$  is affected by the  $b$  components having a scaling term which is small enough so as to affect the accept-reject ratio of the algorithm in the limit. The AOAR is unfortunately not independent of the target distribution anymore, and varies according to the choice of  $f(\cdot)$



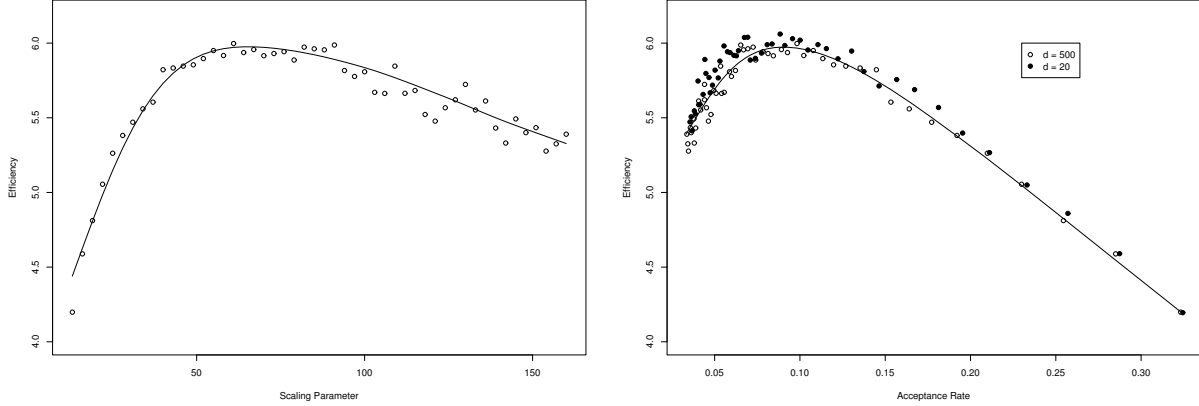


Figure 3: Left graph: efficiency of  $X_3$  versus  $\hat{\ell}^2$ ; the dotted curve represents the results of simulations with  $d = 500$ . Right graph: efficiency of  $X_3$  versus the acceptance rate; the results of simulations with  $d = 20$  and  $500$  are pictured by the dotted curves. In both cases, the theoretical curve  $v(\ell)$  is depicted (solid line).

in (2) and  $\Theta^{-2}(d)$ . It is then simpler to optimize the efficiency of the algorithm by numerically determining  $\hat{\ell}$  from (10) rather than monitoring the acceptance rate, since in any case finding the AOAR implies solving for  $\hat{\ell}$ . As before,  $\hat{\ell}$  is inversely proportional to  $E_R$  but now also depends on  $X_1, \dots, X_b$ , causing the algorithm to reject a greater proportion of moves. This provokes a reduction of  $a(\hat{\ell})$ , resulting in AOARs that are smaller than 0.234 and vary inversely proportionally to  $b$ . Correspondingly, the greater  $b$  is, the smaller  $\hat{\ell}$  will be. For more details on this case, see Bédard (2006b).

The following example illustrates how to solve for the appropriate  $\hat{\ell}$  and AOAR using (10). It presents a situation where tuning the acceptance rate to 0.234 results in an algorithm whose performance is substantially less than when using the correct AOAR.

**Example 5.** Consider the  $d$ -dimensional target density introduced in (1). Consistent with the notation of Sections 2.1 and 2.2, we find  $\Theta^{-2}(d) = (1, 1, 25d, 25d, 25d, \dots)$  and so  $\sigma^2(d) = \ell^2$ . We remark that the first two scaling terms are significantly smaller than the balance since we have  $\lim_{d \rightarrow \infty} (2 + (d - 2)/25d)^{-1} = 25/51$ . Even though  $\theta_1^{-2}(d)$  and  $\theta_2^{-2}(d)$  are significantly small, they still share the responsibility of selecting  $\sigma_\alpha^2(d)$  with the other  $d - 2$  components since  $\lim_{d \rightarrow \infty} c(1, d) \theta_3^2(d) \sigma_\alpha^2(d) = \lim_{d \rightarrow \infty} (d - 2)/25d = 1/25$ . Conditions (8) and (9) being satisfied, we thus use (10) to optimize the efficiency of the algorithm. After having estimated the expectation term in (10) for various values of  $\ell$ , a scan of the vector  $v(\ell)$  produces  $\hat{\ell}^2 = 61$  and  $a(\hat{\ell}) = 0.0981$ . Note that the term  $E_R = 1/75$  causes an increase of  $\hat{\ell}$ , but the components  $X_1$  and  $X_2$  ( $b = 2$ ) act in the opposite direction. This is why  $\hat{\ell}^2 < 424.83$ , which would be the optimal value for  $\ell$  if  $X_1$  and  $X_2$  were ignored.

Figure 3 illustrates the result of 500,000 iterations of a Metropolis algorithm in dimensions 500 for the left graph and in dimensions 20 and 500 for the right one. On both graphs, the maximum occurs close to the theoretical values mentioned previously. We note that the AOAR is now quite far from 0.234, and that tuning the proposal scaling so as to produce this acceptance rate would

contribute to considerably lessen the performance of the method. In particular, this would generate a drop of at least 20% in the efficiency of the algorithm.

In Example 3, it did not matter which of the  $d$  components was selected to compute first order efficiency, as all of them would have yielded similar efficiency curves. In Example 4, the choice of the component became important since  $X_1$  had a scaling term much smaller than the others, resulting in a lengthy convergence to the right optimal acceptance rate. In Examples 5, it is now crucial to choose this component judiciously since  $X_1$  has an asymptotic distribution that remains discrete. The AOAR generated by this sole component is thus specific to the chosen measure of efficiency, which is not representative of the target distribution as a whole.

### 3.3 Inhomogeneous Proposal Distribution: An Alternative

We finally consider the remaining situation where there exist  $b$  components having scaling terms that are of extremely small order (where  $b$  is as in Section 3.2), meaning that they are the only ones to have an impact on the selection of the proposal variance. This is mathematically translated by the satisfaction of Condition (8) along with

$$\lim_{d \rightarrow \infty} c(i, d) \theta_{n+i}^2(d) \sigma_\alpha^2(d) = 0 \quad \forall i \in \{1, \dots, m\}. \quad (11)$$

This means that if we were basing our prognostic for  $\sigma_\alpha^2(d)$  on the last  $d - n$  components only, we would opt for a larger order proposal variance. The  $b$  components thus become the only one to have an impact on the accept/reject ratio as the dimension of the target increases.

In these circumstances, the value  $\hat{\ell}$  is maximizing the equation

$$v(\ell) = 2\ell^2 \mathbf{P}_{\mathbf{Y}^{(b)}, \mathbf{X}^{(b)}} \left( \sum_{j=1}^b \varepsilon(X_j, Y_j) > 0 \right) = \ell^2 a(\ell). \quad (12)$$

Attempting to optimize  $v(\ell)$  leads to an impasse, since this function is unbounded for basically any smooth density  $f(\cdot)$ . That is,  $v(\ell)$  increases with  $\ell$ , resulting in a null AOAR. This phenomenon can be explained by the fact that the scaling of the first  $b$  components are much smaller than the others, determining the form of  $\sigma^2(d)$  as a function of  $d$ . However, the moves generated by a proposal distribution with such a variance will definitely be too small for the other components, forcing the parameter  $\ell$  to increase in order to generate reasonable moves for them. In practice, it is thus impossible to find a proposal variance that is small enough for the first  $b$  components, but at the same time large enough so as to generate moves that are not compromising the convergence speed of the last  $d - b$  components. In Section 3.2, the situation encountered was similar, except that it was possible to achieve an equilibrium between these two constraints. In the current circumstances, the discrepancy between the scaling terms is too large and the disparities are irreconcilable.

In practice, we then face an important difficulty as no choice of  $\ell$  will lead to an efficient algorithm. Therefore, modifying the proposal distribution becomes essential and we turn to inhomogeneous proposal variances. The approach chosen is to maintain the same parameter  $\ell$  for each component and also the same form for the proposal variance of the first  $n$  terms, while adjusting the form of

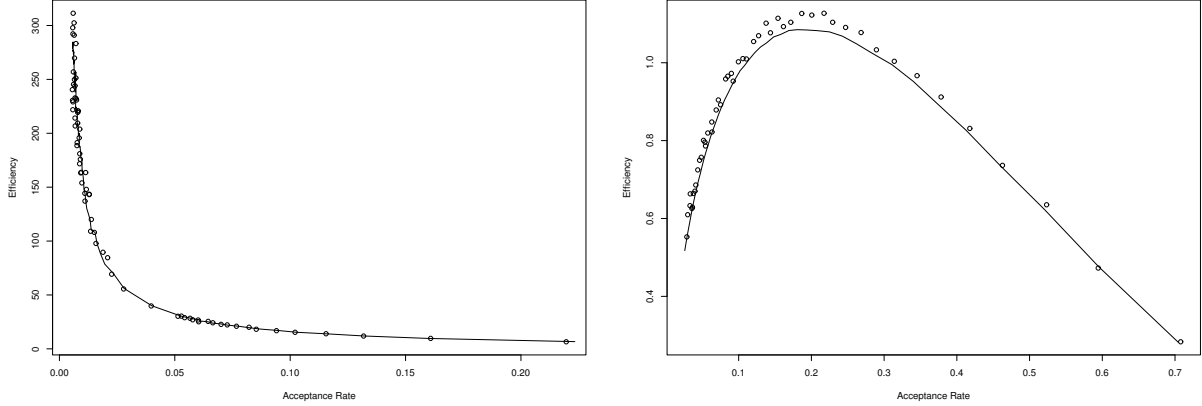


Figure 4: Efficiency of  $X_2$  versus the acceptance rate for homogeneous and inhomogeneous proposal variances respectively. The solid line represents the theoretical curve  $v(\ell)$  and the dotted line has been obtained by running a Metropolis algorithm in dimensions 101.

$\sigma^2(d)$  as a function of  $d$  for the last  $d - n$  components. In particular, we let  $\sigma^2(d) = \ell^2 \sigma_\alpha^2(d)$  with  $\sigma_\alpha^2(d)$  as in (4) for  $j = 1, \dots, n$ . For every component  $j \in \{n + 1, \dots, d\}$  and belonging to the  $i$ -th of the  $m$  groups,  $\sigma_\alpha^2(d)$  is now defined to be the largest order function satisfying

$$0 < \lim_{d \rightarrow \infty} c(i, d) \theta_{n+i}^2(d) \sigma_\alpha^2(d) < \infty. \quad (13)$$

Under the inhomogeneity assumption for the proposal variances, we can then use the results presented in Section 3.2 and determine  $\hat{\ell}$  by maximizing (10). To illustrate such a situation, consider the following example.

**Example 6.** Suppose  $f(\cdot)$  in (2) is the standard normal density and consider the vector of variances  $\Theta^{-2}(d) = (d^{-5}, d^{-1/2}, 3, d^{-1/2}, 3 \dots)$ . The particularity of this setting resides in the fact that  $\theta_1^{-2}(d)$  is extremely small compared to the other scaling terms (so Condition (8) is satisfied). In the present case,  $\sigma^2(d) = \ell^2/d^5$  and the proposal variance is totally governed by  $X_1$ ; indeed,  $\lim_{d \rightarrow \infty} (d - 1)/2d^{4.5} = 0$  and  $\lim_{d \rightarrow \infty} (d - 1)/6d^5 = 0$ , implying that Condition (11) is also verified. We must then use (12) to determine how to optimize the efficiency of the algorithm.

As explained previously and as illustrated by the left graph of Figure 4, the optimal value for  $\ell$  diverges, resulting in an optimal acceptance rate which converges to 0. Obviously, it is impossible to reach a satisfactory level of efficiency in the limit using the prescribed proposal distribution. To overcome this problem, we shall make use of inhomogeneous proposal distributions. The idea is to personalize the proposal variance of the last  $d - 1$  terms. The proposal variance for the first term just stays  $\ell^2/d^5$  and using (13), the vector  $\Theta^{-2}(d)$  becomes  $(\ell^2/d^5, \ell^2/d^{1.5}, \ell^2/d, \dots, \ell^2/d^{1.5}, \ell^2/d)$ . From the results of Section 3.2, we then deduce that  $E_R = \lim_{d \rightarrow \infty} ((d - 3)/2d + (d - 3)/6d) = 2/3$ .

Running the Metropolis algorithm for 100,000 iterations in dimensions 101 yields the curves in Figure 4 (right graph), where the solid line again represents the theoretical curve  $v(\ell)$  in (10). The theoretical values obtained for  $\hat{\ell}^2$  and  $a(\hat{\ell})$  are 6 and 0.1808251 respectively, which agree with the simulations. The inhomogeneous proposal variances have then contributed to decrease  $\hat{\ell}$  while

raising the AOAR. Indeed, large values for  $\hat{\ell}$  are now inappropriate since components with larger scaling terms now possess a proposal variance that is suited to their size, ensuring an reasonable speed of convergence for these components.

## 4 Simulation Studies for some Hierarchical Models

A nice feature of the results presented in this paper is their capability to optimize the efficiency of the Metropolis algorithm when sampling from any multivariate normal target distribution, no matter the correlation structure existing between its components. The normal hierarchical target model considered in Section 4.1 illustrates this property. The last two sections focus on empirically studying the optimal scaling problem for more general hierarchical models, engendering distributions that are not jointly normal.

### 4.1 Normal Hierarchical Model

Consider a model with location parameters  $\mu_1 \sim N(0, 1)$  and  $\mu_2 \sim N(\mu_1, 1)$ . Further suppose, assuming conditional independence, that  $X_i \sim N(\mu_1, 1)$ ,  $i = 1, \dots, 9$  and  $X_i \sim N(\mu_2, 1)$ ,  $i = 10, \dots, 18$ . The joint distribution of  $\mu_1, \mu_2, X_1, \dots, X_{18}$  is multivariate normal with null mean and covariance matrix  $\Sigma_{20}$ . Obtaining the covariances between each pair of components is easily achieved by using conditioning: for the variances, we obtain  $\sigma_1^2 = 1$ ,  $\sigma_i^2 = 2$  for  $i = 2, \dots, 11$  and  $\sigma_i^2 = 3$  for  $i = 12, \dots, 20$ ; for the covariances, we get  $\sigma_{ij} = 2$  for  $i = 2, j = 12, \dots, 20$  (and vice versa) and for  $i = 12, \dots, 20, j = 12, \dots, 20, i \neq j$ ; all the other covariance terms are equal to 1.

A useful property of multivariate normal distributions is their invariance under orthogonal transformations. It is therefore possible to transform  $\Sigma_{20}$  into a diagonal matrix where the diagonal elements consist in the eigenvalues of  $\Sigma_{20}$ . Since the target distribution is still normal but has now independent components, optimizing the efficiency of the Metropolis algorithm can be achieved by using the results presented previously.

In order to determine which one of (6), (10) or (12) should be used for determining  $\hat{\ell}$ , we need to know how the eigenvalues of  $\Sigma_d$  evolve as a function of  $d$ . Obtaining numerical values for the eigenvalues of  $\Sigma_d$  in any dimension is easily achieved with the help of any statistical software; this allows us to deduce that  $d-4$  of the  $d$  eigenvalues are exactly equal to 1. Plots of the four remaining eigenvalues,  $\lambda_i(d)$ ,  $i = 1, \dots, 4$ , clearly show that the two smallest eigenvalues satisfy  $a_i/d = \lambda_i(d)$ ; they also reveal a relation of the form  $a_i d = \lambda_i(d)$  for the two largest eigenvalues. Fitting these linear equations using the eigenvalues of  $\Sigma_{600}$  (say), we obtain fitted values for  $a_i$ ,  $i = 1, 2, 3, 4$ .

Optimizing the efficiency of the algorithm for sampling from this hierarchical model then reduces to optimize a 20-dimensional multivariate normal distribution with independent components, null mean and variances equal to  $(1.983/20, 1.997/20, 0.193(20), 1.311(20), 1, \dots, 1)$ . It is easily verified that this vector satisfies Conditions (8) and (9), and leads to  $\sigma^2(\ell) = \ell^2/d$ . We then turn to equation (10) to optimize the efficiency of the algorithm; using  $E_R = 1$  along with the method described in Example 5, we estimate that  $\hat{\ell}^2 = 3.4$ , for which  $a(\hat{\ell}) = 0.2214$ . The value  $\hat{\ell} = 1.84$  thus differs

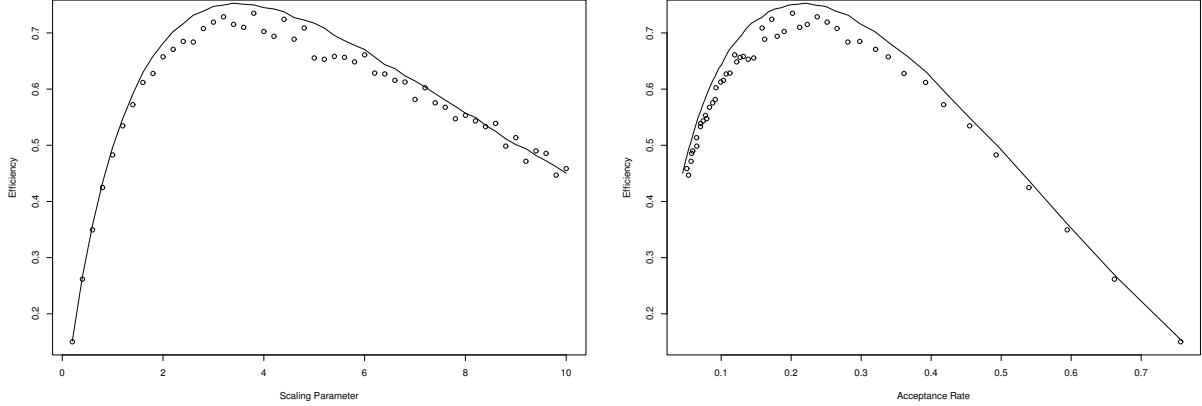


Figure 5: Left graph: efficiency of  $X_3$  versus  $\ell^2$ . Right graph: efficiency of  $X_3$  versus the acceptance rate. The solid line represents the theoretical curve, while the dotted curve is the result of the simulation study.

from the baseline 2.38 in Section 3.1, but still yields an AOAR that is close to 0.234.

Figure 5 presents graphs based on 100,000 iterations of the Metropolis algorithm, depicting how the first order efficiency of  $X_5$  relates to  $\ell^2$  and the acceptance rate respectively. The curves obtained emphasize the rapid convergence of the algorithm in finite dimensions to its asymptotic counterpart, represented by the solid line.

## 4.2 Variance Components Model

The second simulation study focuses on the variance components model. Let  $\mu \sim N(0, 1)$ ,  $\sigma_\theta^2 \sim IG(3, 1)$  and  $\sigma_e^2 \sim IG(2, 1)$ . The means  $\theta_i$  are conditionally IID given  $\mu, \sigma_\theta^2$  and are distributed according to  $\theta_i \sim N(\mu, \sigma_\theta^2)$  for  $i = 1, \dots, 30$ . The 30 groups of data values are conditionally independent given the mean vector  $(\theta_1, \dots, \theta_{30})$  and the variance  $\sigma_e^2$ , while the values within each group are IID. In particular,  $Y_{i,j} \sim N(\theta_i, \sigma_e^2)$  for  $i = 1, \dots, 30$  and  $j = 1, \dots, 10$ .

We are interested in the posterior distribution of  $\mu, \sigma_\theta^2, \sigma_e^2, \theta_1, \dots, \theta_{30}$  given the data  $Y_{i,j}$ ,  $i = 1, \dots, 30$ ,  $j = 1, \dots, 10$ . Since the algorithm does not work well when generating moves from a normal proposal to mimic moves from an inverse gamma distribution, we use inverse transformations and instead deal with gamma distributions. The target is such that

$$\pi(\mu, \varphi_\theta, \varphi_e, \theta_1, \dots, \theta_{30} | \mathbf{Y}) \propto (\varphi_\theta)^{17} (\varphi_e)^{151} \exp\left(-\frac{\mu^2}{2\sigma_0^2} - \varphi_\theta - \varphi_e - \sum_{i=1}^{30} \frac{\varphi_\theta (\theta_i - \mu)^2}{2} - \sum_{i=1}^{30} \sum_{j=1}^{10} \frac{\varphi_e (\theta_i - Y_{i,j})^2}{2}\right), \quad (14)$$

where  $\varphi_\theta = 1/\sigma_\theta^2$  and  $\varphi_e = 1/\sigma_e^2$ .

We run the Metropolis algorithm with a target as in (14), updating the variables  $\mu, \varphi_\theta, \varphi_e, \theta_1, \dots,$

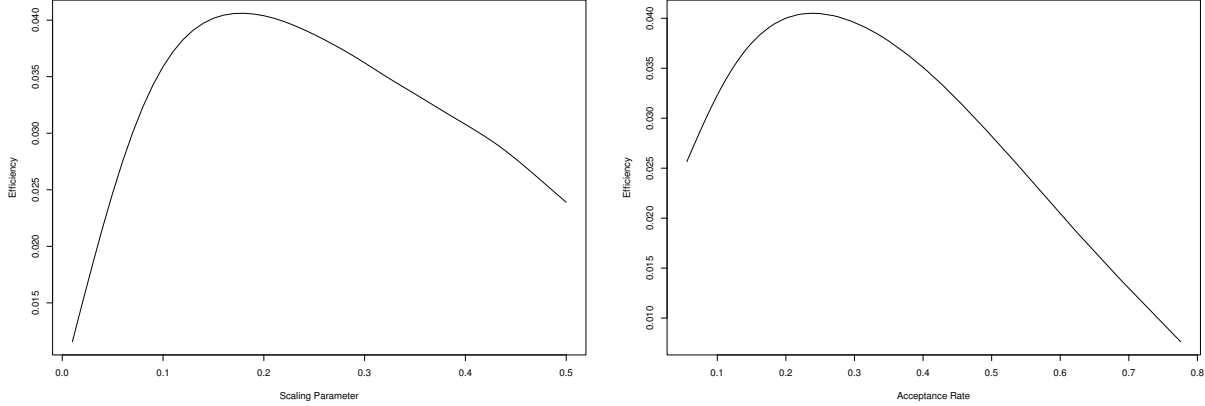


Figure 6: Left graph: efficiency of  $\theta_1$  versus  $\ell^2$ . Right graph: efficiency of  $\theta_1$  versus the acceptance rate.

$\theta_{30}$ . For the sake of the example, the data was simulated from the target. We performed 100,000 iterations and plotted first order efficiency of  $X_4$  versus  $\ell^2$  and the acceptance rate. The maximum is located around 0.17 for  $\hat{\ell}^2$  and choosing an acceptance rate close to 0.2 optimizes the efficiency of the algorithm. Although the AOAR seems to lie close to 0.234, it is hard to tell its exact value from the graph. According to the previous results, we suspect that it might differ from 0.234, which might become clearer when simulating from target distributions possessing a greater number of non-normal components. Although the joint distribution is not normally distributed, it then seems possible to optimize not only hierarchical models where the mean of normally distributed variables is random, but also hierarchical models with more layers and random variances.

### 4.3 Gamma-Gamma Hierarchical Model

Let  $\lambda \sim \Gamma(4, 1)$  and, assuming conditional independence,  $X_i \sim \Gamma(4, \lambda)$  for  $i = 1, \dots, 20$ . The unconditional 21-dimensional target density satisfies

$$\pi(\lambda, x_1, \dots, x_{20}) \propto \lambda^{83} \exp\left(-\lambda\left(1 + \sum_{i=1}^{20} x_i\right)\right) \prod_{i=1}^{20} x_i^3.$$

This time, 10,000,000 iterations of the algorithm were required to reduce Monte Carlo errors and obtain clear curves. Figure 7 shows the existence of a finite value  $\hat{\ell}$  optimizing the efficiency of the method ( $\hat{\ell}^2 = 1.6$ ), resulting in an optimal acceptance rate lying around 0.16. This small AOAR appears to corroborate the discussion at the end of last section. That is, it seems feasible to optimize the efficiency of Metropolis algorithms for general hierarchical target models and this will yield AOARs that are smaller than 0.234.

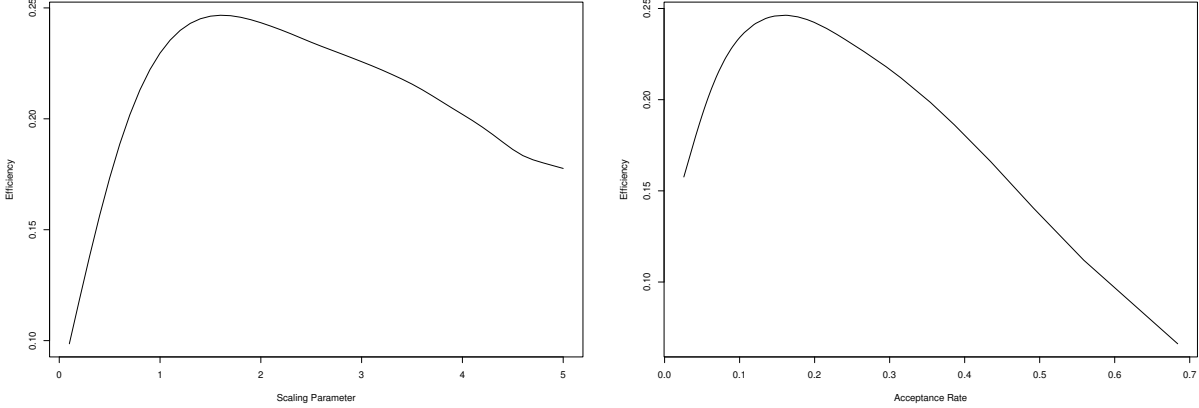


Figure 7: Left graph: efficiency of  $X_1$  versus  $\ell^2$ . Right graph: efficiency of  $X_1$  versus the acceptance rate.

## 5 Discussion

The results presented in this paper permit to optimize the efficiency for sampling from target densities as described in Section 2.1 using Metropolis algorithms with proposal distributions as in Section 2.2. This has been illustrated with numerous examples throughout the paper, which aimed to outline the fact that applying these asymptotic results even to relatively low-dimensional target distributions produced satisfactory conclusions. They also provided evidence that the acceptance rate might considerably differ from 0.234, from where the importance of solving for the correct AOAR. A drastic variation in the AOAR seems to be more common with target distributions that are not normally distributed. In general, AOARs for multivariate normal targets appear to lie close to 0.234, regardless of the correlation structure existing among the components. As discussed in Section 3.3 however, even for the most regular target distributions, an extremely small scaling term causes the algorithm to be inefficient and forces us to resort to inhomogeneous proposal distributions. The AOAR obtained under this method is then not necessarily close to 0.234.

As mentioned previously, since our results can be used to optimize any multivariate normal target distribution, this includes cases where the target is any level normal hierarchical model, meaning that the variances of the distributions are fixed while the mean is random and normally distributed. This raises the question as to whether Metropolis algorithms can be optimized when the variance of the normal is also random, or more generally if similar results can be derived for broader hierarchical models. The examples presented in Sections 4.2 and 4.3 seem to answer this question positively, but AOARs appear to differ from 0.234 with increasing significance as the distribution gets further from normality. The optimization problem for general hierarchical models is presently under investigation (see Bédard, 2006c).

## Acknowledgments

This work is part of my Ph.D. thesis and has been supported by NSERC of Canada. Special thanks are due to my supervisor, Professor Jeffrey S. Rosenthal, without who this work would not have been complete. His expertise, guidance and encouragements have been precious throughout my studies.

## References

Bédard, M. (2006a), "Weak Convergence of Metropolis Algorithms for Non-IID Target Distributions," *in review for publication in Annals of Applied Probability*.

Bédard, M. (2006b), "Optimal Acceptance Rates for Metropolis Algorithms: Moving Beyond 0.234," *submitted for publication in Annals of Statistics*.

Bédard, M. (2006c), "On the Optimization of Metropolis Algorithms for Hierarchical Target Distributions," *in preparation*.

Breyer, L. A., Roberts, G. O. (2000), "From Metropolis to Diffusions: Gibbs States and Optimal Scaling," *Stochastic Processes and their Applications*, **90**, 181-206.

Hastings, W. K. (1970), "Monte Carlo sampling methods using Markov chains and their applications" *Biometrika*, **57**, 97-109.

Metropolis, N., Rosenbluth, A. W., Rosenbluth, M. N., Teller, A. H. & Teller, E. (1953), "Equations of state calculations by fast computing machines," *Journal of Chemical Physics*, **21**, 1087-92.

Neal, P., Roberts, G. O. (2004), "Optimal Scaling for Partially Updating MCMC Algorithms," *to appear in Annals of Applied Probability*.

Pasarica, C., Gelman A. (2003), "Adaptively scaling the Metropolis algorithm using expected squared jumped distance," technical report, Department of Statistics, Columbia University.

Roberts, G. O., Gelman, A., Gilks, W. R. (1997), "Weak Convergence and Optimal Scaling of Random Walk Metropolis Algorithms," *Annals of Applied Probability*, **7**, 110-20.

Roberts, G. O., Rosenthal, J. S. (1998), "Optimal Scaling of Discrete Approximations to Langevin Diffusions," *Journal of the Royal Statistical Society, Series B*, **60**, 255-68.

Roberts, G. O., Rosenthal, J. S. (2001), "Optimal Scaling for various Metropolis-Hastings algorithms," *Statistical Science*, **16**, 351-67.