

The Bouncy Particle Sampler: A Non-Reversible Rejection-Free Markov Chain Monte Carlo Method

Alexandre Bouchard-Côté*, Sebastian J. Vollmer† and Arnaud Doucet†

18th January 2017

*Department of Statistics, University of British Columbia, Canada.

†Department of Statistics, University of Oxford, UK.

1 Direct proof of invariance

Let μ_t be the law of $z(t)$. In the following, we prove invariance by explicitly verifying that the time evolution of the density $\frac{d\mu_t}{dt} = 0$ is zero if the initial distribution μ_0 is given by $\rho(z) = \pi(x)\psi(v)$ in Proposition 1. This is achieved by deriving the forward Kolmogorov equation describing the evolution of the marginal density of the stochastic process. For simplicity, we start by presenting the invariance argument when $\lambda^{\text{ref}} = 0$.

Notation and description of the algorithm. We denote a pair of position and velocity by $z = (x, v) \in \mathbb{R}^d \times \mathbb{R}^d$ and we denote translations by $\Phi_t(z) = (\Phi_t^{\text{pos}}(z), \Phi_t^{\text{dir}}(z)) = (x + vt, v)$. The time of the first bounce coincides with the first arrival T_1 of a PP with intensity $\chi(t) = \lambda(\Phi_t(z))$ where:

$$\lambda(z) = \max\{0, \langle \nabla U(x), v \rangle\}. \quad (1)$$

It follows that the probability of having no bounce in the interval $[0, t]$ is given by:

$$\text{No}_t(z) = \exp\left(-\int_0^t \lambda(\Phi_s(z)) ds\right), \quad (2)$$

and the density of the random variable T_1 is given by:

$$q(t_1; z) = \mathbf{1}[t_1 > 0] \frac{d}{dt_1} (1 - \text{No}_{t_1}(z)) \quad (3)$$

$$= \mathbf{1}[t_1 > 0] \text{No}_{t_1}(z) \lambda(\Phi_{t_1}(z)). \quad (4)$$

If a bounce occurs, then the algorithm follows a translation path for time T_1 , at which point the velocity is updated using a bounce operation $C(z)$, defined as:

$$C(z) = (x, R(x)v) \quad (5)$$

where

$$R(x)v = v - 2 \frac{\langle \nabla U(x), v \rangle \nabla U(x)}{\|\nabla U(x)\|^2}. \quad (6)$$

The algorithm then continues recursively for time $t - T_1$, in the following sense: a second bounce time T_2 is simulated by adding to T_1 a random increment with density $q(\cdot; C \circ \Phi_{t_1}(z))$. If $T_2 > t$, then the output of the algorithm is $\Phi_{t-T_1} \circ C \circ \Phi_{t_1}(z)$, otherwise an additional bounce is simulated,

etc. More generally, given an initial point z and a sequence $\mathbf{t} = (t_1, t_2, \dots)$ of bounce times, the output of the algorithm at time t is given by:

$$\Psi_{\mathbf{t},t}(z) = \begin{cases} \Phi_t(z) & \text{if } t_1 > 0 \text{ or } \mathbf{t} = (), \\ \Psi_{\mathbf{t}',t-t_1}(z) \circ C \circ \Phi_{t_1}(z) & \text{otherwise,} \end{cases} \quad (7)$$

where $()$ denotes the empty list and \mathbf{t}' the suffix of \mathbf{t} : $\mathbf{t}' = (t_2, t_3, \dots)$. As for the bounce times, they are distributed as follows:

$$T_1 \sim q(\cdot; z) \quad (8)$$

$$T_i - T_{i-1} | T_{1:i-1} \sim q\left(\cdot; \underbrace{\Psi_{T_{1:i-1}, T_{i-1}}(z)}_{\text{Pos. after collision } i-1}\right), \quad i \in \{2, 3, 4, \dots\} \quad (9)$$

where $T_{1:i-1} = (T_1, T_2, \dots, T_{i-1})$.

Decomposition by the number of bounces. Let h denote an arbitrary non-negative measurable test function. We show how to decompose expectations of the form $\mathbb{E}[h(\Psi_{\mathbf{T},t}(z))]$ by the number of bounces in the interval $(0, t)$. To do so, we introduce a function $\#\text{Col}_t(\mathbf{t})$, which returns the number of bounces in the interval $(0, t)$:

$$\#\text{Col}_t(\mathbf{t}) = \min\{n \geq 1 : t_n > t\} - 1. \quad (10)$$

From this, we get the following decomposition:

$$\mathbb{E}[h(\Psi_{\mathbf{T},t}(z))] = \mathbb{E}[h(\Psi_{\mathbf{T},t}(z)) \sum_{n=0}^{\infty} \mathbf{1}[\#\text{Col}_t(\mathbf{T}) = n]] \quad (11)$$

$$= \sum_{n=0}^{\infty} \mathbb{E}[h(\Psi_{\mathbf{T},t}(z)) \mathbf{1}[\#\text{Col}_t(\mathbf{T}) = n]]. \quad (12)$$

On the event that no bounce occurs in the interval $[0, t)$, i.e. $\#\text{Col}_t(\mathbf{T}) = 0$, the function $\Psi_{\mathbf{T},t}(z)$ is equal to $\Phi_t(z)$, therefore:

$$\mathbb{E}[h(\Psi_{\mathbf{T},t}(z)) \mathbf{1}[\#\text{Col}_t(\mathbf{T}) = 0]] = h(\Phi_t(z)) \mathbb{P}(\#\text{Col}_t(\mathbf{T}) = 0) \quad (13)$$

$$= h(\Phi_t(z)) \text{No}_t(z). \quad (14)$$

Indeed, on the event that $n \geq 1$ bounces occur, the random variable $h(\Phi_t(z))$ only depends on a finite dimensional random vector, (T_1, T_2, \dots, T_n) , so we can write the expectation as an integral with respect to the density $\tilde{q}(t_{1:n}; t, z)$ of these variables:

$$\mathbb{E}[h(\Psi_{\mathbf{T},t}(z)) \mathbf{1}[\#\text{Col}_t(\mathbf{T}) = n]] \quad (15)$$

$$= \mathbb{E}[h(\Psi_{\mathbf{T},t}(z)) \mathbf{1}[0 < T_1 < \dots < T_n < t < T_{n+1}]]$$

$$= \int \dots \int_{0 < t_1 < \dots < t_n < t < t_{n+1}} h(\Psi_{t_{1:n},t}(z)) q(t_1; z) \prod_{i=2}^{n+1} q(t - t_{i-1}; \Psi_{t_{1:i-1}, t_{i-1}}(z)) dt_{1:n+1}$$

$$= \int \dots \int_{0 < t_1 < \dots < t_n < t} h(\Psi_{t_{1:n},t}(z)) \tilde{q}(t_{1:n}; t, z) dt_{1:n}, \quad (16)$$

where:

$$\tilde{q}(t_{1:n}; t, z) = q(t_1; z) \times \begin{cases} \text{No}_{t-t_1}(\Phi_{t_1}(z)) & \text{if } n = 1 \\ \text{No}_{t-t_n}(\Phi_{t_{1:n},t_n}(z)) \prod_{i=2}^n q(t_i - t_{i-1}; \Psi_{t_{1:i-1}, t_{i-1}}(z)) & \text{if } n \geq 2. \end{cases}$$

To include Equations (14) and (16) under the same notation, we define $t_{1:0}$ to the empty list, $()$, $\tilde{q}((); t, z) = \text{No}_t(z)$, and abuse the integral notation so that for all $n \in \{0, 1, 2, \dots\}$:

$$\mathbb{E}[h(\Psi_{\mathbf{T},t}(z)) \mathbf{1}[\#\text{Col}_t(\mathbf{T}) = n]] = \int \dots \int_{0 < t_1 < \dots < t_n < t} h(\Psi_{t_{1:n},t}(z)) \tilde{q}(t_{1:n}; t, z) dt_{1:n}. \quad (17)$$

Marginal density. Let us fix some arbitrary time $t > 0$. We seek a convenient expression for the marginal density at time t , $\mu_t(z)$, given an initial vector $Z \sim \rho$, where ρ is the hypothesized stationary density $\rho(z) = \pi(x)\psi(v)$ on Z . To do so, we look at the expectation of an arbitrary non-negative measurable test function h :

$$\mathbb{E}[h(\Psi_{\mathbf{T},t}(Z))] = \mathbb{E}\left[\mathbb{E}[h(\Psi_{\mathbf{T},t}(Z))|Z]\right] \quad (18)$$

$$= \sum_{n=0}^{\infty} \mathbb{E}\left[\mathbb{E}[h(\Psi_{\mathbf{T},t}(Z))\mathbf{1}_{\{\#\text{Col}_t(\mathbf{T})=n\}}|Z]\right] \quad (19)$$

$$= \sum_{n=0}^{\infty} \int_{\mathcal{Z}} \rho(z) \int \cdots \int_{0 < t_1 < \cdots < t_n < t} h(\Psi_{t_{1:n},t}(z)) \tilde{q}(t_{1:n}; t, z) dt_{1:n} dz \quad (20)$$

$$= \sum_{n=0}^{\infty} \int \cdots \int_{0 < t_1 < \cdots < t_n < t} \int_{\mathcal{Z}} \rho(z) h(\Psi_{t_{1:n},t}(z)) \tilde{q}(t_{1:n}; t, z) dz dt_{1:n} \quad (21)$$

$$= \sum_{n=0}^{\infty} \int \cdots \int_{0 < t_1 < \cdots < t_n < t} \int_{\mathcal{Z}} \rho(\Psi_{t_{1:n},t}^{-1}(z')) h(z') \tilde{q}(t_{1:n}; t, \Psi_{t_{1:n},t}^{-1}(z')) \left| \det D\Psi_{t_{1:n},t}^{-1} \right| dz' dt_{1:n}$$

$$= \underbrace{\int_{\mathcal{Z}} h(z') \sum_{n=0}^{\infty} \int \cdots \int_{0 < t_1 < \cdots < t_n < t} \rho(\Psi_{t_{1:n},t}^{-1}(z')) \tilde{q}(t_{1:n}; t, \Psi_{t_{1:n},t}^{-1}(z')) dt_{1:n} dz'}_{\mu_t(z')} \quad (22)$$

We used the following in the above derivation successively the law of total expectation, equation (12), equation (18), Tonelli's theorem and the change of variables, $z' = \Psi_{t_{1:n},t}(z)$, justified since for any fixed $0 < t_1 < t_2 < \cdots < t_n < t < t_{n+1}$, $\Psi_{t_{1:n},t}(\cdot)$ is a bijection (being a composition of bijections). Now the absolute value of the determinant is one since $\Psi_{t,t}(z)$ is a composition of unit-Jacobian mappings and, by using Tonelli's theorem again, we obtain that the expression above the brace is necessarily equal to $\mu_t(z')$ since h is arbitrary.

Derivative. Our goal is to show that for all $z' \in \mathcal{Z}$

$$\frac{d\mu_t(z')}{dt} = 0.$$

Since the process is time homogeneous, once we have computed the derivative, it is enough to show that it is equal to zero at $t = 0$. To do so, we decompose the computation according to the terms I_n in Equation (22):

$$\mu_t(z') = \sum_{n=0}^{\infty} I_n(z', t) \quad (23)$$

$$I_n(z', t) = \int \cdots \int_{0 < t_1 < \cdots < t_n < t} \rho(\Psi_{t_{1:n},t}^{-1}(z')) \tilde{q}(t_{1:n}; t, \Psi_{t_{1:n},t}^{-1}(z')) dt_{1:n}. \quad (24)$$

The categories of terms in Equation (23) to consider are:

No bounce: $n = 0$, $\Psi_{t_{1:n},t}(z) = \Phi_t(z)$, or,

Exactly one bounce: $n = 1$, $\Psi_{t_{1:n},t}(z) = F_{t,t_1} := \Phi_{t-t_1} \circ C \circ \Phi_{t_1}(z)$ for some $t_1 \in (0, t)$, or,

Two or more bounces: $n \geq 2$, $\Psi_{t_{1:n},t}(z) = \Psi_{t-t_2} \circ C \circ F_{t_2,t_1}(z)$ for some $0 < t_1 < t_2 < t$

In the following, we show that the derivative of the terms in the third category, $n \geq 2$, are all equal to zero, while the derivative of the first two categories cancel each other.

No bounce in the interval. From Equation (14):

$$I_0(z', t) = \rho(\Phi_{-t}(z')) \text{No}_t(\Phi_{-t}(z')). \quad (25)$$

We now compute the derivative at zero of the above expression:

$$\begin{aligned} \left. \frac{d}{dt} I_0(z', t) \right|_{t=0} &= \text{No}_0(\Phi_0(z')) \left. \frac{d\rho(\Phi_{-t}(z'))}{dt} \right|_{t=0} + \\ &\quad \rho(\Phi_0(z')) \left. \frac{d\text{No}_t(\Phi_{-t}(z'))}{dt} \right|_{t=0} \end{aligned} \quad (26)$$

The first term in the above equation can be simplified as follows:

$$\text{No}_0(\Phi_0(z')) \frac{d\rho(\Phi_{-t}(z'))}{dt} = \frac{d\rho(\Phi_{-t}(z'))}{dt} \quad (27)$$

$$= \left\langle \frac{\partial \rho(\Phi_{-t}(z'))}{\partial \Phi_{-t}^{\text{pos}}(z')}, \frac{d\Phi_{-t}^{\text{pos}}(z')}{dt} \right\rangle + \left\langle \frac{\partial \rho(\Phi_{-t}(z'))}{\partial \Phi_{-t}^{\text{dir}}(z')}, \underbrace{\frac{d\Phi_{-t}^{\text{dir}}(z')}{dt}}_{=0} \right\rangle \quad (28)$$

$$= \left\langle \frac{\partial \rho(z)}{\partial x}, -v' \right\rangle \quad (29)$$

$$= \left\langle \frac{\partial}{\partial x} \frac{1}{Z} \exp(-U(x)) \psi(v), -v' \right\rangle = \rho(\Phi_{-t}(z')) \langle \nabla U(x), v' \rangle, \quad (30)$$

where $x = \Phi_{-t}^{\text{pos}}(z')$. The second term in Equation (26) is equal to:

$$\rho(\Phi_0(z')) \left. \frac{d\text{No}_t(\Phi_{-t}(z'))}{dt} \right|_{t=0} = -\rho(\Phi_0(z')) \text{No}_0(z') \lambda(\Phi_0(z')) \quad (31)$$

$$= -\rho(z') \lambda(z'), \quad (32)$$

using Equation (4). In summary, we have:

$$\left. \frac{d}{dt} I_0(z', t) \right|_{t=0} = \rho(z') \langle \nabla U(x'), v' \rangle - \rho(z') \lambda(z').$$

Exactly one bounce in the interval. From Equation (16), the trajectory consists in a bounce at a time T_1 , occurring with density (expressed as before as a function of the final point z') $q(t_1; F_{t_1}^{-1}(z'))$, followed by no bounce in the interval $(T_1, t]$, an event of probability:

$$\text{No}_{t-t_1}(C \circ \Phi_{t_1}(z)) = \text{No}_{t-t_1}(C \circ \Phi_{t_1} \circ F_{t_1}^{-1}(z')) \quad (33)$$

$$= \text{No}_{t-t_1}(\Phi_{t_1-t}(z')), \quad (34)$$

where we used that $C^{-1} = C$. This yields:

$$I_1(z', t) = \int_0^t q(t_1; F_{t_1}^{-1}(z')) \rho(\Psi_{t_1, t}^{-1}(z')) \text{No}_{t-t_1}(\Phi_{t_1-t}(z')) dt_1. \quad (35)$$

To compute the derivative of the above equation at zero, we use again Leibniz's rule:

$$\left. \frac{d}{dt} I_1(z', t) \right|_{t=0} = \rho(C(z')) \lambda(C(z')).$$

Two or more bounces in the interval. For a number of bounce, we get:

$$I_n(z', t) = \int_0^t \left[\underbrace{\int \cdots \int_{t_2:n:t_1 < t_2 < \cdots < t_n < t} \rho(\Psi_{t_1:n, t}^{-1}(z')) \tilde{q}(t_1:n; t, \Psi_{t_1:n, t}^{-1}(z')) dt_2:n}_{\tilde{I}(t_1, t, z')} \right] dt_1, \quad (36)$$

and hence, using Leibniz's rule on the integral over t_1 :

$$\left. \frac{d}{dt} I_n(z', t) \right|_{t=0} = \tilde{I}(0, 0, z') = 0. \quad (37)$$

Putting all terms together. Putting everything together, we obtain:

$$\left. \frac{d\mu_t(z')}{dt} \right|_{t=0} = \rho(z') \langle \nabla U(x'), v' \rangle - \underbrace{\rho(z') \lambda(z') + \rho(C(z')) \lambda(C(z'))}_{=0}. \quad (38)$$

From the expression of $\lambda(\cdot)$, we can rewrite the two terms above the brace as follows:

$$\begin{aligned}
& -\rho(z')\lambda(z') + \rho(C(z'))\lambda(C(z')) \\
&= -\rho(z')\lambda(z') + \rho(z')\lambda(C(z')) \\
&= -\rho(z')\max\{0, \langle \nabla U(x'), v' \rangle\} + \rho(z')\max\{0, \langle \nabla U(x'), R(x')v' \rangle\} \\
&= -\rho(z')\max\{0, \langle \nabla U(x'), v' \rangle\} + \rho(z')\max\{0, \langle \nabla U(x'), R(x')v' \rangle\} \\
&= -\rho(z')\max\{0, \langle \nabla U(x'), v' \rangle\} + \rho(z')\max\{0, -\langle \nabla U(x'), v' \rangle\} \\
&= -\rho(z')\langle \nabla U(x'), v' \rangle,
\end{aligned}$$

where we used that $\rho(z') = \rho(C(z'))$, $\langle \nabla U(x'), R(x')v' \rangle = -\langle \nabla U(x'), v' \rangle$ and $-\max\{0, f\} + \max\{0, -f\} = -f$ for any function f . Hence we have $\left. \frac{d\mu_t(z')}{dt} \right|_{t=0} = 0$, establishing that the bouncy particle sampler $\lambda^{\text{ref}} = 0$ admits ρ as invariant distribution. The invariance for $\lambda^{\text{ref}} > 0$ then follows from Lemma 1 given below.

Lemma 1. *Suppose P_t is a continuous time Markov kernel and Q is a discrete time Markov kernel which are both invariant with respect to μ . Suppose we construct for $\lambda^{\text{ref}} > 0$ a Markov process \hat{P}_t as follows: at the jump times of an independent PP with intensity λ^{ref} we make a transition with Q and then continue according to P_t , then \hat{P}_t is also μ -invariant.*

Proof. The transition kernel is given by

$$\begin{aligned}
\hat{P}_t &= e^{-\lambda t} P_t + \int_0^t dt_1 \lambda e^{\lambda t_1} e^{-\lambda(t-t_1)} P_{t-t_1} Q P_{t_1} \\
&\quad + \int_0^t dt_1 \int_{t_1}^{t_2} dt_2 \lambda^2 e^{\lambda t_1} e^{\lambda(t_2-t_1)} e^{-\lambda(t-t_2)} P_{t-t_2} Q P_{t_2-t_1} Q P_{t_1} + \dots
\end{aligned}$$

Therefore

$$\begin{aligned}
\mu \hat{P}_t &= \mu \left(e^{-\lambda t} + \lambda t e^{-\lambda t} + \frac{(\lambda t)^2}{2} e^{-\lambda t} \dots \right) \\
&= \mu.
\end{aligned}$$

Hence \hat{P}_t is μ -invariant. □

2 Invariance of the local sampler

The generator of the local BPS is given by

$$\begin{aligned}
\mathcal{L}h(z) &= \langle \nabla_x h(x, v), v \rangle \\
&\quad + \sum_{f \in F} \lambda_f(x, v) \{h(x, R_f(x)v) - h(x, v)\} \\
&\quad + \lambda^{\text{ref}} \int (h(x, v') - h(x, v)) \psi(dv').
\end{aligned} \tag{39}$$

The proof of invariance of the local BPS is very similar to the proof of Proposition 1. We have

$$\int \mathcal{L}h(z) \rho(z) dz = \int \int \langle \nabla_x h(x, v), v \rangle \rho(z) dz \tag{40}$$

$$+ \int \int \sum_{f \in F} \lambda_f(x, v) \{h(x, R_f(x)v) - h(x, v)\} \rho(z) dz \tag{41}$$

$$+ \lambda^{\text{ref}} \int \int \int (h(x, v') - h(x, v)) \psi(dv') \rho(z) dz \tag{42}$$

where the term (42) is straightforwardly equal to 0 while, by integration by parts, the term (40) satisfies

$$\int \int \langle \nabla_x h(x, v), v \rangle \rho(z) dz = \int \int \langle \nabla U(x), v \rangle h(x, v) \rho(z) dz. \quad (43)$$

as h is bounded. Now a change-of-variables shows that for any $f \in F$

$$\int \int \lambda_f(x, v) h(x, R_f(x) v) \rho(z) dz = \int \int \lambda(x, R_f(x) v) h(x, v) \rho(z) dz \quad (44)$$

as $R_f^{-1}(x) v = R(x) v$ and $\|R_f(x) v\| = \|v\|$ implies $\psi(R_f(x) v) = \psi(v)$. So the term (41) satisfies

$$\begin{aligned} & \int \int \sum_{f \in F} \lambda_f(x, v) \{h(x, R_f(x) v) - h(z)\} \rho(z) dz \\ = & \int \int \sum_{f \in F} [\lambda(x, R_f(x) v) - \lambda(x, v)] h(x, v) \rho(z) dz \\ = & \int \int \sum_{f \in F} [\max\{0, \langle \nabla U_f(x), R(x) v \rangle\} - \max\{0, \langle \nabla U_f(x), v \rangle\}] h(x, v) \rho(z) dz \\ = & \int \int \sum_{f \in F} [\max\{0, -\langle \nabla U_f(x), v \rangle\} - \max\{0, \langle \nabla U_f(x), v \rangle\}] h(x, v) \rho(z) dz \\ = & - \int \int \sum_{f \in F} [\langle \nabla U_f(x), v \rangle] h(x, v) \rho(z) dz \\ = & - \int \int \langle \nabla U(x), v \rangle h(x, v) \rho(z) dz, \end{aligned} \quad (45)$$

where we have used $\langle \nabla U_f(x), R_f(x) v \rangle = -\langle \nabla U_f(x), v \rangle$ and $\max\{0, -f\} - \max\{0, f\} = -f$ for any f . Hence, summing (43)-(45)-(42), we obtain $\int \mathcal{L}h(z) \rho(z) dz = 0$ and the result now follows by [1, Proposition 34.7].

3 Calculations in the isotropic normal case

As we do not use refreshment, it follows from the definition of the collision operator that

$$\begin{aligned} \langle x^{(i)}, v^{(i)} \rangle &= \left\langle x^{(i)}, v^{(i-1)} - \frac{2 \langle x^{(i)}, v^{(i-1)} \rangle}{\|x^{(i)}\|^2} x^{(i)} \right\rangle \\ &= -\langle x^{(i)}, v^{(i-1)} \rangle = -\langle x^{(i-1)}, v^{(i-1)} \rangle - \tau_i \\ &= \begin{cases} -\sqrt{-\log V_i} & \text{if } \langle x^{(i-1)}, v^{(i-1)} \rangle \leq 0 \\ -\sqrt{\langle x^{(i-1)}, v^{(i-1)} \rangle^2 - \log V_i} & \text{otherwise} \end{cases}, \end{aligned}$$

and therefore

$$\|x^{(i)}\|^2 = \begin{cases} \|x^{(i-1)}\|^2 - \langle x^{(i-1)}, v^{(i-1)} \rangle^2 - \log V_i & \text{if } \langle x^{(i-1)}, v^{(i-1)} \rangle \leq 0 \\ \|x^{(i-1)}\|^2 - \log V_i & \text{otherwise.} \end{cases}$$

It follows that $\langle x^{(j)}, v^{(j)} \rangle \leq 0$ for $j > 0$ if $\langle x^{(0)}, v^{(0)} \rangle \leq 0$ so, in this case, we have

$$\begin{aligned}
\|x^{(i)}\|^2 &= \|x^{(i-1)}\|^2 - \langle x^{(i-1)}, v^{(i-1)} \rangle^2 - \log V_i \\
&= \|x^{(i-1)}\|^2 + \log V_{i-1} - \log V_i \\
&= \|x^{(i-2)}\|^2 - \langle x^{(i-1)}, v^{(i-1)} \rangle^2 - \log V_{i-1} + \log V_{i-1} - \log V_i \\
&\vdots \\
&= \|x^{(1)}\|^2 - \langle x^{(1)}, v^{(1)} \rangle^2 - \log V_i
\end{aligned}$$

In particular for $x^{(0)} = e_1$ and $v^{(0)} = e_2$ with e_i being elements of standard basis of \mathbb{R}^d , the norm of the position at all points along the trajectory can never be smaller than 1.

4 Supplementary information on the evolutionary parameters inference experiments

4.1 Model

We consider an over-parameterized generalized time reversible rate matrix [2] with $d = 10$ corresponding to 4 unnormalized stationary parameters x_1, \dots, x_4 , and 6 unconstrained substitution parameters $x_{\{i,j\}}$, which are indexed by sets of size 2, i.e. where $i, j \in \{1, 2, 3, 4\}$, $i \neq j$. Off-diagonal entries of Q are obtained via $q_{i,j} = \pi_j \exp(x_{\{i,j\}})$, where

$$\pi_j = \frac{\exp(x_j)}{\sum_{k=1}^4 \exp(x_k)}.$$

We assign independent standard Gaussian priors on the parameters x_i . We assume that a matrix of aligned nucleotides is provided, where rows are species and columns contains nucleotides believed to come from a shared ancestral nucleotide. Given $x = (x_1, \dots, x_4, x_{\{1,2\}}, \dots, x_{\{3,4\}})$, and hence Q , the likelihood is a product of conditionally independent continuous time Markov chains over $\{A, C, G, T\}$, with “time” replaced by a branching process specified by the phylogenetic tree’s topology and branch lengths. The parameter x is unidentifiable, and while this can be addressed by bounded or curved parameterizations, the over-parameterization provides an interesting challenge for sampling methods, which need to cope with the strong induced correlations.

4.2 Baseline

We compare the BPS against a state-of-the-art HMC sampler [3] that uses Bayesian optimization to adapt the the leap-frog stepsize ϵ and trajectory length L of HMC. This sampler was shown in [4] to be comparable or better to other state-of-the-art HMC methods such as NUTS. It also has the advantage of having efficient implementations in several languages. We use the author’s Java implementation to compare to our Java implementation of the BPS. Both methods view the objective function as a black box (concretely, a Java interface supporting pointwise evaluation and gradient calculation). In all experiments, we initialize at the mode and use a burn-in of 100 iterations and no thinning. The HMC auto-tuner yielded $\epsilon = 0.39$ and $L = 100$. For our method, we use the global sampler and the global refreshment scheme.

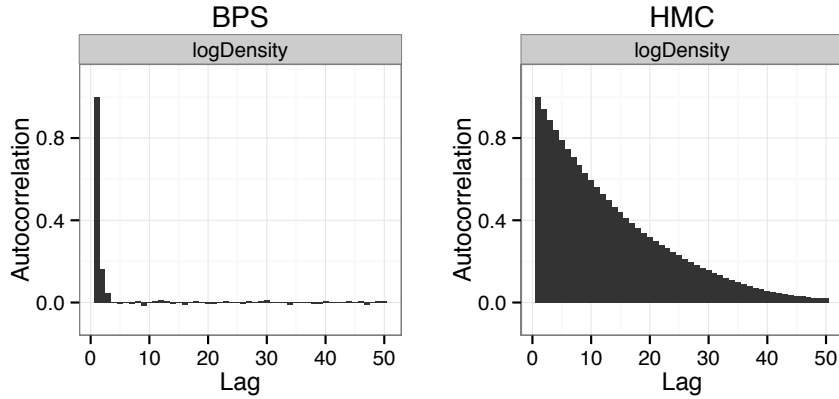


Figure 1: Estimate of the ACF of the log-likelihood statistic for BPS (left) and HMC (right). A similar behavior is observed for the ACF of the other statistics.

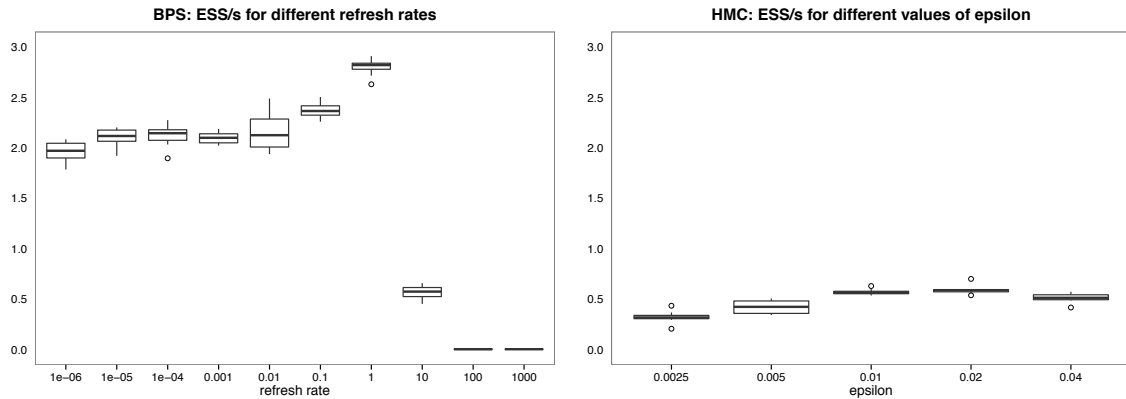


Figure 2: Left: sensitivity of BPS’s ESS/s on the log likelihood statistic. Right: sensitivity of HMC’s ESS/s on the log likelihood statistic. Each setting is replicated 10 times with different algorithmic random seeds.

4.3 Additional experimental results

To ensure that BPS outperforming HMC does not come from a faulty auto-tuning of HMC parameters, we look at the ESS/s for the log-likelihood statistic when varying the stepsize ϵ . The results in Figure 2(right) show that the value selected by the auto-tuner is indeed reasonable, close to the value 0.02 found by brute force maximization. We repeat the experiments with $\epsilon = 0.02$ and obtain the same conclusions. This shows that the problem is genuinely challenging for HMC.

The BPS algorithm also exhibits sensitivity to λ^{ref} . We analyze this dependency in Figure 2(left). We observe an asymmetric dependency, where values higher than 1 result in a significant drop in performance, as they bring the sampler closer to random walk behavior. Values one or more orders of magnitudes lower than 1 have a lower detrimental effect. However for a range of values of λ^{ref} covering six orders of magnitudes, BPS outperforms HMC at its optimal parameters.

References

- [1] M.H.A. Davis. *Markov Models & Optimization*, volume 49. CRC Press, 1993.
- [2] S. Tavaré. Some probabilistic and statistical problems in the analysis of DNA sequences. *Lectures on Mathematics in the Life Sciences*, 17:56–86, 1986.

- [3] Z. Wang, S. Mohamed, and N. de Freitas. Adaptive Hamiltonian and Riemann manifold Monte Carlo. In *Proceedings of the 30th International Conference on Machine Learning*, pages 1462–1470, 2013.
- [4] T Zhao, Z. Wang, A. Cumberworth, J. Gsponer, N. de Freitas, and A. Bouchard-Côté. Bayesian analysis of continuous time Markov chains with application to phylogenetic modelling. *Bayesian Analysis*, 2016. To appear.