



A note on geometric ergodicity and
floating-point roundoff error

by

Gareth O. Roberts

Department of Mathematics and Statistics
Lancaster University

and

Laird Breyer

Department of Mathematics and Statistics
Lancaster University

and

Jeffrey S. Rosenthal
Department of Statistics
University of Toronto

Technical Report No. 2004, August 8, 2000.

TECHNICAL REPORT SERIES

University of Toronto

Department of Statistics

A note on geometric ergodicity and floating-point roundoff error

by

Laird Breyer*, Gareth O. Roberts*, and Jeffrey S. Rosenthal**

(August 8, 2000.)

Abstract. We consider the extent to which Markov chain convergence properties are affected by the presence of computer floating-point roundoff error. Both geometric ergodicity and polynomial ergodicity are considered. This paper extends previous work of Roberts, Rosenthal, and Schwartz (1998); connections between that work and the present paper are discussed.

1. Introduction.

Geometric ergodicity is an important concept in convergence of Markov chains to their stationary distributions. For example, this property is used to justify the applicability of the central limit theorem to ergodic averages along the path of the chain. When run on an actual computer, Markov chains are subject to floating-point roundoff errors. This paper considers the extent to which geometric (and other) ergodicity is affected by small roundoff errors.

A Markov chain on a state space \mathcal{X} , with transition probabilities $P(x, \cdot)$ and stationary distribution $\pi(\cdot)$, is said to be *geometrically ergodic* if there is $\rho < 1$ and $M : \mathcal{X} \rightarrow [0, \infty)$ such that

$$\|P^n(x, \cdot) - \pi(\cdot)\| \leq M(x) \rho^n,$$

where

$$\|P^n(x, \cdot) - \pi(\cdot)\| \equiv \sup_{A \subseteq \mathcal{X}} |P^n(x, A) - \pi(A)|$$

* Department of Mathematics and Statistics, Fylde College, Lancaster University, Lancaster, LA1 4YF, England. Internet: g.o.roberts@lancaster.ac.uk and l.breyer@lancaster.ac.uk.

** Department of Statistics, University of Toronto, Toronto, Ontario, Canada M5S 3G3. Internet: jeff@math.toronto.edu. Supported in part by NSERC of Canada.

is the total variation distance between the law of the Markov chain after n steps (when started at the point $x \in \mathcal{X}$), and the stationary distribution $\pi(\cdot)$.

This is known (cf. Roberts and Rosenthal, 1997) to be equivalent to the existence of $\lambda < 1$, $b < \infty$, and a small set $C \subseteq \mathcal{X}$ such that

$$PV(x) \leq \lambda V(x) + b \mathbf{1}_C(x), \quad x \in \mathcal{X}, \quad (1)$$

where $PV(x) = \int V(y) P(x, dy)$. (Recall that a set is *small* for a Markov chain if there exists a positive integer n_0 , a positive constant ϵ , and a probability measure ν on \mathcal{X} , such that $P^{n_0}(x, \cdot) \geq \epsilon \nu(A)$ for all $x \in C$ and $A \subseteq \mathcal{X}$.)

Roberts, Rosenthal, and Schwartz (1998) considered issues related to running such a Markov chain on a computer, and in particular the effect of various roundoff errors during the simulation. They introduced a summary roundoff function $h : \mathcal{X} \rightarrow \mathcal{X}$, with $h(x)$ “close” to x for all x . This leads to a modified Markov chain \tilde{P} given by

$$\tilde{P}(x, A) = P(x, h^{-1}(A)). \quad (2)$$

In this framework, the assumption of small computer errors can be taken as

$$\|h(x) - x\| \leq \delta, \quad x \in \mathcal{X}, \quad (3)$$

where $\|x\|$ is the norm of $x \in \mathcal{X}$. (We assume throughout that \mathcal{X} is a normed vector space over \mathbf{R} , e.g. $\mathcal{X} \subseteq \mathbf{R}^d$.) It is shown by Roberts et al. (1998) that, if P is geometrically ergodic with drift function V such that $\log V$ is uniformly continuous, and δ is sufficiently small, then \tilde{P} will also be geometrically ergodic. That is, geometric ergodicity is preserved under small perturbations in that case.

A common case not included in the above arises when instead we have merely

$$\|h(x) - x\| \leq \delta \|x\|, \quad x \in \mathcal{X}, \quad (4)$$

as may occur with floating-point computations (cf. Section 2 below). That is, the roundoff errors may have magnitude proportional to the magnitude of x , rather than being uniformly bounded. In this note, we shall show that this case is amenable to a technique similar to that of Roberts et al. (1998).

As motivation, in Section 2, we introduce the IEEE standard floating point representation for real numbers used in modern computers. Section 3 considers the robustness of

geometric ergodicity to perturbations of Markov chain dynamics which are similar to those caused by computer roundoff errors. Section 4 then investigates connections between these robustness properties and those given in Roberts et al. (1998). Some simple examples are described in Section 5, and in Section 6 the methodology is extended to consider robustness of perturbations to polynomial ergodicity.

We note that the results we present are only the beginning of a rigorous analysis of how computer engineering realities affect the dynamics of mathematically specified Markov chains. We hope to pursue such questions more comprehensively in the future.

Remark. Even if a modified Markov chain \tilde{P} is proven to converge as quickly as the original chain, there is still the question of what the new target distribution is. Roberts et al. (1998) investigate this issue in total variation distance and in the weak topology. The methods in our paper could also be extended to consider this issue, but we do not pursue that here.

2. Floating point representations in computers.

IEEE standard 754 (IEEE, 1985) is a specification commonly adhered to for the representation of floating point numbers in computers, using a fixed number B of bits (e.g. $B = 32$ with single precision numbers or $B = 64$ with double precision numbers). Mathematically, numbers x are encoded using $B = M + N + 1$ bits to a finite precision, in the following way (called normalized floating point representation):

$$x := \sigma \cdot (1 + k/2^N) \cdot 2^e,$$

where $\sigma = \pm 1$ is the sign (1 bit), $k \in \{0, \dots, 2^N - 1\}$ is the fractional part (N bits), and $e \in \{-2^{M-1} + 2, \dots, 2^{M-1} - 1\}$ is the exponent (M bits). Single precision numbers use $M = 8$ and $N = 23$, giving an effective range (excluding the sign) of $2^{-126} \approx 10^{-44.85}$ to $(2 - 2^{-23}) \cdot 2^{127} \approx 10^{38.53}$, while double precision is represented by $M = 11$, $N = 52$, with an effective range (excluding sign) of $2^{-1022} \approx 10^{-323.3}$ to $(2 - 2^{-52}) \cdot 2^{1023} \approx 10^{308.3}$. Numbers larger than this are represented by the special symbol **+Infinity**, which is often encoded as $\sigma = 0$, $k = 0$, $e = 2^{M-1}$. Moreover, the number zero is nonunique; more precisely there exist two distinct values $+0$ and -0 which are only equal when compared directly.

Clearly, not all real numbers x can be represented with a fixed number B of bits in this way. Indeed, given a real number x , setting $e = \lfloor \log_2 |x| \rfloor$ and $\sigma = \text{sign}(x)$, the closest

the computer can come to approximating x is as

$$h(x) = \sigma \lfloor \frac{1}{2} + |x|2^{N-1-e} \rfloor 2^{-(N-1-e)}.$$

It follows that

$$|h(x) - x| \leq 2^{-(N-1-e)} \leq 2^{-(N-1)}|x|.$$

We see from the above that the error $|h(x) - x|$ is proportional to $|x|$, thus violating (3). (Strictly speaking, (3) holds for a sufficiently large δ since $|x|$ is bounded by the finite range of the computer. However, in the present paper, we ignore issues related to finite range, and concentrate solely on issues related to finite precision, i.e. to roundoff errors. It is that sense that (3) is violated.) However, the assumption (4) does hold here with $\delta = 2^{-(N-1)}$.

With regard to Markov chain algorithms and their implementations on computer systems, we shall therefore assume that the final error for each update behaves as (4). Of course, this is meant as a convenient summary of the cumulative effect of various complicated roundoff errors introduced at each stage of the update calculation. (For example, a side effect of using floating point representations is that the corresponding arithmetic becomes inexact and non-commutative, e.g. perhaps $(x \cdot y)/y \neq x$ or $x + y \neq y + x$.)

3. Geometric ergodicity under perturbations satisfying (4)

Suppose a Markov chain P is geometrically ergodic, thus satisfying (1) for some function V and small set C . Suppose further that \tilde{P} is obtained via (2), for some roundoff function h satisfying (4) for some $\delta > 0$. Assume also that V satisfies

$$V(y+u) - V(y) \leq \delta K V(y), \quad \|u\| \leq \delta \|y\|, \quad y \in \mathcal{X}, \quad (5)$$

for some $K < \infty$. (Of course, we could subsume the product of δ and K into a single constant, but our notation better emphasizes the dependence upon δ .)

Note that the condition (5) is implied, if \mathcal{X} is finite-dimensional and $V(x)$ is continuously differentiable, by

$$\|\nabla \log V(y)\| \leq K' / \|y\|, \quad (6)$$

where $K' = \delta^{-1} \log(1 + K\delta) \approx K$.

Proposition 1. If (1) and (5) hold, and if \tilde{P} is derived from P via (2), where h satisfies (4), then

$$\tilde{P}V(x) \leq (1 + \delta K)(\lambda V(x) + b \mathbf{1}_C(x)).$$

Proof. We have that

$$\begin{aligned}
\tilde{P}V(x) &= PV(x) + (\tilde{P} - P)V(x) \\
&= PV(x) + \int (V(h(y)) - V(y)) P(x, dy) \\
&\leq \lambda V(x) + b \mathbf{1}_C(x) + \int \delta K V(y) P(x, dy) \\
&\leq \lambda V(x) + b \mathbf{1}_C(x) + \delta K (\lambda V(x) + b \mathbf{1}_C(x)),
\end{aligned}$$

which gives the result. ■

This ensures geometric ergodicity provided that $(1 + \delta K)\lambda < 1$, or equivalently

$$\delta < K^{-1}(\lambda^{-1} - 1). \tag{7}$$

We thus obtain

Theorem 2. Suppose a Markov chain P is geometrically ergodic, satisfying (1) for some V and C . Assume that V satisfies (5) for some $K < \infty$. Suppose further that \tilde{P} is obtained via (2), for some perturbation function h satisfying (4) and (7). Then \tilde{P} is also geometrically ergodic.

This theorem thus proves that geometric ergodicity is preserved, under sufficiently small floating-point-type perturbations, provided that the drift function V satisfies the smoothness condition (5) (or (6)).

4. Connection between perturbations of type (3) and (4).

It is shown by Roberts et al. (1998) that, assuming cumulative roundoff error is governed by (3), geometric ergodicity is preserved for sufficiently small δ provided only that $\log V$ is uniformly continuous.

By contrast, we have shown in Theorem 2 above that, if the cumulative roundoff error is instead governed by (4), then geometric ergodicity is preserved if the gradient of $\log V$ decays sufficiently fast, i.e. (5) or (6) holds.

The stronger condition required when (4) holds instead of (3) is not surprising, since the absolute magnitude of errors is unbounded in this case, thus giving arbitrarily large perturbations of the Markov chain.

Our aim in this section is to connect Theorem 2 above with results of Roberts et al. (1998). We shall need the following lemma, whose proof is straightforward.

Lemma 3. Let $\{X_t\}$ be a geometrically ergodic Markov chain with state space \mathcal{X} , which satisfies (1) for some drift function V and small set C . Let $\varphi : \mathcal{X} \rightarrow \mathcal{X}'$ be a bi-measurable bijection, and define a Markov chain $\{X_t^\varphi\}$ on \mathcal{X}' by $X_t^\varphi = \varphi(X_t)$, with corresponding transition kernel P^φ . Then $\varphi(C)$ is a small set for $\{X_t^\varphi\}$, and furthermore

$$P^\varphi V^\varphi(x') \leq \lambda V^\varphi(x') + b \mathbf{1}_{\varphi(C)}(x'), \quad x' \in \mathcal{X}', \quad (8)$$

where

$$V^\varphi(x') = V(\varphi^{-1}(x')). \quad (9)$$

Hence, $\{X_t^\varphi\}$ is also geometrically ergodic with the same constants. Indeed, for all positive integers n , $\|(P^\varphi)^n(\varphi(x), \cdot) - \pi^\varphi\| = \|P^n(x, \cdot) - \pi\|$, where $\pi^\varphi(dx') = \pi(\varphi^{-1}(dx'))$ is stationary for $\{X_t^\varphi\}$.

For a concrete example of Lemma 3, recall that a *Metropolis-Hastings algorithm* with target density $\pi(x)$ and proposal kernel $q(x, y)$ is a Markov chain X_t which can be constructed by the recurrence

$$X_{t+1} = \begin{cases} Y_t, & \pi(Y_t)q(Y_t, X_t)/\pi(X_t)q(X_t, Y_t) > \xi_t, \\ X_t, & \text{otherwise,} \end{cases} \quad (10)$$

where $Y_t \sim q(X_t, \cdot)$, and where the ξ_t are independently chosen as i.i.d. Uniform $[0, 1]$. The chain's transition kernel is given by

$$P(x, dy) = q(x, y) \left(1 \wedge \frac{\pi(y)q(y, x)}{\pi(x)q(x, y)} \right) dy + r(x)\delta_x(dy),$$

where $r(x)$ is chosen to make $P(x, \cdot)$ into a probability measure. Let X_t be a Metropolis-Hastings chain with state space \mathcal{X} , target density $\pi(x)$ and proposal density $q(x, y)$, and let $\varphi(x) : \mathcal{X} \rightarrow \mathcal{X}'$ be a bijection. Then it is easily verified that $X_t^\varphi := \varphi(X_t)$ is a Metropolis-Hastings chain with state space \mathcal{X}' , target probability measure $\pi^\varphi(dx') = \pi(\varphi^{-1}(dx'))$, and proposal density

$$q^\varphi(x', y') = q(\varphi^{-1}(x'), \varphi^{-1}(y')) \cdot |J_\varphi(\varphi^{-1}(y'))|$$

i.e. if $Y \sim q(x, \cdot)$, then

$$Y^\varphi := \varphi(Y) \sim q^\varphi(\varphi(x'), \cdot).$$

Here, J_φ denotes the Jacobian of the transformation φ .

Using Lemma 3, it is now possible to connect the stability results for perturbations of types (3) and (4) respectively.

Theorem 4. Suppose that \tilde{X} is a perturbation of X with approximation function h . Then $\varphi(\tilde{X})$ is a perturbation of X^φ with approximation function $h^\varphi \equiv \varphi h \varphi^{-1}$. In the special case that $\mathcal{X} = \mathbf{R}$, $\mathcal{X}' = \mathbf{R}^+$, and $\varphi(x) = \exp(x)$, then h satisfies (3) if and only if h^{\exp} satisfies (4). In that case, if P^{\exp} is geometrically ergodic, then it is robust to perturbations of the kind satisfying (4) provided it has a drift function V^{\exp} satisfying that $\log V^{\exp}(\exp(\cdot))$ is uniformly continuous.

Proof. The first statement is straightforward, and the equivalence of (3) for h and (4) for h^{\exp} (with two distinct values δ) is a straightforward computation.

For the final statement, we note that by Proposition 6 of Roberts, Rosenthal and Schwartz (1998), the geometric ergodicity of P is robust to perturbations by h satisfying (3) so long as its drift function V satisfies $\log V$ is uniformly continuous. However, by (9), this is equivalent to $\log V^{\exp}(\exp(\cdot))$ being uniformly continuous. ■

5. Examples.

In this section, we give three examples to illustrate our results.

Example 1. For a first example, suppose that $\mathcal{X} = \mathbf{R}$ and $V(x) = C_1|x|^n + C_2$, with $C_1, C_2 \geq 0$. In that case, for $\delta' \leq \delta$,

$$V(y + \delta'y) - V(y) = C_1(|y + \delta'y|^n - |y|^n) \leq C_1(1 + \delta')^n|y|^n \leq (1 + \delta)^n V(y).$$

Hence, (5) holds with $K = (1 + \delta)^n/\delta$. We thus conclude that, for a Markov chain which is geometrically ergodic with drift function of the form $C_1|x|^n + C_2$, sufficiently small roundoff errors of the form (4) still preserve geometric ergodicity.

Example 2. On the other hand, if the error summary function is of the type (4) but not (3), then many rounded-off Markov chains as in (2) can fail to be geometrically ergodic. For example, consider the Gaussian Random Walk Metropolis algorithm (that is an algorithm satisfying (10) with a random walk proposal $Y_t = X_t + \sigma W_t$ with Gaussian increment $W_t \sim \mathcal{N}(0, 1)$) with Gaussian target distribution (zero mean, unit variance) and Gaussian proposal increment with fixed standard deviation $\sigma \leq 1$.

drift function V which satisfies (14) for some constants $\gamma \leq \min(1, \alpha)$ and $\epsilon > 0$. Define \tilde{P} by (2), and assume that (13) holds for some $\beta \leq \epsilon$ and $c \leq \delta$. Then

$$\tilde{P}V \leq V - aV^\alpha + b' \mathbf{1}_C + cKV^\alpha, \quad (15)$$

for some $b' < \infty$. In particular, if $c < a/K$, then the chain defined by \tilde{P} is also polynomially ergodic, with the same polynomial rate α .

Proof. We compute that

$$\begin{aligned} \tilde{P}V(x) &= PV(x) + (\tilde{P} - P)V(x) \\ &= PV(x) + \int P(x, dy)(V(h(y)) - V(y)) \\ &\leq PV(x) + \delta K(P(V^\gamma))(x) \\ &\leq PV(x) + \delta K(PV)^\gamma(x) \\ &\leq V(x) - aV^\alpha(x) + b \mathbf{1}_C(x) + cK(V(x) - aV^\alpha(x) + b \mathbf{1}_C(x))^\gamma \\ &\leq V(x) - aV^\alpha(x) + b \mathbf{1}_C(x) + cK(V(x) + b \mathbf{1}_C(x))^\gamma \\ &\leq V(x) - aV^\alpha(x) + b \mathbf{1}_C(x) + cK(V^\gamma(x) + b^\gamma \mathbf{1}_C(x)) \\ &\leq V(x) - aV^\alpha(x) + b' \mathbf{1}_C(x) + cKV^\alpha(x), \end{aligned}$$

as claimed. Here the first inequality combines (13) and (14), the second inequality is Jensen's inequality, and $b' = b + b^\gamma$. For the final inequality, we have used that $\gamma \leq \alpha$ and $V \geq 1$ so that $V^\gamma \leq V^\alpha$. For the second-last inequality, we have used the general fact that $(A + B)^\gamma \leq A^\gamma + B^\gamma$ when $A, B, \gamma \geq 0$ and $\gamma \leq 1$; to see that set $f(x) = (A + x)^\gamma - x^\gamma$ and note that $f'(x) \leq 0$ for $x > 0$ so that $f(B) \leq f(0)$. \blacksquare

REFERENCES

IEEE (1985), ANSI/IEEE Standard 754-1985. In *IEEE Standard for Binary Floating-Point Arithmetic*, IEEE, New York.

K.L. Mengersen and Tweedie, R.L. (1996), Rates of convergence of the Hastings and Metropolis algorithms. *Annals of Statistics* **24**, 101–12.

G.O. Roberts (1999), A note on acceptance rate criteria for CLTs for Metropolis-Hastings algorithms. *Journal of Applied Probability* **36**, 1210–1217.

G.O. Roberts and J.S. Rosenthal (1997), Geometric ergodicity and hybrid Markov chains. *Electronic Communications in Probability* **2**, Paper no. 2, 13–25.

G.O. Roberts, J.S. Rosenthal, and P.O. Schwartz (1998), Convergence properties of perturbed Markov chains. *Journal of Applied Probability* **35**, 1–11.

G.O. Roberts and R.L. Tweedie (1996), Geometric convergence and central limit theorems for multidimensional Hastings and Metropolis algorithms. *Biometrika* **83**, no. 1, 95–110.

